

Распределенная обработка «больших» данных в семантических NoSQL-приложениях

Докладчик:
к.т.н, доц., с.н.с.
Щербак Сергей Сергеевич

Рассматриваемые вопросы



- **Big Data**
- **Linked Data**
- **NoSQL**
- **Endpoints**
- **Sparql**
- **Facebook Open Graph**

Buzzword -Big Data

“Большие данные: новый рубеж для инноваций, конкуренции и производительности” — McKinsey Institute

“Большие данные объединяют техники и технологии, которые извлекают смысл из данных на экстремальном пределе практичности” — Консалтинговая компания Forrester

На пути к йотта байтам...

- мега — 10^6 — 1 000 000 (10^6) или 1 048 576 (2^{20}) байт
- гига — 10^9 — 1 000 000 000 (10^9) или $2^{30} = 1 073 741 824$ байт
- тера — 10^{12} — 1 099 511 627 776 байт или 1024 гигабайт
- пета — 10^{15} или 2^{50} байт
- экса — 10^{18} или 2^{60} байт
- зетта — 10^{21} или 2^{70} байт
- йотта — 10^{24} или 2^{80} байт

Big Data - проблема или решение???

- основное их предназначение в том, чтобы помочь *понять* прошлое *предсказать* будущее



Большую часть данных, которая будет произведена в период с 2012 по 2020 годы, сгенерируют не люди, а машины в ходе взаимодействия друг с другом и другими сетями данных.



три «V»



Методики анализа больших данных

8

A/B testing

Association rule learning

Classification

Cluster analysis

Crowdsourcing.

Data fusion and data integration

Data mining

Ensemble learning

Genetic algorithms

Аналитический инструментарий

- 1010data
- Apache Chukwa
- Apache Hadoop
- Apache Hive
- Apache Pig!
- Jaspersoft
- MapReduce



Технологии

1
0



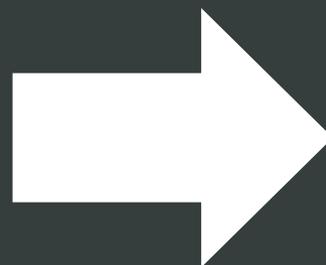
HOW TO WRITE A CV

11

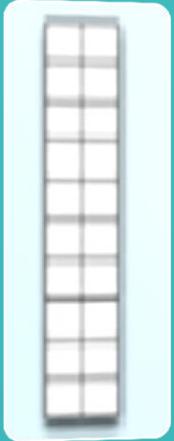
NoSQL-решения



Leverage the NoSQL boom



NO
SQL



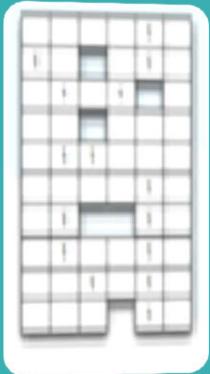
Key-value

- Dynamo
- Voldemo
- Redis
- Riak



Column (tabular)

