Ольга Ляшевская

Цифровые ресурсы для лингвистических исследований

Лекция 2



Курс «Лингвистические данные» • НИУ ВШЭ, ФиКЛ, 1 курс бакалавриата

Форматы документации (лингвистических) данных

Кодировки текстов



Форматы документации (лингвистических) данных Unicode (UTF-8, UTF-16), ASCII, Windows-1251

Кодировки текстов

HTML, YAML, XML, JSON...



Форматы документации

Unicode (UTF-8, UTF-16), ASCII, Windows-1251

Кодировки текстов

HTML, YAML, XML, JSON...

1DF00

1DFFF

	1DF0	1DF1	1DF2	***	1DFF
0	fij	K	dţ		
1	g	ŀ	dţ		
2	Ð	dз"	tł		
3	k	ţ	번		
4	Ł	ŋ,	t O		
5	ß	ત્ર	rd		
6	K	£	ŀ		
7	ũ	tſ,	n		
8	1	3₀	r		
9	f	ф	'n		
Α	Į	į	t		
В	€	b	ďз		
С	£	Ą	ф		
D	J	ત			
E	7	S			
F	G	ďð			

Сериализация

- представление информации в виде более низкоуровневой структуры
 - формат хранения
 - битовая последовательность



Инструменты для разметки

Разметка (аннотация) – добавление информации к первоисточнику

делают эксперты (в том числе студенты и волонтеры) и машины (в том числе нейросетевые модели и скрипты, основанные на правилах)

