



NATIONAL RESEARCH  
UNIVERSITY

# Лингвистические электронные ресурсы для лингвистических исследований (введение)

О.Н.Ляшевская  
школа лингвистики НИУ ВШЭ

# На какие ресурсы нужно опираться при исследовании (неизвестного) языка?

# Что нужно для исследования языка?

Документация языков:

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

# Что нужно для исследования языка?

## 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

# Цифровая революция и лингвистика

книги

радио, телевидение

словари

энциклопедии

учебники

электронные книги

аудиокниги

интернет

+      электронные словари

электронные пособия,  
медиакурсы, тренажеры

автоматические переводчики,  
Skype-обучение ....

# Язык в Интернете

(Ру)нет как фонд текстов на (русском) языке:

- Источник полезной информации: новости, статьи из газет и журналов, электронные версии книг, сценарии кинофильмов, сайты музеев и учебных заведений, обзоры товаров, транскрипты интервью...
  - + аудио- и видеозаписи: радиопрограммы, интервью, аудиокниги, радиоспектакли, песни, youtube (rutube)
- Сети для электронной коммуникации: facebook, вконтакте, Живой журнал и др., форумы, чаты
- Справочные, энциклопедические и образовательные ресурсы
- Поисковые системы и переводчики

# Интернет представляет испорченный язык

(Ру)нет как фонд текстов на (русском) языке:

- Источник полезной информации: новости, статьи из газет и журналов, электронные версии книг, сценарии кинофильмов, сайты музеев и учебных заведений, обзоры товаров, транскрипты интервью...
  - + аудио- и видеозаписи: радиопрограммы, интервью, аудиокниги, радиоспектакли, песни, youtube (rutube)
- Сети для электронной коммуникации: facebook, вконтакте, Живой журнал и др., форумы, чаты
- Справочные, энциклопедические и образовательные ресурсы
- Поисковые системы и переводчики

# Электронные библиотеки (примеры)

## Google Books - [books.google.com](http://books.google.com)

- "Народные" проекты
  - [lib.ru](http://lib.ru) Библиотека Максима Мошкова
  - [lib.aldebaran.ru](http://lib.aldebaran.ru) Библиотека "Альдебаран"
  - [netslova.ru](http://netslova.ru) Сетевая словесность
  - [russ.ru](http://russ.ru) Русский журнал
- Академические проекты
  - [feb-web.ru](http://feb-web.ru) Фундаментальная электронная библиотека "Русская литература и фольклор" -- аннотированные электронные версии классики, включая варианты изданий (там же словари и литературные энциклопедии)
  - [ru.wikipedia.org](http://ru.wikipedia.org) Википедия (архив Википедии как большой текстовый ресурс)
  - Проект Гутенберг - [www.gutenberg.org](http://www.gutenberg.org)

# Электронные корпуса

- Задачам лингвистического исследования лучше всего отвечают не просто **тексты** (архивы текстов),  
**а корпуса**
  - коллекции текстов, снабженные специальной разметкой (информация о текстах в общем, о каждом предложении и слове)

# 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...)

# Классификация ресурсов

- грамматики
  - базы данных
    - структурированные факты по грамматикам
  - корпуса
    - сводные данные по употреблению языковых единиц в корпусе, в т.ч. частотные
  - словари
    - структурированные факты о лексике
- справочные системы
- другие специальные ресурсы (геоинформация по диалектам и т.п.)

# Чем еще пользуются лингвисты?



# Чем еще пользуются лингвисты?

- интуицией:

если являешься носителем языка, можно спросить себя, "можно ли так сказать"?

# Чем еще пользуются лингвисты?

- интуицией:

если являешься носителем языка, можно спросить себя, "можно ли так сказать"?

*К сожалению, интуицию трудно превратить в ресурс!*

# Чем еще пользуются лингвисты?

- интуицией:

если являешься носителем языка, можно спросить себя, "можно ли так сказать?"

- опросами информантов + экспериментами

если сомневаешься, можно спросить носителя языка, "можно ли так сказать?"

