

Построение машинно- читаемого словаря на основе Русского Викисловаря



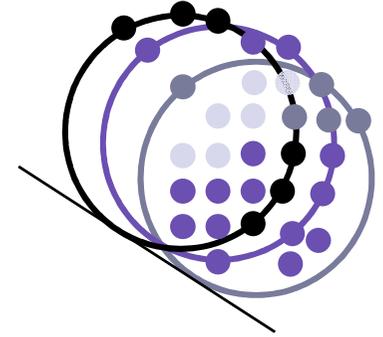
Санкт-Петербургский институт
информатики и автоматизации РАН



Крижановский Андрей (andrew.krizhanovsky  gmail.com)



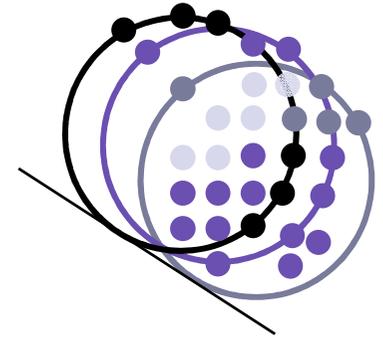
Содержание



- Викисловарь
 - применение
 - достоинства и трудности обработки
- MRD, парсер и сравнение Викисловарей
- Эксперимент
 - Корреляция мер семантической близости
- Результаты



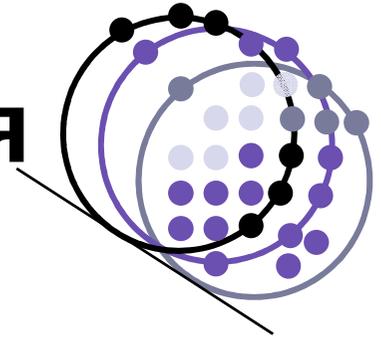
Цель



Сделать возможным
применение данных Викисловаря
(как лингвистического ресурса)
в различных компьютерных программах,
в задачах, связанных с обработкой текста.



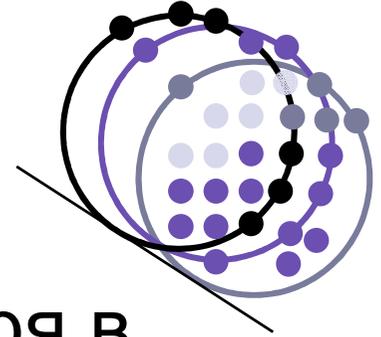
Применение Викисловаря



- В компьютерных программах:
 - текстовые поисковые системы
расширение / переформулировка запросов с помощью тезаурусов
 - запросно-ответные системы
распознавание запроса
- В задачах:
 - определение значения многозначного слова
 - сравнение онтологий (ontology matching)
 - автоматическое создание тезаурусов
 - машинный перевод
 - компьютерные **игры** для изучения языков
Медиа данные (звук + иллюстрации)



Задача



- Преобразования данных Викисловаря в машинную форму, а именно:
машинно-читаемый словарь (MRD).

MRD включает:

- Данные (база данных),
- Алгоритмы и функции (API)

Викисловарь –

МНОГО-

функциональный

многоязычный

словарь и

тезаурус

Wiktionary

Français
Le dictionnaire libre
856 000+ articles

English
The free dictionary
841 000+ articles

Tiếng Việt
Từ điển mở
227 000+ mục từ

Русский
Свободный словарь
137 000+ статей

中文
自由的多语言词典
116 000+ 词条

தமிழ்
கட்டற்ற ஆகரமுதலி
102 000+ கட்டுரைகள்

Türkçe
Özgür sözlük
208 000+ madde

Ido
La libera vortaro
137 000+ artikli

Ελληνικά
Το Ελεύθερο Λεξικό
107 000+ λέξεις

Polski
Wolny słownik
93 000+ stron

a multilingual tree encyclopedia
Wiktionary
[ˈwɪkʃənəri] n.,
a wiki-based Open Content dictionary
Wileo [ˈwɪl kəʀɪ]

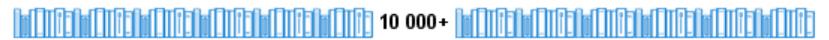
rechercher • search • tìm kiếm • ara • поиск • serchez • 搜索
αναζήτηση • தேடு • szukaj • haku • ricerca • suche • keresés • sök

English

Грамматический
Толковый
Этимологический
Многоязычный



Ελληνικά • English • Français • Ido • Русский • தமிழ் • Türkçe • Tiếng Việt • 中文



Afrikaans • العربية • Български • Brezhoneg • Deutsch • Eesti • Español • فارسی • Galego • 한국어 / 조선어 • Bahasa Indonesia • Íslenska • Italiano • Kurdî / كوردی • Lietuvių • Limburgs • Magyar • 日本語 • Nederlands • Polski • Português • Română • Sicilianu • Српски / Srpski • Suomi • Svenska • తెలుగు • Volapük



Asturiano • Bân-lâm-gú / Hō-ló-oē • Català • Corsu • Český • Dansk • Englisc • Esperanto • Frysk • Gaeilge • ગુજરાતી • हिन्दी • Hornjoserbsce • Hrvatski • Interlingua • עברית • Kalaallisut • Kaszëbsczi • ལྷན་རྒྱུ་གུ་ • Latina • മലയാളം • Bahasa Melayu • Norsk (bokmål) • Occitan • Қазақша • Sesotho • Shqip • Simple English • Slovenčina • Slovenščina • Kiswahili • Tatarça / татарча • Ἰουρῶν • Українська • اردو

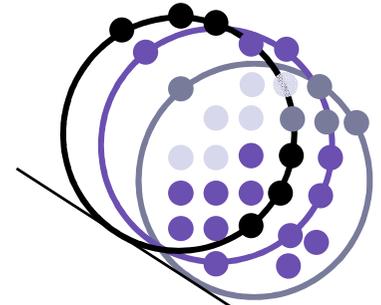


አማርኛ • Aragonés • Avañe'ê • Azərbaycan • Беларуская • Bosanski • Cymraeg • Euskara • Føroyskt • Gàidhlig • ગુજરાતી • Interlingue • ཐིམ་ཕུ་རྫོང་གི་སྐད་ • Kinyarwanda • Кыргызча • Latviešu • Македонски • मराठी • монгол • Nāhuatlahtōlli • पंजाबी • Plattdütsch • Runa Simi • سنڌي • Basa Sunda • Tagalog • བོད་སྐད་ • Gŵy • Xitsonga • ئۇيغۇرچە • Wolof • עברית • isiZulu

Other languages



Викисловарь = вики + ?



? Структура статьи = f (язык, ~часть речи)



? Определена последовательность частей статьи

? Шаблоны:

? структурные шаблоны ({{пример}}, {{морфо}})

? словоизменений, этимологии, родств. слова, пометы...