связана с интерпретацией (анализом) данных



Инструменты для разметки: примеры

Разметка текстов

Praat – разметка звука

ELAN – для аудио и видео

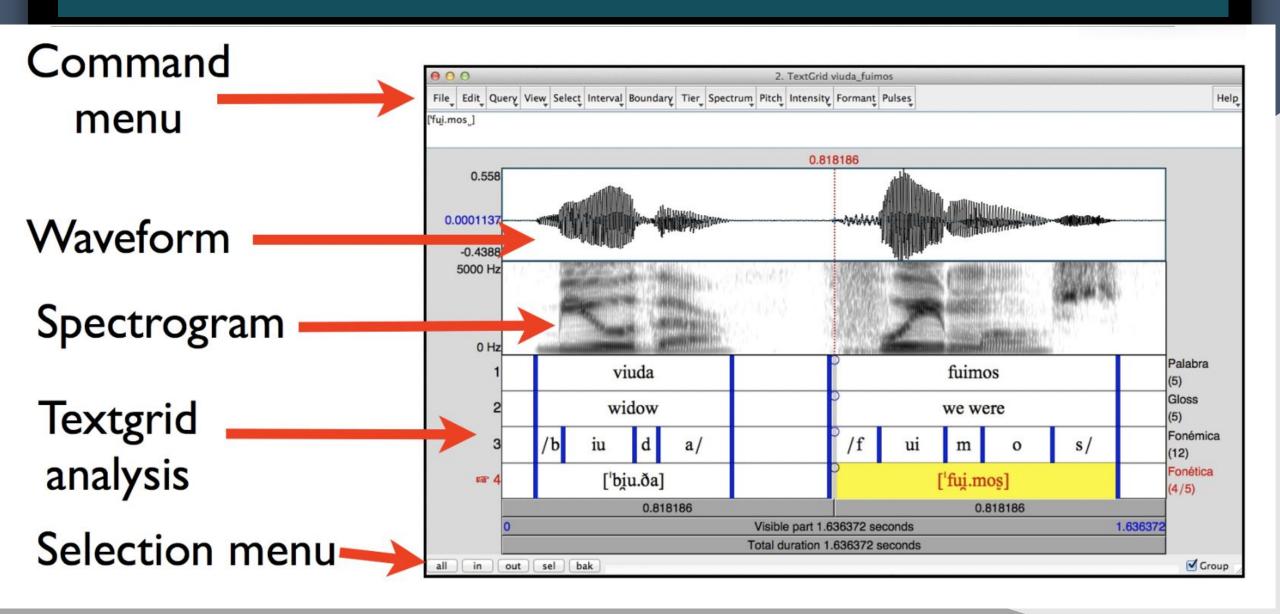
Inception, GATE – для текстов

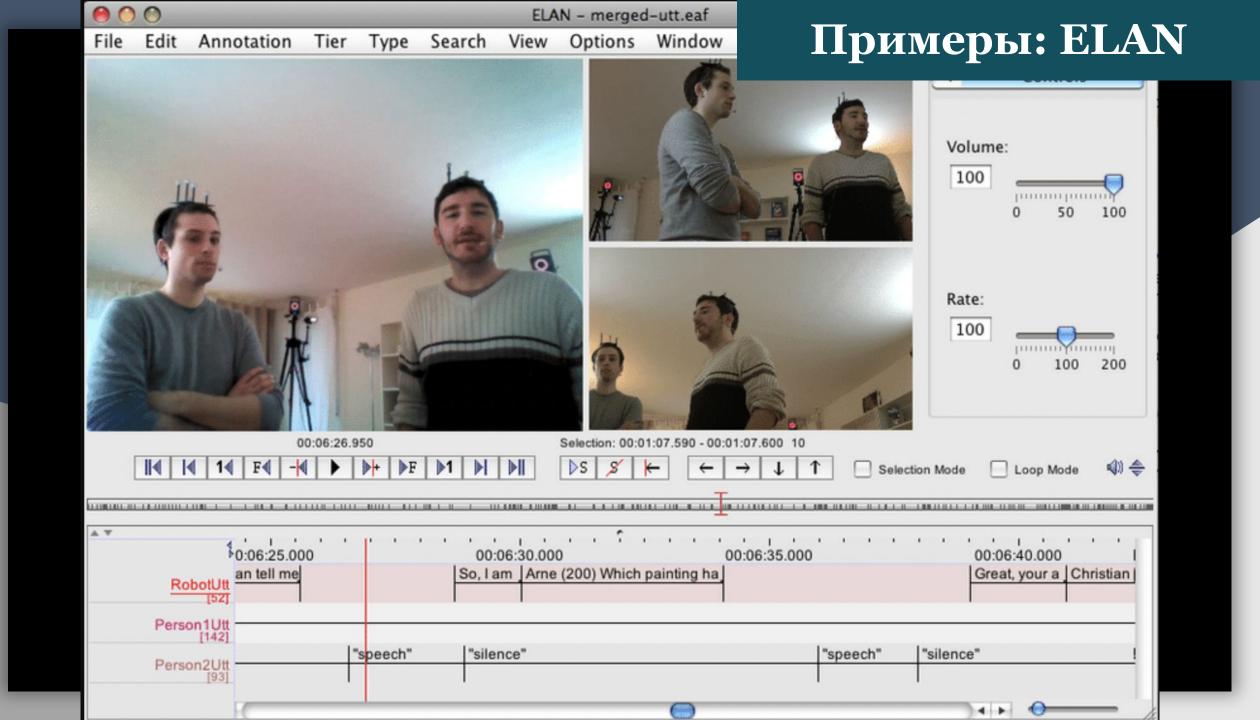
(адаптированы для компьютерно-лингвистических задач)

Fieldworks – для полевых лингвистов



Примеры: Praat





Примеры: ELAN N 3.9.1 File Edit Annotation Tier Type Search View Options Window Help Elan - pear story.eaf Grid Text Subtitles Audio Recognizer Video Recognizer Metadata Controls Recognizer: Tag vowels (volume peaks of voiced timespans) File(s): Pear.wav Selections **Parameters** Add Remove Pitch ceiling [Hz], eg. 300 male / 550 female (150.0 - 1500.0) Add Tier... Intensity change [dB] to start/end a peak (0.5 - 5.0) 2.0 Minimum amplitude (0..1) for pitch analysis (0.01 - 0.3) □ □ □ E Segmentations Progress Create Tier(s)... Report... 00:00:16.840 Selection: 00:00:00:00.000 - 00:00:00.000 0 DS S ← → ↓ ↑ Selection Mode Loop Mode 78.005 32.598 340.3425 00:00:15.000 00:00:16.000 00:00:12.000 00:00:13.000 00:00:14.000 00:00:09.000 00:00:10.000 00:00:11.000 00:00:12.000 00:00:13.000 00:00:14.000 00:00:15.000 00:00:16.000 and he starts picking pears off the tree and so he climbs up a tree and he starts with the ladder and he puts the pears into an apron ause Transcri non-motion non-motion sture # is Hand

Инструменты для разметки: примеры

Разметка текстов

Praat – разметка звука
ELAN – для аудио и видео
Inception, GATE – для текстов
Fieldworks – для полевых

