

# Cross-Language Harmonization of Linguistic Resources

Daniel Zeman

📅 November 6, 2024

## Multilingual:

- not just for English
- not just for resource-rich European languages (English, Czech, ...)
- ideally **all languages** (up to 7000?)
- including
  - endangered
  - unwritten
  - very different from Indo-European

## Multilingual:

- not just for English
- not just for resource-rich European languages (English, Czech, ...)
- ideally all languages (up to 7000?)
- including
  - endangered
  - unwritten
  - very different from Indo-European
- focus on
  - Morphology
  - Syntax

# Morphology: Explaining Word Forms



## English

*do, does, did, done, doing*



## Czech

*dělat, dělali, dělám, děláš, dělá, děláme, děláte, dělají, dělej, dělejme, dělejte, dělal, dělala, dělalo, dělali, dělaly, dělaje, dělajíc, dělajíce, dělán, dělána, děláno, dělání, dělány*



## Finnish

*tehdä, tee, teemme, teen, teenkin, teenkö, teet, teette, tehden, tehdessä, tehdessään, tehdyillä, tehdyissä, tehdyistä, tehdyille, tehdyllä, tehdyn, tehdyssä, tehdystä, tehdyt, tehdään, tehkää, tehkäämme, tehkөөn, tehkөөt, tehnee, tehneen, tehneensä, tehneet, tehneille, tehnumme, tehnen, tehnet, tehnette, tehnevät, tehny, tehnyt, tehnytkin, tehtiin, tehtiinkin, tehty, tehtyjen, tehtyjä, tehtynä, tehtyyn, tehtyä, tehtyäni, tehtyään, tehtäessä, tehtäisi, tehtäisiin, tehtäkö, tehtäkөөn, tehtämän, tehtäne, tehtäneen, tehtävien, tehtäviin, tehtäville, tehtävistä, tehtäviä, tehtävä, tehtävän, teimme, tein, teinkin, teinpä, teit, teitte, tekee, tekeekin, tekeekö, tekemiensä, tekemiin, tekeminen, tekemistä, tekemiä, tekemiäni, tekemiään, tekemä, tekemäisilläni, tekemällä, tekemän, tekemäni, tekemänsä, tekemässä, tekemästä, tekemästään, tekemät, tekemättä, tekemää, tekemään, tekemäänsä, tekevien, tekeville, tekevä, tekeväni, tekevät, tekevää, tekevään, teki, tekikin, tekis, tekisi, tekisimme, tekisin, tekisit, tekisitte, tekisivät, tekivät, tekivätkö*

# Morphology: Explaining Word Forms



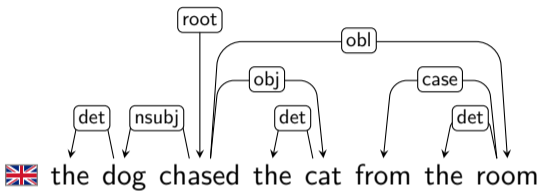
word form: *dělám* “I do”  
lemma: *dělat* “do”  
word category  
(part of speech): **VERB**  
features: **indicative mood**  
**present tense**  
**active voice**  
**first person singular**




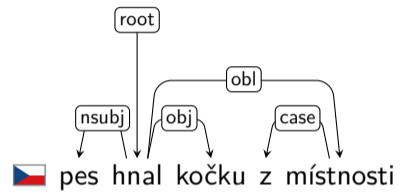
word form: *teen* “I do”  
lemma: *tehdä* “do”  
word category  
(part of speech): **VERB**  
features: **indicative mood**  
**present tense**  
**active voice**  
**first person singular**




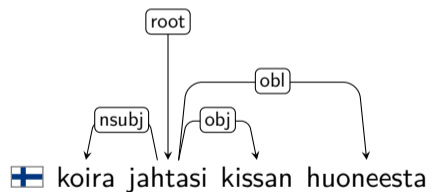
# Syntax: Relations between Words




 the dog chased the cat from the room



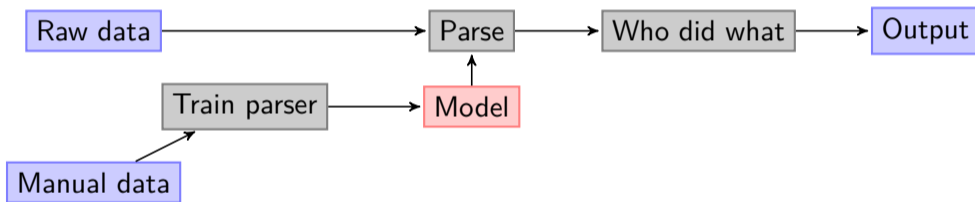
 pes hnal kočku z místnosti



 koira jahtasi kissan huoneesta

# What Is It Good For?

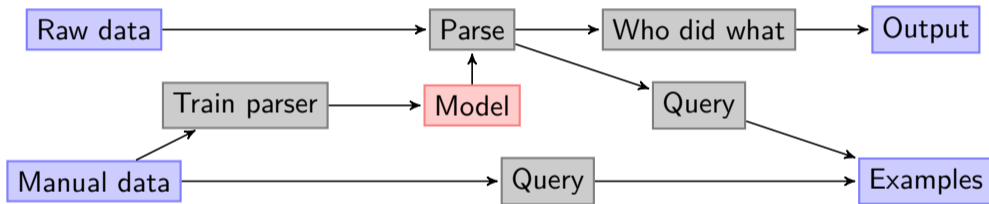
- Downstream **language understanding** technology





# What Is It Good For?

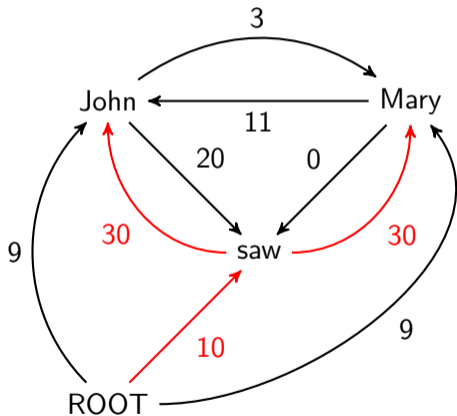
- Downstream **language understanding** technology
- **Linguistic research**



- More broadly: **digital humanities** (contrastive linguistics, typology, literary studies, cultural heritage, works in classical languages)