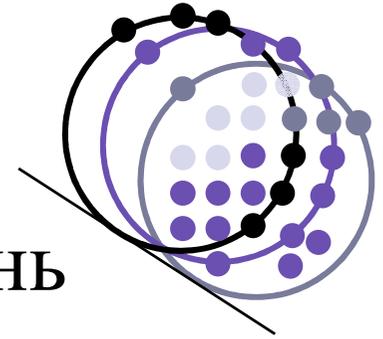
Т.о. жёсткая схема даёт:

+ единообразие, системность

+ возможность автоматически анализировать текст



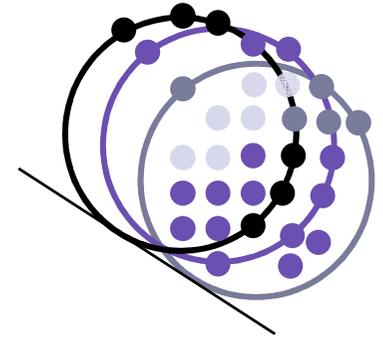
Данные Викисловаря: плюсы и трудности



- + Богатство
 - + тезаурус
(синонимы, антонимы...)
 - + фразеологизмы
 - + этимология
 - + произношение
 - + примеры употреб-ий
 - + переводы
 - + ...
 - + Быстрый рост
 - + Интервики (доп. д.)
 - + Свободная лицензия
- Разная степень стандартизации и формализации (структура статьи) в разных Викисловарях
 - Быстрый рост данных, но *толпа*:
 - Ручной ввод данных =>
 - Ошибки =>
Парсер д.б. устойчив!
 - Омонимия вне страницы (см. даль⁸ше)

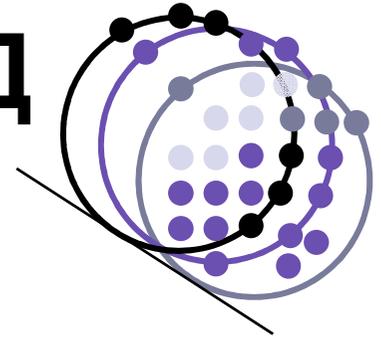


Данные Викисловаря: какие ещё статьи?

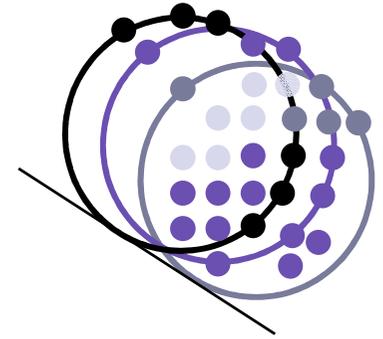


- Слова
- Устойчивые выражения, пословицы, поговорки, крылатые слова, народные приметы, загадки, скороговорки, сокращения
- Отдельные морфемы — корни, суффиксы, приставки и т. д.
- Омонимы, омографы, анаграммы, метаграммы и рифмы

Требования к парсеру, БД и процессу разработки



- Надёжность и устойчивость (lang=zzz, 8)
 - Unit-тесты > 200, визуализация
- Гибкость (раскопки форматов и правил)
 - Тестирование («живая» документация)
- Визуализация (Wiwordik, JavaFX)
- Викисловарь ++. (рост в ширину)
 - парсер = ядро + языкозависимая часть, ru + en
- Инкрементальный подход (рост в глубину)
- ¿Интеграция?



Структура словарной статьи в Русском Викисловаре



Морфологические и синтаксические свойства

[\[прав.\]](#)

ра-бó-чий

Прилагательное (относительное), тип склонения по классификации А. Зализняка — 4а.

Корень: **-раб-**; суффикс: **-оч-**; окончание: **-ий**.

Произношение

МФА: [rɐˈbɔtɕɨ]  Пример произношения

Семантические свойства

Значение

1. относящийся к работе ♦ Запиши мой рабочий телефон.
2. относящийся к рабочему (рабочим) I ♦ Рабочий посёлок. ♦ Она-де горожанка, руки у нее не **рабочие**, она не знает, что делать в огороде, со скотиной не умеет обращаться, доить не умеет, и при этом перечислении матушка заметно оживилась и повеселела. Альфред Хейдок, «Грешница», 1924-1934 г. (цитата из Национального корпуса русского языка, см. Список литературы)
3. находящийся в работе; предварительный, пробный, черновой ♦ Рабочий вариант.
4. способный нормально работать, функционировать ♦ У тебя случайно не заваялся какое-нибудь телефон, пусть старенький, но чтобы рабочий. ♦ А что, этот план вполне рабочий!

Синонимы

[\[прав.\]](#)

1. частичн.: трудовой

падеж	ед. ч.			мн. ч.
	м.	с.	ж.	
Им.	рабо́чий	рабо́чее	рабо́чая	рабо́чие
Р.	рабо́чего	рабо́чего	рабо́чей	рабо́чих
Д.	рабо́чему	рабо́чему	рабо́чей	рабо́чим
В. ^(одуш./неодуш.)	рабо́чего рабо́чий	рабо́чее	рабо́чую	рабо́чих рабо́чие
Тв.	рабо́чим	рабо́чим	рабо́чей рабо́чею	рабо́чими
Пр.	рабо́чем	рабо́чем	рабо́чей	рабо́чих

[\[прав.\]](#)[\[прав.\]](#)[\[прав.\]](#)

Структура статьи и БД

Морфологические и син

ра-бó-чий

Прилагательное (относительное),
классификации А. Зализняка — 4а.

Корень: **-раб-**; суффикс: **-оч-**; оконч

Произношение

МФА: [rɐ'botɕɨ]  Пример произ

Семантические свойства

Значение

1. относящийся к работе ♦ Зап телефон.
2. относящийся к рабочему (рабочим) I ♦ Рабочий посёлок. ♦ Она-де горожанка, руки у нее не **рабочие**, она не знает, что делать в огороде, со скотиной не умеет обращаться, доить не умеет, и при этом перечислении матушка заметно оживилась и повеселела. Альфред Хейдок, «Грешница», 1924-1934 г. (цитата из Национального корпуса русского языка, см. Список литературы)
- ★ 3. находящийся в работе; предварительный, пробный, черновой ♦ Рабочий вариант.
4. способный нормально работать, функционировать ♦ У тебя случайно не заваялся какое-нибудь телефон, пусть старенький, но чтобы рабочий. ♦ А что, этот план вполне рабочий!

Синонимы

1. частичн.: трудовой

page	
*id=1	INT(10)
°page_title="рабочий"	VARCHAR(255)

lang_pos	
*id=2	
*page_id=1	INT(10)
°lang_id=390	SMALLINT
*pos_id=3	TINYINT
*etymology_n=1	TINYINT

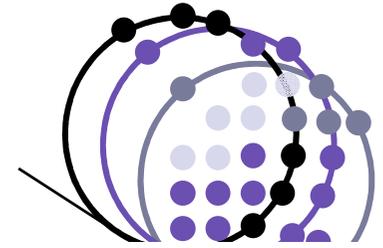
lang	
*id=390	SMALLINT
°name="Russian"	VARCHAR(255)
°code="ru"	VARCHAR(12)

