# DGT-UD: a Parallel 23-language Parsebank
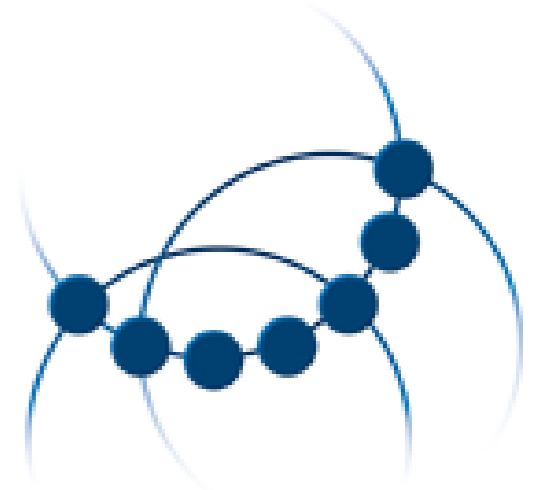
Nikola Ljubešić, Tomaž Erjavec

Department of Knowledge Technologies, Jožef Stefan Institute
SI-1000 Ljubljana, Slovenia
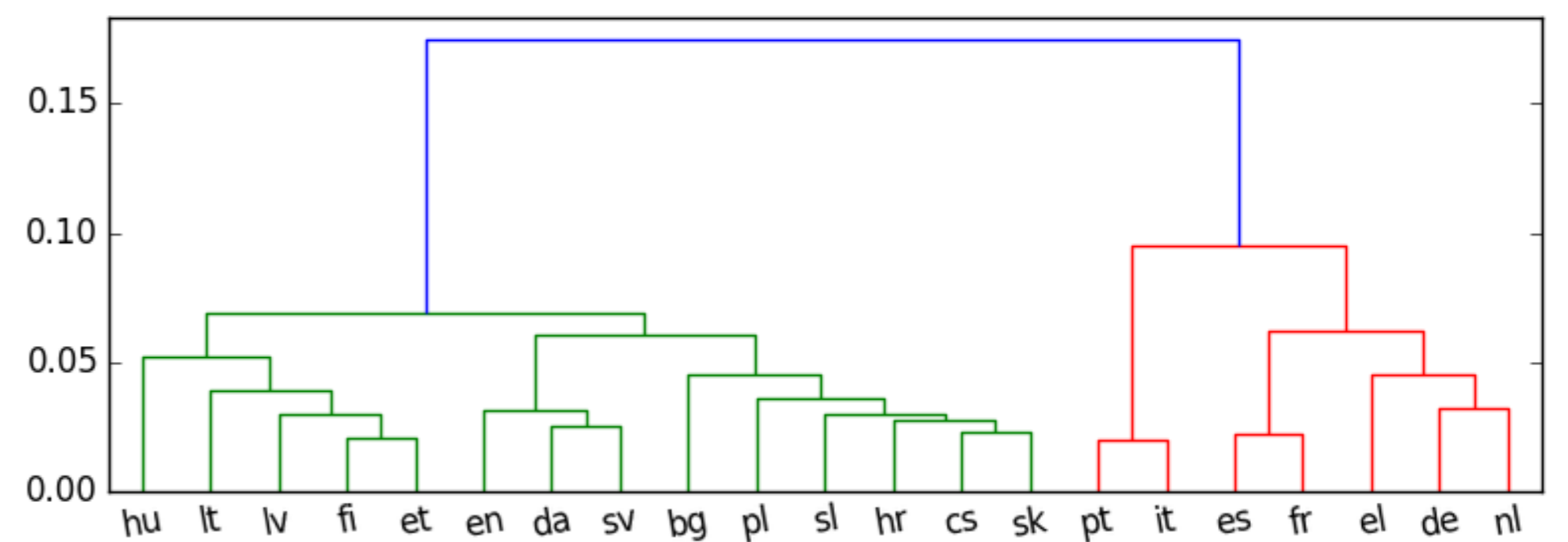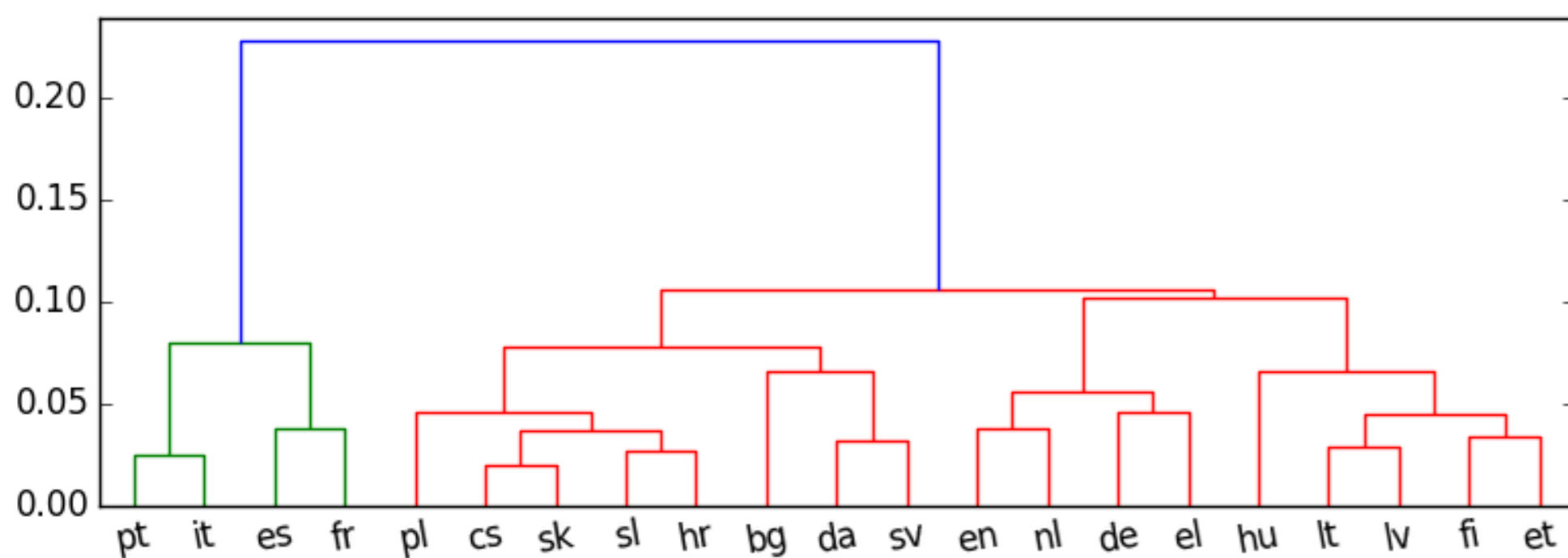{nikola.ljubesic,tomaz.erjavec}@ijs.si