- HBase
- Hypertable
- Cassandra



Graph

- Neo4J
- Virtuoso
- AllegroGraph
- CloudGRAPH

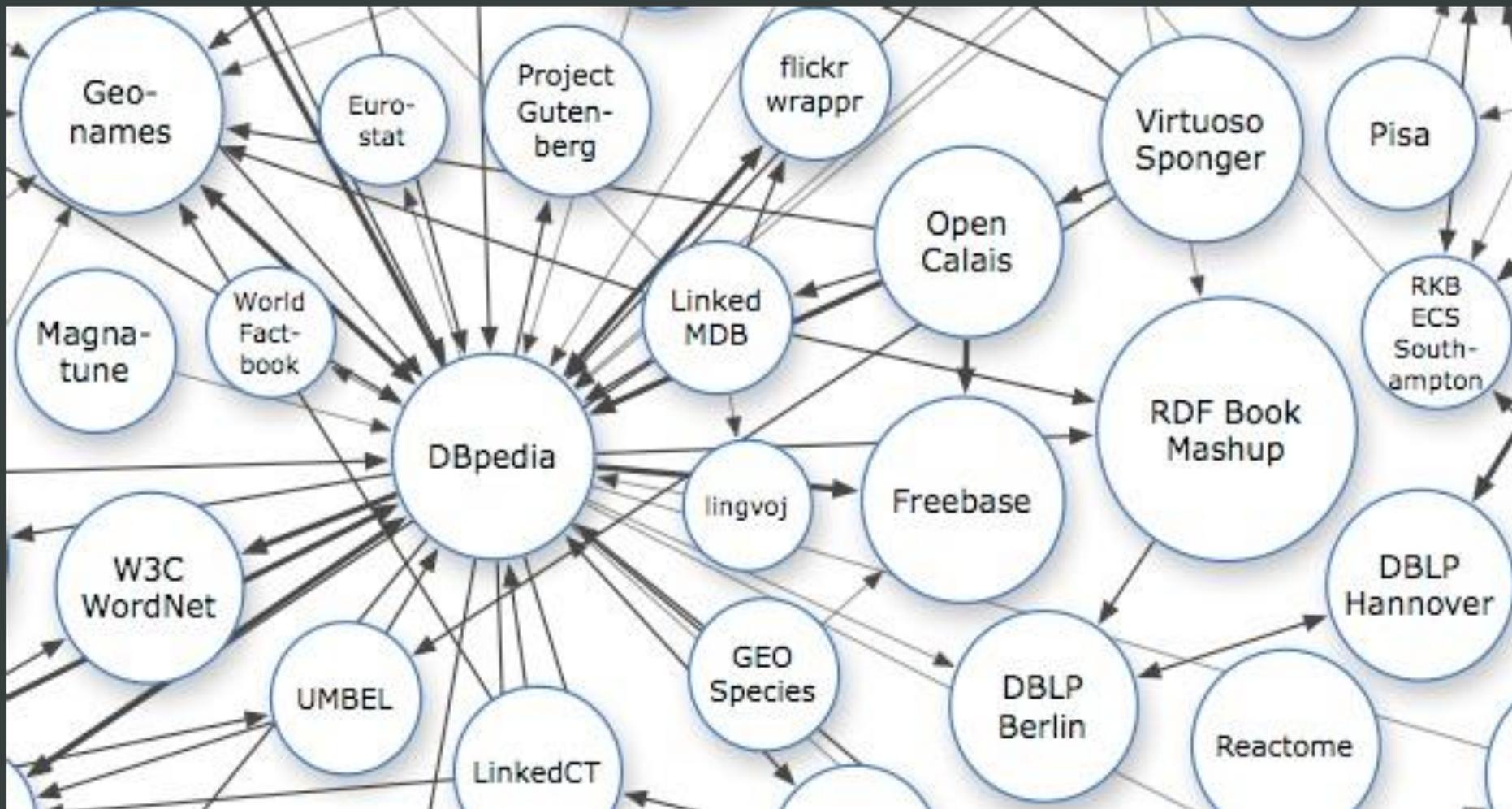
semantic



Document store

- CouchDB
- MongoDB

Множество источников данных, доступных в Internet



Datasets



DBpedia



FOAF



CKAN



GeoNames



Freebase

Интерфейсы доступа к ресурсам Data Sets

ODBC

- JDBC

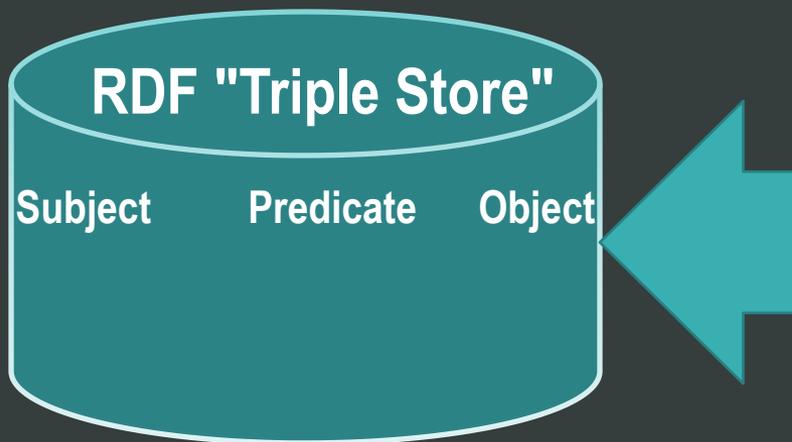
ADO.NET

- OLE DB

ADO.NET Entity Framework

Обеспечение унифицированного программного API для доступа к ресурсам Graph Triple-, Quad Store

Примеры клиентов для доступа к ресурсам Dataset



Interactive SQL

isql

Basic | Advanced

Server-side script

WebDAV source Local file

```
sparql SELECT ?s ?p ?o
FROM <http://shcherbak.net/User>
WHERE
{ ?s ?p ?o }
```

Show no more than rows

Sesame Windows Client

File Help

Connection | Namespaces | SPARQL EVALUATE | SPARQL CONSTRUCT | Browse Hierarchy | Export | Add data | Remove data | Add/Remove Repository | Community

Namespaces mode: User supplied Repository defined

```
SELECT ?s ?p ?o
FROM <http://shcherbak.net/User>
WHERE
{ ?s ?p ?o }
```

Result mode: current window new window Save result to file

o	s	p
"1" <http://www.w3.org/2001/XMLSchema...>	<http://shcherbak.net/log1>	<http://shcherbak.net/id_us>