можно (в ходе эксперимента) спровоцировать носителя языка произнести или не произнести интересующую меня языковую единицу

К сожалению, результаты опросов информантов и экспериментов малодоступны и еще не стали общеспецифичными

# Примеры (just a few...)

## Словари в электронном формате

- **dic.academic.ru**
  - словари и энциклопедии на русском языке
- **slovvari.ru** под эгидой Института русского языка им. В.В. Виноградова
  - академические словари русского языка + грамматики русского языка
- **feb-web.ru/feb/feb/dict.htm** - словари Фундаментальной е-библиотеки
- **etymolog.ruslang.ru** - этимологические словари
- **dict.ruslang.ru** - словари на основе НКРЯ
- **ru.wiktionary.org** - Вики-словарь (сделай словарь сам!)
- **<http://www.onelook.com>** - поиск по словарям английского языка
- **~~slovvari.yandex.ru~~** на портале Яндекса

и др. словари русского языка – двуязычные словари –  
энциклопедии онлайн

# Электронные корпуса русского языка

- **ruscorpora.ru** - Национальный корпус русского языка
- **ruTenTen** (sketch engine)
- **корпуса Сергея Шарова** (Лидс, Англия)
- **Упсальский корпус**
- **Тюбингенский корпус**
- **ХАНКО** Хельсинкский аннотированный корпус русских текстов
- **Компьютерный корпус текстов русских газет конца XX века**  
(МГУ)
- **Корпус русского литературного языка** (С.-Петербург)
- **Регенсбургский диахронический корпус русского языка**  
(древнерусские тексты)
- Рукописные памятники Древней Руси: берестяные грамоты, летописи, рукописная книга
- Параллельный корпус переводов «Слова о полку Игореве»
- Корпус русских публицистических текстов второй половины XIX века

# Электронные корпуса других языков

- **BNC** - Британский национальный корпус (английский язык)
- [corpus.byu.edu/coca](http://corpus.byu.edu/coca) - корпус COCA (American English)
- **COHA** - исторический корпус американского английского
- **COSMAS** - Мангеймские корпуса немецкого языка
- **Lexicum** - корпус канадского французского
- **VISL** - корпуса Университета Южной Дании (Э.Бик)
- **Leeds corpora** - интернет-корпуса Университета Лидса (С.Шаров)
- [web-corpora.net](http://web-corpora.net) - корпуса осетинского, калмыцкого, бурятского, монгольского, албанского и многих других "небольших" языков
- **JRC-Acquis** - параллельный корпус документов Евросоюза на 23 языках
- см. больше ссылок на [studiorum.ruscorpora.ru](http://studiorum.ruscorpora.ru) и [linguistilist.org](http://linguistilist.org)