# Parsing

Total graph + edge scores  $\rightarrow$  **maximum spanning tree**




$G = \langle V, E \rangle =$

$$A^* = \underset{\substack{A \subseteq G \\ A \text{ is tree}}}{\operatorname{argmax}} \sum_{e \in A} \text{score}(e)$$

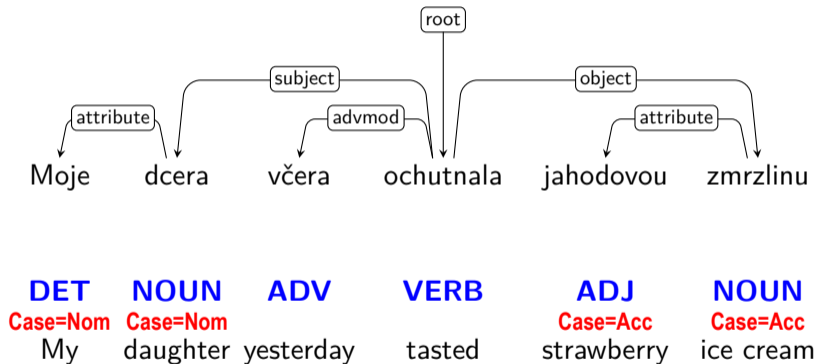
# Parsing Low-Resource Languages

- Machine learning:
- Manual annotation needed to train a parser
  - Thousands of sentences – hard and expensive
- Available for a few “lucky” languages
- But there are thousands of languages – **what about the less lucky ones?**

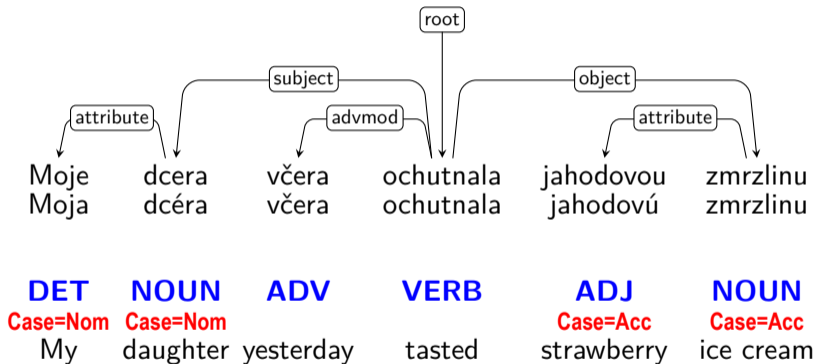
---

*Daniel Zeman, Philip Resnik (2008): Cross-Language Parser Adaptation between Related Languages. In: IJCNLP 2008 Workshop on NLP for Less Privileged Languages, pp. 35-42, Hyderabad, India  265 cit.*

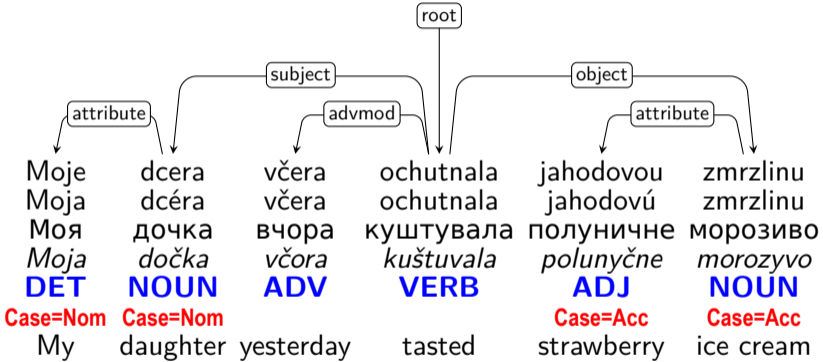
# Delexicalized Parsing



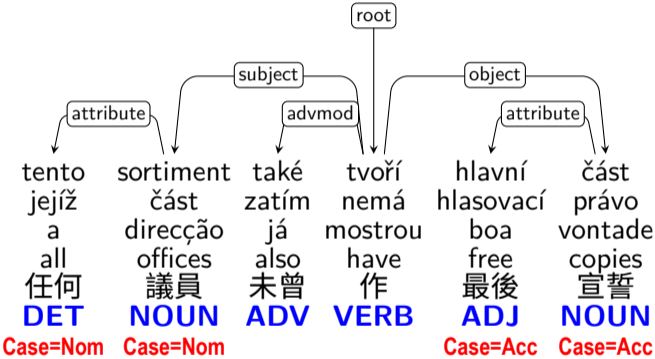
# Delexicalized Parsing



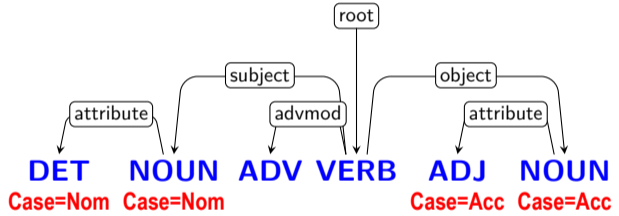
# Delexicalized Parsing



# Delexicalized Parsing




# Delexicalized Parsing





## Evaluation: Does It Work?

- Manually annotated **test data** needed  
⇒ Simulate on languages for which we have test data
- Train parser on  Czech  
(large training data available)



*truepos* ... correctly  
identified dependencies

$$P = \frac{\textit{truepos}}{\textit{truepos} + \textit{falsepos}}$$

$$R = \frac{\textit{truepos}}{\textit{truepos} + \textit{falseneg}}$$

$$LAS = F_1 = \frac{2 \times P \times R}{P + R}$$

## Evaluation: Does It Work?

- Manually annotated **test data** needed  
⇒ Simulate on languages for which we have test data
- Train parser on  Czech  
(large training data available)
- Evaluate on  Slovak  
(some test data available)



*truepos* ... correctly  
identified dependencies

$$P = \frac{\textit{truepos}}{\textit{truepos} + \textit{falsepos}}$$

$$R = \frac{\textit{truepos}}{\textit{truepos} + \textit{falseneg}}$$

$$LAS = F_1 = \frac{2 \times P \times R}{P + R}$$

## Evaluation: Does It Work?

- Manually annotated **test data** needed  
⇒ Simulate on languages for which we have test data
- Train parser on  Czech  
(large training data available)
- Evaluate on  Slovak  
(some test data available)
- If OK, assume it also works for  Lower Sorbian  
(no data available)

*truepos* ... correctly  
identified dependencies

$$P = \frac{\textit{truepos}}{\textit{truepos} + \textit{falsepos}}$$

$$R = \frac{\textit{truepos}}{\textit{truepos} + \textit{falseneg}}$$

$$LAS = F_1 = \frac{2 \times P \times R}{P + R}$$


# PROBLEM:

The annotations must be compatible across languages!

