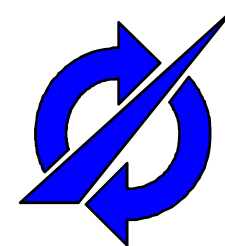


# Машинная обработка Русского Викисловаря



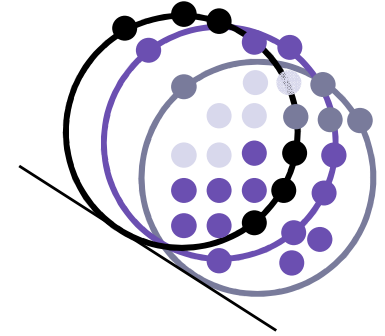
Санкт-Петербургский институт  
информатики и автоматизации РАН



Крижановский Андрей (andrew.krizhanovsky  gmail.com)



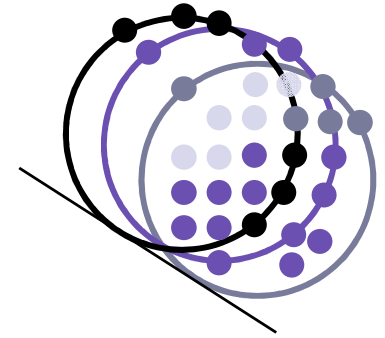
# Содержание



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



# Цель



Сделать возможным

применение данных Викисловаря

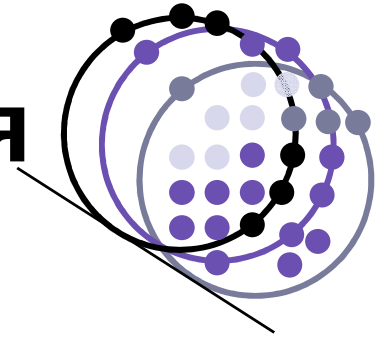
(как лингвистического ресурса)

в различных компьютерных программах,

в задачах, связанных с обработкой текста.



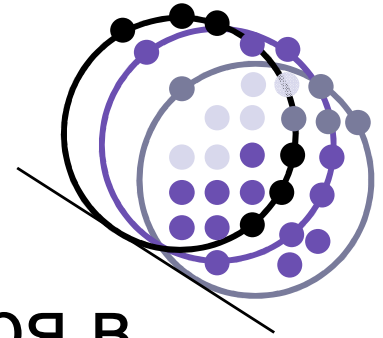
# Применение Викисловаря



- 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)