"2" <http://www.w3.org/2001/XMLSchema...>	<http://shcherbak.net/log2>	<http://shcherbak.net/id_us>
"3" <http://www.w3.org/2001/XMLSchema...>	<http://shcherbak.net/log3>	<http://shcherbak.net/id_us>
"4" <http://www.w3.org/2001/XMLSchema...>	<http://shcherbak.net/log4>	<http://shcherbak.net/id_us>
"5" <http://www.w3.org/2001/XMLSchema...>	<http://shcherbak.net/log5>	<http://shcherbak.net/id_us>
"6" <http://www.w3.org/2001/XMLSchema...>	<http://shcherbak.net/log6>	<http://shcherbak.net/id_us>
"7" <http://www.w3.org/2001/XMLSchema...>	<http://shcherbak.net/log7>	<http://shcherbak.net/id_us>
"8" <http://www.w3.org/2001/XMLSchema...>	<http://shcherbak.net/log8>	<http://shcherbak.net/id_us>
"9" <http://www.w3.org/2001/XMLSchema...>	<http://shcherbak.net/log9>	<http://shcherbak.net/id_us>

2550 query results (took 322 ms)

SPARQL-точка доступа DBPedia

17

<http://dbpedia.org/sparql/>

Virtuoso SPARQL Query Editor

[About](#) | [Namespace Prefixes](#) | [Inference rules](#) | [iSPARQL](#)

Default Data Set Name (Graph IRI)

Query Text

```
select distinct ?Concept where {[] a ?Concept} LIMIT 100
```

(Security restrictions of this server do not allow you to retrieve remote RDF data, see [details](#).)

