

# Основы обработки текстов

Лекция 10

Лексическая семантика  
Извлечение информации

# Возможные взгляды на семантику

- Лексическая семантика
  - значение индивидуальных слов
- Композиционная семантика
  - как значения комбинируются и определяют новые значения для словосочетаний
- Дискурс или прагматика
  - как значения комбинируются между собой и другими знаниями, чтобы задать значение текста или дискурс

# План

- Основные понятия
  - слова и отношения между ними
  - словари и тезаурусы
- Вычислительная семантика
  - Разрешение лексической многозначности
  - Привязка к базам знаний

# ОСНОВНЫЕ ПОНЯТИЯ

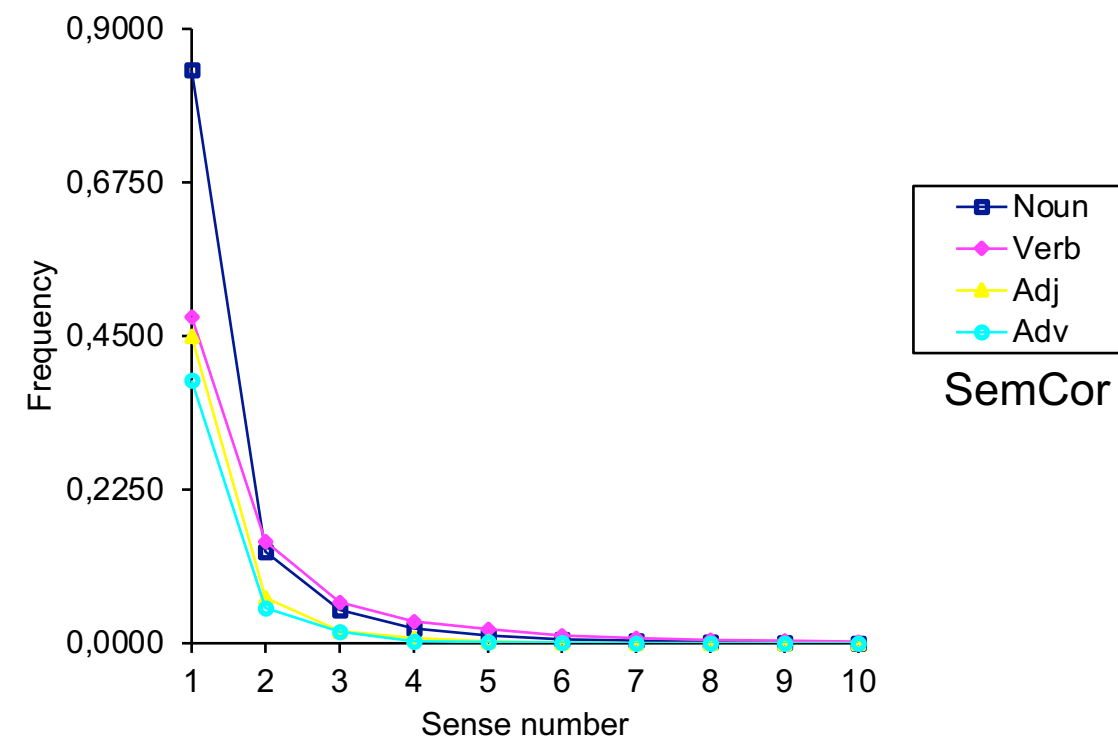
- Значение слова и многозначность
- Омонимия VS многозначность
  - ключ
  - платформа
- Метонимия
  - Я три *тарелки* съел
- Зевгма
  - За окном шел снег и рота красноармейцев
- Типы омонимов
  - омофоны (луг-лук, плод-плот)
  - омографы (м'ука - мук'а, гв'оздик-гвозд'ик)

# Отношения между словами

- Синонимия
  - Машина / автомобиль
- Антонимия
  - большой / маленький, вверх / вниз, ложь / истина
- Обобщение и детализация (hyponym and hypernym/superordinate)
  - машина - транспортное средство
  - яблоко - фрукт
- Меронимы (партонимы) и холонимы
  - колесо - машина

# Многозначность на практике

- Text-to-Speech  
–омографы
- Информационный поиск
- Извлечение информации
- Машинный перевод
- Закон Ципфа (Zipf law)



# Основные вопросы

- Что такое значение?
  - Сколько значений у слова “платформа”?
- Что нужно для того, чтобы понять значение?
  - Что такое контекст?
  - Как определить связь значения с контекстом?
    - Разреженность языка

# Базы знаний

- (В широком смысле) **База знаний** — база данных, содержащая правила вывода и информацию о человеческом опыте и знаниях в некоторой предметной области
- (В узком смысле) **База знаний** — некоторое структурированное описание предметной области



# WordNet

- База лексических отношений
  - содержит иерархии
  - сочетает в себе тезаурус и словарь
  - доступен on-line
  - разрабатываются версии для языков кроме английского (в т.ч. для русского)

| Категория       | Уникальных форм |
|-----------------|-----------------|
| Существительные | 117,097         |
| Глаголы         | 11,488          |
| Прилагательные  | 22,141          |
| Наречия         | 4,601           |

❏ <http://http://wordnet.princeton.edu/>

❏ <http://wordnet.ru/>

# Формат WordNet

The noun “bass” has 8 senses in WordNet.

1. bass<sup>1</sup> - (the lowest part of the musical range)
2. bass<sup>2</sup>, bass part<sup>1</sup> - (the lowest part in polyphonic music)
3. bass<sup>3</sup>, basso<sup>1</sup> - (an adult male singer with the lowest voice)
4. sea bass<sup>1</sup>, bass<sup>4</sup> - (the lean flesh of a saltwater fish of the family Serranidae)
5. freshwater bass<sup>1</sup>, bass<sup>5</sup> - (any of various North American freshwater fish with lean flesh (especially of the genus Micropterus))
6. bass<sup>6</sup>, bass voice<sup>1</sup>, basso<sup>2</sup> - (the lowest adult male singing voice)
7. bass<sup>7</sup> - (the member with the lowest range of a family of musical instruments)
8. bass<sup>8</sup> - (nontechnical name for any of numerous edible marine and freshwater spiny-finned fishes)

The adjective “bass” has 1 sense in WordNet.

1. bass<sup>1</sup>, deep<sup>6</sup> - (having or denoting a low vocal or instrumental range)  
*”a deep voice”; ”a bass voice is lower than a baritone voice”;  
”a bass clarinet”*

# WordNet: отношения между словами

| Relation       | Also called   | Definition                                | Example   |
|----------------|---------------|---|---|
| Hypernym       | Superordinate | From concepts to superordinates           | <i>breakfast</i> <sup>1</sup> → <i>meal</i> <sup>1</sup>    |
| Hyponym        | Subordinate   | From concepts to subtypes                 | <i>meal</i> <sup>1</sup> → <i>lunch</i> <sup>1</sup>        |
| Member Meronym | Has-Member    | From groups to their members              | <i>faculty</i> <sup>2</sup> → <i>professor</i> <sup>1</sup> |
| Has-Instance   |               | From concepts to instances of the concept | <i>composer</i> <sup>1</sup> → <i>Bach</i> <sup>1</sup>     |
| Instance       |               | From instances to their concepts          | <i>Austen</i> <sup>1</sup> → <i>author</i> <sup>1</sup>     |
| Member Holonym | Member-Of     | From members to their groups              | <i>copilot</i> <sup>1</sup> → <i>crew</i> <sup>1</sup>      |
| Part Meronym   | Has-Part      | From wholes to parts                      | <i>table</i> <sup>2</sup> → <i>leg</i> <sup>3</sup>         |
| Part Holonym   | Part-Of       | From parts to wholes                      | <i>course</i> <sup>7</sup> → <i>meal</i> <sup>1</sup>       |
| Antonym        |               | Opposites                                 | <i>leader</i> <sup>1</sup> → <i>follower</i> <sup>1</sup>   |

| Relation | Definition  | Example   |
|----------|---|---|
| Hypernym | From events to superordinate events                               | <i>fly</i> <sup>9</sup> → <i>travel</i> <sup>9</sup>        |
| Troponym | From a verb (event) to a specific manner elaboration of that verb | <i>walk</i> <sup>1</sup> → <i>stroll</i> <sup>1</sup>       |
| Entails  | From verbs (events) to the verbs (events) they entail             | <i>snore</i> <sup>1</sup> → <i>sleep</i> <sup>1</sup>       |
| Antonym  | Opposites   | <i>increase</i> <sup>1</sup> ⇔ <i>decrease</i> <sup>1</sup> |

# Иерархии WordNet

```
Sense 3
bass, basso --
(an adult male singer with the lowest voice)
=> singer, vocalist, vocalizer, vocaliser
    => musician, instrumentalist, player
        => performer, performing artist
            => entertainer
                => person, individual, someone...
                    => organism, being
                        => living thing, animate thing,
                            => whole, unit
                                => object, physical object
                                    => physical entity
                                        => entity
=> causal agent, cause, causal agency
    => physical entity
        => entity
```

```
Sense 7
bass --
(the member with the lowest range of a family of
musical instruments)
=> musical instrument, instrument
    => device
        => instrumentality, instrumentation
            => artifact, artefact
                => whole, unit
                    => object, physical object
                        => physical entity
                            => entity
```