Разработка словарей

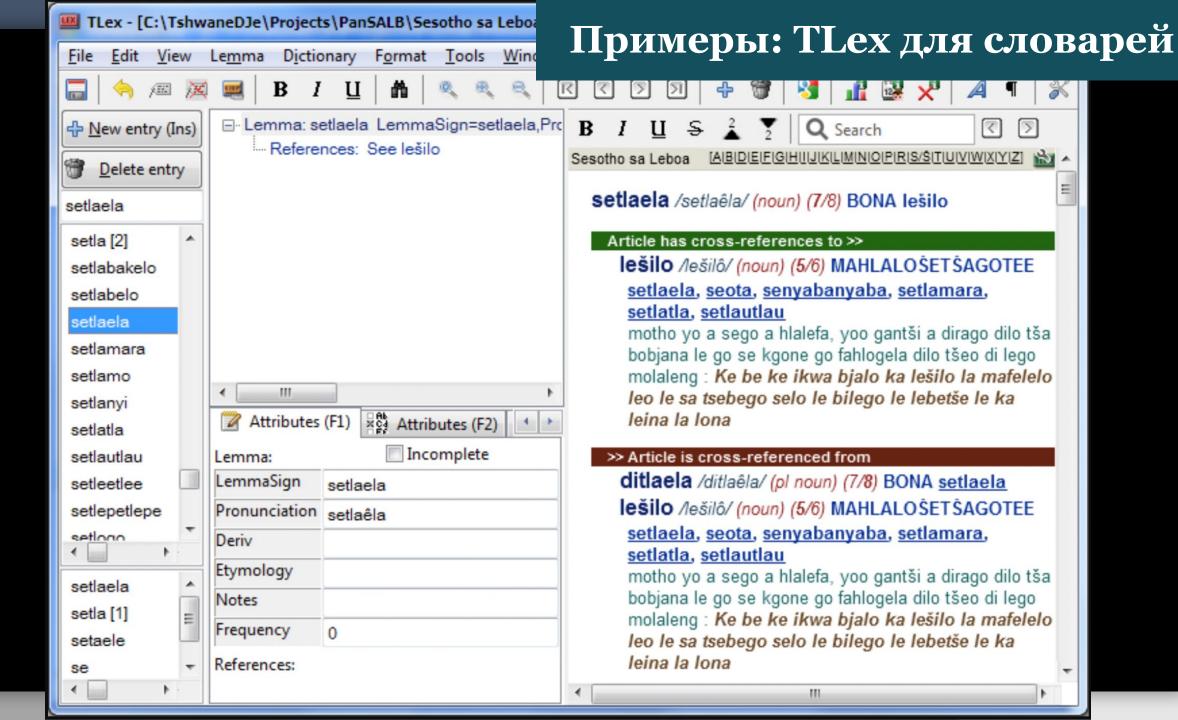
iLex – разработка

<u>TLex</u> – разработка с AI-

поддержкой

<u>Lexique PRO</u> – для публикации Tickbox Lexicography (TBL)





Инструменты для разметки: примеры

Разметка текстов

Praat – разметка звука
ELAN – для аудио и видео
Inception, GATE – для текстов
Fieldworks – для полевых

Разработка корпусов KonText, Corpus Workbench конкордансеры

+ текстовые редакторы, редакторы баз данных...

Разработка словарей

iLex – разработка

<u>TLex</u> – разработка с AI-

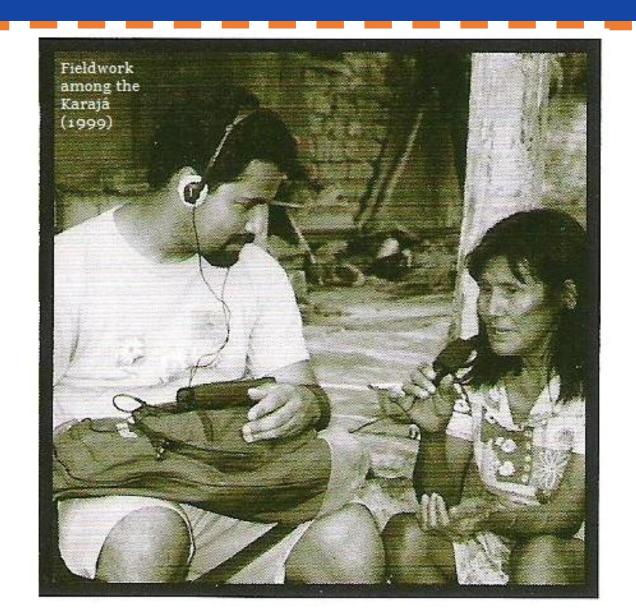
поддержкой

Lexique PRO – для публикации Tickbox Lexicography (TBL)

Перевод

Translation memory tools

Документация языков



Сценарий 1

На какие ресурсы следует опираться при исследовании (неизвестного) языка?

Сценарий 1

Документация языков:

можно ли написать текст на незнакомом языке, прочитав его грамматику и словарь?

Сценарий 1 - документация языков

Можно ли написать текст на незнакомом языке, прочитав его **грамматику** и **словарь**?

Нет - нужно прочитать много **текстов**, а еще лучше, **пообщаться** с носителями языка.

- язык средство общения
- в языке всегда есть много вариантов выражения мысли в зависимости от намерений говорящего и коммуникативной ситуации
- язык живой, языковые средства могут меняться

Откуда тексты?

Цифровая революция и лингвистика

Интернет как фонд текстов на языке N:

- новости, тексты из газет и журналов, электронные версии книг, сценарии кинофильмов, сайты музеев и учебных заведений, обзоры товаров, транскрипты интервью, реклама...
- аудио- и видеозаписи: радиопрограммы, интервью, аудиокниги, радиоспектакли, песни, youtube/rutube..
- социальные сети, форумы, чаты
- справочные, энциклопедические и образовательные ресурсы
- поисковые системы и переводчики

Откуда тексты?

Цифровая революция и лингвистика

стало

книги радио, телевидение

электронные книги аудиокниги интернет-тексты

словари + энциклопедии

электронные словари

учебники



электронные пособия, медиакурсы, тренажеры записи онлайн-обучения (Skype)

ресурсы для переводчиков

Откуда тексты?