=====

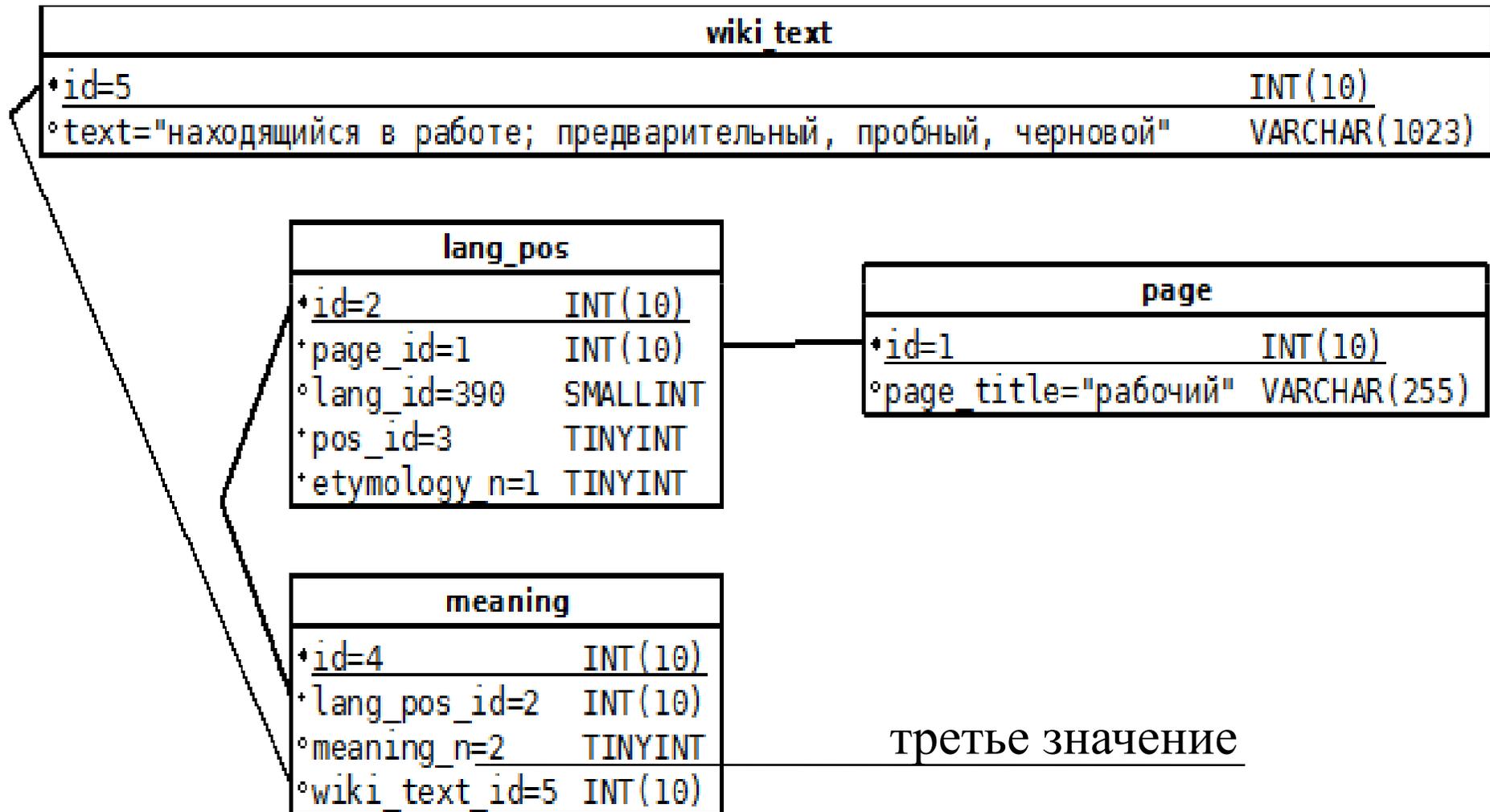
part of speech	
*id=3	TINYINT
*name="adjective"	VARCHAR(255)



Структура (толкование)

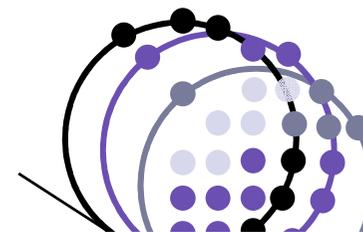


3. находящийся в работе; предварительный, пробный, черновой ♦ Рабочий вариант.





Внутренние ссылки (1)



3. находящийся в работе; предварительный, пробный, черновой ♦ Рабочий вариант.

находящийся в [[работа|работе]];
предварительный, пробный, черновой

wiki_text	
* <u>id=5</u>	INT(10)
°text="находящийся в работе; предварительный, пробный, черновой"	VARCHAR(1023)

wiki_text_words	
* <u>id</u>	INT(10)
°wiki_text_id=5	INT(10)
°page_id=10	INT(10)
°page_inflection_id=7	INT(10)

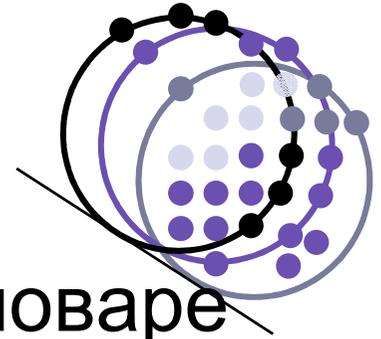
page_inflection	
* <u>id=7</u>	INT(10)
°page_id=10	INT(10)
°inflection_id=6	INT(10)
°term_freq=1	INT(6)

page	
* <u>id=10</u>	INT(10)
°page_title="работа"	VARCHAR(255)

inflection	
* <u>id=6</u>	INT(10)
°freq	INT(11)
°inflected_form="работе"	VARCHAR(255)



Внутренние ссылки (2): ?



- Частота конкретных форм слова в словаре
 - `page_inflection . term_freq`
- Информация о **ссылках / ключевых словах** толкования на другие словарные статьи
 - *в поиск-х алг. (поиск синонимов)*
- Слова, для которых есть **ссылки**, но нет словарных статей – всё равно добавляются в таблицу «**page**».

Значение

Пр.	рабо́чем	рабо́чей
-----	----------	----------

1. относящийся к работе ◆ Запиши мой рабочий телефон.
2. относящийся к рабочему (рабочим) I ◆ Рабочий посёлок. ◆ Она-де горожанка, руки у нее не знает, что делать в огороде, со скотиной не умеет обращаться, доить не умеет, и при этом пер; заметно оживилась и повеселела. *Альфред Хейдок, «Грешница», 1924-1934 г.* (цитата из Национального корпуса литературы)
3. находящийся в работе; предварительный, пробный, черновой ◆ Рабочий вариант.
4. способный нормально работать, функционировать ◆ У тебя случайно не заваялся какой-нибудь старенький, но чтобы рабочий. ◆ А что, этот план вполне рабочий!

Синонимы

1. частичн.: трудовой
2. -
3. предварительный, пробный, черновой
4. работоспособный

Антонимы

Гиперонимы

Гипонимы

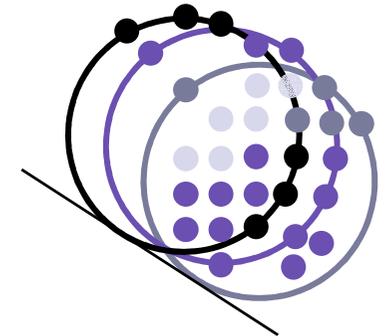
1. офисный
2. пролетарский

Семантические отношения

Значение

Пф

1. относящийся к работе ◆ Запиши мой рабочий телефон.
2. относящийся к рабочему (рабочим) | ◆ Рабочий посёлок. ◆ С знает, что делать в огороде, со скотиной не умеет обращаться заметно оживилась и повеселела. Альфред Хейдок, «Грешница», 19 литературы)