# Как “значение” определяется в WordNet

- Множество синонимов называется **синсет**
- Пример

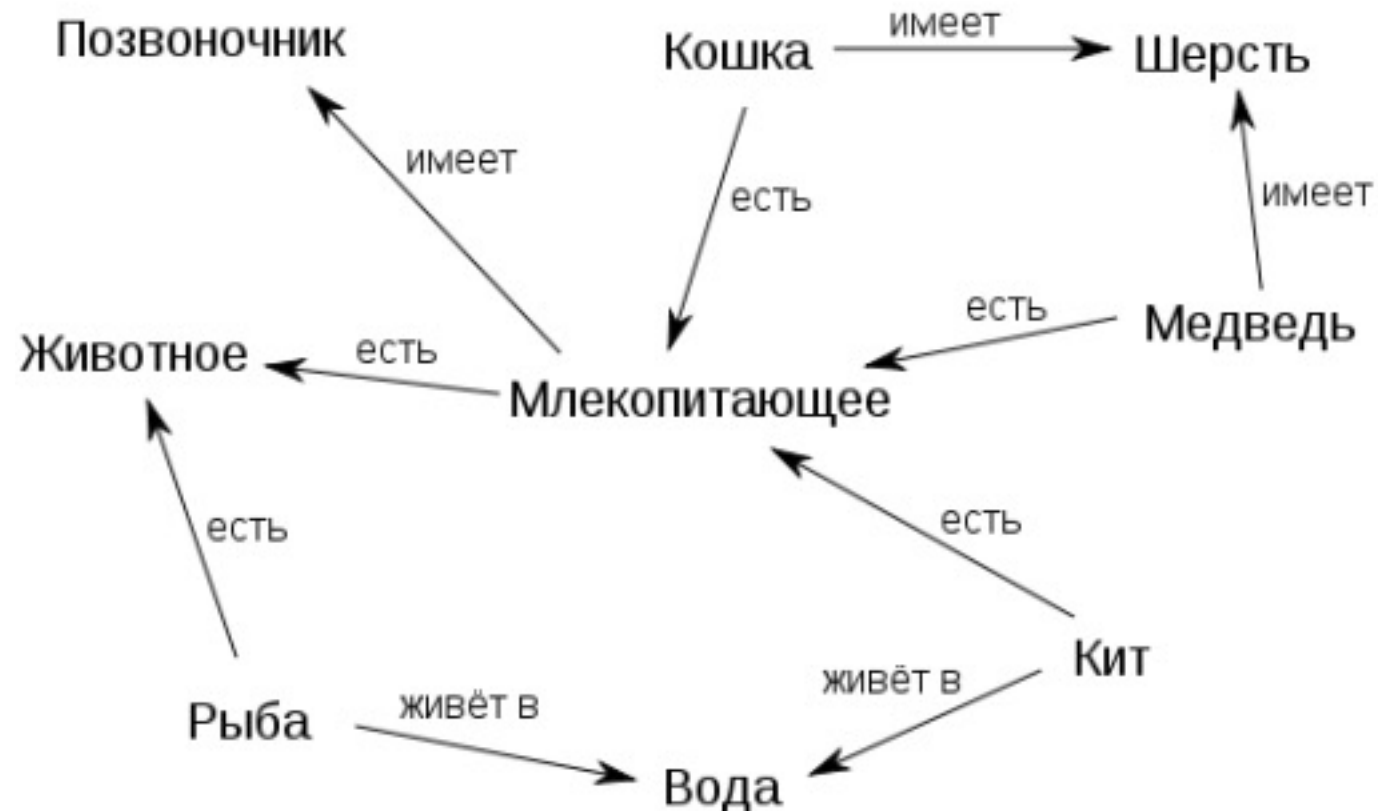
```
from nltk.corpus import wordnet
for synset in wordnet.synsets('platform'):
    print(synset.definition)
    print([lemma.name for lemma in synset.lemmas])
```

```
A document stating the aims and principles of a political party
['platform', 'political platform', 'political program', 'program']
Any military structure or vehicle bearing weapons
[platform, weapons platform]
...
```

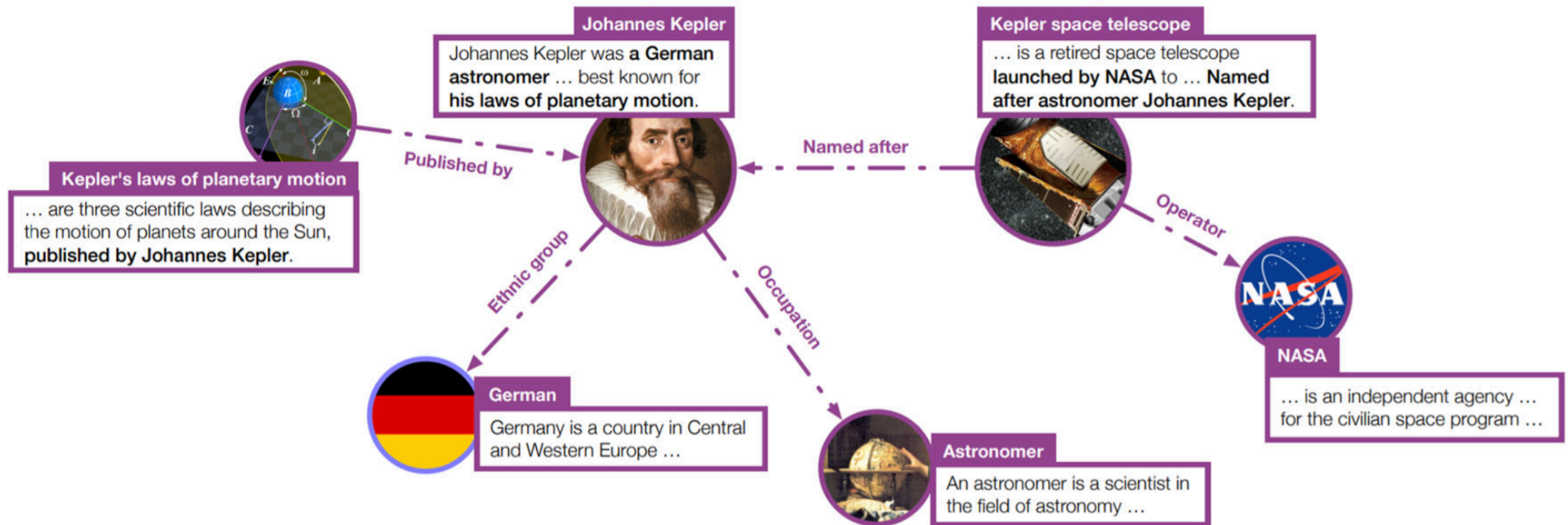


# Семантическая сеть

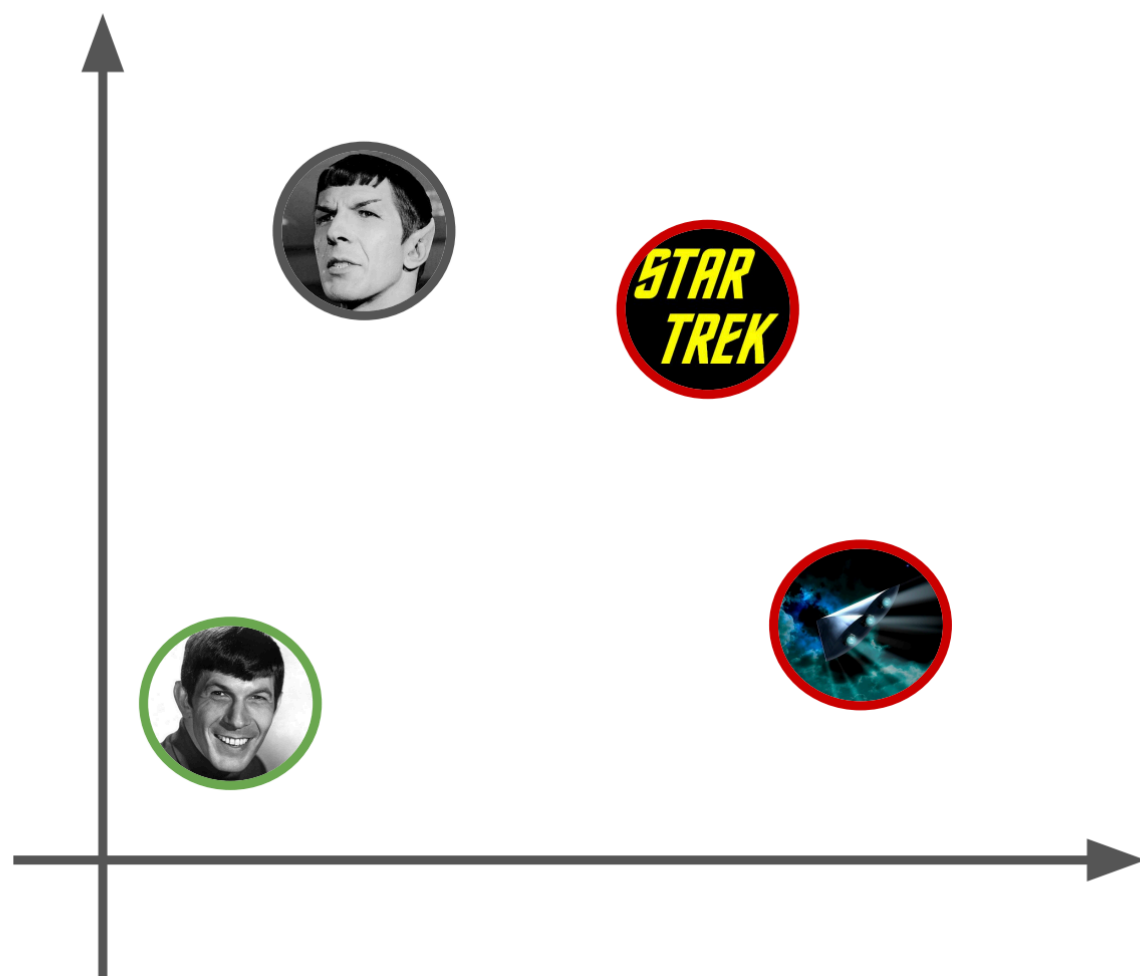
- **Семантическая сеть** — информационная модель предметной области, имеющая вид ориентированного графа, вершины которого соответствуют объектам предметной области, а дуги (рёбра) задают отношения между ними.



