

# زمانه وانیی کۆمپیوتەری

بابەتی 2: سەرچاوه زمانه وانیه کان

سینا ئەحمەدی

[sinaahmadi.github.io/KurdishCL](https://sinaahmadi.github.io/KurdishCL)



1 دهسپیک

2 پیدراوهی زمان

3 پیدراوهکان بهپیی فۆرماتهکان

4 بهکارهینانی سه رچاوه زمانیهکان

5 بهشی پراکتیک

# دەسپیک





© Sina Ahmadi



# سەرچاوه زمانه وانیه کان بو زمانه وانیه کو مپیوتهری

- پیدراوهی زمانه وانیه ← یه کهم ههنگاو بو پرۆسهس کردنی زمان
- به شیککی زۆری زمانه کانیه دنیا سەرچاوهی زمان و زمانه وانیه کان بو کهمه [1]
- زمانی کوردی زمانیککی کهم سەرچاوه کهیه (less-resourced)

زمانه سەرچاوه دهوله مهنده کان  
سه دان ملیۆن به لگه ی سهرهیل و پیدراوه و زانیاری  
وهک ئینگلیسی

زمانه مامناوه ندییه کان  
ملیۆنان به لگه ی سهرهیل و پیدراوه و  
زانیاری زۆربه ی زمانه ئه ورووپیه کان

زمانه هه ژاره کان  
بی یان کهم/سەرچاوه  
زۆربه ی زمانه کانیه دنیا



# به کارهییانی سه رچاوه کانی زمان

- ناسینه وهی به شه کانی وشه (part-of-speech tagging)
- راست کردنه وهی هه لهی رینووس (spelling error correction)
- لیك دانه وهی واتای وشه کان: «داریکی به دهسته وه بوو.» (word-sense disambiguation)
- ناوانی رۆله ماناداره کان (semantic role labeling)
- په یوهندی دانی چه مکه کان له رسته دا (entity linking)
- وپرای پیشکه وتنه کان له بواری ژیری دهستکردا، سه رچاوه کانی زمان هیشتاش رۆلیکی گرینگیان له تهنۆلۆژیای زماندا

پا ریس / [https://ckb.wikipedia.org/wiki/پا\\_ریس](https://ckb.wikipedia.org/wiki/پا_ریس)

پاریس پایتهختی ولاتی فهره نسا په.

فهره نسا / [https://ckb.wikipedia.org/wiki/فهره\\_نسا](https://ckb.wikipedia.org/wiki/فهره_نسا)

## سەرچاوهی زمانه وانیه/زمانی

سەرچاوهی زمانه وانیه بریتیه له پیکهاتهی بهرهمه زمانه وانیه و زمانیهکان که بو چی کردن، باشر کردن یان هه لسه نگاندهی ته کنولۆژیای زمان و تووژینه وهی سهر زمان به کار دههیندرین.

- 1 پیدراوهی زمانه وانیه: ئەو سەرچاوه یانیهی زمانیک له ریگای پیدراوه وه ده ناسین، وهک فهرهنگ، لیستی وشه یان ته نانهت دهنگی شیوه زاریک
- 2 ئەو ئامیرانهی که بو به کارهینان، نیشان دان یان دهست تی وهردانی پیدراوهکان پیوستن
- 3 میتاداتا و وشه سازیهکان که بو مامه له کردن له گه ل پاراستنی سەرچاوهکان به کار دین، به تایه تی بهرده وامی (sustainability)، بهرده ستیوون و ئه رشیفکردنی درێژخایه ن

# پیداوہی زمان





هه‌ندیك له گرنه‌تیرین جوهره‌كانی پیداوه زمانه‌وانیه‌كان ئه‌مانه‌ن:

- 1 وه‌سفی زمان وه‌ك ریزمان (grammar) و مودیل کردنی فۆرمال
- 2 کۆرپس (corpora / corpus به‌ کۆ): پیکهاته‌یه‌ک له ده‌ق یان ده‌نگ به‌ شیوه‌ی دیجیتال
- 3 سه‌رچاوه‌ی واتاناسی-و‌شه‌ناسی (lexical-semantic): ریک خستی سه‌رچاوه‌ی زمان به‌ پیی بو‌چوون یان تیورییه‌ک له واتاناسی یان و‌شه‌ناسیدا
- 4 ئۆنتۆلۆژی (ontology)، زاراوه‌ناسی (terminology) و فه‌ره‌ه‌نگۆک (glossary)
- 5 مۆدیلی زمان (language model) که له سه‌ر کۆرپسه‌ گه‌وره‌کاندا فی‌رده‌کری
- 6 گرافی زانیاری (knowledge graph) و بنکه‌ی زانیاری (knowledge base)