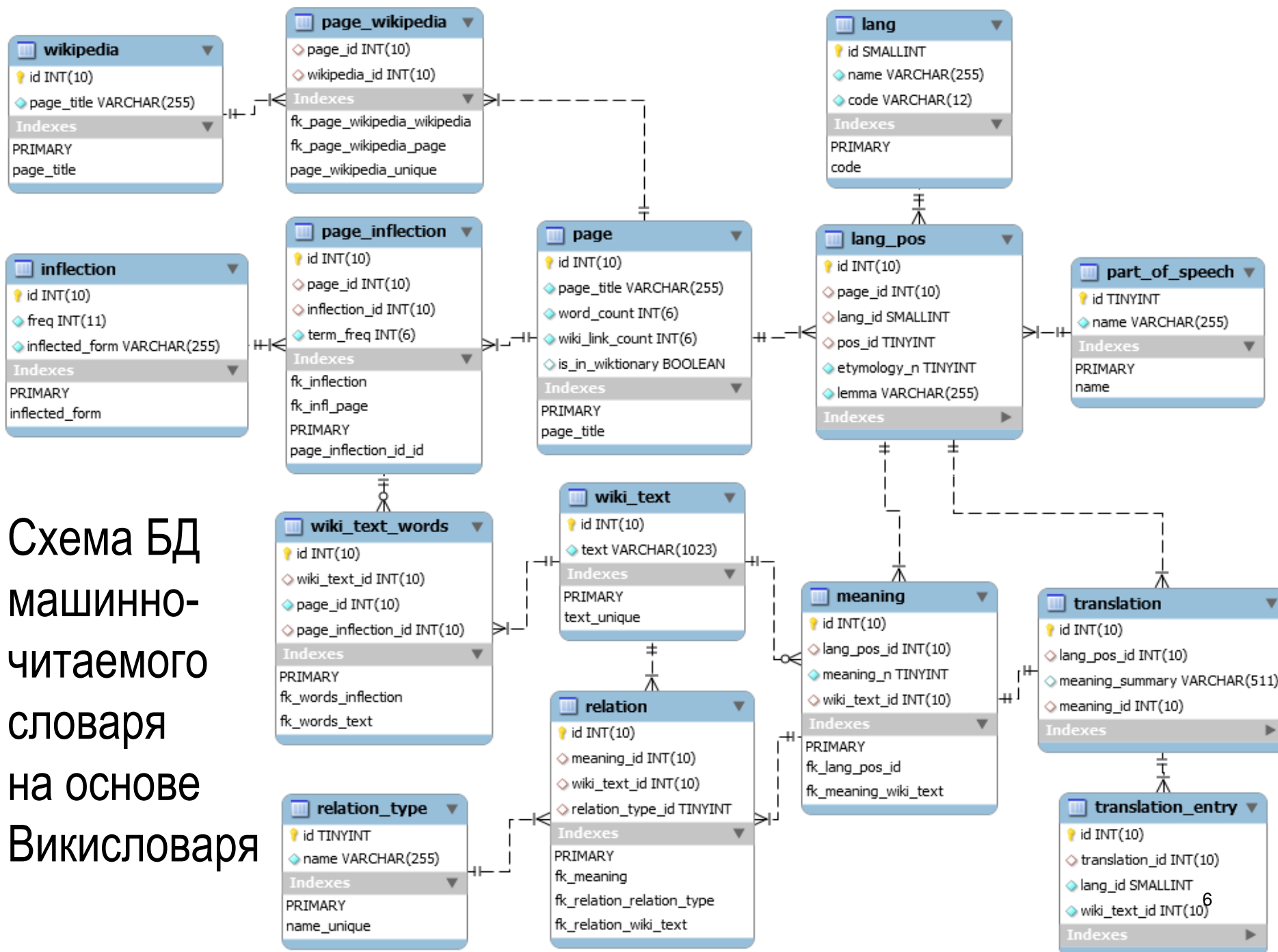
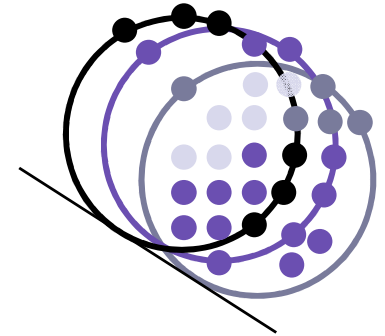


Схема БД  
машинно-  
читаемого  
словаря  
на основе  
Викисловаря



# 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

# 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

a multilingual tree  
encyclopedia

**Wiktionary**  
[ˈwɪkʃənəri] *n.*,  
a wiki-based Open  
Content dictionary

Wikeo [ˈwɪl kəʊ]

**Türkçe**

*Özgür sözlük*  
208 000+ madde

**Ido**

*La libera vortaro*  
137 000+ artikli

**Русский**

*Свободный словарь*  
137 000+ статей

**中文**

*自由的多语言词典*  
116 000+ 词条

**Ελληνικά**

*Το Ελεύθερο Λεξικό*  
107 000+ λέξεις

**தமிழ்**

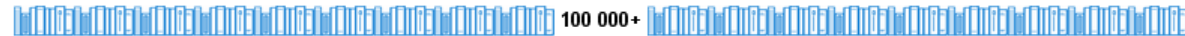
*கட்டற்ற அகரமுதலி*  
102 000+ கட்டுரைகள்