# Граф знаний

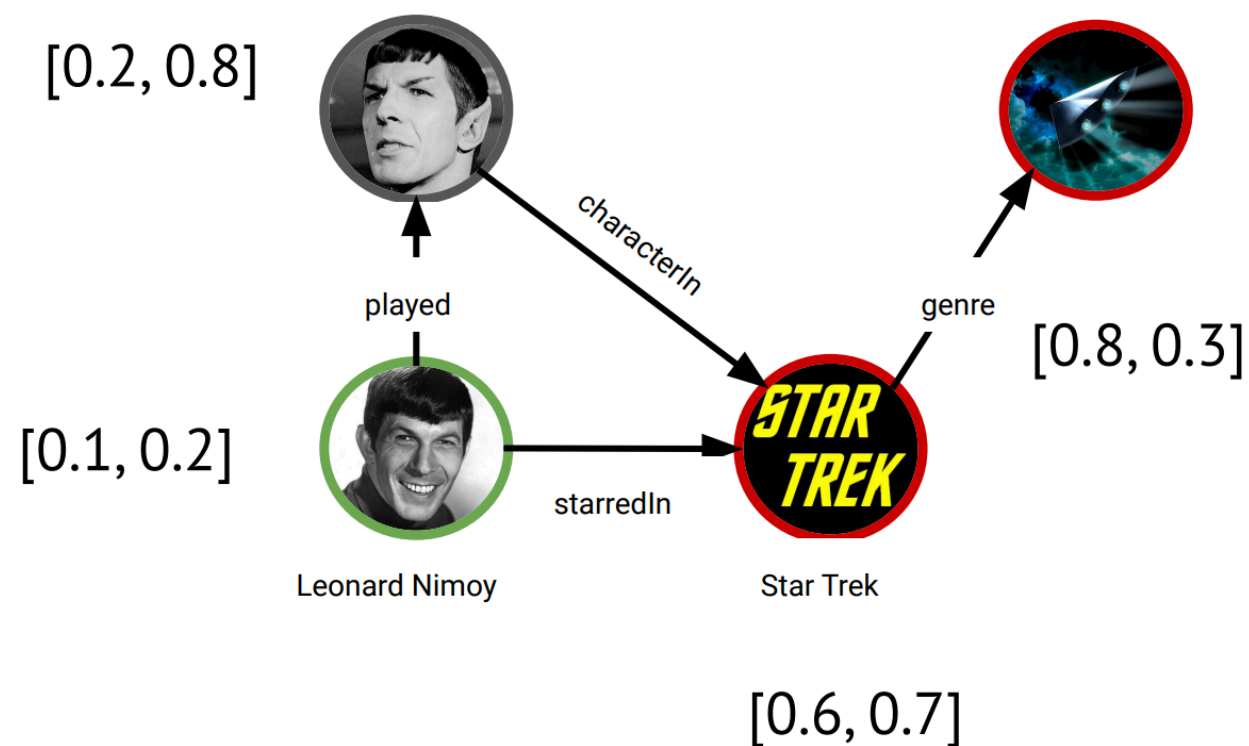


# Векторное представление



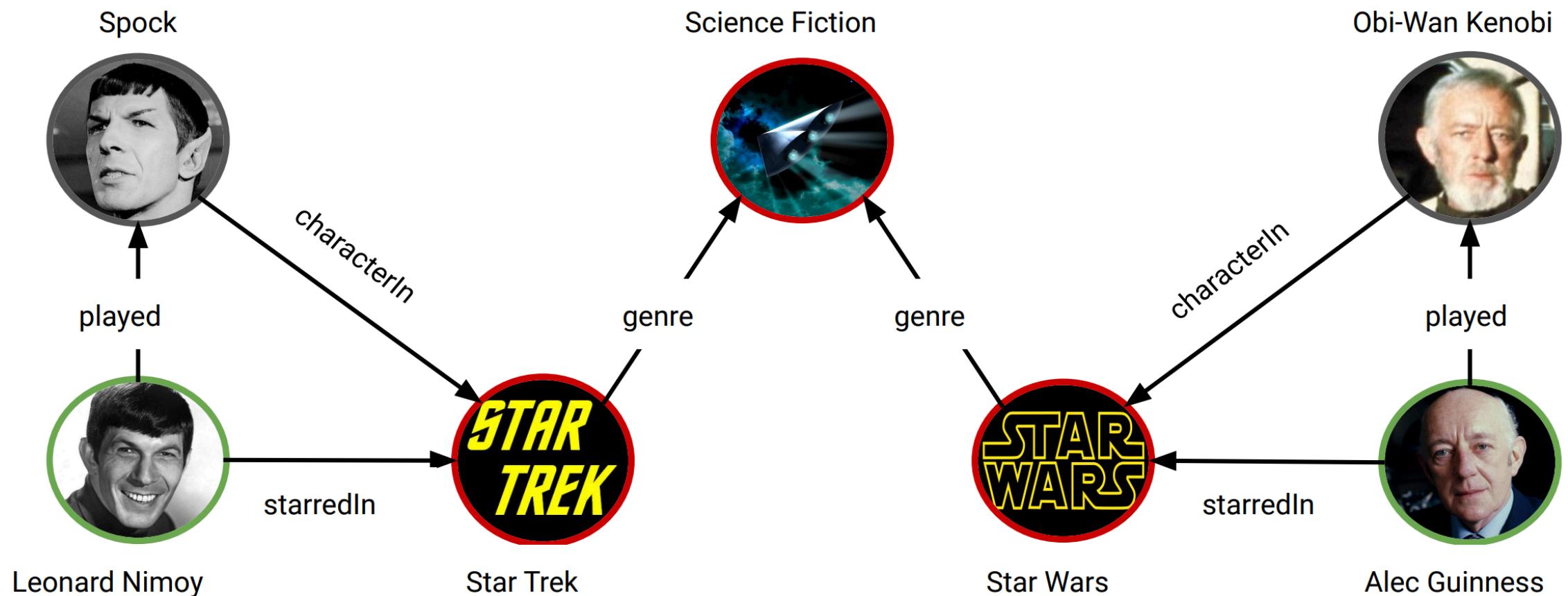
$$E \in \mathbb{R}^{N_e \times d}$$

$$R \in \mathbb{R}^{N_r \times d}$$





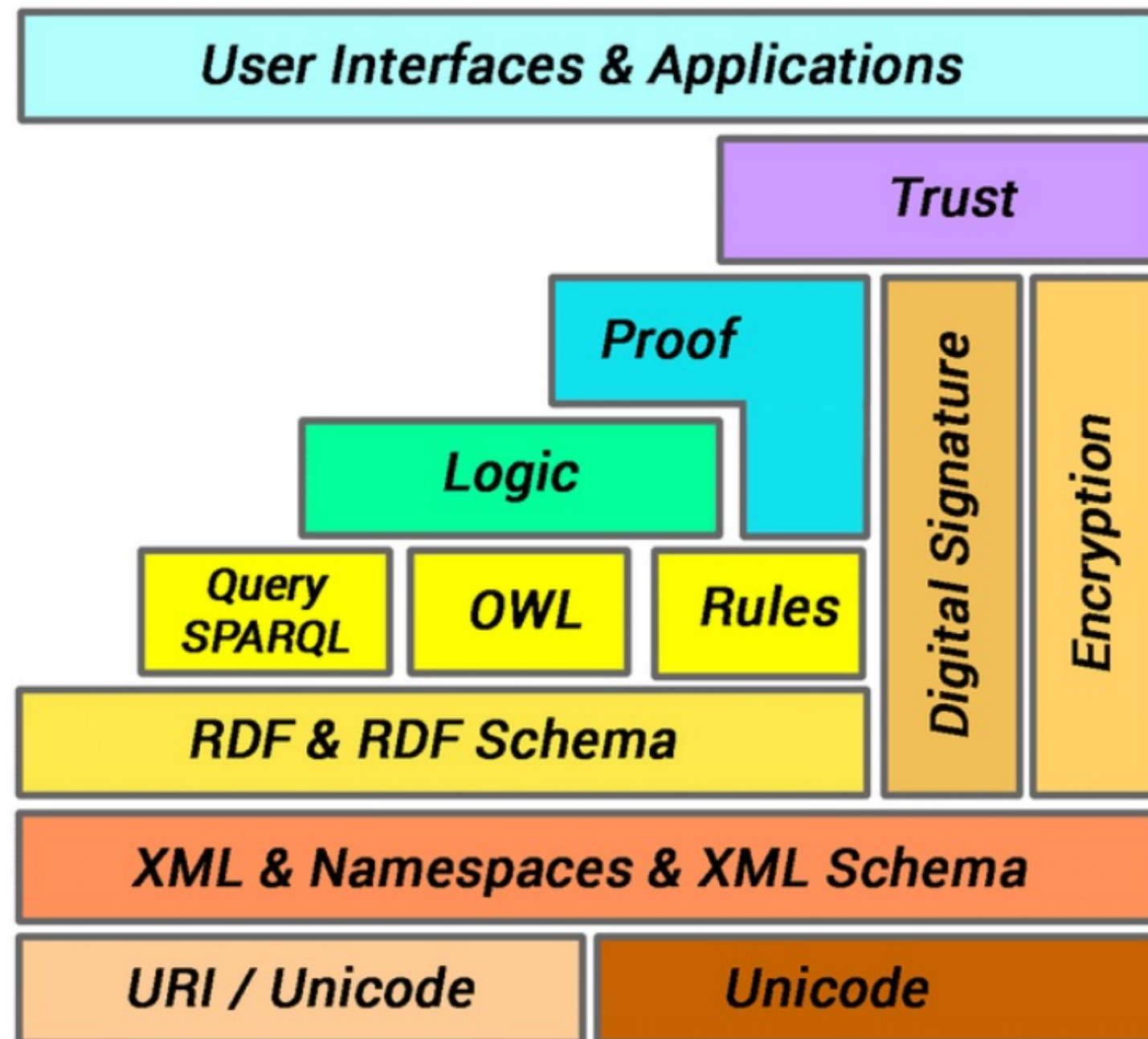
# Символьное представление



Leonard Nimoy played Spock  
 Leonard Nimoy starredIn Star Trek  
 Spock characterIn Star Trek  
 Star Trek genre Science Fiction