# Базы данных

## О языках (типологические БД)

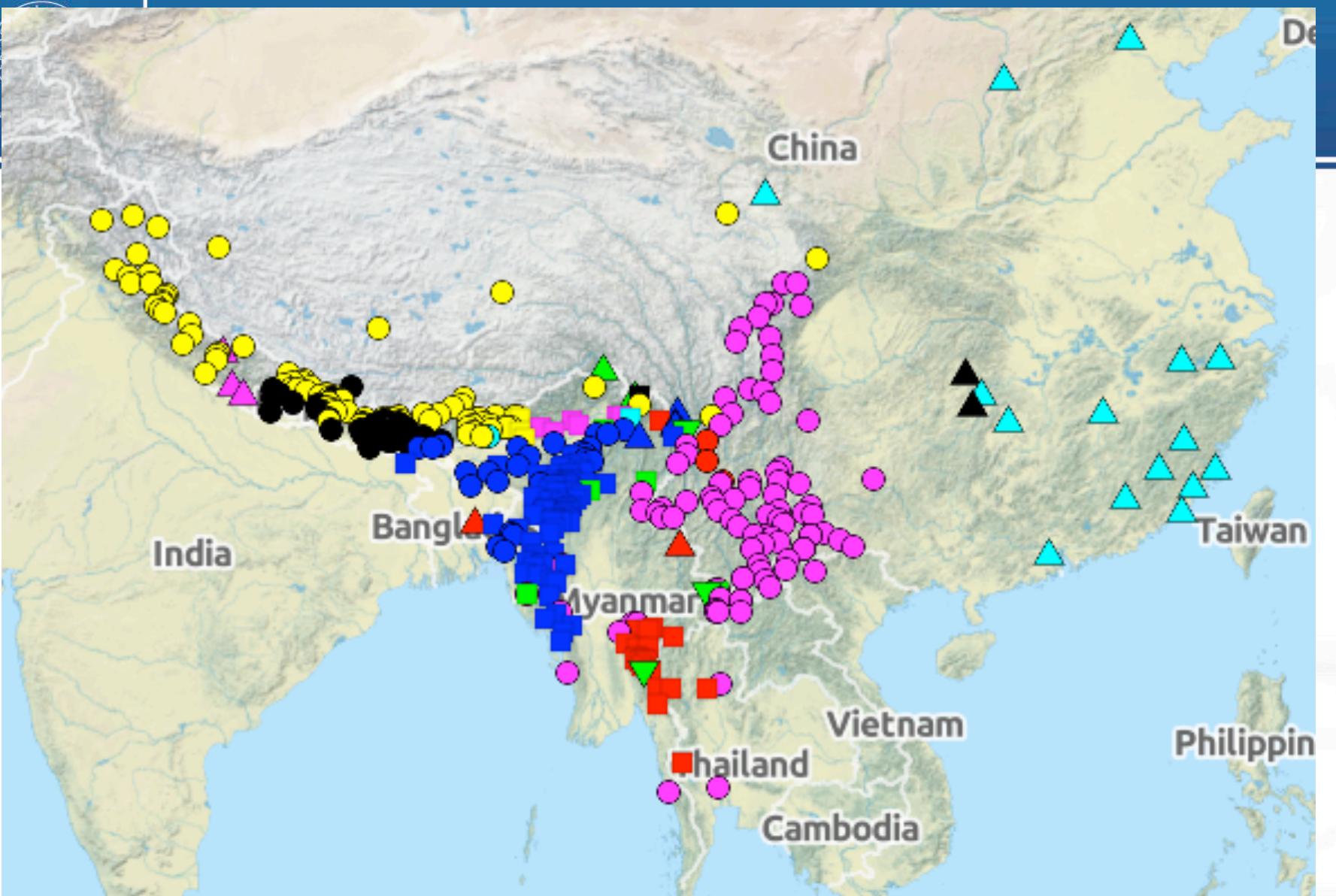
- **Ethnologue** [www.ethnologue.com](http://www.ethnologue.com) - база данных языков мира (семья, численность, ареал, живой/вымирающий, карты)
- **Glottolog** [glottolog.org](http://glottolog.org) - генеалогическая классификация + библиография
- **WALS** [wals.info](http://wals.info) - The World Atlas of Language Structures, + типологические свойства языков, типологические очерки

## О лексике

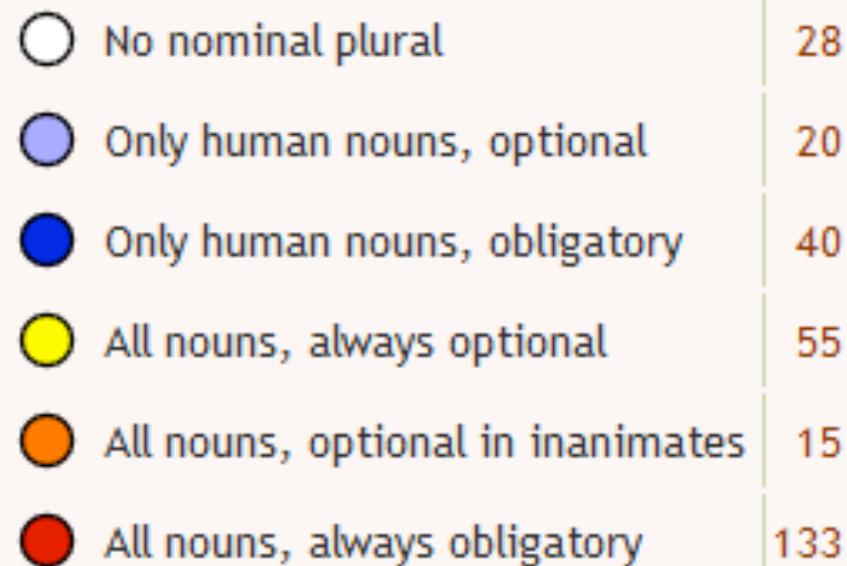
- **The Tower of Babel** [starling.rinet.ru](http://starling.rinet.ru) - этимологическая база данных
- **WordNet** - семантическая сеть лексики (синонимы, гипонимы...)
- **visuwords.com** - online graphical dictionary, связи между словами

## О синтаксисе и сочетаемости

- **SketchEngine** - помощник лексикографа
- **FrameNet** - конструкции глаголов



Glottolog: ареал распространения сино-тибетских языков



WALS: языки с разными системами грамматического числа

**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), viis (T), vuiss (Kld.), vuass (Not.) 'klein'

**Mordovian:** vaz (E M), vazne (E), vaznä (M) 'Kalb'

**Mansi (Vogul):** (wēsəj KM, wēsəy P, wāsiy So. 'Elchkalb' - rejected)

Вавилонская башня: этимологически связанное гнездо (фино-угорск.)