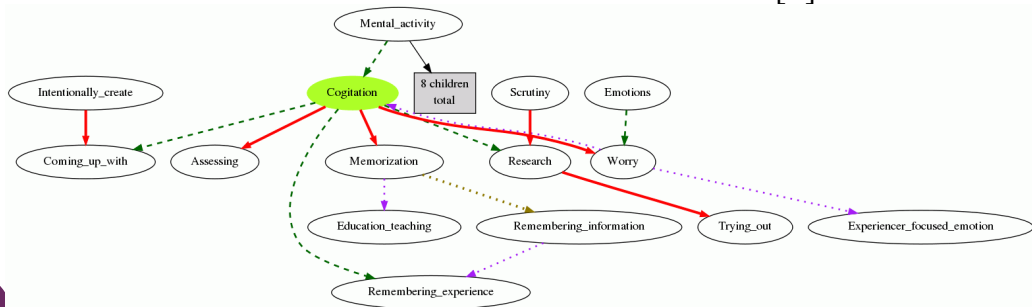
- ریزمانی زمان وهک ریزمانی کوردیی تاکستۆن [2]
- فورمالیسم و یاساکان

$$\left[ \begin{array}{l} /XY(\mathbf{in})/_{V} \\ x (= \text{چاوغ}) \\ \text{کردن} (kirdin) \end{array} \right] \leftrightarrow \left[ \begin{array}{l} /XY/_{V} \\ x \text{ رهگی رابردووی} \\ \text{کرد} (kird) \end{array} \right] \leftrightarrow \left[ \begin{array}{l} /XZ/_{V} \\ x \text{ رهگی ئیستای} \\ \text{که} (ke) \end{array} \right]$$

که تییدا  $X$  رهگی کرداره که یه و  $Y$  و  $Z$  یش ههردووکیان پاشگری دارپژانن.

# پیداوهی زمان: سه‌رچاوهی واتاناسی-وشه‌ناسی

- له سه‌ر بنه‌مایه‌کی واتاناسی یان وشه‌ناسیدا به شیوه‌یه‌کی وه‌سفکه‌ر و پیکهاته‌دار ری‌ک خراون.
- هه‌ندی‌ک له گرینگترینی ئه‌و سه‌رچاوه‌یانه ئه‌مانه‌ن:
- فره‌هنگی وشه وه‌ک هه‌نبانه بو‌رینه یان [3] Ferhenga Birûskî
- وو‌ردنی‌ت (WordNet) [4] (به‌سته‌ر) وه‌ک کوردنی‌ت [5]
- [6] FrameNet



- زاراوه یان زار (terminology) به وشه‌یه‌ک ده‌کوتری که تایبته به بواریکه وه‌ک زاراوه‌ی پزیشکی یان زاراوه‌ی ته‌کنۆلۆژیا
- سه‌رجه‌م زاره‌کان زۆر جار پییان ده‌کوتری فه‌ره‌ه‌نگۆک (glossary)
- هاوشیوه‌ی فه‌ره‌ه‌نگ له سه‌روه‌ه‌ پیک هاتووه و هه‌ندی‌ک جار زانیاریی زمانه‌وانیسی له‌گه‌ل‌دایه.
- تی‌رمینۆلۆژیی یه‌کیه‌تیی ئه‌ورووپا (به‌سته‌ر)
- ئامانجی زاراوه‌ریک خستن و پۆلاندنی زانیارییه‌ بو‌ سازکردنی پۆلینه‌ناسی (taxonomy)
- له‌ ته‌کنۆلۆژیای زمان زۆر به‌سوودن بو‌ زۆرکردنه‌وه‌ی ئاستی ناسینه‌وه‌ی وشه‌کان
- هه‌ندی‌ک بنکه‌ی زاری کوردیی به‌سوود
- وشه‌نامه‌ی هوژین: <https://hojan.org/dic/>
- یاپراخ: <https://dict.linux.krd/>

## کۆرپس

کۆمهلیک پیدراوهی زمانه وانیهکان که به شیوهی دهق یان دهنگ کۆده کریتته وه و وهکوو خالی دهستییک بۆ سهلماندنی گریمانیهکان سهبارته به زمانیک به کار دههیندری.