Alec Guinness starredIn Star Wars  
 Alec Guinness played Obi-Wan Kenobi  
 Obi-Wan Kenobi characterIn Star Wars  
 Star Wars genre Science Fiction

# Semantic Web



[Search Wikidata](#)

## Configure

| Language | Label  | Description                              | Also known as   |
|----------|--------|--|---|
| English  | Moscow | capital and most populous city of Russia | Moskva<br>Moscow, Russia<br>Moskva Federal City, Russia<br>Moscow, USSR<br>Moskva, Russia<br>City of Moscow<br>Moscow, Russian Federation<br>Moscow, Soviet Union<br>Moscow, Russian SFSR<br>Muscovite<br>Moscovite |
| Russian  | Москва | столица и крупнейший город России        | Первопрестольная<br>Порт пяти морей<br>Москва (город)<br>Москва, Россия<br>Москва (Россия)<br>Москва Златоглавая  |

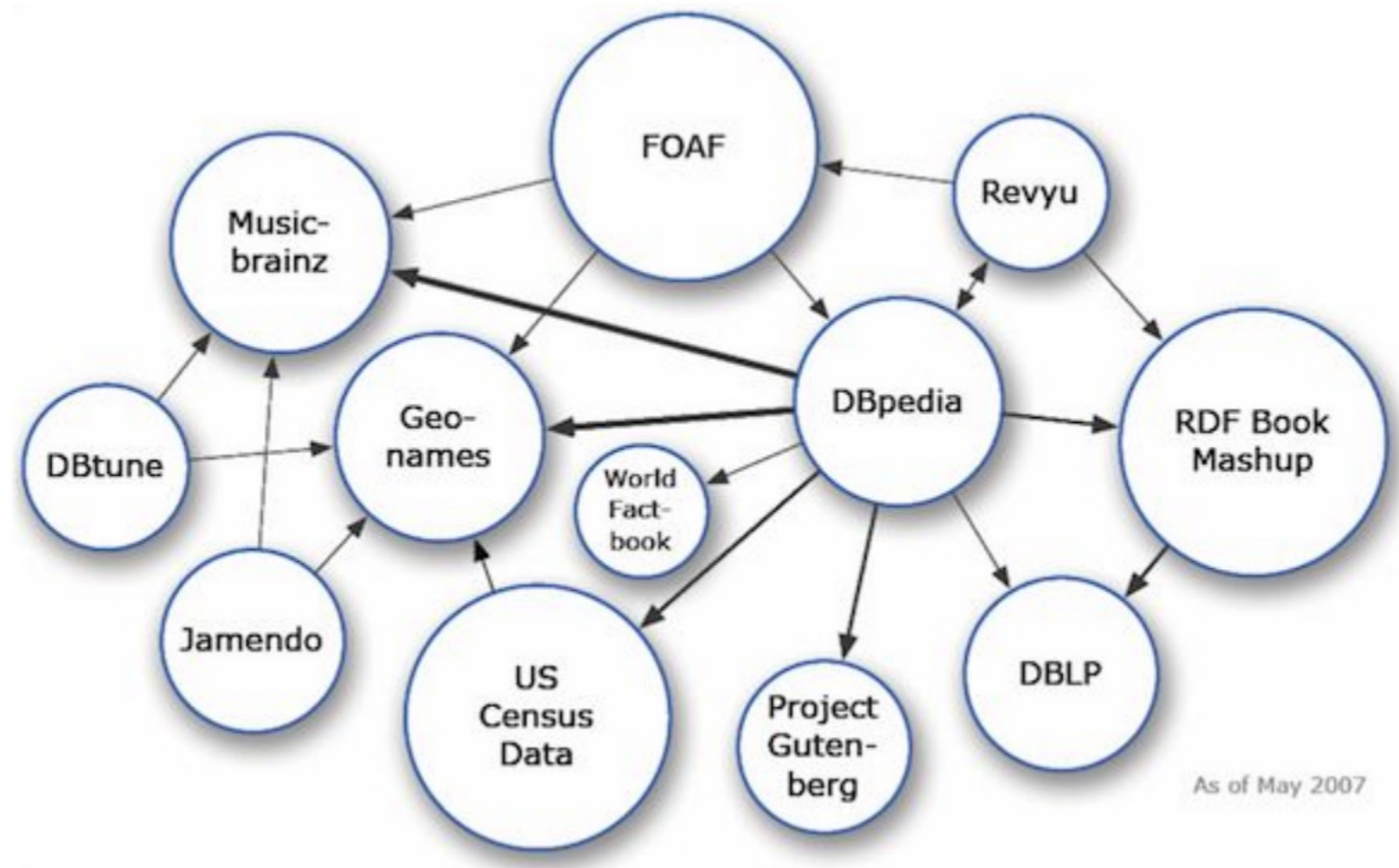
instance of

|   |                        |      |
|---|------------------------|------|
|   | start time             | 1918 |
|   | ▸ 2 references         |      |
| <hr/>   |                        |      |
|  | federal city of Russia |      |
|   | start time             | 1991 |
|   | ▸ 1 reference          |      |

|     |        |
|-----|--------|
| ab  | Москва |
| ace | Moskōw |
| ady | Москва |
| af  | Moskou |
| ak  | Moscow |
| als | Moskau |
| alt | Москва |
| am  | ሞስኮ    |
| ang | Moscow |
| an  | Moscú  |
| arc | ܡܫܟܘܐ  |
| ar  | موسكو  |
| ary | موسكو  |
| arz | موسكو  |
| ast | Moscú  |
| avk | Moskva |
| av  | Москва |
| awa | मस्को  |
| ay  | Mosku  |