Results Format: (The CXML output is disabled, see [details](#))

Execution timeout: milliseconds (values less than 1000 are ignored)

Options: Strict checking of void variables

(The result can only be sent back to browser, not saved on the server, see [details](#))

SPARQL запрос для получения перечня 18 названий первых десяти фильмов с указанием lang-тега

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX dbpedia-owl: <http://dbpedia.org/ontology/>
SELECT ?film
WHERE
{
  ?a rdf:type dbpedia-owl:Film;
  rdfs:label ?film.
  FILTER( lang(?film) = "ru" )
}
limit 10
```

SPARQL запрос для подключения к удаленному источнику DBpedia

19

```
sparql
SELECT ?s WHERE
{
  SERVICE <http://dbpedia.org/sparql>
  { graph ?g
    {
      ?a rdf:type  dbpedia-owl:Film;
        rdfs:label ?film.
      FILTER( lang(?film) = "ru" )
    }
  }
}
```

Пример страницы, посвященной фильму «Неудержимые 2»

About: [Неудержимые 2](#)

An Entity of Type : [movie](#), from Named Graph :
<http://dbpedia.org>, within Data Space : [dbpedia.org](#)



«Неудержимые 2» — сиквел вышедшего в 2010 году боевика «Неудержимые», снятый режиссёром Саймоном Уэстом по сценарию Сильвестра Сталлоне, Ричарда Уэнка, Кена Кауфмана и Дэвида Агосто. Прокатом в США занялась компания Lionsgate, в России — Universal Pictures. Сюжет фильма разворачивается вокруг отряда наёмников «Неудержимые», которые в ходе выполнения очередного задания сталкиваются с убийством одного из своих товарищей.

Property	Value
dbpedia-owl:Work/runtime	<ul style="list-style-type: none"> 103.0
dbpedia-owl:abstract	<ul style="list-style-type: none"> The Expendables 2 ist ein US-amerikanischer Actionfilm von Richard Wenk, Ken Kaufman, David Agosto und Fortsetzung zu The Expendables aus dem Jahr 2012. Statham, Jet Li, Dolph Lundgren, Chuck Norris, F. Williams und Arnold Schwarzenegger über ein All-Star Team in Deutschland der 30. August 2012. The