Цифровая революция и лингвистика

стало

книги радио, телевидение

электронные книги аудиокниги интернет-тексты

словари + энциклопедии

электронные словари

учебники

электронные пособия, медиакурсы, тренажеры записи онлайн-обучения (Skype)

было

ресурсы для переводчиков

Электронные библиотеки (примеры)

Google Books - <u>books.google.com</u>

- Universal Digital <u>Library</u>
- Проект Гутенберг www.gutenberg.org
 Internet Archive.org
- •"Народные" проекты
 - <u>lib.ru</u> Библиотека Максима Мошкова
 - netslova.ru Сетевая словесность
 - <u>russ.ru</u> Русский журнал, <u>stihi.ru</u> Стихи.ру...
- Академические проекты
- <u>feb-web.ru</u> Фундаментальная электронная библиотека "Русская

литература и фольклор" — аннотированные электронные версии классики, включая варианты изданий (там же словари и литературные энциклопедии)

• wikipedia.org Википедия (архив как большой текстовый ресурс)



Сценарий 1 - документация языков

Корпус - лучше, чем текстовый архив

Корпус – коллекция текстов, снабженная специальной разметкой (информация о самих текстах, о каждом предложении и слове)

Электронные корпуса



Типичные вопросы, на которые отвечают корпуса:

- отличается ли речь авторов-женщин от авторов-мужчин?
- когда впервые появилось в языке слово слямзить?
 (NB! не появилось, а задокументировано)

• отличается ли сочетаемость слов *хотеть* и *стремиться*? (ср. [?]я стремился, чтобы...)

Электронные корпуса

Типичные вопросы, на которые отвечают корпуса:

- отличается ли речь авторов-женщин от авторов-мужчин?
 - <все тексты должны иметь помету "пол автора">
- когда впервые появилось в языке слово слямзить?
 (NB! не появилось, а задокументировано)
 - <все тексты должны иметь помету "дата создания">
 - <корпус должен уметь находить слово во всех формах разметка лексем>
- отличается ли сочетаемость слов хотеть и стремиться? (ср. [?]я стремился, чтобы...)
 - <синтаксическая разметка связь и расстояние между словами>

Электронные корпуса

- Собираются на основе существующих текстов, но включают добавленное знание (added knowledge)
 - лингвистическое
 про каждый звук, букву/иероглиф, слог, слово, предложение,
 абзац и так далее в зависимости от той или иной теории
 - энциклопедическое разметка топонимов, годы жизни авторов и т.п.
 - стиховедческое / литературоведческое
 метрика, виды рифмы, типы элегий и т.п.
 - историческое, Computer Science/NLP/AI, психологическое...

Классификация ресурсов

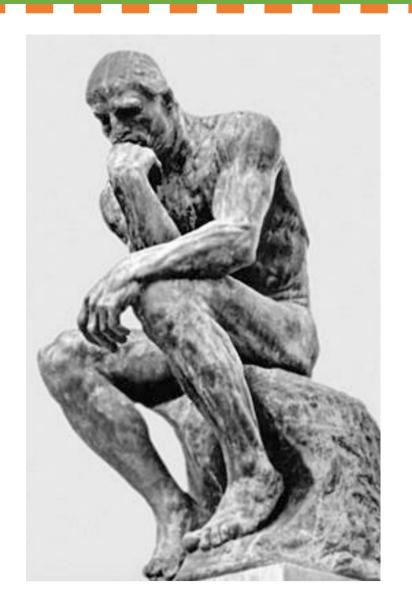
- грамматики
- корпуса
- словари

базы данных

- грамматические БД структурированные факты по грамматикам
- надкорпусные БД, в т.ч. частотные
- лексикографические БД структурированные факты о лексике
- справочные системы (Грамота.ру)
- другие специальные ресурсы (GIS-системы для диалектных исследований, <u>common voice</u> и т.п.)



Что еще?



Чем еще пользуются лингвисты?



Что еще?

- языковая интуиция?
 - если исследователь сам носитель языка, интуиция поможет ему ответить на вопрос, существует та или иная конструкция/единица

- если нет... можем спросить других носителей!
- собрать непрямые свидетельства
- и какого типа у нас получатся ресурсы?



Ресурсы для/как результат социо-, психо-, нейролингвистических, полевых и т.п. исследований

- унифицированные анкеты/опросники
- сбалансированные базы стимулов
- наборы данных (ответы, измерения времени реакции, ЭЭГ...)
 корпусные наборы данных (выборки)
- библиографии / mind maps
 - воспроизводимость исследований



Примеры





Корпуса

- Национальные корпуса ВNС, НКРЯ...
- Интернет-корпуса EnTenTen, RuTenTen, <u>Wacky</u>, <u>Aranea</u>...
- Мониторинговые газетный корпус НКРЯ
- Диахронические <u>СОНА</u>
- Устные корпуса
- Мультимедийные The Longitudinal Corpus of Finnish Spoken in Helsinki (1970s, 1990s and 2010s),
- Параллельные <u>OPUS</u> project, параллельные корпуса НКРЯ
- Сопоставимые корпуса <u>CHILDES</u> TalkBank, Wikipedia как корпус
- Диалектные корпуса HSE <u>LingConLab</u> согрога
- Учебные & эритажные корпуса LINDSEI (Louvain International Database of Spoken English Interlanguage), REALEC



Map data @2018 Google, INEGI, ORION-ME Terms of Use

Kingdom

Google

Nordic Dialect Corpus: Search result for the word ikke 'not' (a map view).