2007

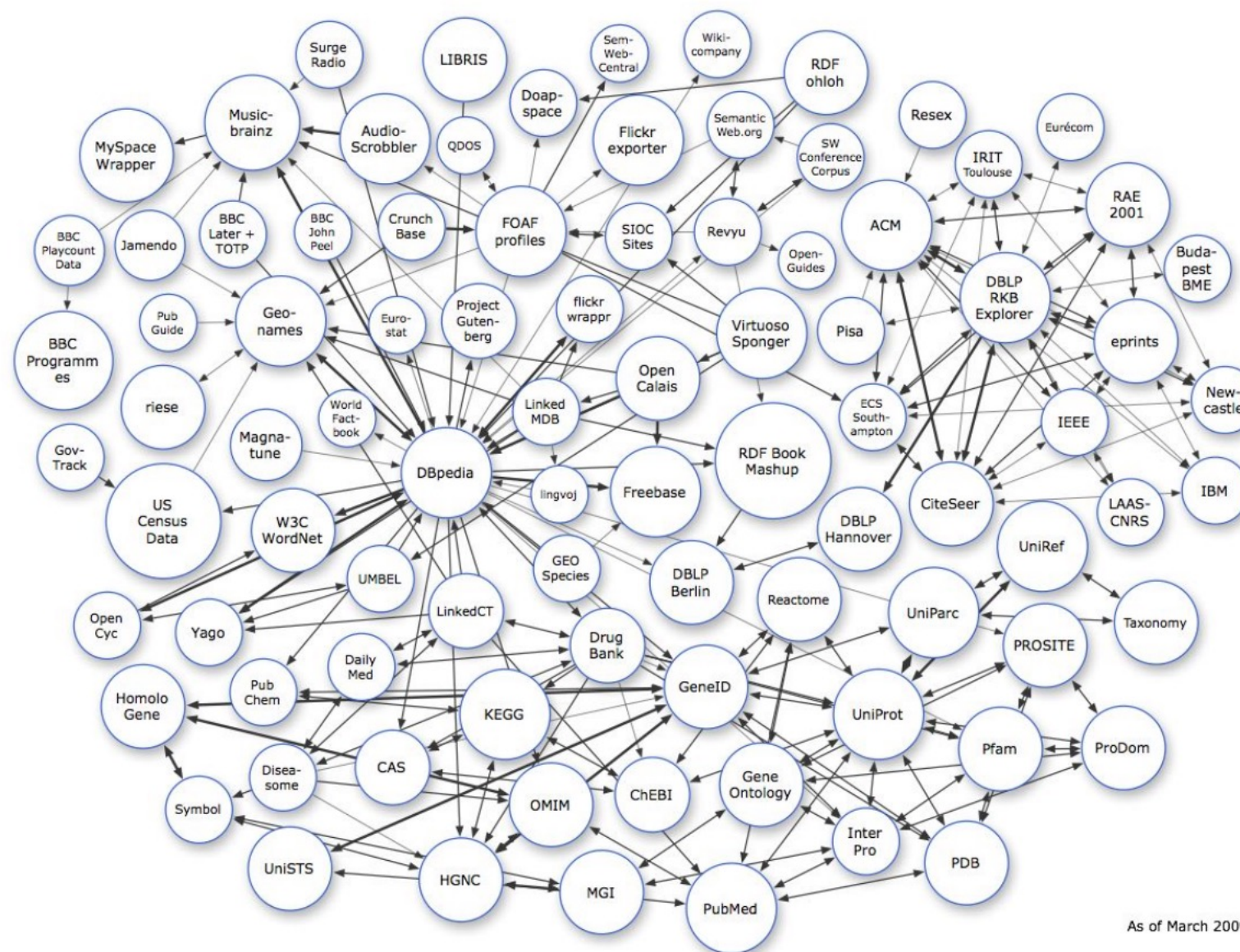
# Link Open Data





2009

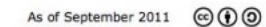
# Link Open Data



As of March 2009

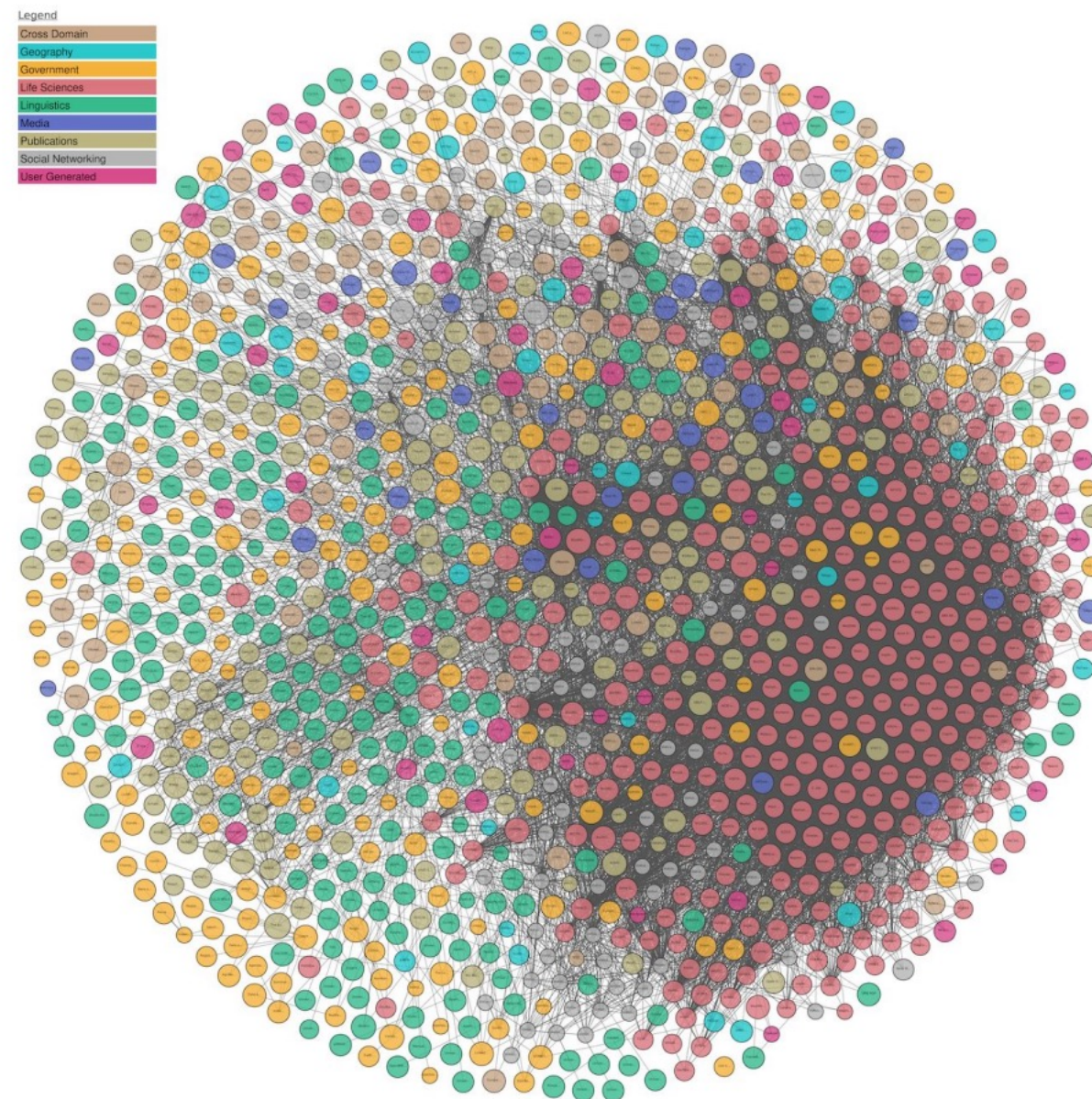


## 2011





# Link Open Data<sup>2020</sup>

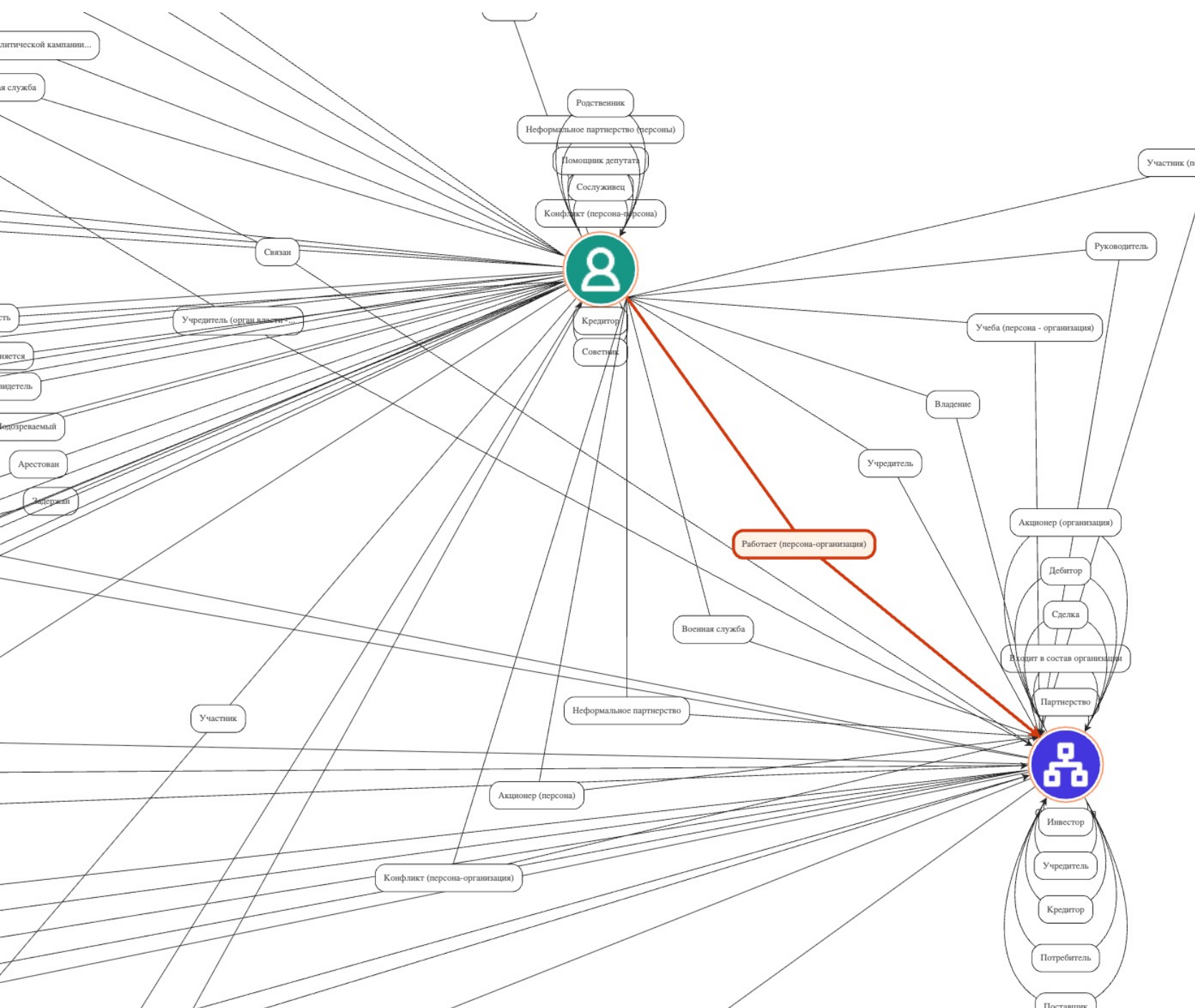


<https://lod-cloud.net>





# Типизация графа



Тип связи № 786 Работает (персона-организация)

 6

Тип сущности (из)













Персона

Тип сущности (в)

Организация

☒ Направленная☐ Иерархическая

### ✓ Типы характеристик (6)

|                       |        |   |   |
|-----------------------|--------|---|---|
| ID факта              | Строка |  |  |
| Комментарий           | Строка |  |  |
| Характер деятельности | Строка |  |  |
| Дата начала           | Дата   |  |  |
| Дата окончания        | Дата   |  |  |
| Должность             | Работа |  |  |