**Polski**

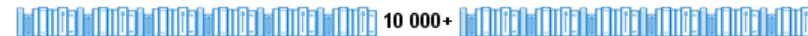
*Wolny słownik*  
93 000+ stron

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

English >



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 • କଞ୍ଚକାଠ • Latina • മലയാളം • Bahasa Melayu • Norsk (bokmål) • Occitan • Қазақша • Sesotho • Shqip • Simple English • Slovenčina • Slovenščina • Kiswahili • Tatarça / татарча • Ἰου • Українська • اردو



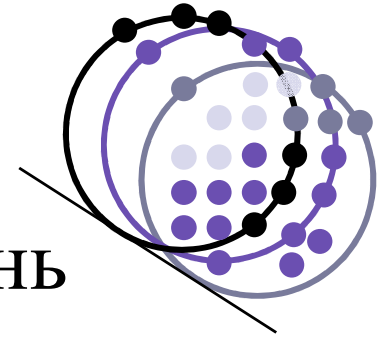
አማርኛ • Aragonés • Avañe'ê • Azərbaycan • Беларуская • Bosanski • Cymraeg • Euskara • Føroyskt • Gàidhlig • ગુજરાતી • Interlingue • ଶାସ୍ତ୍ରୀୟ • சமீப • Kinyarwanda • Кыргызча • Latviešu • Македонски • मराठी • монгол • Nāhuatlāhtōlli • पञ्जाबी • Plattdütsch • Runa Simi • سنڌي • Basa Sunda • Tagalog • ཇལ་ཅེ • GŵY • Xitsonga • ئۇيغۇرچە • Wolof • עברית • isiZulu

Other languages





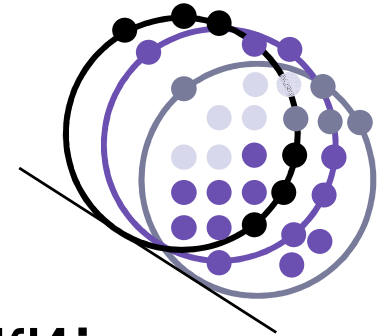
# Данные Викисловаря: плюсы и трудности



- + Богатство
    - + тезаурус  
(синонимы, антонимы...)
    - + словосочетания
    - + 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. даль<sup>9</sup>ше)