Базы данных

• Типологические

7,000 living languages: families, location & maps, population size, dialects

Ethnologue

WALS

World Atlas of Language Structures

languages & languoids, incl. sign languages and artificial languages

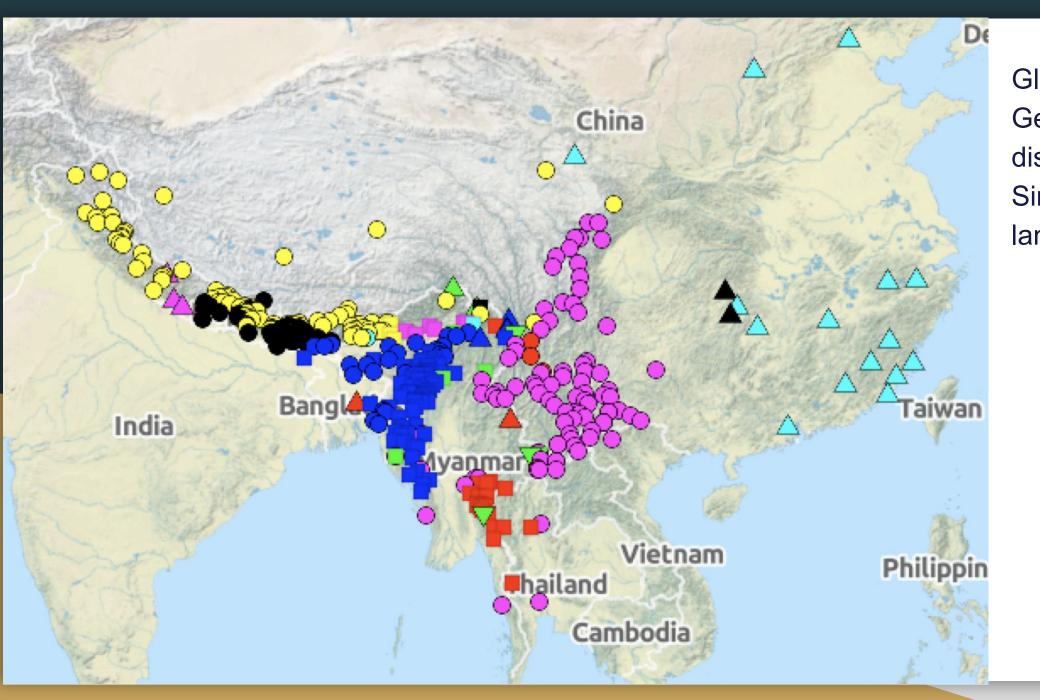
Glottolog

PHOIBLE

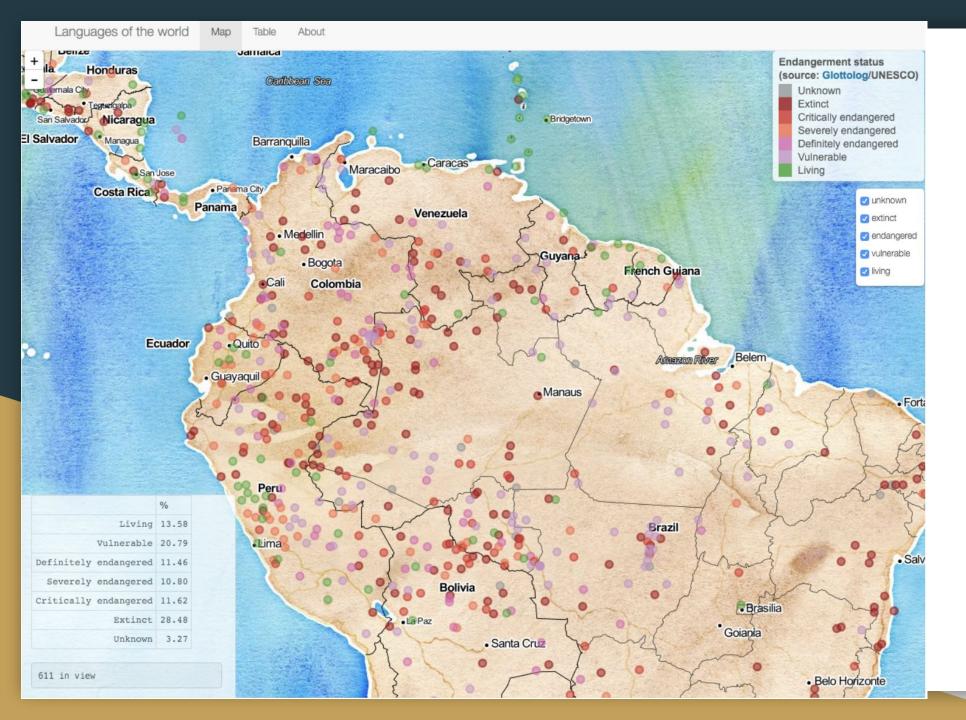