# Сценарии использования графов знаний

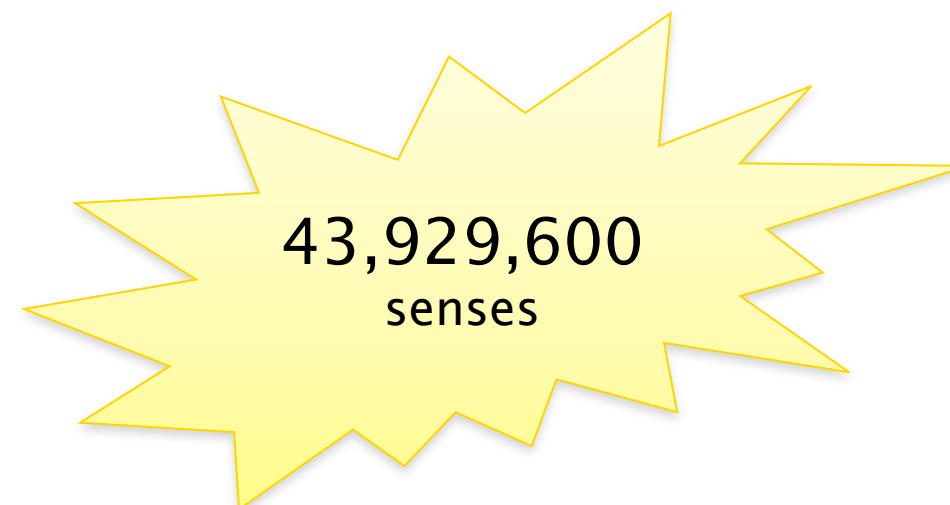
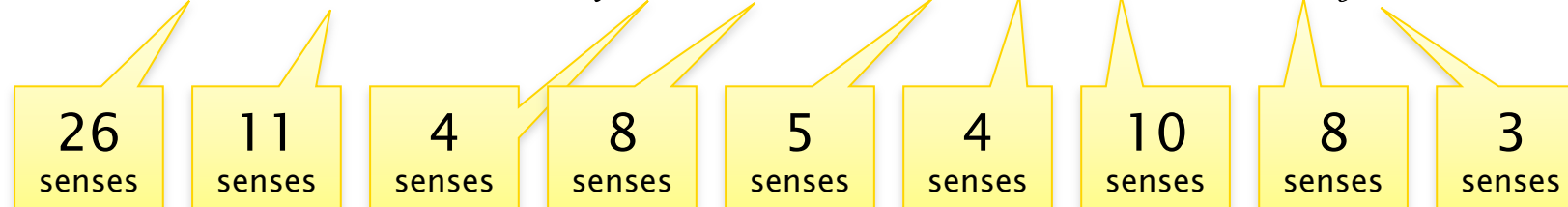
- Информационный поиск, расширение запросов, GraphRAG
- Проверка фактов, в том числе поиск противоречий (например, несовместимость лекарств)
- Интеграция структурированных и неструктурированных данных (корпоративные системы)
- Рекомендательные и вопросно-ответные системы
- Гибридный и объяснимый ИИ
- Аналитика сложных сетей и прогнозирование (поиск аномалий, вывод неочевидных закономерностей)

# Вычислительная лексическая семантика

- Разрешение лексической многозначности
- Привязка к базам знаний
- Наполнение графов знаний

# Трудность разрешения лексической многозначности

*I saw a man who is 98 years old and can still walk and tell jokes*



# Разрешение лексической многозначности

- Word Sense Disambiguation (WSD)
  - определение значения слова в контексте
  - обычно предполагается фиксированный список значений (например WordNet)
- Сводится к задаче классификации
- Отличается от задачи разграничения значений (word sense discrimination)

# Разрешение лексической многозначности: варианты

- Определение значений только заранее выбранных слов (lexical sample task)
  - line - hard - serve; interest
  - Ранние работы
  - Обучение с учителем
- Определение значений всех слов (all-word task)
  - Проблема разреженности данных
  - Невозможно натренировать отдельный классификатор для каждого слова

# Признаки

- Должны описывать **контекст**
- Предварительная обработка текста
  - параграфы, предложения, части речи, леммы, синтаксический разбор?
- Признаки в словосочетаниях с позициями
- Множества соседей
- Проблема разреженности языка
  - Использовать *семантическую близость*.  
(длина пути в иерархии WordNet, косинус между векторами слов:  $word2vec$  и т.п.)

# Пример

*An electric guitar and **bass** player stand off to one side, not really part of the scene, just as a sort of nod to gringo expectations perhaps.*

| Collocational features |          |
|------------------------|----------|
| word_L3                | electric |
| POS_L3                 | JJ       |
| word_L2                | guitar   |
| POS_L2                 | NN       |
| word_L1                | and      |
| POS_L1                 | CC       |
| word_R1                | player   |
| POS_R1                 | NN       |
| word_R2                | stand    |
| POS_R2                 | VB       |
| word_R3                | off      |
| POS_R3                 | RB       |

| Bag-of-words features |   |
|-----------------------|---|
| fishing               | 0 |
| big                   | 0 |
| sound                 | 0 |
| player                | 1 |
| fly                   | 0 |
| rod                   | 0 |
| pound                 | 0 |
| double                | 0 |
| runs                  | 0 |
| playing               | 0 |
| guitar                | 1 |
| band                  | 0 |



# Алгоритмы

- Любые методы классификации

# Вопрос на засыпку

- Как сделать классификатор для задачи определения значений всех слов (all-word task)?

# Методы оценки

- Внешние (in vivo)
  - Машинный перевод с/без WSD
- Внутренние (in vitro)
  - Применение к размеченным данным (SemCor, SENSEVAL, SEMEVAL)
  - Измерение точности и полноты в сравнении со стандартными значениями
- Нижняя граница
  - Выбор случайных значений работает плохо
  - Более сильные границы: наиболее частое значение, алгоритм Леска
- Верхняя граница: согласие экспертов
  - 75-80 для задачи определения значений всех слов со значениями из WordNet
  - до 90% с менее гранулированными значениями

# Методы основанные на словарях и тезаурусах

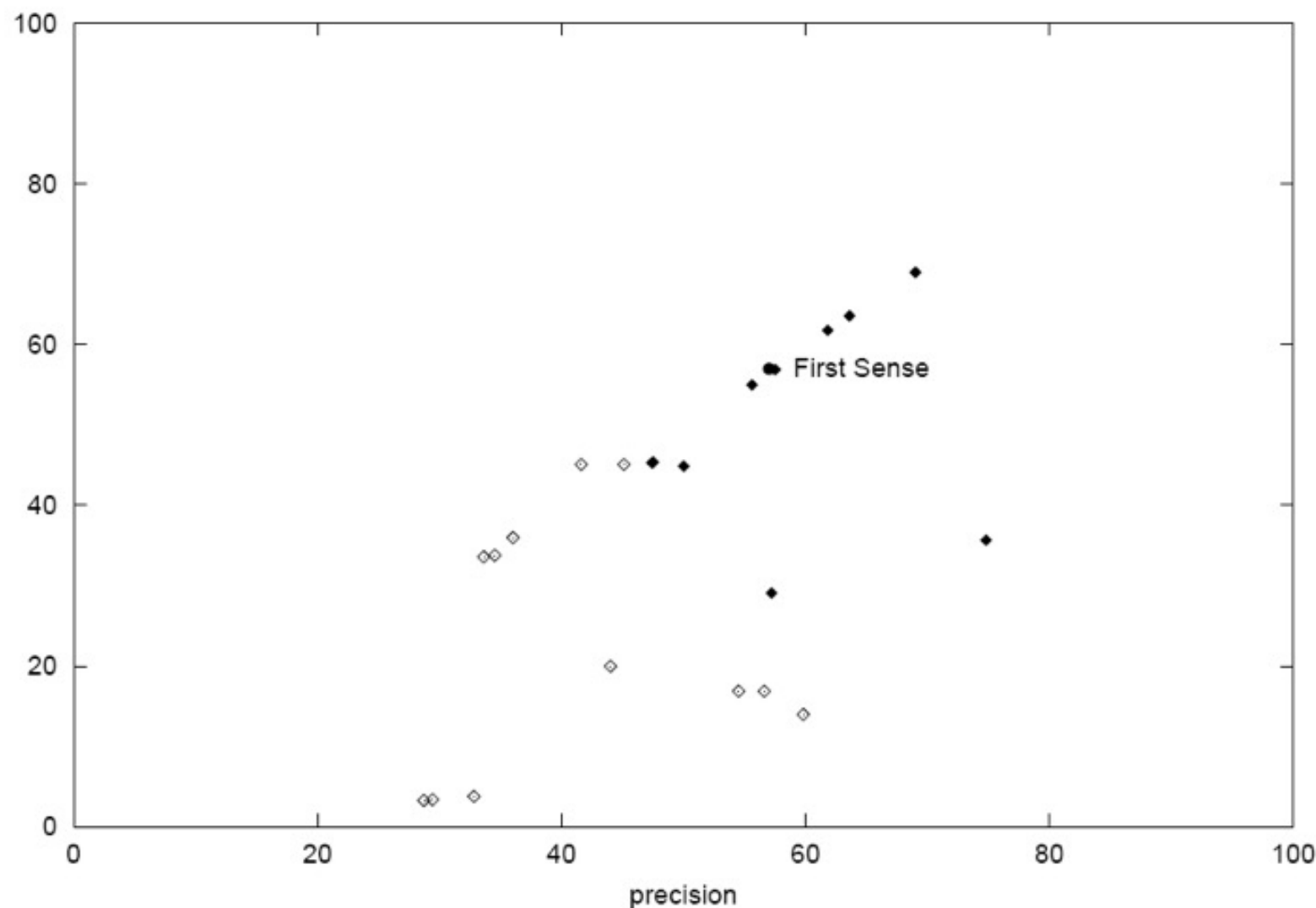
- Алгоритм Леска (1986)
  - Взять все определения целевого слова из словаря
  - Сравнить с определениями слов в контексте
  - Выбрать значение с максимальным пересечением
- Пример
  - *pine*
    1. a kind of **evergreen tree** with needle-shaped leaves
    2. to waste away through sorrow or illness
  - *cone*
    1. A solid body which narrows to a point
    2. Something of this shape, whether solid or hollow
    3. Fruit of certain **evergreen trees**
  - Определить значение: *pine cone*