- زۆربهی کۆرپسهکان لانی کهم چهند ملیۆن وشه یان تیدایه.
- کۆرپس ته نیا چیژیک له زمانه کهیه، نهک زمانه که بۆ خوی.
- به رای چۆمسکی، کۆرپس قهت ناتوانی هه موو دیاردهکانی زمانیک بینیتته سهلماندن [7].
- بۆ ههر بابتهیک، وهک زانست یان رامیاری، و ههر پیناسهیهکی زمان، وهک بنزاراوه یان ههرهتی قسه کردن، دهتوانین کۆرپسمان هه بی.
- کۆرپس ده بی نمونه، هاوسهنگ، نوینه ر و سروشتی بی [8].
- مژاری سه رهکیی «زمانه وانیهی کۆرپس» (corpus linguistics) تیگه یشتن له کۆرپسه.
- کۆرپسهکان به بهربلاوی له پرۆسهس کردنی زمانیشدا به کار دههیندرین.



# سەرچاوه زمانه وانیه کان: کۆرپس (دریژه)

نیوی کۆرپسه که	سەرچاوه	بنزاراوه	قهباره (وشه #)
کۆرپسی پیوان [9]	هه و آل	سۆرانی	18M
		کورمانجی	4M
کۆرپسی دهقی فۆلکلۆر [10]	گۆرانی، بهند و بهیت	سۆرانی	49K
کۆرپسی ئاسۆسافت [11]	هه و آل	سۆرانی	80M
کۆرپسی زازا-گۆرانی [12]	هه و آل	زازاکی	10M
		گۆرانی	2M
کۆرپسی پهروه ده [13]	کتیبی پهروه ده	سۆرانی	0.6M
کۆرپسی مه نچستر [14]	کۆکردنه وهی مهیدانی	سۆرانی	
		کورمانجی	

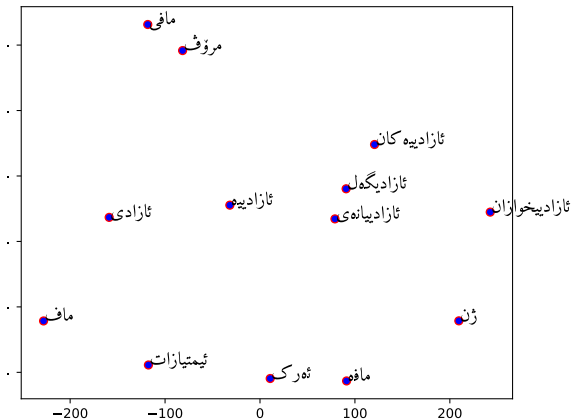
لیستی هه موو کۆرپسه کان بۆ کوردی لیڤه دایه:

<https://github.com/sinaahmadi/awesome-kurdish>



# سەرچاوه زمانه وانیه کان: مۆدیلی زمان

- مۆدیل کردنی ئاماریی زمان که تایبه تمه نندییه کانی زمانیک «فیر دهبی» رۆلی مۆرفۆلۆژی و سینتاکس، په یوه نندی نیوان وشه کان و هتد
- به سوود له مۆدیل کردنی زۆر دیاردهی زمان و بهربلاوه له ئه ورۆی ته کۆنۆلۆژیای زماندا ئه و وشانهی له باری واتاوه له وشه «ماف» نزیکیه تیان ههیه ↓





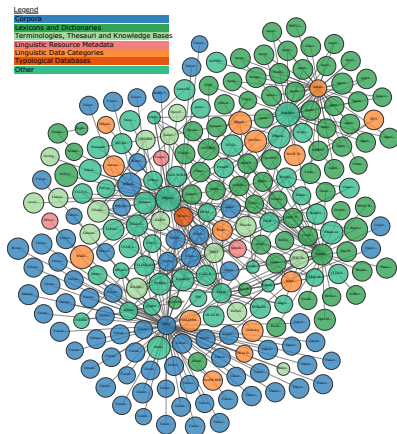
## ئۆنتۆلۆژی (ontology)

به پیکهاتهی ئه و یاسا و زانیاریانی ده کوترئ که بو ناساندنی چه مکه کانی بواریک و په یوهندی نیوانیان به شیوهی فۆرمال به کارده هیندرین.