# goal

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

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

SketchEngine: типичные контексты слова *goal*

# clever/intelligent

ukWaC freqs = 20589/26115

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

and/or	4955	10062	2.2	3.6	modifier	4950	3168	0.9	0.5	modifies	10948	16081	2.0	2.4
perceptive	0	<u>34</u>	0.0	6.4	emotionally	0	<u>111</u>	0.0	8.6	being	0	<u>208</u>	0.0	6.1
thought-provoking	0	<u>32</u>	0.0	6.2	artificially	0	<u>52</u>	0.0	7.9	robot	0	<u>77</u>	0.0	6.1
informed	0	<u>66</u>	0.0	6.2	fiercely	0	<u>26</u>	0.0	7.0	agent	<u>9</u>	<u>455</u>	0.4	6.0
autonomous	0	<u>46</u>	0.0	6.2	highly	0	<u>570</u>	0.0	6.9	guess	0	<u>35</u>	0.0	5.5
adaptive	0	<u>39</u>	0.0	6.1	ferociously	0	<u>8</u>	0.0	6.2	routing	0	<u>27</u>	0.0	5.3
well-informed	0	<u>24</u>	0.0	6.0	supposedly	0	<u>28</u>	0.0	6.2	layman	0	<u>22</u>	0.0	5.3
literate	0	<u>26</u>	0.0	5.9	averagely	0	<u>7</u>	0.0	6.1	conversation	0	<u>88</u>	0.0	5.1
compassionate	0	<u>27</u>	0.0	5.9	moderately	0	<u>11</u>	0.0	5.7	creature	<u>11</u>	<u>137</u>	2.4	5.9
well-educated	0	<u>17</u>	0.0	5.7	reasonably	0	<u>54</u>	0.0	5.7	lyric	<u>81</u>	<u>80</u>	5.8	5.7
cultured	0	<u>19</u>	0.0	5.7	computationally	0	<u>6</u>	0.0	5.6	fellow	<u>52</u>	<u>14</u>	5.1	3.1
rational	0	<u>46</u>	0.0	5.6	supremely	0	<u>7</u>	0.0	5.5	pass	<u>67</u>	<u>9</u>	5.2	2.2
playful	0	<u>22</u>	0.0	5.6	culturally	0	<u>12</u>	0.0	5.5	stuff	<u>146</u>	<u>6</u>	5.1	0.4
sensitive	<u>8</u>	<u>134</u>	2.0	5.9	exceptionally	<u>29</u>	<u>25</u>	6.0	5.9	gimmick	<u>15</u>	0	5.1	0.0
thoughtful	<u>14</u>	<u>121</u>	5.0	7.7	remarkably	<u>24</u>	<u>11</u>	5.7	4.8	satire	<u>19</u>	0	5.1	0.0
affectionate	<u>6</u>	<u>31</u>	4.5	6.2	amazingly	<u>17</u>	<u>7</u>	5.9	5.0	flick	<u>21</u>	0	5.2	0.0
sophisticated	<u>23</u>	<u>75</u>	4.2	5.7	wonderfully	<u>20</u>	<u>9</u>	5.4	4.5	lob	<u>15</u>	0	5.3	0.0
charming	<u>21</u>	<u>50</u>	4.9	5.9	very	<u>1707</u>	<u>596</u>	5.6	4.0	pun	<u>19</u>	0	5.3	0.0
insightful	<u>11</u>	<u>31</u>	5.2	6.1	too	<u>476</u>	<u>76</u>	5.4	2.8	ruse	<u>17</u>	0	5.5	0.0
resourceful	<u>12</u>	<u>29</u>	5.8	6.3	damn	<u>12</u>	0	5.6	0.0	eh	<u>24</u>	0	5.8	0.0
witty	<u>132</u>	<u>166</u>	8.3	8.2	dead	<u>16</u>	0	5.8	0.0	wordplay	<u>21</u>	0	5.8	0.0
inventive	<u>22</u>	<u>24</u>	5.9	5.5	diabolically	<u>9</u>	0	5.9	0.0	chap	<u>47</u>	0	5.9	0.0
clever	<u>54</u>	<u>30</u>	5.8	4.8	awfully	<u>15</u>	0	6.1	0.0	twist	<u>94</u>	0	6.5	0.0
funny	<u>233</u>	<u>103</u>	7.0	5.7	terribly	<u>25</u>	0	6.2	0.0	trick	<u>166</u>	0	6.7	0.0
cunning	<u>16</u>	<u>7</u>	5.9	4.0	devilishly	<u>17</u>	0	6.8	0.0	clog	<u>50</u>	0	7.0	0.0
catchy	<u>19</u>	0	5.8	0.0	fiendishly	<u>45</u>	0	8.1	0.0	ploy	<u>68</u>	0	7.2	0.0