# Morphology Example: SynTagRus Tags vs. PDT Tags

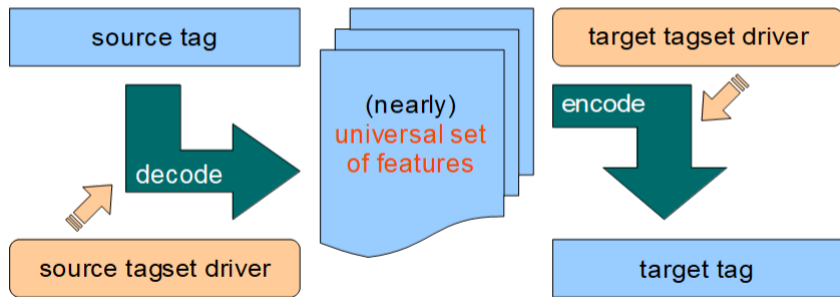
SynTagRus annotation	PDT annotation	Meaning
S ЕД МУЖ ИМ	NNMS1-----A----	noun masculine singular nominative
S МН РОД ОД	PSXXXXP3-----	pronoun possessive plural 3 <sup>rd</sup> person
A МН ИМ	AAXP1----1A----	adjective positive plural nominative
NUM ВИН	C1XX4-----	numeral cardinal accusative
V НЕСОВ ИЗЪЯВ...	VB-P---3P-AA---	verb imperfective present indicative ...
... НЕПРОШ МН 3-Л		... 3 <sup>rd</sup> person plural
ADV СРАВ	Dg-----2A----	adverb comparative
PR	RR--6-----	preposition
CONJ	J^-----	coordinating conjunction
PART	TT-----	particle
INTJ	II-----	interjection

---

*Daniel Zeman (2008): Reusable Tagset Conversion Using Tagset Drivers. In: Proceedings of LREC, pp. 213–218, ELRA, Marrakech, Morocco  234 cit.*

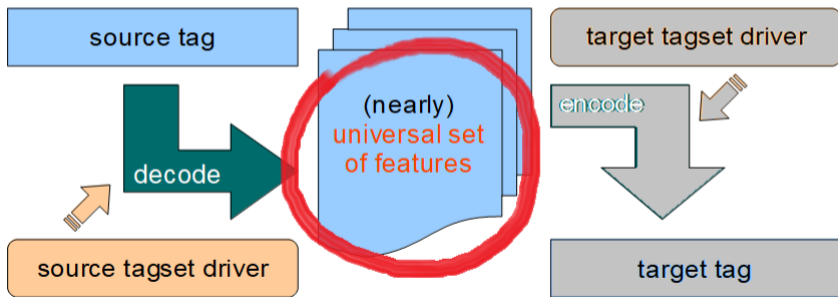
# Tagset Drivers

- A module with the following functions:
  - `decode()` ... converts a tag to Interaset
  - `encode()` ... generates a tag from Interaset
  - `list()` ... lists known tags in the tagset (optional)



# Tagset Drivers

- A module with the following functions:
  - `decode()` ... converts a tag to Interaset
  - `encode()` ... generates a tag from Interaset
  - `list()` ... lists known tags in the tagset (optional)









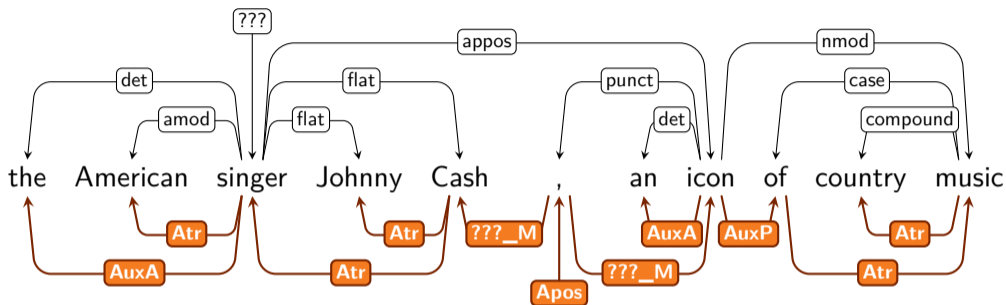








# Syntax: Stanford vs. Prague Dependencies



- Change relation labels
- Bottom-up **tree transformations**

*Daniel Zeman, Ondřej Dušek, David Mareček, Martin Popel, Loganathan Ramasamy, Jan Štěpánek, Zdeněk Žabokrtský, Jan Hajič (2014): HamleDT: Harmonized Multi-Language Dependency Treebank. In: Language Resources and Evaluation, ISSN 1574-020X, vol. 48, no. 4, pp. 601–637 [88 cit.](#)*




# Universal Dependencies


- Joined our (Prague) forces with similar efforts at Stanford, Google, and Uppsala
- Defined universally-applicable annotation standard
  - Morphology from Intersect
  - Syntax adapted from “Stanford Dependencies”



---

Marie-Catherine de Marneffe, Christopher Manning Joakim Nivre, **Daniel Zeman** (2021): *Universal Dependencies*. In: *Computational Linguistics*, vol. 47, no. 2, pp. 255–308  610 cit. *WoS Q1*

---

Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, ... **Daniel Zeman** (2016): *Universal Dependencies v1: A Multilingual Treebank Collection*. In: *Proceedings of LREC*, pp. 1659–1666  1694 cit.

# Universal Dependencies

- Joined our (Prague) forces with similar efforts at Stanford, Google, and Uppsala
- Defined universally-applicable annotation standard
  - Morphology from Intersect
  - Syntax adapted from “Stanford Dependencies”
- Converted existing datasets to UD



---

Marie-Catherine de Marneffe, Christopher Manning Joakim Nivre, **Daniel Zeman** (2021): *Universal Dependencies*. In: *Computational Linguistics*, vol. 47, no. 2, pp. 255–308 📖 610 cit. WoS Q1

---

Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, ... **Daniel Zeman** (2020): *Universal Dependencies v2: An Evergrowing Multilingual Treebank Collection*. In: *Proceedings of LREC*, pp. 4034–4043 📖 574 cit.