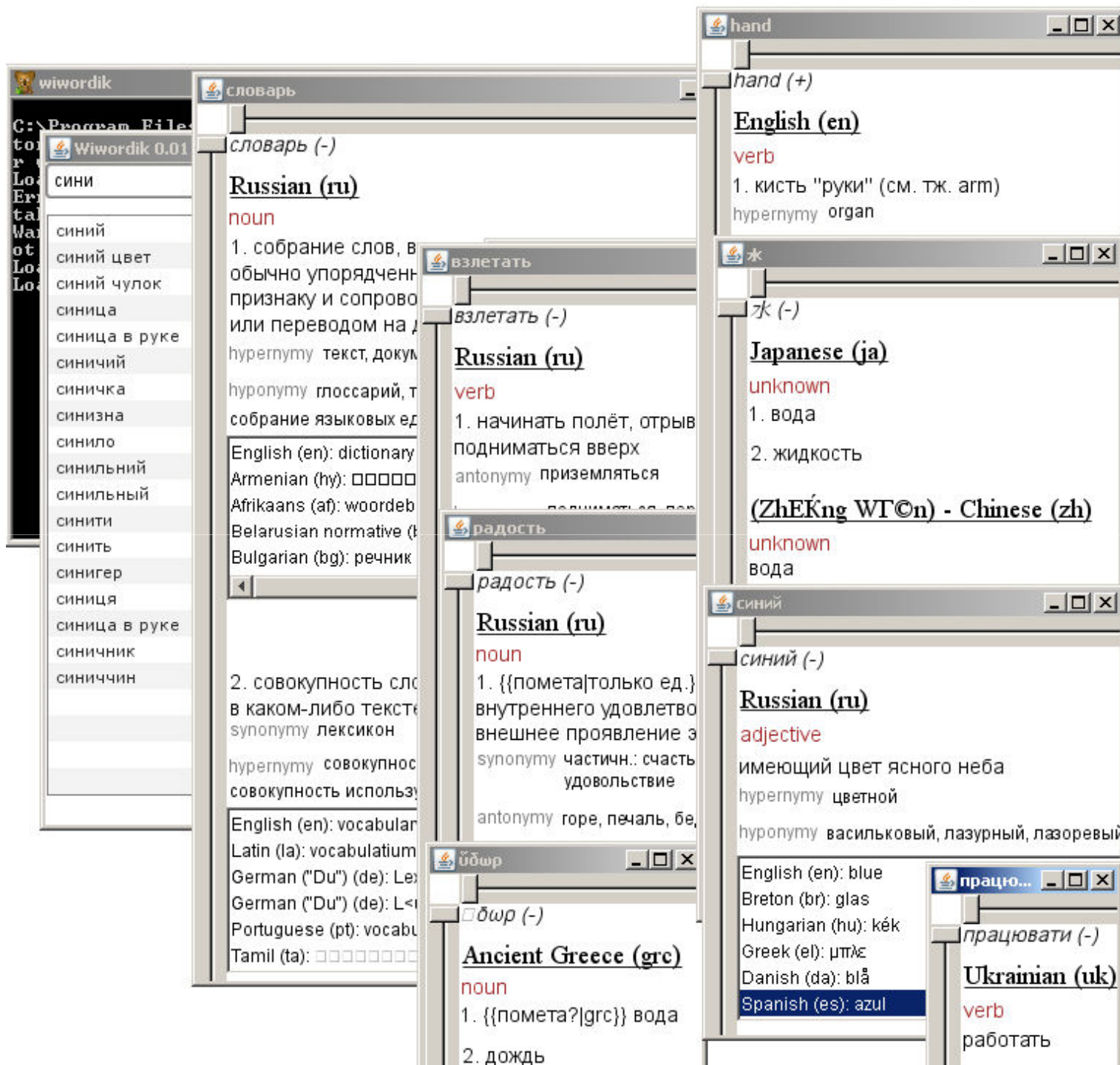


# Реализация 1



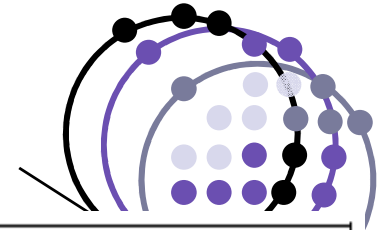
- Программный код включает наработки:
  - synarcher – поиск синонимов в Википедии
  - wikidf – индексирование текстов Википедии
- Java
- База данных:
  - MySQL - для разработки и тестирования
  - SQLite – в скачиваемом приложении
- JUnit тестирование

