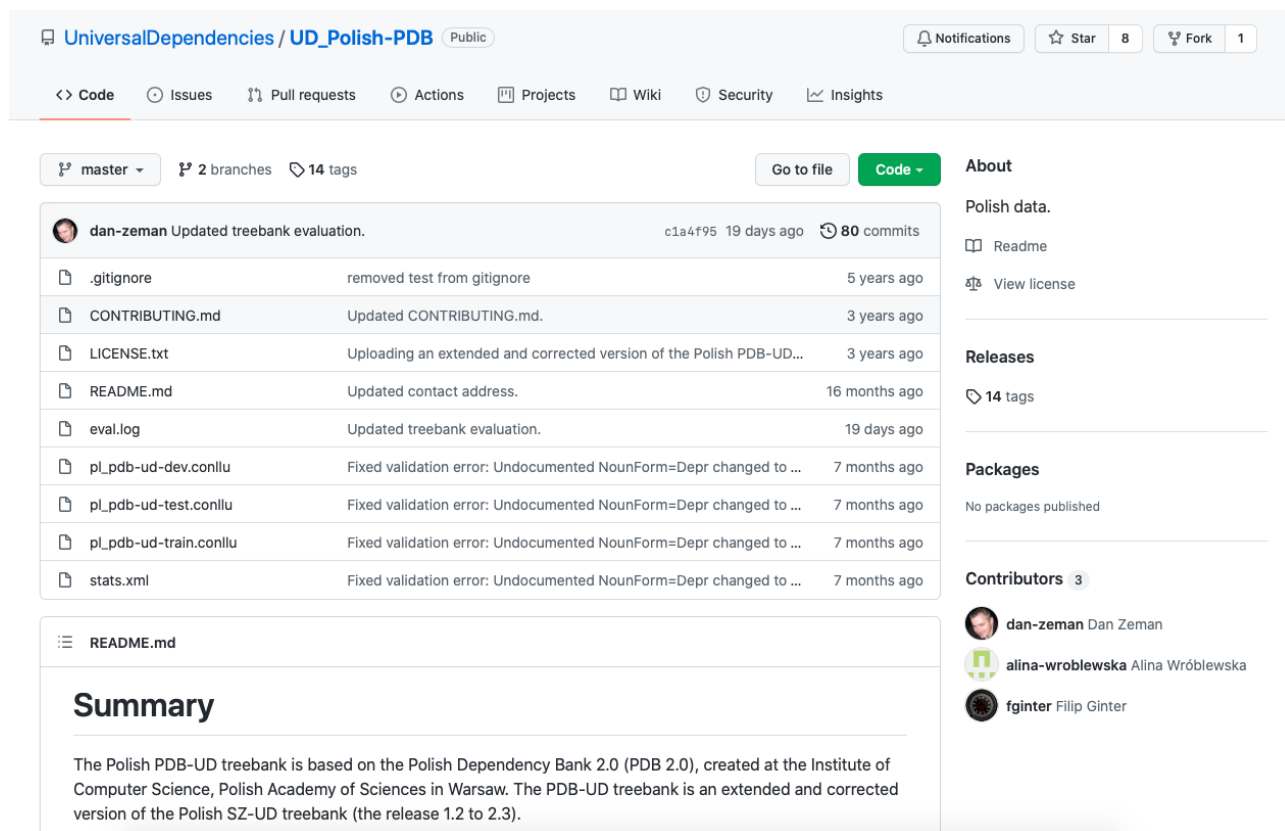


Universal Dependencies to międzynarodowa inicjatywa mająca na celu opracowanie uniwersalnego schematu anotacji drzew zależnościowych i zbudowanie ogromnej wielojęzycznej kolekcji banków drzew zależnościowych anotowanych zgodnie z tym schematem. Największym polskim bankiem drzew w tej kolekcji jest [UD\\_Polish-PDB](#) (Wróblewska, 2018) zawierający 22 tys. grafów zależności i 352 tys. segmentów. Duża część zdań leżących u podstaw banku drzew, tj. 14 tys. zdań, pochodzi z NKJP. Oprócz morfosyntaktycznych anotacji zgodnych ze schematem UD, UD\_Polish-PDB zawiera warstwę anotacji ze specyficznymi dla języka polskiego tagami wywodzącymi się z NKJP. Warto również podkreślić, że segmentacja całego banku PDB-UD jest również zbieżna ze schematem segmentacji opracowanym na potrzeby NKJP.