# Universal Dependencies


- Joined our (Prague) forces with similar efforts at Stanford, Google, and Uppsala
- Defined universally-applicable annotation standard
  - Morphology from Intersect
  - Syntax adapted from “Stanford Dependencies”
- Converted existing datasets to UD
- Collecting new datasets for new languages
- Building infrastructure to lower the entry barrier for new data contributors
- Organizing shared tasks (evaluation campaigns for parsers)



---

Marie-Catherine de Marneffe, Christopher Manning Joakim Nivre, **Daniel Zeman** (2021): *Universal Dependencies*. In: *Computational Linguistics*, vol. 47, no. 2, pp. 255–308  610 cit. *WoS Q1*

---

**Daniel Zeman**, Jan Hajič, Martin Popel, ... Slav Petrov (2018): *CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*. In: *Proceedings of CoNLL*. pp. 1–21  671 cit.

**Daniel Zeman, Joakim Nivre, Mitchell Abrams, Elia Ackermann, Noëmi Aeppli, Hamid Aghaei, Željko Agić, Amir Ahmadi, Lars Ahrenberg, Chika Kennedy Ajede, Salih Furkan Akkurt, Gabriél Aleksandravičiūtė, Ika Alfina, Avner Algom, Khalid Alnajjar, Chiara Alzetta, Erik Andersen, Lene Antonsen, Tatsuya Aoyama, Katya Aponova, Angelina Aquino, Carolina Aragon, Glyd Aranes, Maria Jesus Aranzabe, Bilge Nas Arcan, Pórunn Arnardóttir, Gashaw Arutie, Jessica Naraiswari Arwidarasti, Masayuki Asahara, Katla Ásgeirsdóttir, Deniz Baran Aslan, Cengiz Asmazoğlu, Luma Ateyah, Furkan Atmaca, Mohammed Attia, Aitziber Atutxa, Liesbeth Augustinus, Mariana Avelás, Elena Badmaeva, Keerthana Balasubramani, Miguel Ballesteros, Esha Banerjee, Sebastian Bank, Verginica Barbu Mititelu, Starkardur Barkason, Rodolfo Basile, Victoria Basmov, Colin Batchelor, John Bauer, Seyyit Talha Bedir, Shabnam Behzad, Juan Belieni, Kepa Bengoetxea, Ibrahim Benli, Yifat Ben Moshe, Ansu Berg, Gözde Berk, Riyaz Ahmad Bhat, Erica Biagetti, Eckhard Bick, Agnė Bielinskiėnė, Esma Fatima Bilgin Taşdemir, Kristin Bjarnadóttir, Verena Blaschke, Rogier Blokland, Victor Bobicew, Loïc Boizou, Johnatan Bonilla, Emanuel Borges Völker, Carl Börstell, Cristina Bosco, Gosse Bouma, Sam Bowman, Adriane Boyd, Anouck Braggarr, António Branco, Kristina Brookaitė, Aljoscha Burchardt, Marisa Campos, Marie Candito, Bernard Caron, Gauthier Caron, Catarina Carvalho, Rita Carvalho, Lauren Cassidy, Maria Clara Castro, Sérgio Castro, Tatiana Cavalcanti, Gülşen Cebirođlu Eryiđit, Flavio Massimiliano Cecchini, Giuseppe G. A. Celano, Slavomir Čepłó, Neslihan Cesur, Savas Cetin, Özlem Çetinođlu, Fabricio Chalub, Lityanage Chamila, Shweta Chauhan, Yifei Chen, Ethan Chi, Taishi Chika, Yongseok Cho, Jinho Choi, Bermet Chontayeva, Jayeol Chun, Juyeon Chung, Alessandra T. Cignarella, David Cinková, Aurélie Collomb, Çađır Çöltekin, Miriam Connor, Claudia Corbetta, Daniela Corbetta, Francisco Costa, Marine Courtin, Benoît Crabbé, Mihaela Cristescu, Vladimir Cvetkosić, Ingerid Løyning Dale, Philemon Daniel, Elizabeth Davidson, Leonel Figueiredo de Alencar, Mathieu Dehouck, Martina de Laurentiis, Marie-Catherine de Marneffe, Valeria de Paiva, Mehmet Oguz Derin, Elvis de Souza, Arantza Diaz de Ilarraza, Rønning Antonio Díaz Hernández, Carly Dickerson, Arawinda Dinakaramani, Elisa Di Nuovo, Bamba Dione, Peter Dirix, Hoa Do, Kaja Dobrovoljc, Caroline Döhmer, Adrian Doyle, Timothy Dozat, Kira Droganova, Magali Sanches Duran, Puneet Dwivedi, Christian Ebert, Hanne Eckhoff, Masaki Eguchi, Sandra Eiche, Ronald Eiselein, Marhaba Eli, Ali Elkahky, Binyam Ephrem, Olga Erina, Tomaž Erjavec, Soudabeh Esлами, Farah Essaidi, Aline Etienne, Wograine Evelyn, Sidney Facundes, Richard Farkas, Federica Favero, Jannatul Ferdousi, Marilía Fernanda, Hector Fernandez Alcalde, Amal Fethi, Jennifer Foster, Theodor Fransen, Cláudia Freitas, Kazunori Fujita, Katarína Gajdošová, Daniel Galbraith, Edith Galy, Federica Gamba, Marcos Garcia, Moe Gårdenfors, Tanja Gaustad, Efe Eren Genç, Fabricio Ferraz Gerardi, Kim Gerdes, Luke Gessler, Filip Ginter, Gustavo Godoy, Iakes Goenaga, Koldo Gojenola, Memduh Gökürmak, Yoav Goldberg, Xavier Gómez Guinovart, Berta González Saavedra, Bernadeta Gričiūtė, Matias Groni, Loic Grobel, Normunds Grūzītis, Bruno Guillaume, Kirian Guiller, Céline Guillot-Barbance, Tunga Güngör, Nizar Habash, Henrik Hafsteinsson, Jan Hajić, Jan Hajić, jr., Mika Hämäläinen, Linh Hà My, Na-Rae Han, Muhammad Yudistira Hanifmuti, Takahiro Harada, Sam Hardwick, Kim Harris, Naima Hassert, Dag Haug, Johannes Heinecke, Oliver Hellwig, Felix Hennig, Barbara Hladká, Jaroslava Hlaváčová, Florinel Hociung, Diana Hoefels, Petter Hohle, Yidi Huang, Marivel Huerta Mendez, Jena Huang, Takumi Ikeda, Inessa Iliadova, Anton Kar Ingason, Radu Ion, Elena Irimia, Oljādiė Ishola, Artan Islamaj, Kaoru Ito, Federica Iurescia, Sandra Jagodzinska, Siratun Jannat, Tomáš Jelinek, Apoorva Jha, Katharine Jiang, Mayank Jobanputra, Anders Johanness, Hildur Jónsdóttir, Fredrik Jørgensen, Markus Juutinen, Hüner Kaşıkara, Nadezhda Kabacheva, Sylvain Kahane, Hiroshi Kanayama, Jenna Kanerva, Neslihan Kara, Ritvān Karahöđe, Andre Käsen, Tolga Kayadelen, Sarveswaran Kengathariyer, Václava Kettnerová, Liilit Kharatyan, Jesse Kirchner, Elena Klementyeva, Elena Klyachko, Petr Kocharov, Arne Köhn, Abdullatif Köksal, Kamil Kopicewicz, Timo Korhikangas, Mehmet Köse, Alexey Koshevoy, Natalia Kotsyba, Barbara Kováčik, Jolanta Kovalcskaić, Simon Kreh, Parameswari Krishnamurthy, Sandra Kübler, Adrian Kuqi, Öđuzhan Kuyrukçu, Aslı Kuzgun, Sookyoung Kwak, Kris Kyle, Kābi Laan, Veronika Laippala, Lorenzo Lambertino, Tatiana Lando, Septina Diana Larasati, Alexei Lavrentiev, John Lee, Phuong Lê Hông, Alessandro Lenci, Saran