Р  
е  
а  
л  
и  
з  
а  
ц  
и  
я



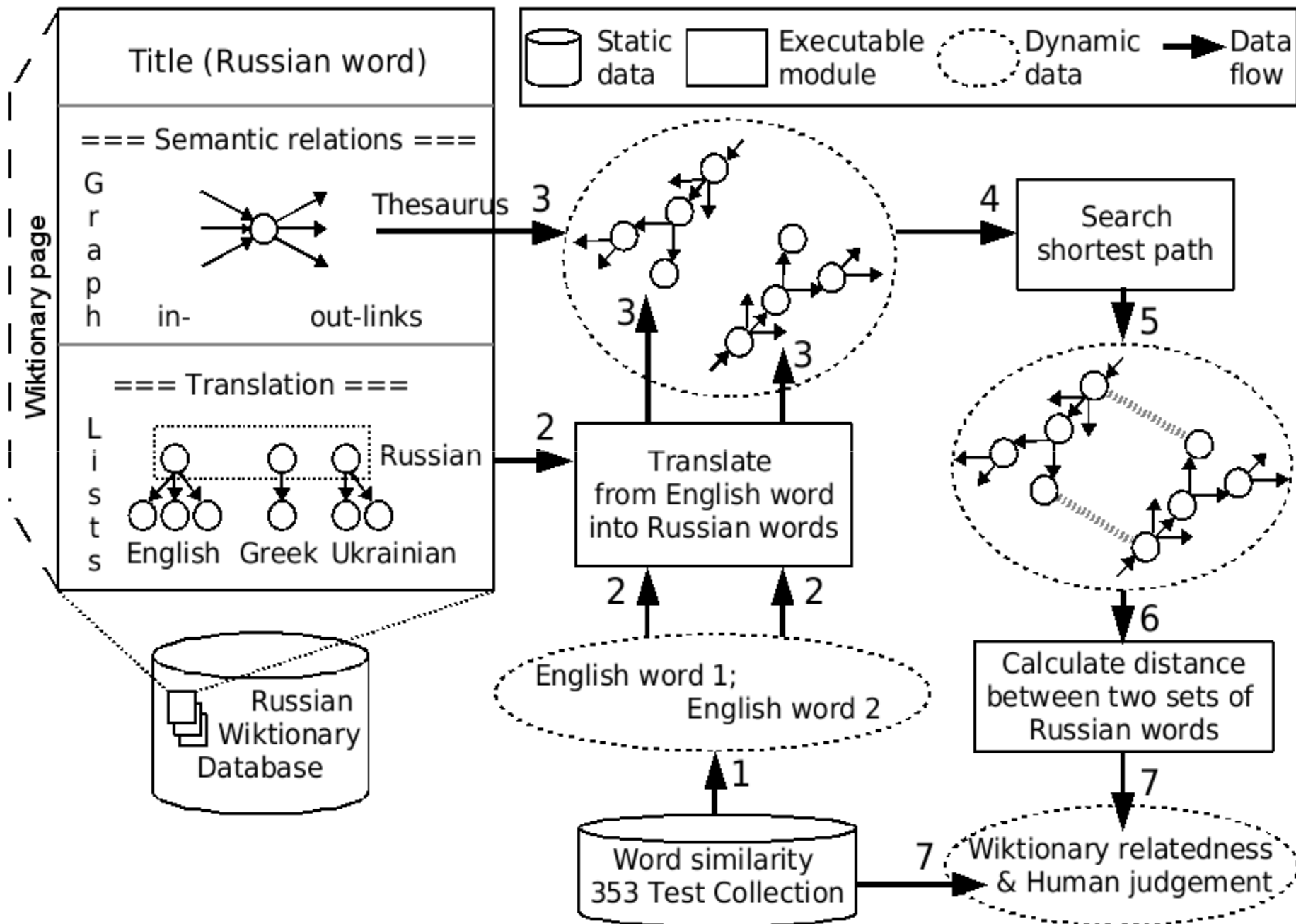


# Размеры Викисловарей



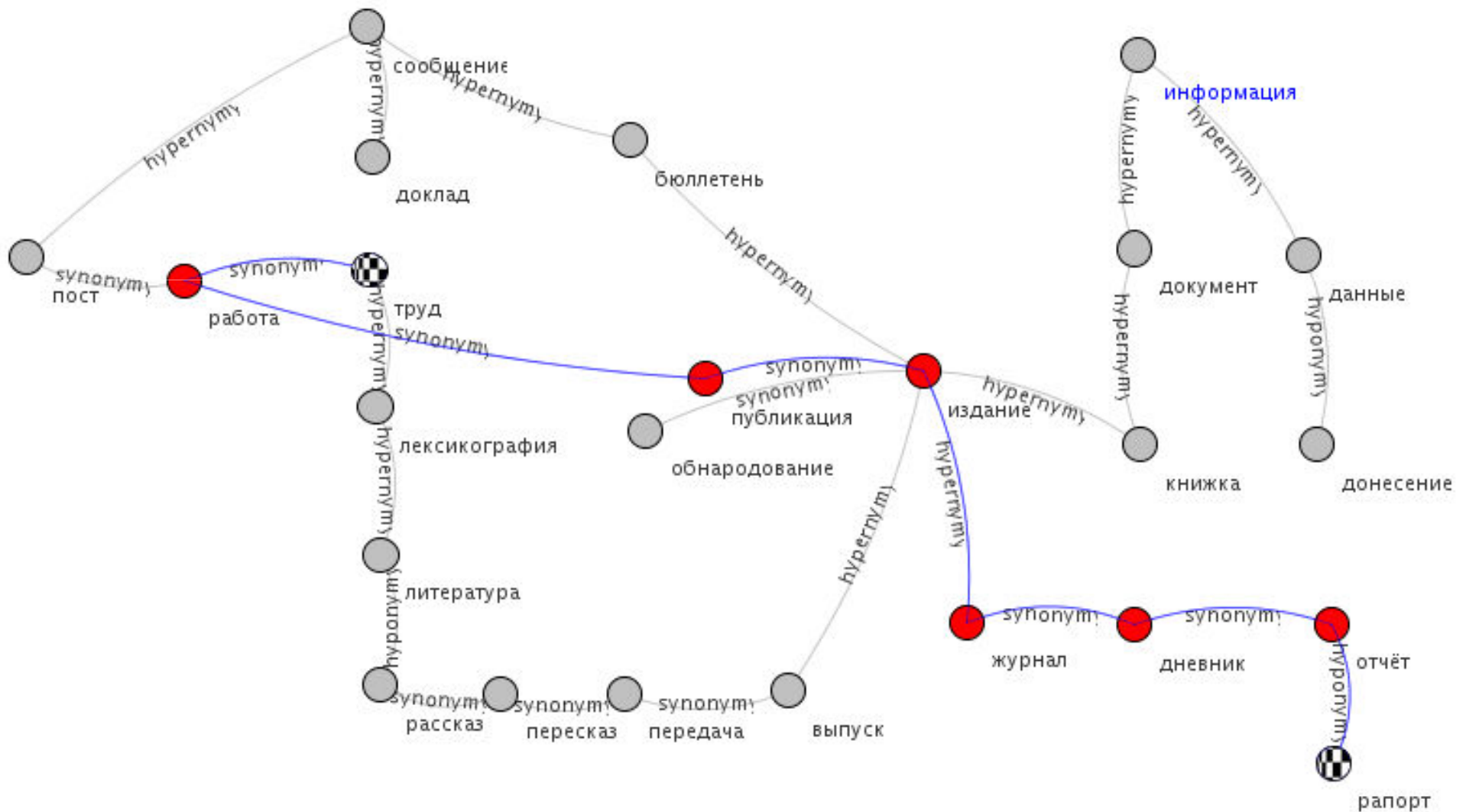
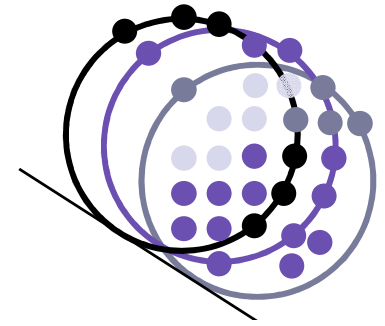
	Wiktionary editions as of September 2007, from [24]				A part of Wiktionary extracted by the parser. Wiktionary edition as of January 2009.				
	English Wiktionary		German Wiktionary		Russian Wiktionary				
	English	German	English	German	Total <sup>5</sup>	English	German	Russian	Ukrainian
Entries	176,410	10,487	3,231	20,557	247,580	2,813 <sup>6</sup>	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 <sup>7</sup>	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	<b>28,718</b>	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	–	–	–	–	<b>67,296</b>	2,212	1,848	58,528	491