CLARIN.SI

## The JRC DGT corpus

- JRC DGT corpus contains the translation memory of Acquis Communautaire (European Union law)

- Made publicly available in 2007 by the European Commission's Directorate-General for Translation (DGT) and the Joint Research Centre (JRC) of the European Commission

- Contains texts in 24 languages aligned on the sentence level and is distributed as a set of TMX files identified by the EUR-Lex number of the underlying documents.

- Proof-read texts, carefully checked translations, no duplicate sentences, is highly multilingual, large, available under a very permissive licence, and regularly updated.

- To create the DGT-UD corpus we took all the releases of JRC DGT, up to 2017.

## UDPipe

- UDPipe 1.2 – tool developed by UFAL to be trained on the Universal Dependencies (UD) data, used as baseline in shared tasks run on UD data

- Performs tokenization, sentence segmentation, morphosyntactic tagging, lemmatization and parsing, all learned from UD data

- Covers 68 models of 50 languages (UD v2.0), 23 out of 24 DGT languages, Maltese the only unsupported language

- LAS scores: Bulgarian (84.8), Croatian (77.9), Czech (83.2), Danish (74.7), Dutch (69.6), English (77.2), Estonian (65.6), Finnish (76.9), French (80.7), German (68.6), Greek (80.4), Hungarian (68.1), Irish (62.4), Italian (86.1), Latvian (62.8), Lithuanian (32.4), Polish (80.7), Portuguese (82.6), Romanian (80.2), Slovak (75.2), Slovenian (81.6), Spanish (81.4), Swedish (77.8)

## Comparative linguistic analysis

- **Question**: Is the automatically annotated DGT-UD better suited for comparative analyses than the manually annotated UD data?

- **Method**: Each CLARIN language represented as probability distribution of UD label trigrams, hierarchical clustering of languages.

- **Figure**: Dendrograms representing clustering results on DGT-UD (left) and UD data (right)



## Structure

The corpus consists of 1 alignment file, 23 vertical files and registry files (as used by the Manatee corpus manager back-end)
Structural attributes of vertical files:

- *Text* with id, year of publication & release e.g.
  `<text id="213A10" year="-2003" release="DGT-TM-2007">`

- *Anonymous block*, the unit of alignment, e.g. `<ab n="0">`

- *Sentence*, as annotated by UD-Pipe, i.e. `<s>`

Positional attributes of vertical files:

- *word*: the token, e.g. `förbindelser`

- *lempos*: the lemma of the token with added part-of-speech, e.g. `förbindelse-n;`

- *tag*: the "PoS tag" of the token, given for convenience, e.g. `ADJ_NumSing_GenMasc_DegPos_Def_Nom`

- *pos*: the UD part-of-speech of the token, e.g. `ADJ`

- *feats*: UD features of the token, e.g. `Case=Nom | Definite=Def | Degree=Pos | Gender=Masc | Number=Sing`

- *deprel*: the UD dependency relation, e.g. `nmod`

- *head_word*, *head_lempos*, *head_tag*, *head_pos*, *head_feats*: same as above, but for the token that is the head of the current token

- *id*, *head_id*: the index of the token and its head

## Corpus @ CLARIN.SI

### Explore with KonText



### Explore with noSketch Engine



### Download corpus