Lertrpadit, Herman Leung, Maria Levina, Lauren Levine, Cheuk Ying Li, Josie Li, Keying Li, Yixuan Li, Yuan Li, KyungTae Lim, Bruna Lima Padovani, Yi-Yu Jessica Lin, Kristér Lindén, Yang Janet Liu, Nikola Lubješić, Irina Lobzhanidze, Olga Loginova, Lucelene Lopes, Stefano Lusito, Anne-Marie Lutgen, Andry Luthfi, Mikko Luukko, Olga Lyashevskaya, Teresa Lynn, Vivien Macketanz, Menel Mahamdi, Jean Maillard, Ilya Makarchuk, Aibek Makazhanov, Francesco Mambrini, Michael Mandl, Christopher Manning, Ruli Manurung, Büsra Marşan, Cătălina Măranduc, David Mareček, Katrin Marheinecke, Stella Markantonatou, Héctor Martínez Alonso, Lorena Martín Rodríguez, André Martins, Cláudia Martins, Jan Mašek, Hiroshi Matsuda, Yuji Matsumoto, Alessandro Mazzei, Ryan McDonald, Sarah McGuinness, Maitrey Mehta, Pierre André Ménard, Gustavo Mendonça, Tatiana Merzhevich, Paul Meurer, Niko Miekka, Emilia Milano, Aaron Miller, Karina Mischenkova, Anna Missilá, Cătălin Mititelu, Maria Mitrofan, Yusuke Miyao, AmirHossein Mojiri Ferozshani, Judit Molnár, Amirsaeed Moloodi, Simonetta Montemagni, Amir More, Laura Moreno Romero, Giovanni Moretti, Shinsuke Mori, Tomohiko Morioka, Shigeki Moro, Bjartur Mortensen, Bohdan Moskalevskiy, Kadri Muischnek, Robert Munro, Yugo Murawaki, Kaili Müürisep, Pinkey Nainwani, Mariam Nakhlé, Juan Ignacio Navarro Horrićacek, Anna Nedoluzhko, Gunta Nešpore-Bėrzkalne, Manuela Nevaci, Lương Nguyễn Thị, Huýn Nguyễn Thị Minh, Yoshihiro Nikaido, Vitaly Nikolaev, Rattima Nitisarog, Victor Norman, Alireza Nourian, Maria das Graças Volpe Nunes, Hanna Nurmi, Stina Ojala, Atul Kr. Ojha, Hulda Óladóttir, Adedayó Olóluòkun, Mai Omura, Emeka Onwuegbuzia, Noam Ordan, Petya Osenova, Robert Ostling, Annika Ott, Liija Övrelid, Šaziye Betül Özates, Merve Özçelik, Arzucan Özgür, Balkız Öztürk Başaran, Teresa Paccosi, Alessio Palermò Aprosio, Anastasia Panova, Thiago Alexandre Salgueiro Pardo, Hyunji Hayley Park, Niko Partanen, Elena Pascual, Marco Passarotti, Agnieszka Patejuk, Guilherme Paulino-Passos, Giulia Pedoneso, Angelika Peljak-Lapińska, Siyao Peng, Siyao Logan Peng, Rita Pereira, Sílvia Pereira, Cene-Augusto Perez, Natalia Perková, Guy Perrier, Slav Petrov, Daria Petrowa, Andrea Peverelli, Jason Phelan, Claudel Pierre-Louis, Jussi Piitulainen, Yuval Pinter, Clara Pinto, Rodrigo Pintucci, Tommi A Pirinen, Emily Pitler, Magdalena Plamada, Barbara Plank, Alistair Plum, Thierry Poibeau, Larisa Ponomareva, Martin Popel, Lauma Pretkalniņa, Rigardt Pretorius, Sophie Prévost, Prokopis Prokopidis, Adam Przepiórkowski, Robert Pugh, Tiina Puolakainen, Christoph Puschke, Sampo Pyysalo, Peng Qi, Andrea Querido, Andriela Rääbis, Alexandre Rademacher, Mizanur Rahaman, Taraka Rama, Loganathan Ramasamy, Carlos Ramisch, Joana Ramos, Fan Rashel, Mohammad Sadeq Rasooli, Vinit Ravishanker, Livy Real, Petru Rebeja, Siva Reddy, Mathilde Regnault, Georg Rehm, Arij Riabi, Ivan Riabov, Michael Rießler, Erika Rimkute, Larissa Rinaldi, Laura Rituma, Putri Rizqiyah, Luisa Rocha, Eirikur Rögnvaldsson, Ivan Roksandic, Mykhailo Romanenko, Rudolf Rosa, Valentin Roşca, Davide Rovati, Ben Rozonoyer, Olga Rudina, Jack Rueter, Paolo Ruffolo, Kristján Rúnarsson, Shoval Sadde, Pegah Safari, Aleksí Sahala, Shadi Saleh, Alessio Salomoni, Tanja Samardžić, Stephanie Samson, Xulia Sánchez-Rodríguez, Manuela Sanchis Sanguinetti, Ezgi Sanıyur, Dage Särg, Marta Sartor, Albina Sarymskova, Mitsuya Sasaki, Baiba Saulite, Agata Savary, Yanin Sawanakunanon, Shefali Saxena, Kevin Scannell, Salvatore Scarlata, Emmanuel Schang, Nathan Schneider, Sebastian Schuster, Lena Schwartz, Djameđ Seddah, Wolfgang Seeker, Sven Sellmer, Moigan Seraji, Syeda Shahzadi, Mo Shen, Atsuko Shimada, Hiroyuki Shirasu, Yana Shishkina, Muh Shohibussirri, Maria Shvedova, Janine Siewert, Einar Frey Sigurðsson, João Silva, Aline Silveira, Natalia Silveira, Sara Silveira, Maria Simi, Radu Simionescu, Katalin Simkó, Mária Šimková, Haukur Barri Símonarson, Kiril Simov, Dmitri Sitchinava, Ted Sither, Aaron Smith, Isabela Soares-Bastos, Per Erik Solberg, Barbara Sonnenhauser, Shafi Sourov, Rachele Sprugnoli, Vivian Stamu, Steinþór Steingrímsson, Antonio Stella, Abishek Stephen, Milan Straka, Emmett Strickland, Jana Strnadová, Alane Suhr, Yogi Lesmana Sulestio, Umut Sulubacak, Shingo Suzuki, Daniel Swanson, Zolt Szántó, Chihiro Taguchi, Dima Taji, Fabio Tamburini, Mary Ann C. Tan, Takaaki Tanaka, Dipta Tanaya, Mirko Tavoni, Samson Tella, Isabelle Tellier, Marinella Testori, Guillaume Thomas, Tarik Emre Tıraş, Sara Tonelli, Liisi Torga, Marsida Toska, Trost Trosterud, Anna Trukhina, René Tšarfaty, Utku Türk, Francis Tyers, Sveinbjörn Þórðarson, Vilhjálmur Þorsteinsson, Sumire Uematsu, Roman Untilov, Zdeňka Urešová, Larraitz Uriá, Hans Uszkoreit, Andrius Utkas, Elena Vagnoni, Sommya Vajjala, Socrates Vak, Rob van der Goot, Martine Vanhove, Daniel van Niekerk, Gertjan van Noord, Viktor Varga, Uliana Vedenina, Giulia Venturi, Eric Villemonte de la Clergerie, Veronika Vincze, Anishka Vissamsetty, Natalia Vlasova, Eleni Vligouridou, Aya Wakasa, Joel C. Wallenberg, Lars Wallin, Abigail Walsh, John Wang, Jonathan North Washington, Maximilian Wendt, Paul Widmer, Shira Wigderson, Sri Hartati Wijono, Vanessa Berwanger Wille, Seyi Williams, Mats Wirén, Christian Wittern, Tsegay Woldemariam, Tak-sum Wong, Alina Wróblewska, Qishen Wu, Mary Yako, Kayo Yamashita, Naoki Yamazaki, Chunxiao Yan, Koichi Yasuoka, Marat M. Yavrumyan, Arife Betül Yenice, Enes Yilandiođlu, Olcay Taner Yildiz, Zhuoran Yu, Arisa Yuliawati, Zdeněk Žabokrtský, Shorouq Zahra, Amir Zeldes, He Zhou, Hanzhi Zhu, Yilun Zhu, Anna Zhuravleva, Rayan Ziane**