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



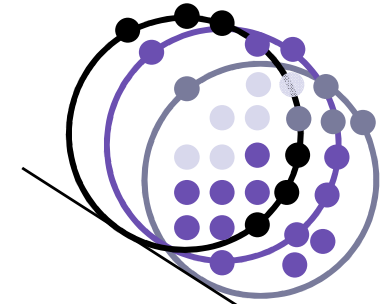


# Кратчайший путь в Русском Викисловаре





# Корреляция мер семантической близости

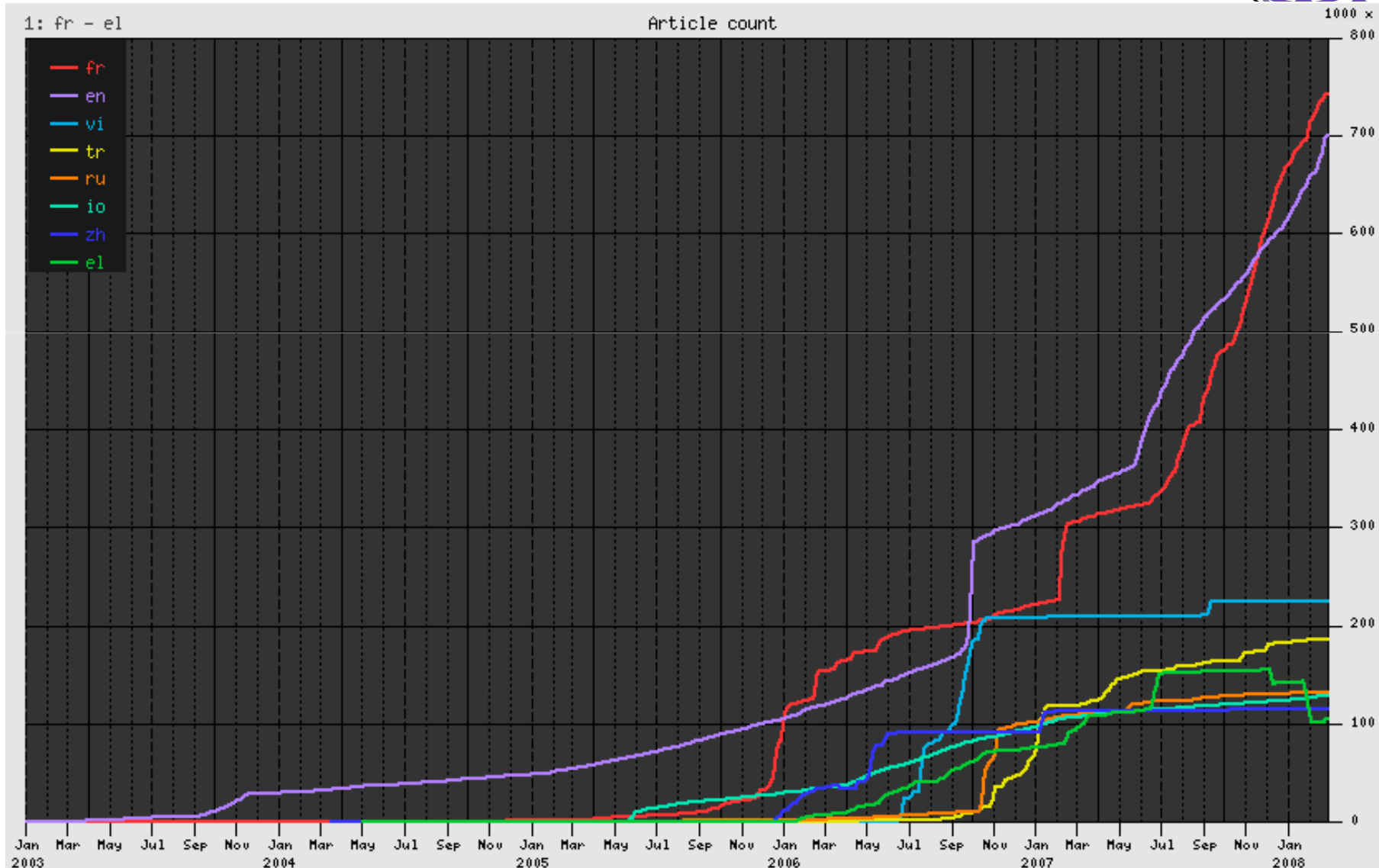
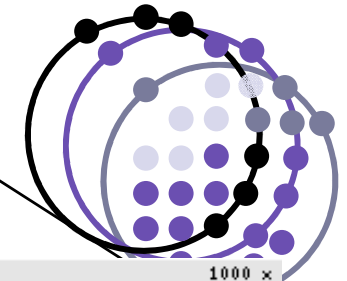