- ئۆنتۆلۆژییه کان به توندی فهرمین و به زمانیکی مهنتیقی دیاری کراون
- «مانگ» یهک وشهیه له فهرهنگیکدا، به لام دهتوانئ به چه ند شیوه پیناسه بکری:
- مانگی ههتاوی و مانگی هیجری
- سەرچاوه زمانه وانیه کان به پئی ئاستی فۆرمال بوونیان له ئۆنتۆلۆژی ده چن.
- ئۆنتۆلۆژی له بواری مهنتیق و کۆمپیوتهردا جیاوازی ههیه. چاو بکه له [15]
- ئەندازیاریی ئۆنتۆلۆژی بابهتی سهرهکیی بواری «وییی سیمانتیک» (semantic web) ه.

# سەرچاوه زمانه وانیه کان: ئۆنتۆلۆژی (دریژه)

- وشه سازی و مۆدیلیکی پیدراوه بۆ سه رچاوه کانی وشه، وه ک فهرههنگه کان [16]
- مۆدیلیکی پیدراوه بۆ نیشان دانی په یوهندی واتایی و بنه مالهی نیوان وشه کان [17]
- هه وری سه رچاوه په یوهندیاره کان (به سته ر)



# سەرچاوه زمانه وانیهکان: گرافی زانیاری و بنکهی زانیاری

- نواندنی زانیاری و پیدراوهکان به شیوازی گراف (graph)
- پشت بهستن به بناغهی ویی سیمانتیک و مۆدیل کردنی پیدراوهکان
- لیک کارکردنی یه کجار زۆری زانیاری له سهر وییدا ← چاره سه ریکی باش بو ئیستای «پیدراوه په یوه ندیدار» هکان (linked data)
- به سوودیشه بو زانیاریی زمانه وانیه
- ئەژماری وشه ئینگلیسییهکان به پیی یه کهم پیتیان (بهسته ر)
- کۆمه لی گرافهکانی زانیاری ده بی به بنکه ی زانیاری (knowledge base)
- لیستی پشکنه رهکانی ئاسمان به یه ک کرته (بهسته ر)
- بنکه به سوودهکان:
- Wikidata (سهیری کوردی بکه: بهسته ر)
- DBpedia (سهیری کوردی بکه: بهسته ر)
- ConceptNet (سهیری کوردی بکه: بهسته ر)

# پیداوه‌کان به‌پیی فۆرماتەکان



SEROKÊ IRAQÊ ŞPÊWAZÎYA ENDAMÊ SEROKATÎYA ENCÛMENA  
 BILINDA ISLAMÎYA IRAQÊ KIR [10:30] 2010/Oct  
 /25 PNA Serokê komara Iraqê Celal Talebanî,  
 şpewaziya endamê serokatîya Encûmena  
 Bilinda Islamiya Iraqê Adil Abdulmehdî û  
 şandeya pê re kir. Herdu aliyan daxwaz  
 kirin, ku divê hukûmeta nû ya Iraqê zû pêk  
 bihê da ku pirs û şarêyên hevvelatîyan  
 çareser bibin û kar û barên gel bi rê ve  
 biçin. Di evê didarê de ya ku doh êvarê hat  
 sazdan, herdu aliyan bir û nerîn di bareya  
 pêkanîna hukûmeta nû ya Iraqê de bi hev  
 guherandin û behsa encamên diyaloga navbera  
 aliyên siyasî kirin.