SketchEngine: синонимы и сочетаемость слов *clever* и *intelligent*

# Справочно-информационные ресурсы

- Справочные порталы (на примере русского)
- [gramota.ru](http://gramota.ru) - Грамота.ру, портал "Русский язык"
  - Справочная служба русского языка
  - словари, статьи, интерактивные диктанты, игры
- [gramma.ru](http://gramma.ru) "Культура письменной речи"
  - академические словари русского языка
  - грамматики русского языка
- [pishu-pravilno.livejournal.com](http://pishu-pravilno.livejournal.com) "Пишу правильно", сообщество в Живом журнале
- Порталы для исследователей
- [studiorum.ruscorpora.ru](http://studiorum.ruscorpora.ru) - справочная система в помощь пользователям корпусов
- [linguistlist.org](http://linguistlist.org) - информация о конференциях, журналах, исследователях и институтах, ресурсах и т.д.

# Ресурсы компьютерной лингвистики

- Яндекс.ru, Google.com, Bing.com, Baidu.com - работают на гигантских размеченных (индексированных) архивах текстов и на специальных словарях
- Системы проверки орфографии - используют словари и базы данных
- Системы автоматического перевода ([translate.google.com](https://translate.google.com), [multitran](https://multitran.org) и другие) - используют параллельные корпуса и словари
- Системы классификации новостей - словари + базы знаний
- Системы анализа мнения о товарах (opinion mining) и т.д.