(2024): Universal Dependencies 2.14. LINDAT/CLARIAH Digital Repository.




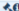




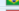
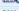

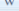













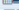








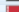

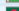




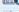









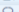



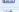

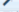








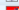


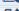
# The UD Repository

UD v2.14:

- 31 language families
- 161 languages
- 283 treebanks
- 616 contributors
- 1.9 million sentences
- 32 million words
- 200 thousand downloads (all versions)

## Current UD Languages

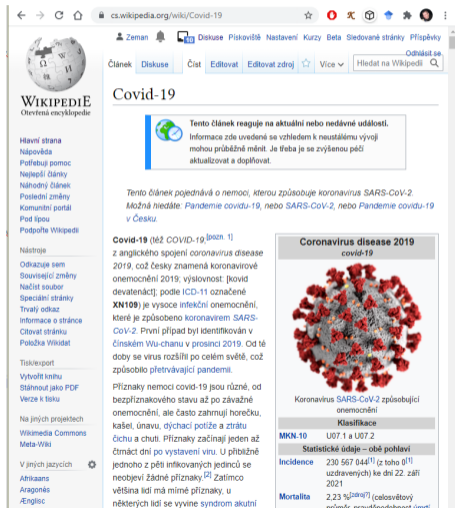
Information about language families (and genera for families with multiple branches) is mostly taken from [WALS Online](#) (IE = Indo-European).

▶		Abaza	1	<1K		Northwest Caucasian
▶		Afrikaans	1	49K		IE, Germanic
▶		Akkadian	2	25K		Afro-Asiatic, Semitic
▶		Akuntsu	1	1K		Tupian, Tupari
▶		Albanian	1	<1K		IE, Albanian
▶		Amharic	1	10K		Afro-Asiatic, Semitic
▶		Ancient Greek	2	416K		IE, Greek
▶		Ancient Hebrew	1	39K		Afro-Asiatic, Semitic
▶		Apurina	1	<1K		Arawakan
▶		Arabic	3	1,042K		Afro-Asiatic, Semitic
▶		Armenian	2	94K		IE, Armenian
▶		Assyrian	1	<1K		Afro-Asiatic, Semitic
▶		Bambara	1	13K		Mande
▶		Basque	1	121K		Basque
▶		Beja	1	<1K		Afro-Asiatic, Cushitic
▶		Belarusian	1	305K		IE, Slavic
▶		Bengali	1	<1K		IE, Indic
▶		Bhojpuri	1	6K		IE, Indic
▶		Bororo	1	<1K		Bororoan
▶		Breton	1	10K		IE, Celtic
▶		Bulgarian	1	156K		IE, Slavic
▶		Buryat	1	10K		Mongolic
▶		Cantonese	1	13K		Sino-Tibetan
▶		Catalan	1	553K		IE, Romance
▶		Cebuano	1	1K		Austronesian, Central Philippine
▶		Chinese	6	287K		Sino-Tibetan
▶		Chukchi	1	6K		Chukotko-Kamchatkan
▶		Classical Chinese	1	433K		Sino-Tibetan
▶		Coptic	1	55K		Afro-Asiatic, Egyptian
▶		Croatian	1	199K		IE, Slavic
▶		Czech	5	2,247K		IE, Slavic
▶		Danish	1	100K		IE, Germanic
▶		Dutch	2	306K		IE, Germanic
▶		English	10	726K		IE, Germanic
▶		Erzya	1	20K		Uralic, Mordvin