- ناسراوترین پیناسه‌ی فۆرمات: .txt
- ✓ ساکار، خیرا و ئاسان
- × پیدراوه‌ی به‌بێ دارپێژران یان وه‌سف





```

{
  "firstName": "John",
  "lastName": "Smith",
  "isAlive": true,
  "age": 27,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021-3100"
  },
  "phoneNumbers": [
    {
      "type": "home",
      "number": "212 555-1234"
    }
  ]
}

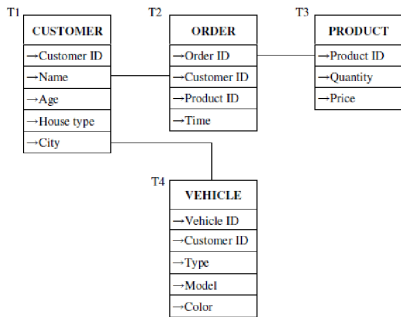
```

## JavaScript Object Notation (.json)

- ✓ خۆوهسفکەر و دارپژراو
- ✓ باش بۆ گواستنهوه و دهست وهردان له پیدراوه
- × کۆل و کورتتر له XML
- JSON و XML بربرهه
- هه لگواستنهوهی پیدراوهن له سهه و ییدا.

## بنکه دراوه (database)

سیستهمی نهرمه کالاً بۆ پیناسه کردن، چۆ کردن، پاراستن و دهست وهردان له پیدراوه به کار دههیندری [18].



- یه کجار خیرا و جی متمانهیه
- ✓ زۆر بهربلاو له ئه ورۆی ویدا
- × خووه سفکه ر نین و پیکهاته کان به پیی
- مهنتقی په یوه ندییه وه له خشته دا ریک ده خرین
- به نیو بانگترین بنکه دراوه کان  
په یوه ندیدارن (relational)  
PostgreSQL، MySQL  
و Oracle Database



```

:lex_kmr_0278943883 a ontolex:
  LexicalEntry, ontolex:Word ;
ontolex:canonicalForm :
  form_kmr_0278943883 ;
rdfs:label "bend"@kmr-latn .

:form_kmr_0278943883 a ontolex:Form ;
dct:language <www.lexvo.org/page/
  iso639-3/kmr> ;
ontolex:writtenRep "bend"@kmr-latn ;
lexinfo:partOfSpeech lexinfo:noun ;
lexinfo:gender lexinfo:feminine ;
lexinfo:number lexinfo:singular ;
ontolex:sense :kmr_9120343779_sense.

```

- چوارچێۆه‌ی وه‌سفکردنی سه‌رچاوه‌ (Resource) (Description Framework (.ttl))
- ✓ خۆوه‌سفکه‌ر و دارپێژاو
- ✓ زمانی ویب و پیدراوه‌ په‌یوه‌ندی‌داره‌کان
- ✓ خیرا و به‌ربلاو
- × هه‌ندی‌ک جار ئه‌ندازه‌ی پیدراوه‌کان زۆر ده‌کاته‌وه‌
- ✓ پێشنیاری W3C (به‌سته‌ر)
- لێره‌دا ده‌توانن له‌گه‌ڵ RDF دا زۆرتر ئاشنا بن: به‌سته‌ر.

# به کارهینانی سه رچاوه زمانیه کان



- هاوشیوهی بنه ماکانی بهرپوه بردنی پیدراوه زانستییه کان [19]، سهرچاوه یه کی زمانه وانی باش ده بی:
  - **F**indable: بدۆزریته وه
  - **A**ccessible: بهرده ست بی
  - **I**nteroperable: کاری له گه ل بکری
  - **R**eusable: به که لک بیته وه
- 
- گرینگایه تی پیدراوه زمانه وانیه په یوه ندیداره کان (linguistic linked data) [20]
  - پیدراوه بناغه ی تیگه یشتن له زمان و گرینگایه تیشیان له زانسته مروییه دیجیتالییه کاند (digital humanities) هه یه.

## بەشی پراکتیک



ئەم بەشە لە خولەكە، بەشیکى پراکتیکیشی هەیه که خویندکار دەتوانی لە رینگای مالپەری  
خولەكەوه شوینی بکا.