3. находящийся в ра
4. способный нормал старенький, но что

Синонимы

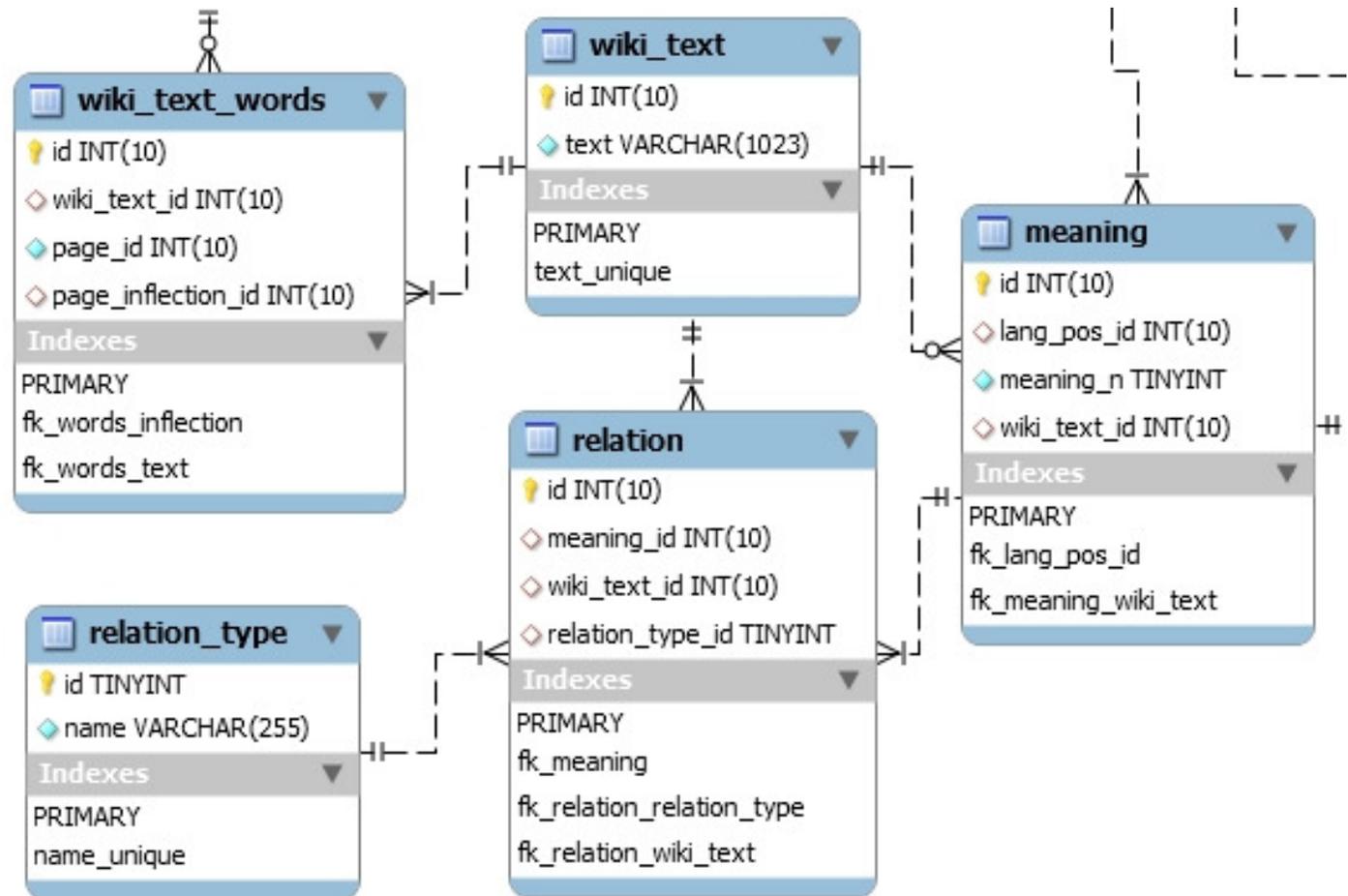
1. частичн.: трудовой
2. -
3. предварительный,
4. работоспособный

Антонимы

Гиперонимы

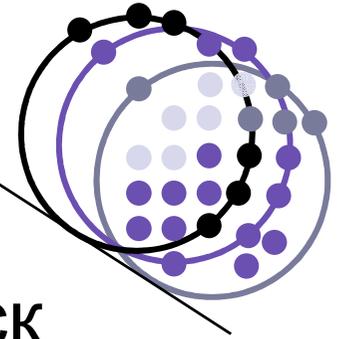
Гипонимы

1. офисный
2. пролетарский





[[Категория:Имя категории]]



Цель: Автоматизация оглавления, поиск

- Грамматические категории

- Часть речи
- Тип словоизменения
- Одушевлённость
- Грамматический род

- Стилистические свойства

- Служебные

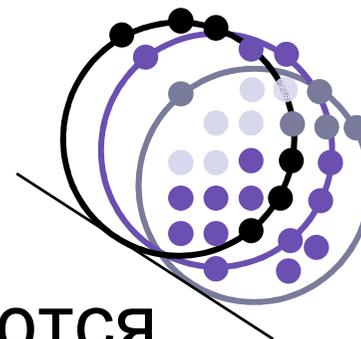
- Семантика

- {{categ|Работа и труд|Рабочие|lang=}}

Вшиты
в
шаблоны



{{Шаблонизация}} всей страны!

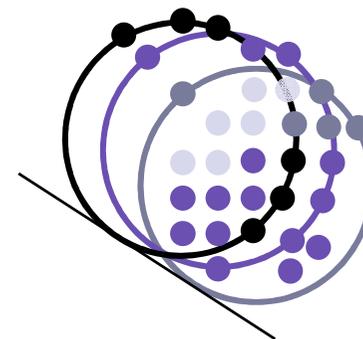


- + Шаблоны автоматически проставляются ботами при создании статьи
- + Централизованная смена внешнего вида сразу многих статей
- + Автоматизация редактирования (ботами, парсером), т.к. есть разметка спец-ми конструкциями
- + Автоматизация категоризации
 - × {{сущ ru m ina}} → категории «Мужской род» и «Неодушевлённые»
- Сложность освоения ☺



{{Шаблонизация}}

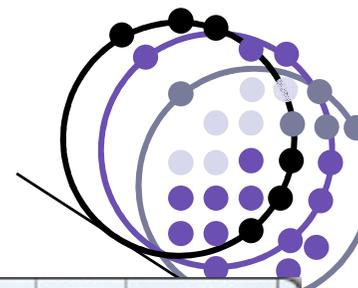
Примеры



- × {{-ru-}}, {{пример|}}
- × Фонетические: {{transcriptions|}}, {{медиа}}
- × {{морфо|под|вод|и|ть|ся}}
- × Морфологические:
 - × {{сущ ru}}, {{прил ru}}, {{сущ ru m ina}}, {{adv ru}}
 - × {{сущ eo}}, {{прил eo}}, {{adv eo}}, {{гл eo}}
- × Шаблоны библиографии
 - × {{НКРЯ}}, {{Ушаков1940}}, {{Эпитет1913}}
- × Технические (из Википедии):
 - × {{За}}, {{wikify}}, «вавилонские шаблоны»



Быстрый поиск на заданном языке (1)



Wiwordik 0.03 (ruwikt20091228_parsed)

l??kk

Source language

Meaning Semantic Relation Translation

Translation language

- aallokko
- kolpakko
- kynsileikkuri
- laukka
- leikki
- linkku
- luokka
- alpakka -> альпака
- balalaikka -> балалайка**
- maljakka -> ваза
- kukkamaljakko -> ваза
- ylipäällikkö -> главнокомандующ...
- leikkuuruumuri -> зерноуборочн...
- leikki -> игра
- kalmukki -> калмыцкий язык
- päällikkö -> командир
- sotapäällikkö -> полководец
- leikkuuruumuri -> хлебоуборочн...
- matkalaukku -> чемодан

балалайка

балалайка (-)