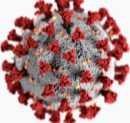
# Linguists Can Parse and Search New Data

<https://lindat.mff.cuni.cz/services/udpipe/>



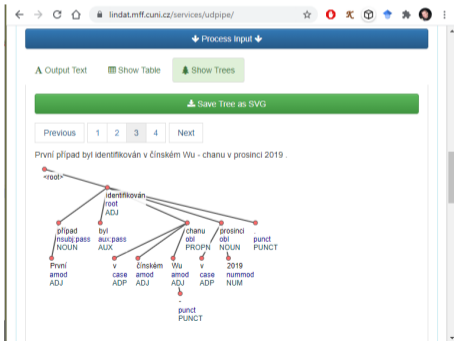
The screenshot shows the Wikipedia page for "Covid-19" in Czech. The page title is "Covid-19". A notice at the top states: "Tento článek reaguje na aktuální nebo nedávné události. Informace zde uvedené se vzhledem k neustálému vývoji mohou průběžně měnit. Je třeba je se zvýšenou péčí aktualizovat a doplňovat." Below this, a summary paragraph reads: "Tento článek pojednává o nemoci, kterou způsobuje koronavirus SARS-CoV-2. Možná hledáte: *Pandemie covidu-19*, nebo *SARS-CoV-2*, nebo *Pandemie covidu-19 v Česku*." The main text begins with "Covid-19 (též COVID-19<sup>[pozn. 1]</sup>) z anglického spojení *coronavirus disease 2019*, což česky znamená koronavirové onemocnění 2019; výslovnost [kvoíd devatenáct], podle ICD-11 označené **XN109**) je vysoce infekční onemocnění, které je způsobeno koronavirem SARS-CoV-2. První případ byl identifikován v čínském Wu-chanu v prosinci 2019. Od té doby se virus rozšířil po celém světě, což způsobilo pletnáváající pandemii. Příznaky nemoci covid-19 jsou různé, od bezpříznakového stavu až po závažné onemocnění, ale často zahrnují horečku, kašel, únavu, dýchací potíže a ztrátu čichu a chuti. Příznaky začínají jeden až čtrnáct dní po vystavení viru. U přibližně jednoho z pěti infikovaných jedinců se neobjeví žádné příznaky.<sup>[2]</sup> Zatímco většina lidí má mírné příznaky, u některých lidí se vyvine syndrom akutní

**Coronavirus disease 2019 covid-19**



Koronavirus SARS-CoV-2 způsobující onemocnění

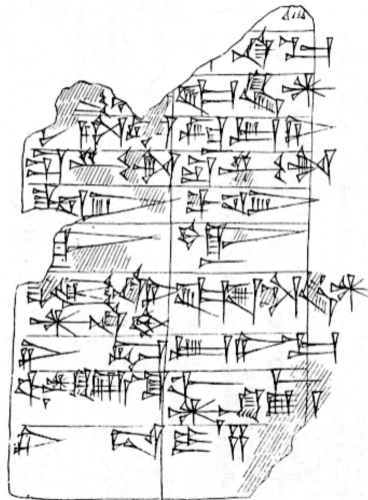
Klasifikace	
<b>MKN-10</b>	U07.1 a U07.2
Statistické údaje – obě pohlaví	
<b>Incidence</b>	230 567 044 <sup>[1]</sup> (z toho 0 <sup>[1]</sup> uzdravených) ke dni 22. září 2021
<b>Mortalita</b>	2,23 % <sup>[2007]</sup> (celosvětový průměr: pravděpodobnost úmrtí



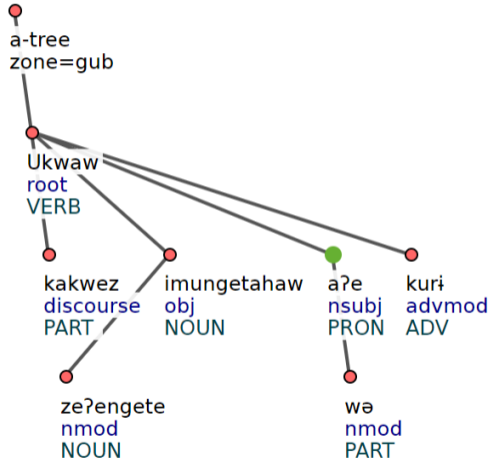
The screenshot shows the UDpipe web interface. At the top, there is a "Process Input" button. Below it, there are buttons for "Output Text", "Show Table", and "Show Trees". A "Save Tree as SVG" button is also present. A navigation bar shows "Previous", "1", "2", "3", "4", and "Next". The main content area displays the sentence "První případ byl identifikován v čínském Wu-chanu v prosinci 2019." and a corresponding parse tree. The tree root is "identifikován", which branches into "případ", "byl", "chanu", and "prosinci". "případ" branches into "případ", "nsubj", and "pass". "byl" branches into "aux" and "pass". "chanu" branches into "obl" and "PROPN". "prosinci" branches into "obl" and "NOUN". "2019" branches into "v", "case", "ADP", "čínském", "amod", "ADJ", "Wu", "amod", "ADJ", "v", "case", "ADP", "2019", "nummod", and "NUM". "punct" branches into "punct" and "PUNCT".

# Historical Linguistics, Classical Languages

- Old Turkish
- Classical Chinese
- Sanskrit
- Akkadian
- Ancient Hebrew
- Coptic
- Ancient Greek
- Latin
- Old French
- Old Irish
- Gothic
- Old Church Slavonic



# Documentation of Endangered Languages



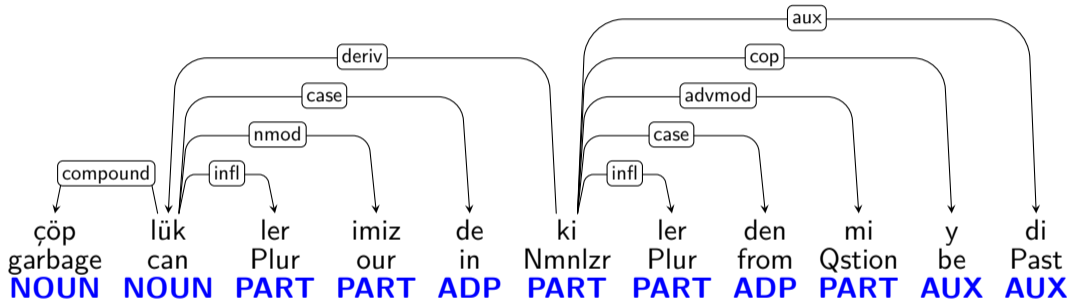




**What's Next?**



# Subword Relations: Agglutinative Languages



“was it from those that were in our garbage cans?”