UniversalDependencies / UD\_Polish-PDB Public

Notifications Star 8 Fork 1

<> Code Issues Pull requests Actions Projects Wiki Security Insights

master 2 branches 14 tags Go to file Code

dan-zeman Updated treebank evaluation. c1a4f95 19 days ago 80 commits

File	Commit Message	Time Ago
.gitignore	removed test from gitignore	5 years ago
CONTRIBUTING.md	Updated CONTRIBUTING.md.	3 years ago
LICENSE.txt	Uploading an extended and corrected version of the Polish PDB-UD...	3 years ago
README.md	Updated contact address.	16 months ago
eval.log	Updated treebank evaluation.	19 days ago
pl_pdb-ud-dev.conllu	Fixed validation error: Undocumented NounForm=Depr changed to ...	7 months ago
pl_pdb-ud-test.conllu	Fixed validation error: Undocumented NounForm=Depr changed to ...	7 months ago
pl_pdb-ud-train.conllu	Fixed validation error: Undocumented NounForm=Depr changed to ...	7 months ago
stats.xml	Fixed validation error: Undocumented NounForm=Depr changed to ...	7 months ago

README.md

## Summary

The Polish PDB-UD treebank is based on the Polish Dependency Bank 2.0 (PDB 2.0), created at the Institute of Computer Science, Polish Academy of Sciences in Warsaw. The PDB-UD treebank is an extended and corrected version of the Polish SZ-UD treebank (the release 1.2 to 2.3).

**About**

Polish data.

Readme

View license

**Releases**

14 tags

**Packages**

No packages published

**Contributors** 3

- dan-zeman Dan Zeman
- alina-wroblewska Alina Wróblewska
- fginter Filip Ginter

Drugi pod względem wielkości polski bank drzew w repozytorium UD – [UD\\_Polish-LFG](#) (Patejuk i Przepiórkowski, 2018) – zawiera 17 tyś. grafów zależnościowych (130 tys. segmentów), których podstawę stanowią wyłącznie zdania pochodzące z NKJP.

Polskie banki drzew zależnościowych z repozytorium UD są wykorzystywane w badaniach naukowych. W badaniach Culbertson i in. (2020) z pogranicza lingwistyki i kognitywistyki opublikowanych w prestiżowym czasopiśmie *Language*, autorzy sprawdzali, czy struktura frazy rzeczownikowej jest determinowana przez własności obiektów istniejących w świecie rzeczywistym oraz czy na podstawie cech statystycznych tych obiektów można się nauczyć hipotetycznych struktur fraz rzeczownikowych.

FROM THE WORLD TO WORD ORDER: DERIVING BIASES IN  
NOUN PHRASE ORDER FROM STATISTICAL PROPERTIES OF THE WORLD

JENNIFER CULBERTSON	MARIEKE SCHOUWSTRA	SIMON KIRBY
<i>Centre for Language Evolution, University of Edinburgh</i>	<i>Centre for Language Evolution, University of Edinburgh</i>	<i>Centre for Language Evolution, University of Edinburgh</i>