# Инструменты для разметки текстов и создания ресурсов

## Аннотация текстов:

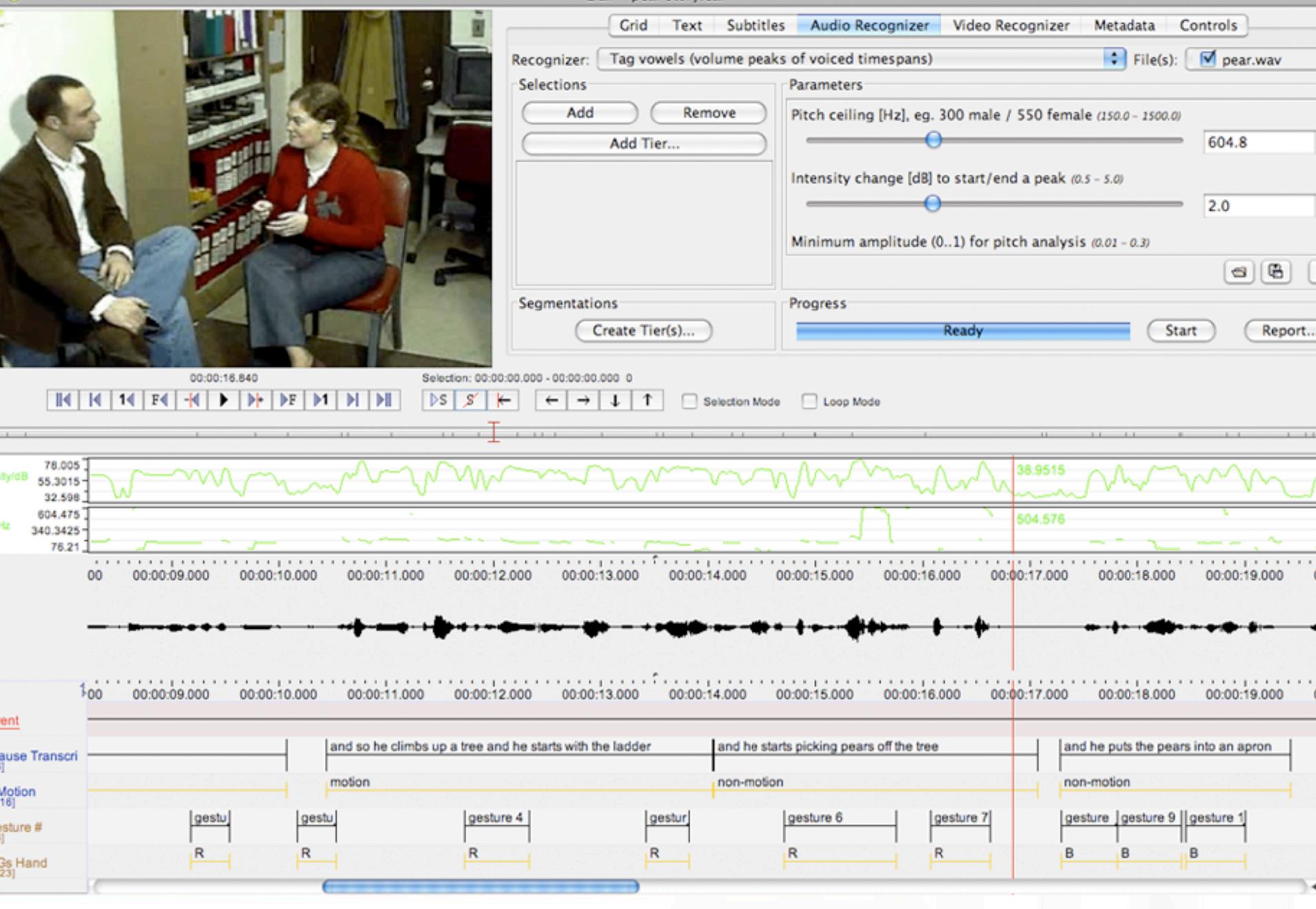
- ELAN [tla.mpi.nl/tools/tla-tools/elan/](http://tla.mpi.nl/tools/tla-tools/elan/) - аннотация аудио и видео
- Praat [www.fon.hum.uva.nl/praat/](http://www.fon.hum.uva.nl/praat/) - для работы с фонетикой
- GATE [GATE.ac.uk](http://GATE.ac.uk) - профессиональные инструменты
- UIMA [uima.apache.org](http://uima.apache.org) компьютерной лингвистики

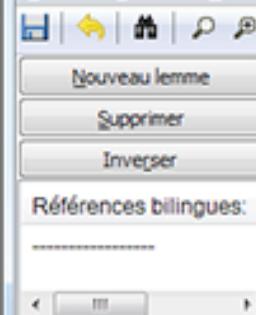
## Создание словарей:

- Lexus [tla.mpi.nl/.../lexus/](http://tla.mpi.nl/.../lexus/)
- iLex [www.emp.dk](http://www.emp.dk)
- IDM [idm.fr](http://idm.fr)
- TshwaneLex [tshwanedje.com/tshwanelex/](http://tshwanedje.com/tshwanelex/)
- Lexique Pro [lexiquepro.com](http://lexiquepro.com)
- ABBYY Lingvo Content
- + текстовые редакторы, базы данных, конкордансеры и т.д

## Корпус-менеджеры:

- WordSmithTools
- Bonito
- Corpus Workbench





sans

sanctuaire (\*)  
sandale (\*)  
sandwich (\*)  
sang (\*)  
sangle (\*)  
sangler (\*)  
sang-mêlé (\*)  
sangsue (\*)

sani

sans (\*)  
sans-coeur (\*)  
sans-joie (\*)  
Santa Claus (\*)  
santé (\*)  
saoul  
saper [1] (\*)  
saper [2]  
sapré (\*)

sans (\*)  
sans-coeur (\*)  
sani

Lemma: sans LemmaSign=sans,Modified=2009-02-23 20  
 Pronunciation: text: 'sɑ̃'  
 POSGroup: AutoNumber=1,PartOfSpeech=prep.  
 Sense: 1 AutoNumber=1  
 TE: TE=without  
 Example: Example=C'est bon quand tu peux danser sans musique. It's good when you can dance without music. (EV) \*On peut faire sans travailler le dimanche. We can do it without working on Sunday. (SL, An94)  
 sans cesse endless, ceaseless <Da84>  
 sans connaissance unconscious <Da84>  
 sans doute no doubt, without a doubt <Da84>  
 sans (que) a unless \*Et on veillait le mort, bien sûr. On aurait jamais laissé le mort sans que quelqu'un soit là. And we waked the body, of course. We would've never left the body unless someone was there. (TB) b without \*T'auras pas battu dans la salle sans il te fuit dehors. You wouldn't have fought in the dance hall without him throwing you out. (LA, An94) <LA, TB, An94, Da84>  
 ça va sans dire it goes without saying <Da84>  
 <Loc: AV, EV, IB, IV, LA, LF, SL, TB, VM, An94, Da84, Gu00, Hi02, Wh83>