Russian (ru)

noun

1. {{илл|TenorBalalaika1.jp

Portuguese (pt): balalaica
Slovenian (sl): balalajka
Finnish (fi): balalaikka
French (fr): balalaïka <i>f</i>
Swedish (sv): balalajka
Esperanto (eo): balalajko

leikki (-)

Finnish (fi)

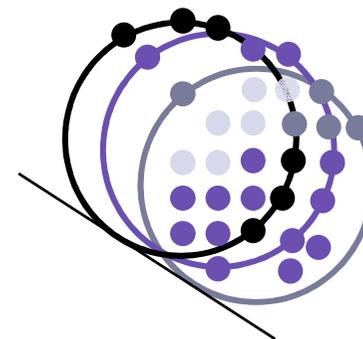
noun

игра

2. русский народный трехструнный щипковый музыкальный инструмент



Быстрый поиск на заданном языке (2)



index_native	
🔑	page_id INT(10)
💎	page_title VARCHAR(255)
💎	has_relation TINYINT(1)
Indexes ▶	

1 таблица

index_uk	
🔑	id INT(10)
💎	foreign_word VARCHAR(255)
💎	foreign_has_definition TINYINT(1)
💎	native_page_title VARCHAR(255)
Indexes ▶	

+ ещё 561 таблица

Список кодов языков: <http://ru.wiktionary.org/wiki/шаблон:перев-блок>

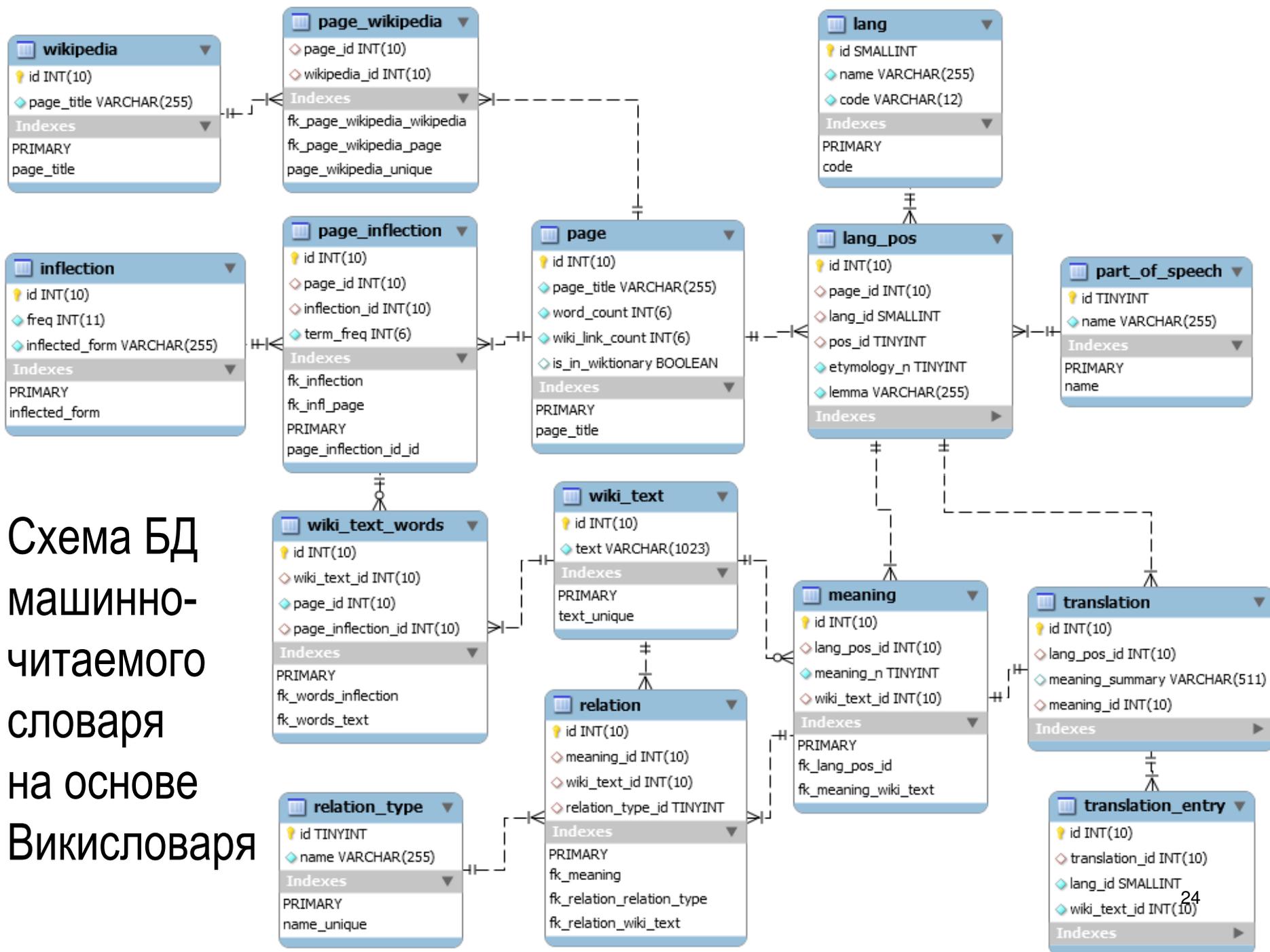
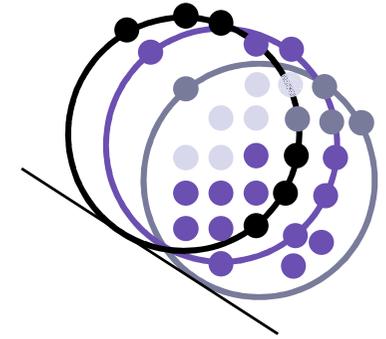
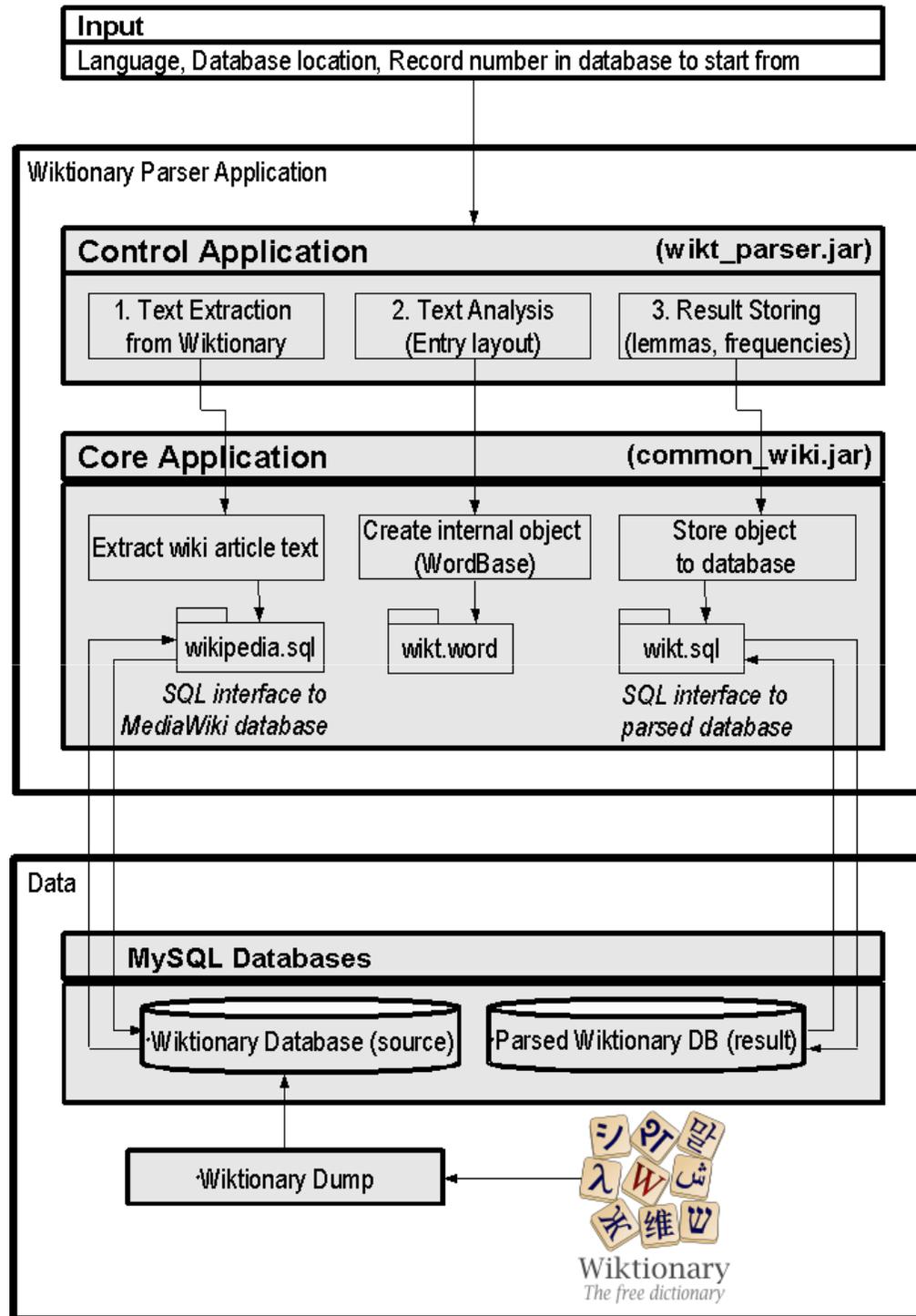