phonological inventories



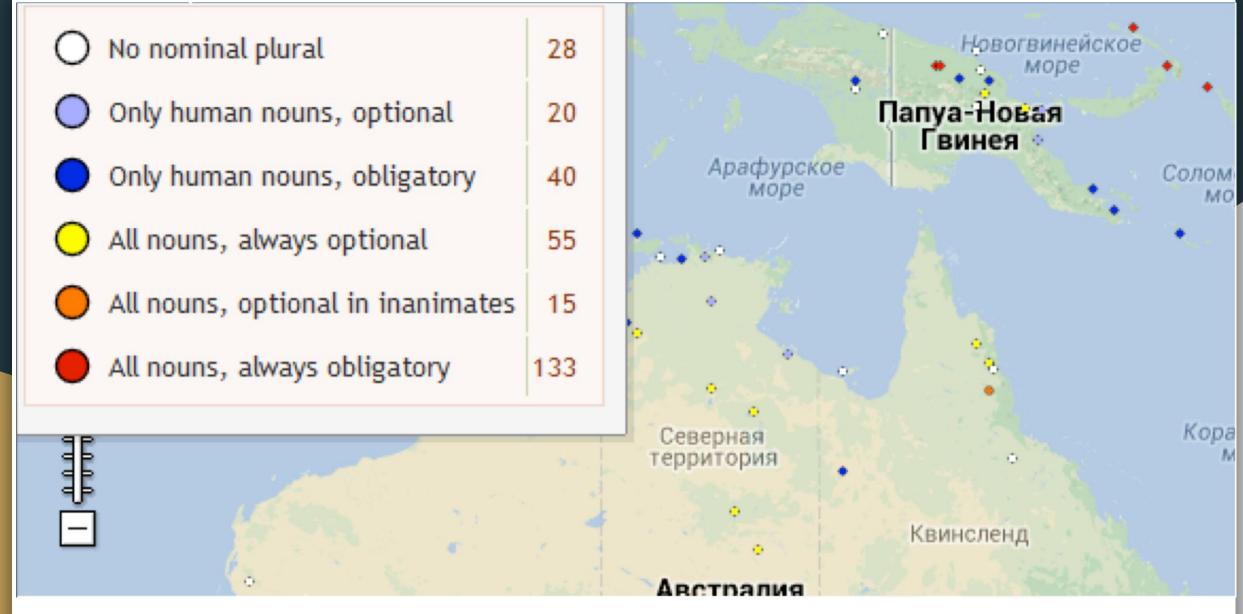




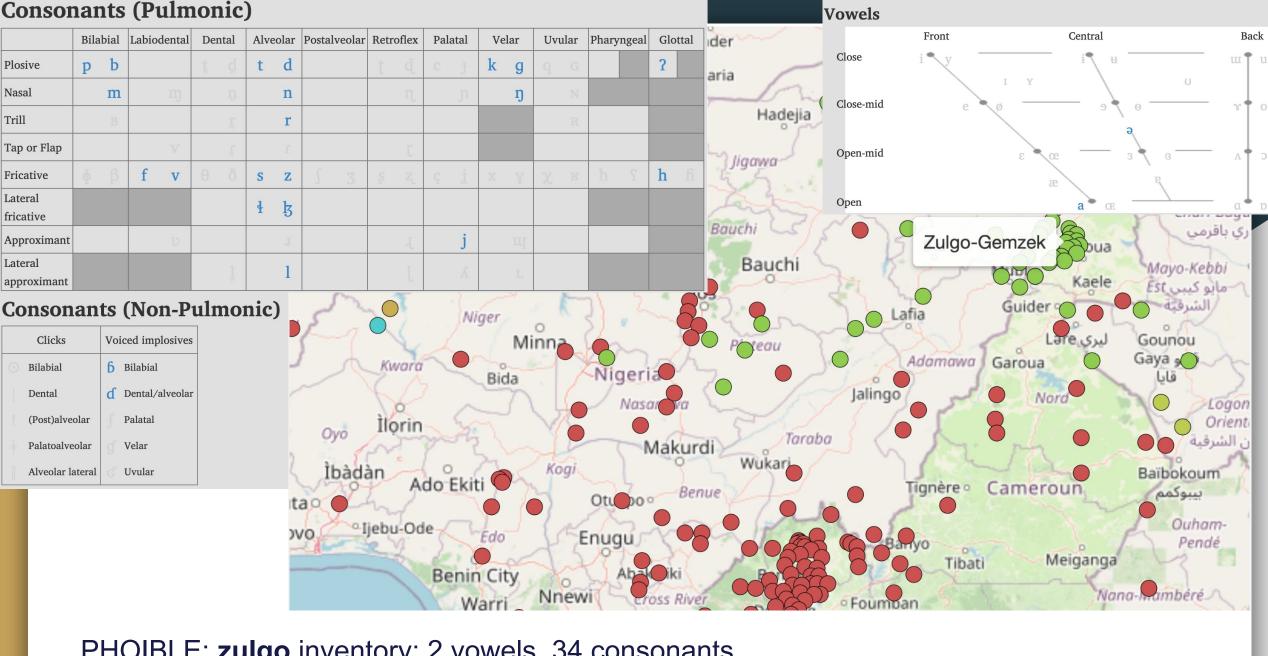
Glottolog: Geographical distribution of Sino-Tibetan languages



Glottolog data
explorer map:
The northern region
of South America
including the
Amazon rainforest
(Caines et al. 2016)