Attributs (F1) Attributs (F2) Rechercher (F3)

Lemma:	sans	Incomplete
LemmaSign	sans	
Comma		
Brackets		
Frequency	0	
Notes		
Pronunciation:		
Audio		Parcourir...
Speaker		
[PCDATA]	sɑ̃	
POSGroup:		
LemmaSign		
PartOfSpeech	prep.	

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z ()

**sans** [sɑ̃] prep.

- 1 without • *C'est bon quand tu peux danser sans musique.* It's good when you can dance without music. (EV) \**On peut faire sans travailler le dimanche.* We can do it without working on Sunday. (SL, An94) ■ **sans cesse** endless, ceaseless <Da84> ■ **sans connaissance** unconscious <Da84> ■ **sans doute** no doubt, without a doubt <Da84> ■ **sans (que) a** unless \**Et on veillait le mort, bien sûr.* On aurait jamais laissé le mort sans que quelqu'un soit là. And we waked the body, of course. We would've never left the body unless someone was there. (TB) **b** without \**T'auras pas battu dans la salle sans il te fuit dehors.* You wouldn't have fought in the dance hall without him throwing you out. (LA, An94) <LA, TB, An94, Da84> ■ **ça va sans dire** it goes without saying <Da84> <Loc: AV, EV, IB, IV, LA, LF, SL, TB, VM, An94, Da84, Gu00, Hi02, Wh83> [Admin]

**sans-cœur** [sɑ̃kõr] n.

- 1 heartless, cruel, pitiless person • *Tu es rien qu'un sans-cœur.* You're nothing but a cruel man. (SB) <Loc: SB, Da84, Di32> [Admin]

**sans-joie** [sɑ̃ʒwa] n.m.

- 1 great blue heron <Loc: Lv88, Re31> [Admin]

**Santa Claus** [sàtaklɒz, sèteklɒz] n.prop.

- 1 Santa Claus <Loc: AC, EV, IB, Lv88, Ph36> [Admin]

**santé** [sàt̪e] n.f.

- 1 health • *J'ai pas pu m'empêcher de marcher à lui. Je dis, "Il y a une question j'aimerais te demander. Quoi c'est tu fais pour ta santé?" Il dit, "Je vas au bal proche tous les soirs."* I couldn't help but walk over to him. I said, "There's a question I'd like to ask you. What do you do for your health?" He said, "I go to the dance almost every night." (ch: *La neige sur la couverture*) ■ **à votre santé** to your health <Da84> ■ **en bonne santé** in good health <Da84> ■ **en mauvaise santé** in bad health <Da84>

# Инструменты для разметки текстов и создания ресурсов

## Аннотация текстов:

- ELAN [tla.mpi.nl/tools/tla-tools/elan/](http://tla.mpi.nl/tools/tla-tools/elan/) - аннотация аудио и видео
- Praat [www.fon.hum.uva.nl/praat/](http://www.fon.hum.uva.nl/praat/) - для работы с фонетикой
- GATE [GATE.ac.uk](http://GATE.ac.uk) - профессиональные инструменты
- UIMA [uima.apache.org](http://uima.apache.org) компьютерной лингвистики

## Создание словарей:

- Lexus [tla.mpi.nl/.../lexus/](http://tla.mpi.nl/.../lexus/)
- iLex [www.emp.dk](http://www.emp.dk)
- IDM [idm.fr](http://idm.fr)
- TshwaneLex [tshwanedje.com/tshwanelex/](http://tshwanedje.com/tshwanelex/)
- Lexique Pro [lexiquepro.com](http://lexiquepro.com)
- ABBYY Lingvo Content
- + текстовые редакторы, базы данных, конкордансеры и т.д

## Корпус-менеджеры:

- WordSmithTools
- Bonito
- Corpus Workbench