Схема БД
машинно-
читаемого
словаря
на основе
Викисловаря

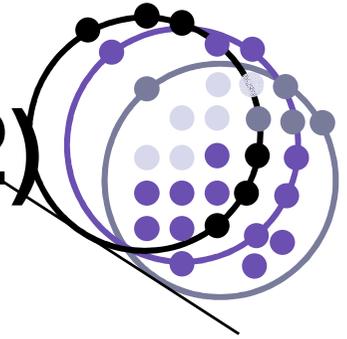


А р х и т е к т у р а





Регулярные выражения (2)



1. `====?\s*Значение\s*====?\s*\n`

2. `(?m)^\s*([\^=]+?)\s*^\s*`

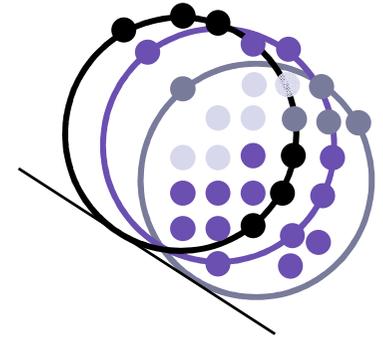
- `==`рабочий `|``==`, `==` Существительное `|` `==`

3. `#(REDIRECT|ПЕРЕНАПРАВЛЕНИЕ) \[\[\[(.+?) \]\]\]`

`#ПЕРЕНАПРАВЛЕНИЕ [[нелётный]] -> нелётный`



Перенаправления



1. Указание основной формы слова

- маня -> манить

2. Подсказка об ошибке

всё-равно -> всё равно

всё-равно



Такого слова не существует!

Вы, возможно, имели в виду **всё равно?**

3. Диакритики

- зверье -> зверьё, соеур -> соеур

4. Неточная кодировка

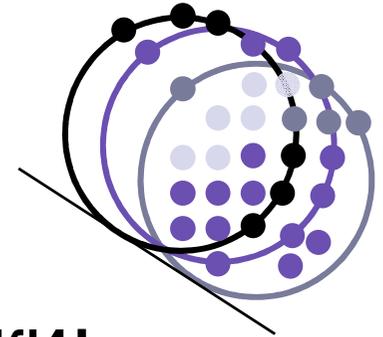
- ліс -> ліс (І латиница, І кириллица)