WALS: Feature 34A: Occurrence of Nominal Plurality

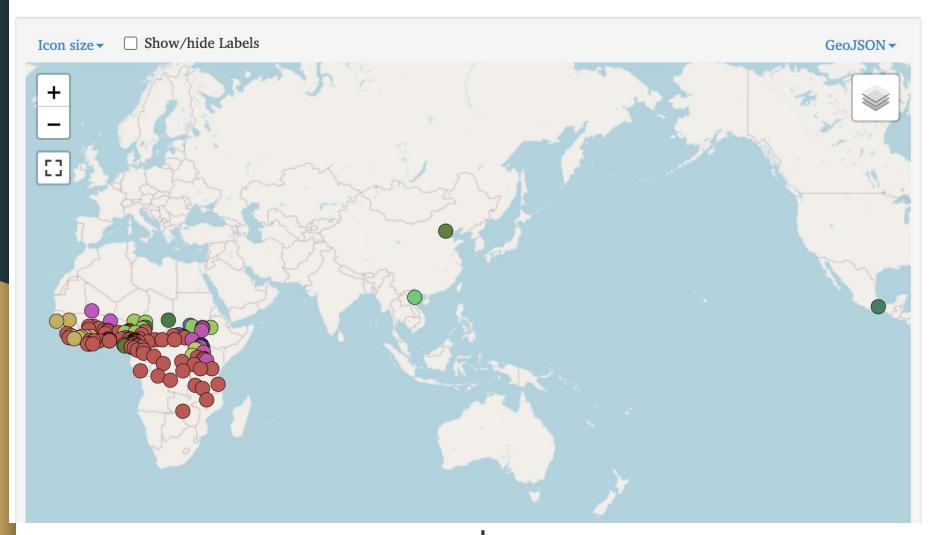


PHOIBLE: zulgo inventory: 2 vowels, 34 consonants





MODIFIER LETTER HIGH TONE BAR - MODIFIER LETTER LOW TONE BAR



PHOIBLE: Languages with tone \ - MODIFIER LETTER HIGH TONE BAR - MODIFIER LETTER LOW TONE BAR

Базы данных

• Лексические

etymological database

The Tower of Babel

<u>VisuWords</u>

lexicon as a net: visual dictionary, visual thesaurus, interactive lexicon

117,000 synsets (synonyms, hyperonyms & hyponyms, meronyms, antonyms, etc.)

WordNet

Common Voice

pronunciation crowdsourcing project





Number: 1713

Proto: *wasa

English meaning: calf, deer calf

German meaning: Kalb, Renkalb

Finnish: vasa 'Kalb, einjähriges Renkalb', vasikka 'Kalb'

Estonian: vasik, vasikas (gen. vasika)

Saam (Lapp): vyesi (I), viisse (T), vūiss (Kld.), vuaiss (Not.) kleines Rentierkalb, bis es

um den Peterstag neues Haar bekommt'

Mordovian: vaz (E M), vazńe (E), vazńä (M) 'Kalb'

Mansi (Vogul): (wēsəj KM, wēsəy P, wāsiy So. 'Elchkalb' - rejected by Redei as "eine vom

FW unabhängige iran. Entlehnung")

StarLing database for Uralic etymology (The Tower of Babel): Entry 1713: *wasa

Гибридные ресурсы

frames as predicate & arguments structures: semantic and syntactic roles, adjuncts, frame scenarios

FrameNet

Constructicon

multiword constructions (core grammatical structures & idiomatic, non-transparent form-meaning pairings)

Syntax, semantics, collocations

> automatic corpus-based summary of a word's grammatical and collocational behaviour

Sketch Engine

Trados

computer-assis ted translation (CAT) tool: translation memory & term database





goa (noun) ukWaC freq = 168345 (107.5 per million)

object of	58924	3.2	subject of	<u>25451</u>	2.4	modifier	67879	1.6	modifies	11026	0.3
score	8390	11.28	score	903	8.59	ultimate	1911	9.27	scorer	389	9.39
achieve	9422	9.9	disallow	223	8.04	long-term	875	7.66	kick	634	8.86
concede	1421	9.39	concede	204	7.53	league	638	7.38	tally	129	7.9
accomplish	<u>585</u>	7.97	gape	76	6.5	winning	401	7.33	keeper	204	7.31
reach	1924	7.66	come	1316	5.44	primary	993	7.24	scramble	<u>50</u>	6.75
net	337	7.42	kick	76	5.44	second	2000	7.19	drought	<u>78</u>	6.65
pursue	648	7.41	rule	<u>61</u>	5.24	common	1529	7.17	difference	676	6.28
attain	400	7.35	orientate	34	5.06	strategic	645	7.1	cushion	<u>53</u>	6.26
grab	406	7.34	arrive	90	4.43	realistic	422	7.05	lead	267	6.24
set	2413	7.01	cap	20	4.38	achievable	290	6.97	setting	<u>405</u>	6.14
pull	<u>501</u>	6.88	beat	<u>53</u>	4.31	stated	<u>259</u>	6.8	kicker	<u>25</u>	6.04
disallow	<u>190</u>	6.67	direct	<u>53</u>	4.22	score	<u>611</u>	6.75	post	<u>482</u>	5.91