# Варианты алгоритма Леска

- Упрощенный (Simplified Lesk)
  - Взять все определения целевого слова из словаря
  - Сравнить со ~~определениями~~ словами в контексте
  - Выбрать значение с максимальным пересечением
- Корпусный (Corpus Lesk)
  - Включить предложения из размеченного корпуса в сигнатуру каждого значения
  - Взвесить слова через IDF
  - $IDF(w) = -\log P(w)$
  - Показывает лучшие результаты
  - Использовался как нижняя граница на SENSEVAL

# Наиболее частое значение

- Сравнение методов на SENSEVAL-2



- McCarthy et. al. 2004 ACL - поиск наиболее частого значения по неразмеченному корпусу

# Привязка к базам знаний

# Связывание (именованных) сущностей

- Вход: неструктурированный текст
- Выход: Значение для каждой (именованной) сущности из базы знаний





# База знаний. Википедия



**Википедия**  
Свободная энциклопедия

[Заглавная страница](#)  
[Рубрикация](#)  
[Указатель А—Я](#)  
[Избранные статьи](#)  
[Случайная статья](#)  
[Текущие события](#)

[Участие](#)  
[Сообщить об ошибке](#)  
[Сообщество](#)  
[Форум](#)  
[Свежие правки](#)  
[Новые страницы](#)  
[Справка](#)  
[Пожертвовать](#)

**Инструменты**  
[Ссылки сюда](#)  
[Связанные правки](#)  
[Служебные страницы](#)  
[Постоянная ссылка](#)  
[Сведения о странице](#)  
[Цитировать страницу](#)

Вы не представились системе [Обсуждение](#) [Вклад](#) [Создать учётную запись](#) [Войти](#)

Статья [Обсуждение](#) [Читать](#) Текущая версия [Править](#) [Править код](#) [История](#)

55°42′11″ с. ш. 37°31′50″ в. д. НГЯО

## Московский государственный университет

Материал из Википедии — свободной энциклопедии

[\[ править \]](#) [\[ править код \]](#)

Текущая версия страницы пока не проверялась опытными участниками и может значительно отличаться от версии, проверенной 21 октября 2018; проверки требуют 7 правок.

*Запрос «МГУ» перенаправляется сюда; см. также другие значения.*

**Моско́вский госуда́рственный университе́т и́мени М. В. Ломоно́сова** — один из старейших<sup>[4][5]</sup> и крупнейших<sup>[6][7]</sup> классических университетов России, один из центров отечественной науки и культуры, расположенный в Москве.

С 1940 года носит имя **Михаила Васильевича Ломоносова**.

Полное наименование — Федеральное государственное бюджетное образовательное учреждение высшего образования «Московский государственный университет имени М. В. Ломоносова». Широко используется аббревиатура «МГУ».

Университет включает в себя 15 научно-исследовательских институтов<sup>[8]</sup>, 43 факультета<sup>[9][10]</sup>, более 300 кафедр и 6 филиалов (в их числе пять зарубежных — все в странах СНГ)<sup>[11]</sup>.

С 1992 года ректором МГУ является академик Виктор Антонович Садовничий.

### Содержание [скрыть]

- История становления и развития Московского университета
  - Императорский Московский университет 1755—1917
    - Основание Московского университета в 1755 году
    - Московский университет в XVIII веке
    - Московский университет в XIX веке