# Universal Semantics?

Kira Droганova



Universal Semantic Roles

Dima Taji



Coreference in Deep UD

Federica Gamba



Uniform Meaning Representation  
for Latin

Minoo Nassajian



Uniform Meaning Representation  
for Persian

# Universal Semantics?

Kira Droганova



Universal Semantic Roles

Dima Taji



Coreference in Deep UD

Federica Gamba



Uniform Meaning Representation  
for Latin

Minoo Nassajian



Uniform Meaning Representation  
for Persian

Diego Alves (defended 2023)



Computational Typological Analysis  
of Syntactic Structures

# Parsing Accuracy (for Discussion)

TREEBANK	MODEL	UPOS	FEATS	LEM	UAS	LAS
Czech PDT (cs_pdt)	UDPipe	99.18	97.23	<b>99.02</b>	93.33	91.31
	Lang	99.18	96.87	98.72	94.35	92.41
	UDify	99.18	96.85	98.56	94.73	92.88
	UDify+Lang	<b>99.24</b>	<b>97.44</b>	98.93	<b>95.07</b>	<b>93.38</b>
German GSD (de_gsd)	UDPipe	94.48	90.68	<b>96.80</b>	85.53	81.07
	Lang	94.77	91.73	96.34	87.54	83.39
	UDify	94.55	90.65	94.82	87.81	83.59
	UDify+Lang	<b>95.29</b>	<b>91.94</b>	96.74	<b>88.11</b>	<b>84.13</b>
English EWT (en_ewt)	UDPipe	96.29	97.10	<b>98.25</b>	89.63	86.97
	Lang	<b>96.82</b>	<b>97.27</b>	97.97	<b>91.70</b>	<b>89.38</b>
	UDify	96.21	96.17	97.35	90.96	88.50
	UDify+Lang	96.57	96.96	97.90	91.55	89.06
Spanish AnCora (es_ancora)	UDPipe	<b>98.91</b>	<b>98.49</b>	<b>99.17</b>	92.34	90.26
	Lang	98.60	98.14	98.52	92.82	90.52
	UDify	98.53	97.84	98.09	92.99	90.50
	UDify+Lang	98.68	98.25	98.68	<b>93.35</b>	<b>91.28</b>
French GSD (fr_gsd)	UDPipe	97.63	<b>97.13</b>	<b>98.35</b>	90.65	88.06
	Lang	<b>98.05</b>	96.26	97.96	92.77	90.61
	UDify	97.83	96.59	97.48	<b>93.60</b>	<b>91.45</b>
	UDify+Lang	97.96	96.73	98.17	93.56	91.45
Russian SynTagRus (ru_syntagrus)	UDPipe	<b>99.12</b>	<b>97.57</b>	<b>98.53</b>	93.80	92.32
	Lang	98.90	96.58	95.16	94.40	92.72
	UDify	98.97	96.35	94.43	94.83	93.13
	UDify+Lang	99.08	97.22	96.58	<b>95.13</b>	<b>93.70</b>
Belarusian HSE (be_hse)	UDPipe	93.63	73.30	87.34	78.58	72.72
	Lang	95.88	76.12	84.52	83.94	79.02
	UDify	<b>97.54</b>	<b>89.36</b>	85.46	<b>91.82</b>	<b>87.19</b>
	UDify+Lang	97.25	85.02	<b>88.71</b>	90.67	86.98
Buryat BDT (bxr_bdt)	UDPipe	40.34	32.40	58.17	32.60	18.83
	Lang	52.54	37.03	54.64	29.63	15.82
	UDify	<b>61.73</b>	<b>47.86</b>	<b>61.06</b>	<b>48.43</b>	<b>26.28</b>
	UDify+Lang	61.73	42.79	58.20	33.06	18.65
Upper Sorbian UFAL (hsb_ufal)	UDPipe	62.93	41.10	68.68	45.58	34.54
	Lang	73.70	46.28	58.02	39.02	28.70
	UDify	84.87	48.63	<b>72.73</b>	<b>71.55</b>	<b>62.82</b>
	UDify+Lang	<b>87.58</b>	<b>53.19</b>	71.88	71.40	60.65
Kazakh KTB (kk_ktb)	UDPipe	55.84	40.40	63.96	53.30	33.38
	Lang	73.52	46.60	57.84	50.38	32.61
	UDify	<b>85.59</b>	<b>65.14</b>	<b>77.40</b>	<b>74.77</b>	<b>63.66</b>
	UDify+Lang	81.32	60.50	67.30	69.16	53.14
Lithuanian HSE (lt_hse)	UDPipe	81.70	60.47	<b>76.89</b>	51.98	42.17
	Lang	83.40	54.34	58.77	51.23	38.96
	UDify	<b>90.47</b>	<b>68.96</b>	67.83	<b>79.06</b>	<b>69.34</b>
	UDify+Lang	84.53	56.98	58.21	58.40	39.91

TREEBANK		UPOS	FEATS	LEM	UAS	LAS
<b>Breton KEB</b>	<b>br_keb</b>	63.67	46.75	53.15	63.97	40.19
<b>Tagalog TRG</b>	<b>tl_trg</b>	61.64	35.27	75.00	64.73	39.38
Faroese OFT	fo_oft	77.86	35.71	53.82	69.28	61.03
Naija NSC	pcm_nsc	56.59	52.75	97.52	47.13	33.43
Sanskrit UFAL	sa_ufal	40.21	18.45	37.60	41.73	19.80

Table 4: Test set results for zero-shot learning, i.e., no UD training annotations available. Languages that are pretrained with BERT are bolded.

Dan Kondratyuk, Milan Straka (2019): 75 Languages, 1 Model: Parsing Universal Dependencies Universally. In: Proceedings of EMNLP, Hong Kong