Dataset	WN	WP	WT	Others
<b>Metric or Algorithm</b>	I. Path based measures (in taxonomy)			
wup	0.3	0.47	–	–
lch	<b>0.34</b>	<b>0.48</b>	–	–
res <sub>hypo</sub>	–	0.25-0.37 <sup>8</sup>	–	–
jarmasz	–	–	–	<b>0.539 RT<sup>9</sup></b>
path <sup>max</sup> <sub>len</sub>	–	–	<b>0.24</b>	–
	II. Words frequency in corpus			
jaccard	–	–	–	Google 0.18
res	0.34	–	–	–
LSA	–	–	–	IntelliZap 0.56
ESA	–	<b>0.75</b>	–	–
	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



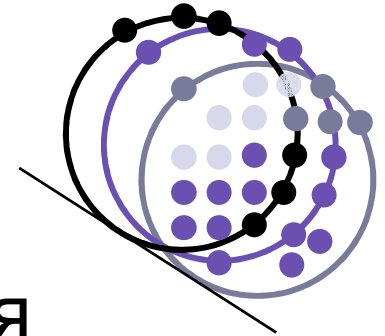
# Восемь самых больших Викисловарей (Март 2008)







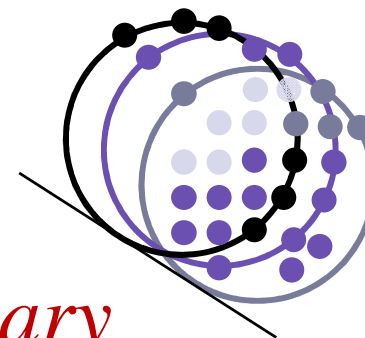
# Результаты



- 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)



## Русский Викисловарь

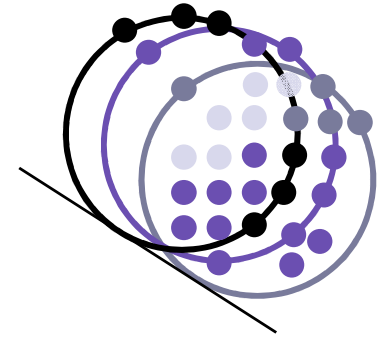
- Извлекаются (RE)
  - Толкование
    - определение
    - помета, цитата, картинка
  - Отношение (синонимы..., помета)
  - Перевод
  - *Фонетика*
    - Транскрипция, Аудио
  - *Этимология*
  - ...
- *API Базы данных*

## *English Wiktionary*

- *Извлечь*
  - Толкование
  - *Отношение* (синонимы...)
  - Перевод
  - ...
- *API Базы данных*



# Планы



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

**Спасибо за внимание!**