**Московский государственный университет имени М. В. Ломоносова**  
(МГУ имени М. В. Ломоносова)

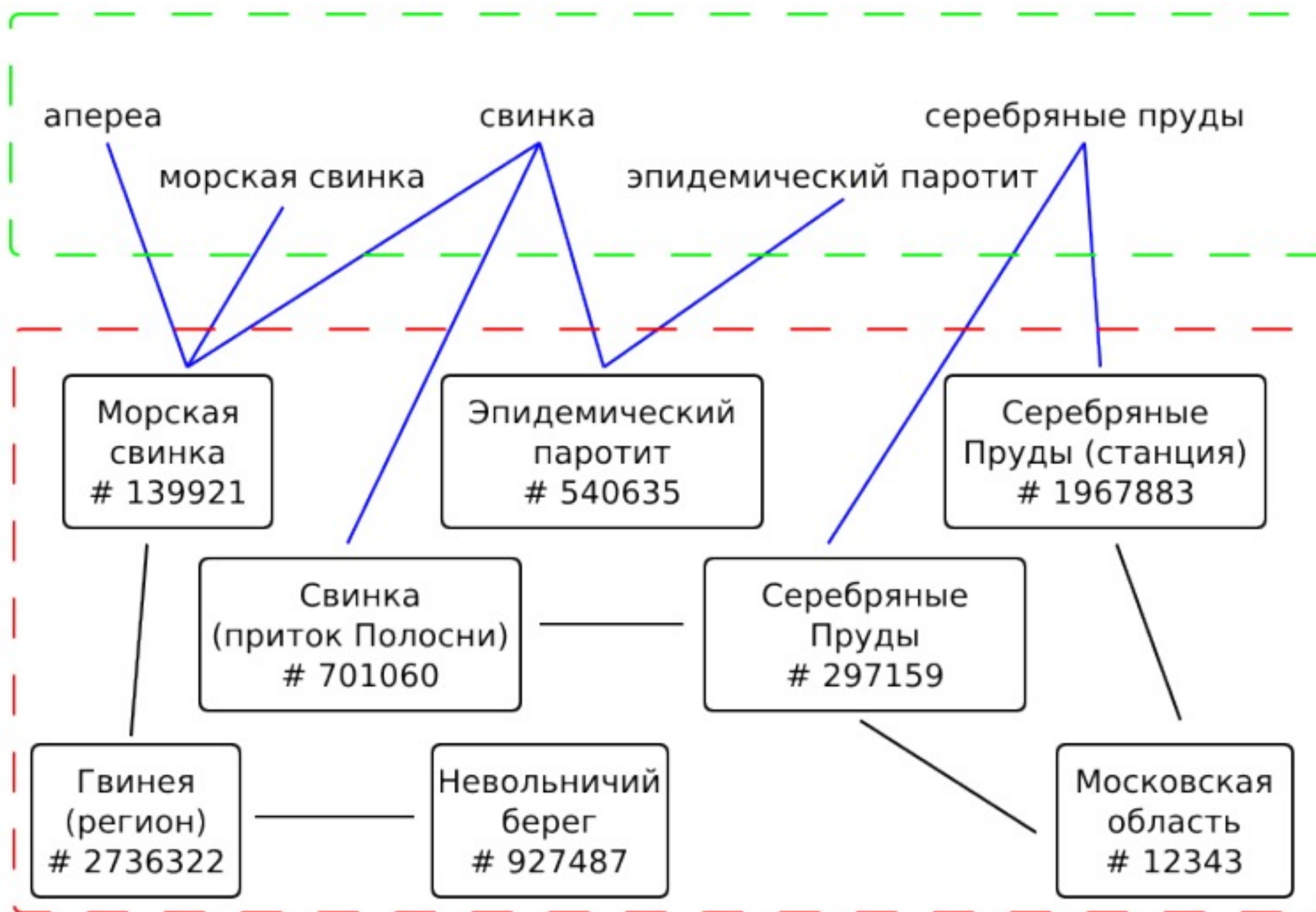


Главное здание МГУ имени М. В. Ломоносова,  
31 мая 2015 года





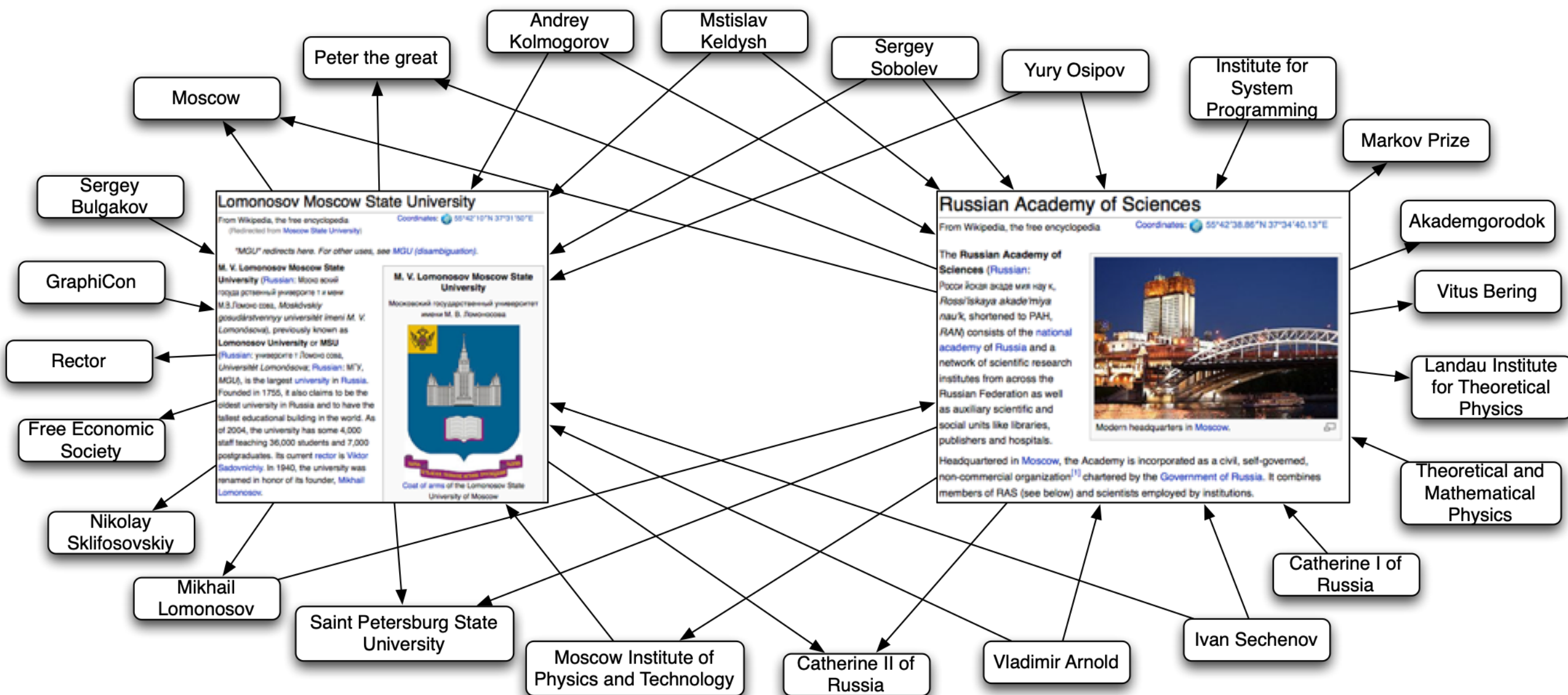
# База знаний. Википедия





# Использование Википедии: семантическая близость

- Нормализованное количество общих соседей



- Близкие концепты чаще встречаются вместе

# Основные проблемы

- Синонимия (вариации имен): New York, NY, Big Apple
- **Многозначность**: NY ? New York / New Year
- Отсутствие значения
- Скорость обработки
- Постоянное изменение базы знаний
- Много языков

# Алгоритм Milne-Witten

- Иcпoльзует oднoзнaчнoе тeрминoвoе кoнтeкcтa

## Depth-first search

From Wikipedia, the free encyclopedia

**Depth-first search (DFS)** is an **algorithm** for traversing or searching a **tree** **tree structure** or **graph**. One starts at the root (selecting some node as the root in the graph case) and explores as far as possible along each branch before **backtracking**.

Formally, DFS is an **uninformed search** that progresses by expanding the first child node of the search **tree** that appears and thus going deeper and deeper until a goal node is found, or until it hits a node that has no children. Then the search **backtracks**, returning to the most recent node it hadn't finished exploring. In a non-recursive implementation, all freshly expanded nodes are added to a **LIFO stack** for exploration.

| sense                        | commonness   | relatedness   |
|------------------------------|--------------|---------------|
| Tree                         | 92.82%       | 15.97%        |
| Tree (graph theory)          | 2.94%        | 59.91%        |
| <b>Tree (data structure)</b> | <b>2.57%</b> | <b>63.26%</b> |
| Tree (set theory)            | 0.15%        | 34.04%        |
| Phylogenetic tree            | 0.07%        | 20.33%        |
| Christmas tree               | 0.07%        | 0.0%          |
| Binary tree                  | 0.04%        | 62.43%        |
| Family tree                  | 0.04%        | 16.31%        |
| ...                          |              |               |

# Алгоритм Milne-Witten

- Commonness (популярность значения)

$$commonness(e, m) = \frac{count(e, m)}{count(e)}$$

- Relatedness (расстояние до контекста)

$$relatedness(a, b) = \frac{\log \max(|A|, |B|) - \log |A \cap B|}{\log |W| - \log \min(|A|, |B|)}$$

$a, b$  - концепты;

$A, B$  - множества концептов, ссылающихся на  $e$  и  $s$  соответственно;

$W$  - множество всех концептов.

- Вероятность термина быть ссылкой

$$link\_prob(t) = \frac{count(t \text{ is link})}{count(t)}$$

# Алгоритм Milne-Witten

- Семантическое расстояние от концепта  $a$  до контекста  $S$ :

$$distance(a, S) = \frac{\sum_{s \in S} w_s \times relatedness(a, s)}{\sum_{s \in S} w_s}$$

- Вес  $w_s$  концепта  $s$ :

$$w_s = \frac{1}{2}(link\_prob(term(s))) + \frac{1}{|S|} \sum_{c \in S} relatedness(s, c)$$

- Качество контекста  $S$ :

$$quality(S) = \sum_{s \in S} w_s$$



# Алгоритм Milne-Witten

- Признаки:
  - популярность значения -  $\text{commonness}(e, m)$ ;
  - расстояние до контекста -  $\text{distance}(e, S)$ ;
  - когерентность контекста -  $\text{quality}(S)$ .
- Алгоритмы машинного обучения:
  - Naïve Bayes;
  - дерево решений C4.5;
  - метод опорных векторов;
  - дерево решений C4.5 с бэггингом.
- Обучение / применение:
  - положительный пример — правильное значение термина;
  - отрицательные примеры — все остальные возможные значения термина;
  - на этапе применения выбирается концепт, с максимальной уверенностью классифицированный как правильный.

# Извлечение информации

| Параметр                                     | Классический IE   | Open IE   |
|--|---|---|
| Схема отношений                              | Фиксированная, заранее определённая (например: «родился в», «работает в»)   | Не фиксирована, извлекает любые отношения из текста                                   |
| Гибкость                                     | Низкая, работает только для конкретных доменов и заранее заданных отношений | Высокая, может извлекать новые, неожиданные факты и связи                             |
| Необходимость разметки данных                | Требует ручной разметки для обучения моделей                                | Может работать с минимальной или автоматической разметкой                             |
| Тип входного текста                          | Обычно структурированные или доменные тексты                                | Любой текст, в том числе новостные статьи, веб-страницы, соцсети                      |
| Формат вывода                                | Обычно фиксированные шаблонные отношения (relation, entity1, entity2)       | Тройки (subject, relation, object), где relation формируется автоматически из текста  |
| Применение                                   | Узкоспециализированные задачи, корпоративные системы, биомедицина, финансы  | Построение графов знаний, QA системы, семантический поиск, аналитика больших корпусов |
| Обработка неопределённых или новых отношений | Сложно, требует добавления новых правил                                     | Естественно извлекает новые отношения без дополнительной разметки                     |
| Сложность реализации                         | Средняя, требует настройки под домен  | Высокая, требует NLP-парсинга, синтаксического анализа и фильтрации шумов             |
| Примеры инструментов                         | GATE, IBM Watson IE, LingPipe   | Stanford OpenIE, OpenIE5, ClausIE, MINIE  |

# Извлечение знаний: LLM

В воскресенье, 20 июля, компания Microsoft выпустила срочные патчи, закрывающие две уязвимости в ПО для совместной работы Microsoft SharePoint. Речь идет о корпоративных инсталляциях SharePoint, облачная версия SharePoint Online, входящая в пакет Microsoft 365, не подвержена проблеме. Уязвимости с идентификаторами [CVE-2025-53770](#) и [CVE-2025-53771](#) являются вариантами проблем, обнаруженных еще в мае этого года и пропатченных (как выяснилось, не до конца) в июльском наборе обновлений от Microsoft.

*Извлеки сущности и связи из следующего текста. Верни структурированные тройки <s,r,o> для сохранения в граф знаний neo4j.*

|                                   |                       |                                   |
|-----------------------------------|-----------------------|-----------------------------------|
| <Microsoft                        | выпустила             | срочные патчи>                    |
| <срочные патчи                    | закрывают             | две уязвимости>                   |
| <уязвимости                       | обнаружены_в          | Microsoft SharePoint>             |
| <Microsoft SharePoint             | является              | ПО для совместной работы>         |
| <SharePoint Online                | входит_в              | Microsoft 365>                    |
| <SharePoint Online                | не_подвержен          | уязвимостям>                      |
| <уязвимости                       | имеют_идентификатор   | CVE-2025-53770>                   |
| <уязвимости                       | имеют_идентификатор   | CVE-2025-53771>                   |
| <CVE-2025-53770                   | является_вариантом    | проблемы_обнаруженной_в_мае_2025> |
| <CVE-2025-53771                   | является_вариантом    | проблемы_обнаруженной_в_мае_2025> |
| <проблемы_обнаруженные_в_мае_2025 | были_пропатчены       | июльским_набором_обновлений>      |
| <июльский_набор_обновлений        | выпущен               | Microsoft>                        |
| <июльский_набор_обновлений        | не_устранил_полностью | уязвимости>                       |
| <патчи                            | выпущены_дата         | 2025-07-20>                       |
| <патчи                            | выпущены_день_недели  | воскресенье>                      |

# Понимание слов предложения ≠ понимание предложения

- Во время интервью:
  - Ваша главная слабость?
  - Правильно интерпретирую семантику вопроса, но игнорирую его суть.
  - Не могли бы вы привести пример?
  - Мог бы.

# Ссылки

- <https://ods.ai/tracks/kgcourse2022>
- <https://dl.acm.org/doi/abs/10.1145/1458082.1458150>

## Learning to link with wikipedia

[D Milne, IH Witten](#) - Proceedings of the 17th ACM conference on ..., 2008 - dl.acm.org

... unstructured text, and enrich it with **links** to the appropriate **Wikipedia** articles. The resulting

... **-learning** approach to disambiguation that uses the **links** found within **Wikipedia** articles for ...

★ Save  Cite Cited by 1733 Related articles All 11 versions

# Следующая лекция

- Машинный перевод