The world's languages exhibit striking diversity. At the same time, recurring linguistic patterns suggest the possibility that this diversity is shaped by features of human cognition. One well-studied example is word order in complex noun phrases (like *these two red vases*). While many orders of these elements are possible, a subset appear to be preferred. It has been argued that this ordering reflects a single underlying representation of noun phrase structure, from which preferred orders are straightforwardly derived (e.g. Cinque 2005). Building on previous experimental evidence using artificial language learning (Culbertson & Adger 2014), we show that these preferred orders arise not only in existing languages, but also in improvised sequences of gestures produced by English speakers. We then use corpus data from a wide range of languages to argue that the hypothesized underlying structure of the noun phrase might be learnable from statistical features relating objects and their properties conceptually. Using an information-theoretic measure of strength of association, we find that adjectival properties (e.g. *red*) are on average more closely related to the objects they modify (e.g. *wine*) than numerosities are (e.g. *two*), which are in turn more closely related to the objects they modify than demonstratives are (e.g. *this*). It is exactly those orders which transparently reflect this—by placing adjectives closest to the noun, and demonstratives farthest away—that are more common across languages and preferred in our silent gesture experiments. These results suggest that our experience with objects in the world, combined with a preference for transparent mappings from conceptual structure to linear order, can explain constraints on noun phrase order.\*

*Keywords:* word order, typology, silent gesture, corpora, information theory

**1. INTRODUCTION.** One of the oldest debates in linguistics concerns whether the languages of the world share a set of core invariant properties reflecting universal features of human cognition. At the center of this debate is a tension between the diversity we see when we look across languages and the similarities that crop up when they are analyzed under a certain lens. This tension, between linguistic diversity on the one hand and universal organizing principles on the other, is on full display in one of the simplest linguistic structures we use: the noun phrase. Given just a noun (e.g. *vases*) and three common categories of words that modify it—a demonstrative (e.g. *these*), a numeral (e.g. *two*), and an adjective (e.g. *blue*)—there are already twenty-four possible ways of ordering the words to make a phrase, almost all of which are found in some language. For example, the English order is *these two blue vases*; in Thai, it would be the equivalent of *vases blue two these*; in Vietnamese, it would be *these two vases blue*; in Basque, it would be *two vases blue these*; and so on. Yet there remains a small subset of orders that no language appears to use systematically. For example, we currently know of no language that systematically uses the equivalent of *blue two these vases* or *blue these vases two*.

Linguists have argued that these missing patterns offer evidence of universal organizing principles underlying how noun phrases are built (Cinque 2005, Steddy & Samek-Lodovici 2011, Abels & Neeleman 2012, Dryer 2018, Steedman 2018). As careful

\* We would like to thank Roger Levy and the referees for their comments on previous versions of this work. This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 Research and Innovation Programme (grant agreement No. 757643).

W dziedzinach lingwistyki obliczeniowej i przetwarzania języka naturalnego banki drzew są wykorzystywane głównie do trenowania modeli umożliwiających analizę morfosyntaktyczną danych w języku naturalnym dla współczesnych systemów wstępnego przetwarzania języka naturalnego. Przykładowo UD\_Polish-PDB został wykorzystany do wytrenowania modeli dla systemów NLP:

- COMBO (Klimaszewski i Wróblewska, 2021): *polish-herbert-large|polish-herbert-base*,
- spaCy (Honnibal i in., 2020): *pl\_core\_news\_sm|md|lg*,
- spaCy v3: *pl\_core\_news\_sm|md|lg*,
- UDPipe (Straka, 2018): *polish-pdb-ud-2.5-191206*,
- Stanza (Qi i in., 2020): *pl*.

Model COMBO dla polskiego wytrenowany na UD\_Polish-PDB został udostępniony publicznie do celów badawczych w [aplikacji webowej](#).



The screenshot displays the COMBO web application interface. On the left, there is a sidebar with the application logo (PAN and CLARIN) and a language selection menu showing 'English' and 'Polish'. The main content area is titled 'Wstępne przetwarzanie języka naturalnego' and includes a description of the system. Below the description, there are input fields for selecting an example or providing a custom sentence. The example sentence 'Mądrość czyni z wiedzy właściwy, korzystny dla człowieka użytek.' is entered and processed. The result is a dependency parse tree for the sentence, with the root node 'czyni' (verb) connected to 'Mądrość' (subject), 'z wiedzy' (oblique argument), and 'właściwy, korzystny dla człowieka użytek' (object). The parse tree nodes are color-coded: green for nouns, pink for verbs, and light blue for oblique arguments and punctuation.