SketchEngine: word sketches of the noun GOAL

clever/intelligent ukwac freqs

ukWaC freqs = 20589/2611

clever 6.0 4.0 2.0 0 -2.0 -4.0 -6.0 intelligent

and/or	4955	10062	2.2	3.6
perceptive	0	<u>34</u>	0.0	6.4
thought-provoking	0	<u>32</u>	0.0	6.2
informed	0	<u>66</u>	0.0	6.2
autonomous	0	<u>46</u>	0.0	6.2
adaptive	0	<u>39</u>	0.0	6.1
well-informed	0	<u>24</u>	0.0	6.0
literate	0	<u>26</u>	0.0	5.9
compassionate	0	<u>27</u>	0.0	5.9
well-educated	0	<u>17</u>	0.0	5.7
cultured	0	<u>19</u>	0.0	5.7
rational	0	<u>46</u>	0.0	5.6
playful	0	22	0.0	5.6
sensitive	8	134	2.0	5.9
thoughtful	<u>14</u>	121	5.0	7.7
affectionate	<u>6</u>	<u>31</u>	4.5	6.2
sophisticated	23	<u>75</u>	4.2	5.7
charming	21	<u>50</u>	4.9	5.9
insightful	11	31	5.2	6.1
resourceful	12	29	5.8	6.3
witty	132	166	8.3	8.2
inventive	22	24	5.9	5.5
clever	<u>54</u>	30	5.8	4.8
funny	233	103	7.0	5.7
cunning	<u>16</u>	<u>7</u>	5.9	4.0
catchy	<u>19</u>	0	5.8	0.0

modifier	4950	3168	0.9	0.5
emotionally	0	111	0.0	8.6
artificially	0	<u>52</u>	0.0	7.9
fiercely	0	26	0.0	7.0
highly	0	570	0.0	6.9
ferociously	0	<u>8</u>	0.0	6.2
supposedly	0	28	0.0	6.2
averagely	0	7	0.0	6.1
moderately	0	11	0.0	5.7
reasonably	0	<u>54</u>	0.0	5.7
computationally	0	<u>6</u>	0.0	5.6
supremely	0	7	0.0	5.5
culturally	0	12	0.0	5.5
exceptionally	29	25	6.0	5.9
remarkably	24	11	5.7	4.8
amazingly	<u>17</u>	7	5.9	5.0
wonderfully	20	9	5.4	4.5
very	1707	<u>596</u>	5.6	4.0
too	<u>476</u>	<u>76</u>	5.4	2.8
damn	12	0	5.6	0.0
dead	<u>16</u>	0	5.8	0.0
diabolically	<u>9</u>	0	5.9	0.0
awfully	<u>15</u>	0	6.1	0.0
terribly	<u>25</u>	0	6.2	0.0
devilishly	<u>17</u>	0	6.8	0.0
fiendishly	<u>45</u>	0	8.1	0.0

modifies	10948	16081	2.0	2.4
being	0	208	0.0	6.1
robot	0	<u>77</u>	0.0	6.1
agent	<u>9</u>	<u>455</u>	0.4	6.0
guess	0	<u>35</u>	0.0	5.5
routing	0	<u>27</u>	0.0	5.3
layman	0	22	0.0	5.3
conversation	0	88	0.0	5.1
creature	11	<u>137</u>	2.4	5.9
lyric	<u>81</u>	<u>80</u>	5.8	5.7
fellow	<u>52</u>	<u>14</u>	5.1	3.1
pass	<u>67</u>	9	5.2	2.2
stuff	<u>146</u>	<u>6</u>	5.1	0.4
gimmick	<u>15</u>	0	5.1	0.0
satire	<u>19</u>	0	5.1	0.0
flick	<u>21</u>	0	5.2	0.0
lob	<u>15</u>	0	5.3	0.0
pun	<u>19</u>	0	5.3	0.0
ruse	<u>17</u>	0	5.5	0.0
eh	24	0	5.8	0.0
wordplay	<u>21</u>	0	5.8	0.0
chap	<u>47</u>	0	5.9	0.0
twist	<u>94</u>	0	6.5	0.0
trick	<u>166</u>	0	6.7	0.0
clog	<u>50</u>	0	7.0	0.0
ploy	<u>68</u>	0	7.2	0.0

SketchEngine: contrastive sketch profiles for the adjectives *clever* and *intelligent*

Словники

Swadesh lists for historical studies of languages

Concepticon

Questionnaires

for linguistic field work (e.g. emotions, body parts, spacial relations)

frequency-balanced word lists for experimental studies

Naming tests





А есть ли у вас идеи?





Заключение

- "The contents of a corpus [and other resources!] should be selected without regard for the language they contain, but according to their communicative function in the community in which they arise."

 John Sinclair 2004
- Конвертация (традиционных) словарей и грамматик в формат структурированных баз данных, статистика и визуализации позволяют исследователям, учащимся и др. быстрее знакомиться с информацией и релевантными исследованиям
- Все ресурсы должны строиться согласно тщательно выработанному дизайну (выборки, баланс и т.п.) и подробно документироваться!
- Cross-Linguistic Linked (Open) Data movement облегчают доступ к ресурсам, распределенным по всему миру, и их интеграцию_____