5. Со строчной буквы на прописную москва -> Москва

Категория: Ашипки

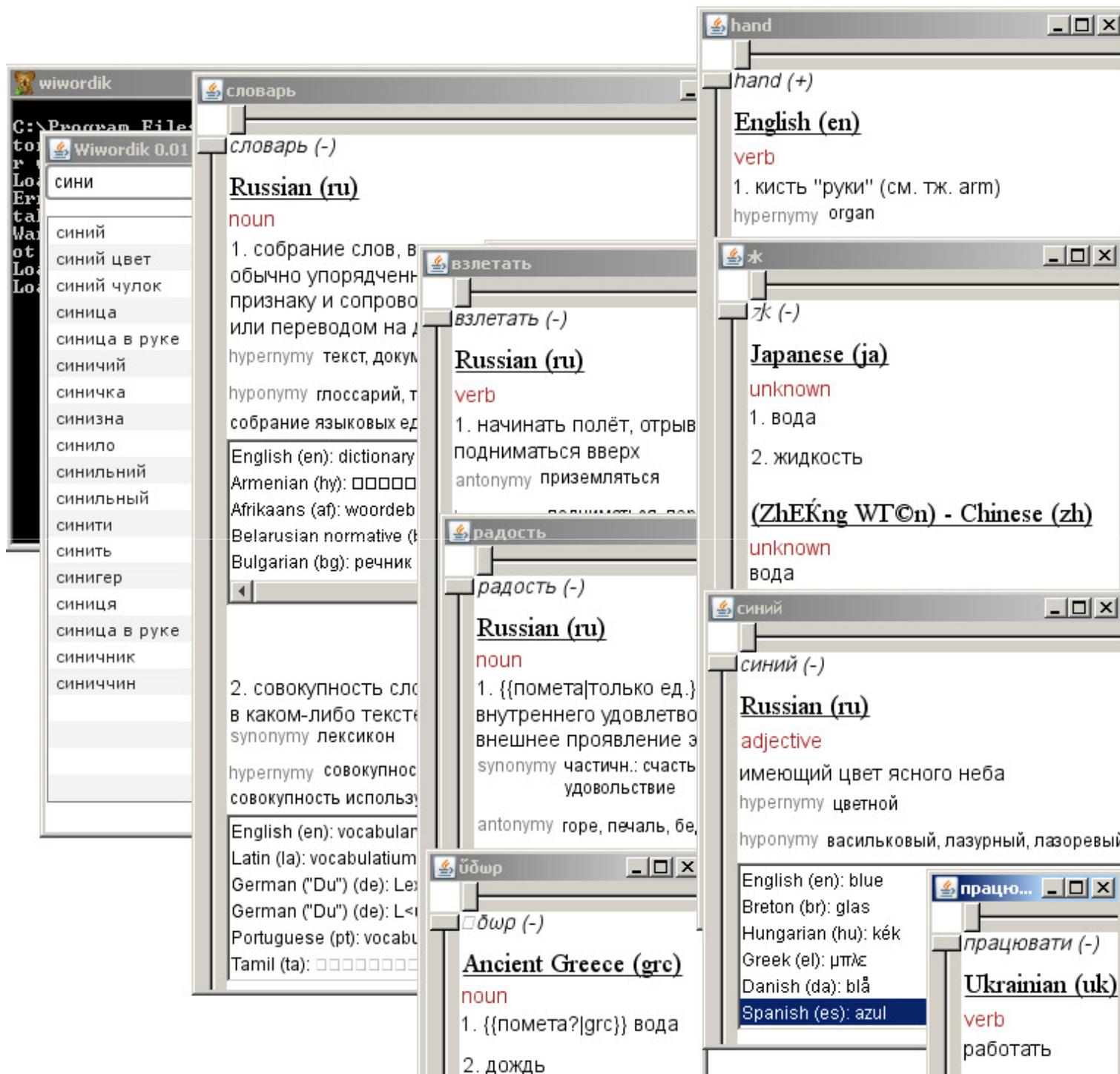


Реализация 1



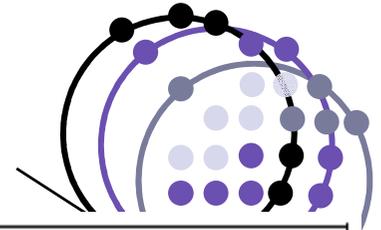
- Программный код включает наработки:
 - synarcher – поиск синонимов в Википедии
 - wikidf – индексирование текстов Википедии
- Java
- База данных:
 - MySQL - для разработки и тестирования
 - SQLite – в скачиваемом приложении
- JUnit тестирование