Expendables 2 is a 2012 American ensemble action film directed by Simon West, written by Richard Wenk and Sylvester Stallone and based on a story by Ken Kaufman, David Agosto and Wenk. Brian Tyler returned to score the film. It is a sequel to the 2010 action film <i>The Expendables</i>, and stars Sylvester Stallone, Jason Statham, Liam Hemsworth, Jean-Claude Van Damme, Mickey Rourke, Randy Couture, Jason Statham, Jet Li, Dolph Lundgren, Chuck Norris, F. Williams and Arnold Schwarzenegger.



WIKIPEDIA
The Free Encyclopedia

[Main page](#)
[Contents](#)
[Featured content](#)
[Current events](#)
[Random article](#)
[Donate to Wikipedia](#)
[Wikimedia Shop](#)

Interaction
[Help](#)
[About Wikipedia](#)
[Community portal](#)
[Recent changes](#)
[Contact page](#)

[Tools](#)

[Print/export](#)

[Languages](#)

Article [Talk](#)

[Read](#) [View source](#)

Search

From Wikipedia, the free encyclopedia

The Expendables

The Expendables 2 is a 2012 American ensemble action film directed by Simon West, written by Richard Wenk and Sylvester Stallone and based on a story by Ken Kaufman, David Agosto and Wenk. Brian Tyler returned to score the film. It is a sequel to the 2010 action film *The Expendables*, and stars Sylvester Stallone, Jason

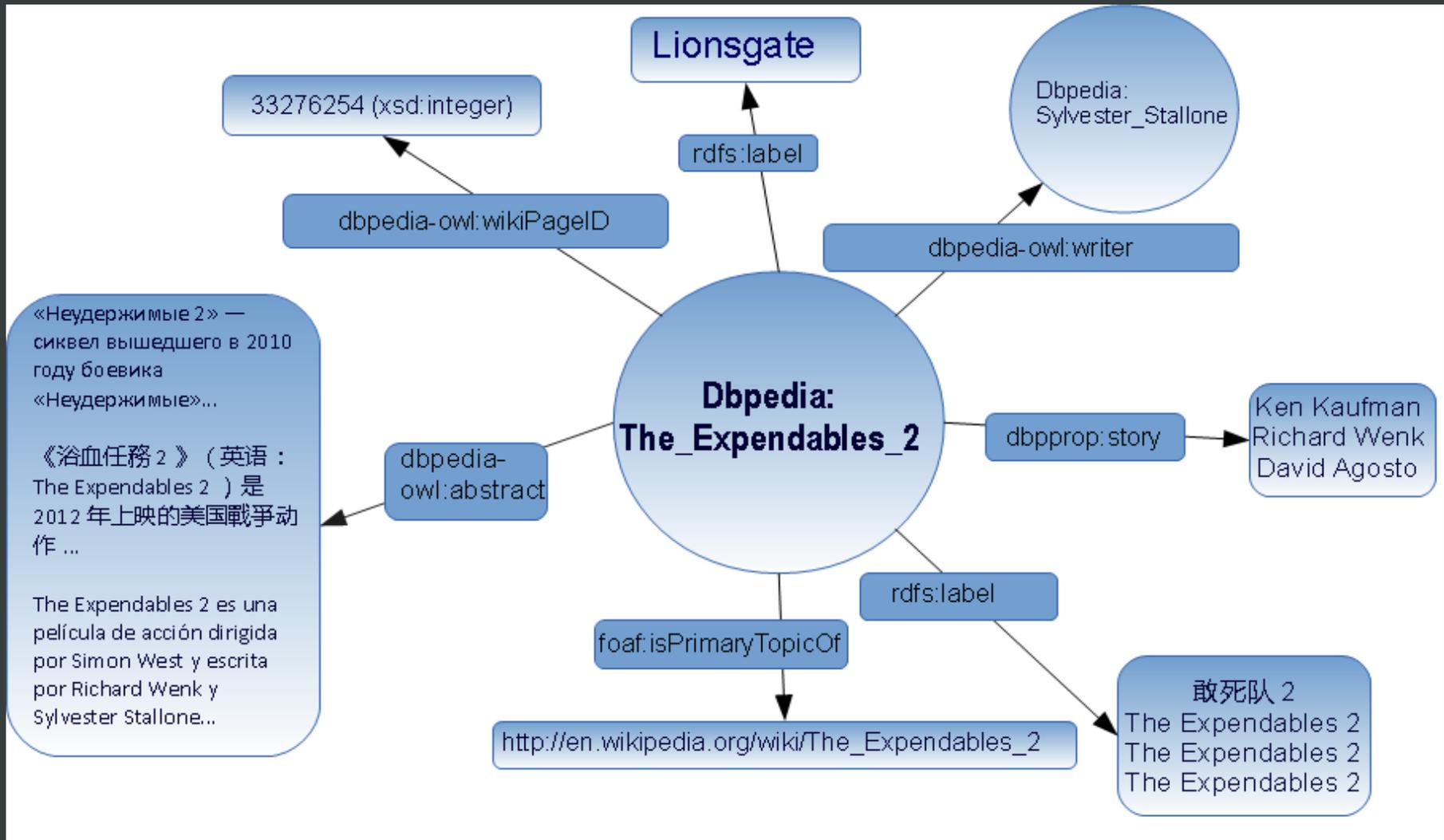
[Create account](#) [Log in](#)

The Expendables 2



Графическое представление данных

21



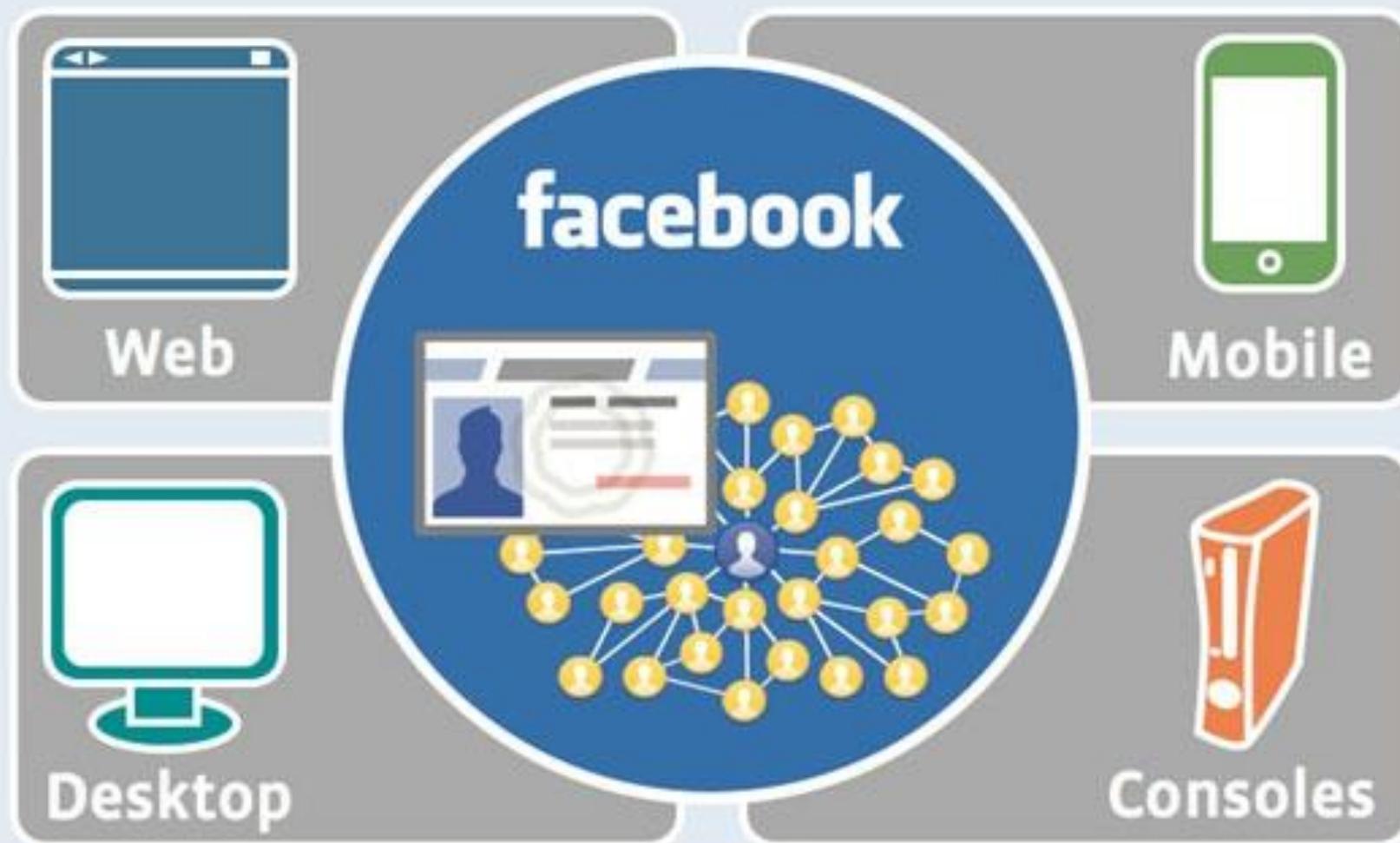
Хранилище триплетов Virtuoso

22

**Virtuoso - Enterprise-
решение на основе
RDBMS Middleware
компании
OpenLink.**

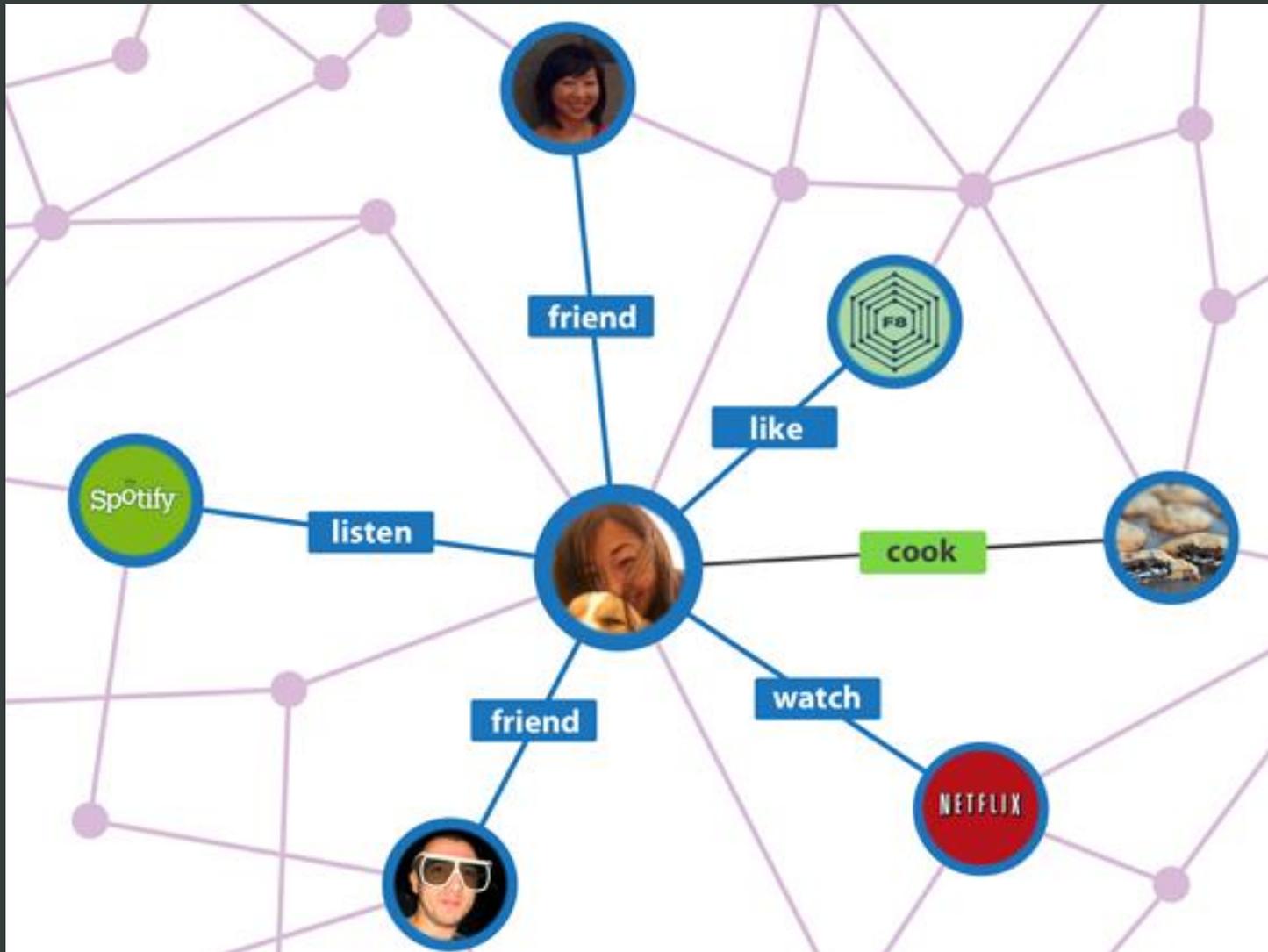


Open Graph



Пример Open Graph

24



Выводы

...

Wolfram Alpha:

Computational Knowledge Engine

IBM Watson

IBM Big Data Platform

Google Knowledge Graph

И многое другое из мира Big Data формирует будущее пользовательских приложений и сервисов.

Уже сегодня

Вопросы???

Спасибо за внимание

<http://shcherbak.net>