- [1] Unsupervised cross-lingual representation learning.  
<https://ruder.io/unsupervised-cross-lingual-learning>.  
 Accessed: 2022-07-30.
- [2] Wheeler M. Thackston.  
*Sorani Kurdish—A Reference Grammar with Selected Readings*.  
 Harvard University, 2006.
- [3] Michael L Chyet and Martin Schwartz.  
*Kurdish-English Dictionary (Kurmanji)*.  
 Yale University Press, 2003.
- [4] George Miller, Christiane Fellbaum, Judy Kegl, and Katherine Miller.  
 Wordnet: An electronic lexical reference system based on theories of lexical memory.  
*Revue quebecoise de linguistique*, 17(2):181–212, 1988.
- [5] Purya Aliabadi, Mohammad Sina Ahmadi, Shahin Salavati, and Kyumars Sheykh Esmaili.  
 Towards building Kurdnet, the Kurdish Wordnet.  
 In *Proceedings of the Seventh Global Wordnet Conference*, pages 1–6, 2014.
- [6] Collin F Baker, Charles J Fillmore, and John B Lowe.  
 The Berkeley Framenet project.  
 In *Proceedings of the 17th international conference on Computational linguistics-Volume 1*, pages 86–90.  
 Association for Computational Linguistics, 1998.

- [7] Jacqueline Léon.  
Claimed and unclaimed sources of corpus linguistics.  
*Henry Sweet Society for the History of Linguistic Ideas Bulletin*, 44(1):36–50, 2005.
- [8] Guillaume Desagulier, Guillaume Desagulier, and Amboy.  
*Corpus linguistics and statistics with R*.  
Springer, 2017.
- [9] Kyumars Sheykh Esmaili, Donya Eliassi, Shahin Salavati, Purya Aliabadi, Asrin Mohammadi, Somayeh Yosefi, and Shownem Hakimi.  
Building a test collection for Sorani Kurdish.  
In *2013 ACS International Conference on Computer Systems and Applications (AICCSA)*, pages 1–7. IEEE, 2013.
- [10] Sina Ahmadi, Hossein Hassani, and Kamaladdin Abedi.  
A Corpus of the Sorani Kurdish Folkloric Lyrics.  
In *Proceedings of the 1st Joint Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL) Workshop at the 12th International Conference on Language Resources and Evaluation (LREC)*, Marseille, France, 2020.
- [11] Hadi Veisi, Mohammad MohammadAmini, and Hawre Hosseini.  
Toward Kurdish language processing: Experiments in collecting and processing the AsoSoft text corpus.  
*Digital Scholarship in the Humanities*, 35(1):176–193, 2020.



- [12] Sina Ahmadi.  
Building a corpus for the zaza–gorani language family.  
*In Proceedings of the Seventh Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2020)*, 2020.
- [13] Roshna Omer Abdulrahman, Hossein Hassani, and Sina Ahmadi.  
A Rule-Based Kurdish Text Transliteration System.  
2019.
- [14] Yaron Matras.  
Revisiting Kurdish dialect geography: findings from the Manchester Database.  
*Current issues in Kurdish linguistics*, 1:225, 2019.
- [15] Nicola Guarino and Pierdaniele Giaretta.  
Ontologies and knowledge bases.  
*Towards very large knowledge bases*, pages 1–2, 1995.
- [16] John P McCrae, Julia Bosque-Gil, Jorge Gracia, Paul Buitelaar, and Philipp Cimiano.  
The ontalex-lemon model: development and applications.  
*In Proceedings of eLex 2017 conference*, pages 19–21, 2017.
- [17] Mustafa Jarrar.  
The Arabic ontology—an Arabic wordnet with ontologically clean content.  
*Applied ontology*, (Preprint):1–26, 2021.





- [18] Thomas M Connolly and Carolyn E Begg.  
*Database systems: a practical approach to design, implementation, and management.*  
 Pearson Education, 2005.
- [19] Mark D Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E Bourne, et al.  
 The fair guiding principles for scientific data management and stewardship.  
*Scientific data*, 3(1):1–9, 2016.
- [20] Fahad Khan, Christian Chiarcos, Thierry Declerck, Daniela Gifu, Elena González-Blanco García, Jorge Gracia, Max Ionov, Penny Labropoulou, Francesco Mambrini, John McCrae, Émilie Pagé-Perron, Marco Passarotti, Salvador Ros, and Ciprian-Octavian Truica.  
 When Linguistics Meets Web Technologies. Recent advances in Modelling Linguistic Linked Open Data.  
<http://www.semantic-web-journal.net/content/when-linguistics-meets-web-technologies-recent-advances-modelling-linguistic-linked-open-0>, 2021.
- [21] Steven Bird and Gary Simons.  
 Extending dublin core metadata to support the description and discovery of language resources.  
*Computers and the Humanities*, 37(4):375–388, 2003.