Р е а л и з а ц и я



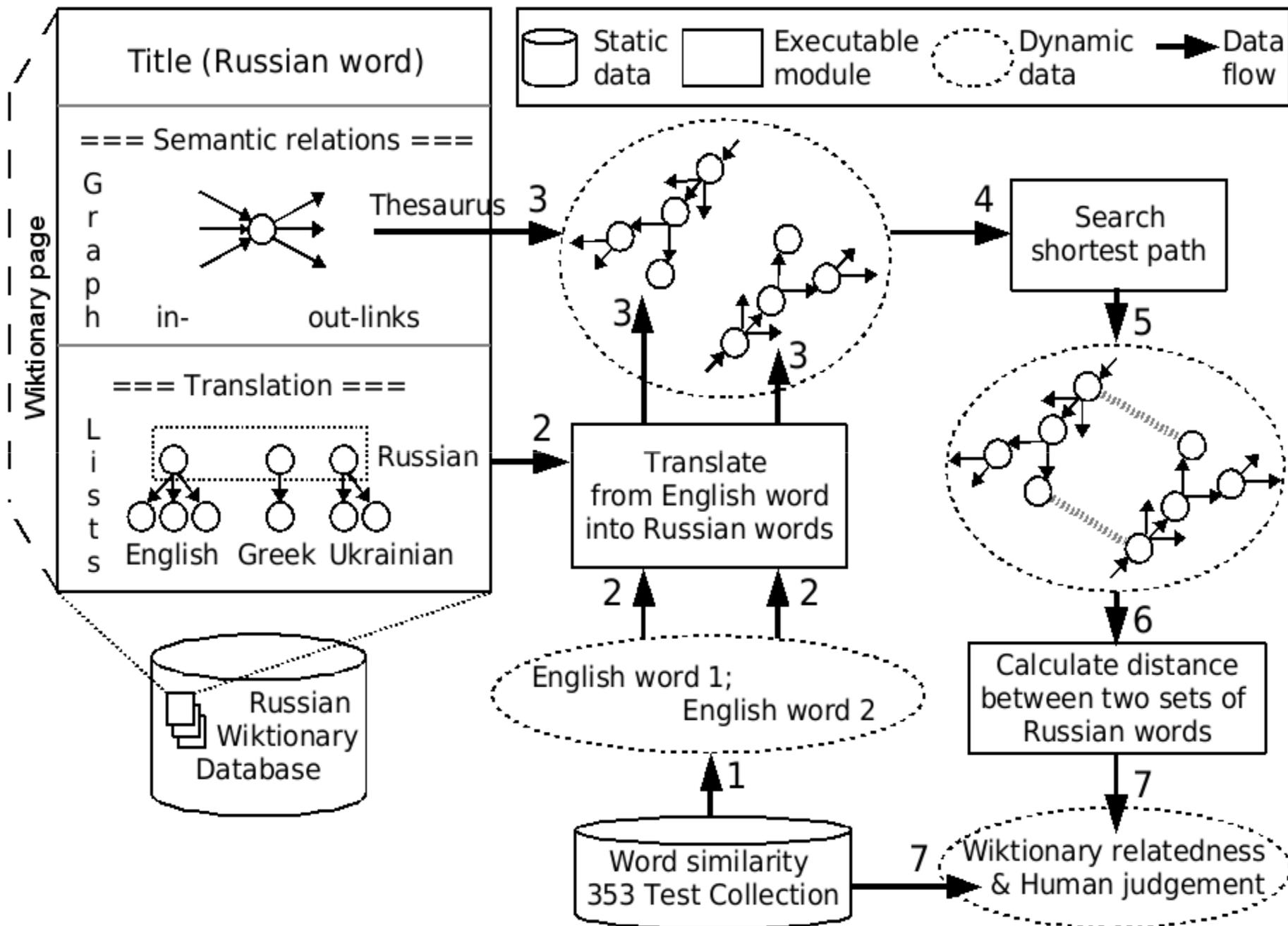


Размеры Викисловарей



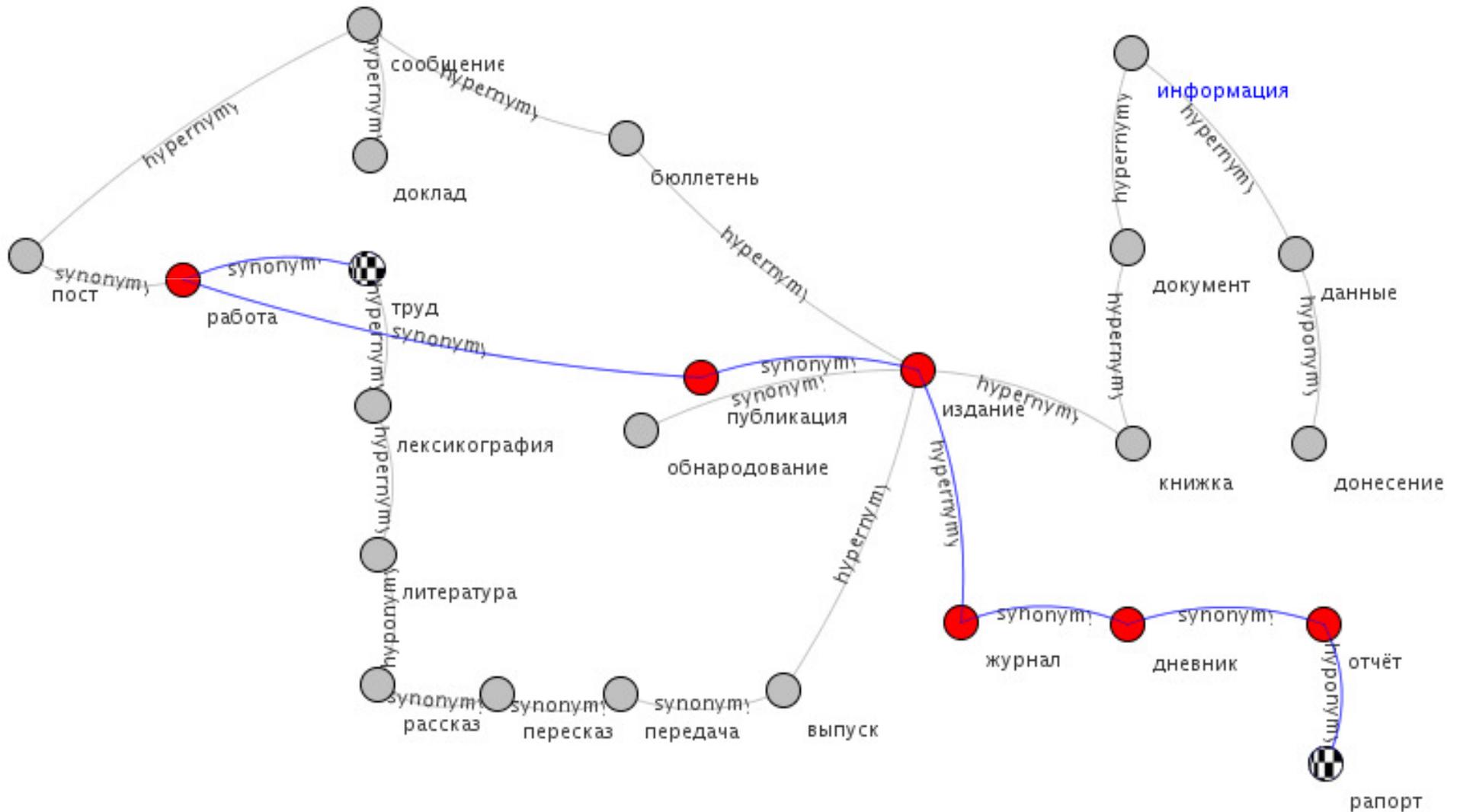
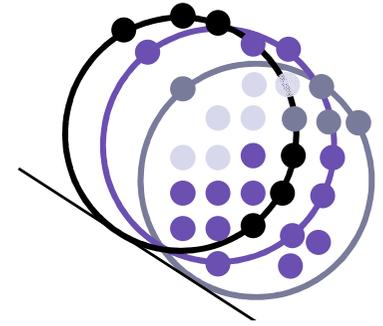
	Wiktionary editions as of September 2007, from [24]				A part of Wiktionary extracted by the parser. Wiktionary edition as of January 2009.				
	English Wiktionary		German Wiktionary		Russian Wiktionary				
	English	German	English	German	Total ⁵	English	German	Russian	Ukrainian
Entries	176,410	10,487	3,231	20,557	247,580	2,813 ⁶	13,072	124,301	88,575
Part of speech (POS)									
Nouns	99,456	6,759	2,116	13,977	108,448	935	336	58,843	40,607
Verbs	31,164	1,257	378	1,872	26,290	342	49	356 ⁷	24,096
Adjectives	23,041	1,117	357	2,261	26,864	184	18	2,168	23,536
Unknown	POS which were not recognized by the parser				80,293	1,321	12,648	57,573	331
Semantic relations									
Synonyms	29,703	1,916	2,651	34,488	28,718	1,345	665	24,338	310
Antonyms	4,305	238	283	10,902	10,480	238	234	9,062	54
Hypernyms	42	0	336	17,286	18,975	444	474	17,033	115
Hyponyms	94	0	390	17,103	8,585	176	473	7,574	12
Holonyms	–	–	–	–	216	1	0	215	0
Meronyms	–	–	–	–	322	8	2	306	0
Total	–	–	–	–	67,296	2,212	1,848	58,528	491

WordNet (2006): 150,000 слов, 115,000 синсетов (наборов синонимов)



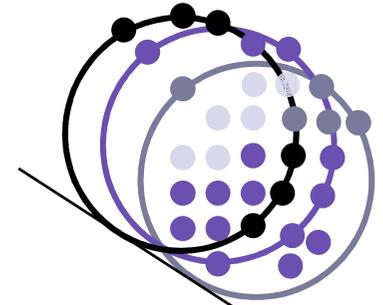


Кратчайший путь в Русском Викисловаре





Корреляция мер семантической близости

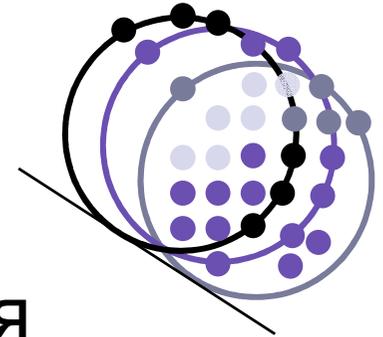


Dataset	WN	WP	WT	Others
Metric or Algorithm	I. Path based measures (in taxonomy)			
wup	0.3	0.47	–	–
lch	0.34	0.48	–	–
res _{hypo}	–	0.25-0.37 ⁸	–	–
jarmasz	–	–	–	0.539 RT⁹
path ^{max} _{len}	–	–	0.24	–
	II. Words frequency in corpus			
jaccard	–	–	–	Google 0.18
res	0.34	–	–	–
LSA	–	–	–	IntelliZap 0.56
ESA	–	0.75	–	–
	III. Text overlapping			
lesk	0.21	0.2	–	–
text	–	0.19	–	–

Корреляция мер семантической близости слов:
1) значения экспертов (набор 353-ТС),
2) значения вычислены автоматически на основе WordNet, Английской Википедии, Русского Викисловаря



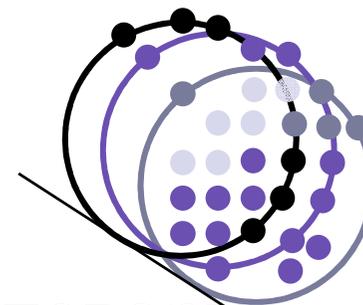
Результаты



- Создан парсер Русского Викисловаря
 - Спроектирована схема БД
 - Реализован доступ к БД (API, Java)
- Выполнено сравнение результатов поиска семантически близких слов на основе Викисловаря и тезауруса WordNet
- Сайт проекта (Wiki tool kit)
 - <http://code.google.com/p/wikokit/>

Сделано и *ещё делать*

(схема БД, парсер)



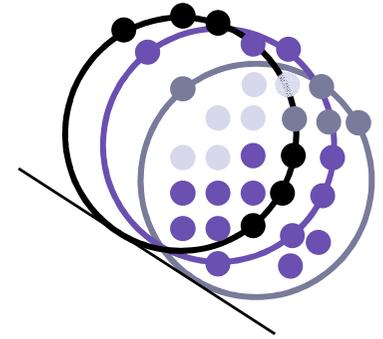
- Извлекаются (RE)
 - Толкование
 - определение
 - *помета, цитата, картинка*
 - Отношение
(синонимы..., *помета*)
 - Перевод
 - *Фонетика*
 - *Транскрипция, Аудио*
 - *Этимология*
 - *Фразеологизмы, поговорки, пословицы*
 - ...

Русский Викисловарь
English Wiktionary

- *Уровни*
 - *Схема БД (+ table)*
 - *API Базы данных*
 - *Код (+ class, RE)*



Планы



- Продолжить разработку MRD
 - Нарращивание функц-ти парсера, отладка
 - + **English Wiktionary**
- Визуализация (JavaFX)
 - MRD браузер
 - Игры и тесты (изучение иностранных языков)

Спасибо за внимание!



<http://ru.wiktionary.org/>