

**빅데이터로 지능적 처리에
날개를 달다**

Overview

- Big data
- Text classification
- How to handle dark data?

Big data

Big data $\approx 3V$

[Volume]

parallel computing

- data partitioning (*****)
- process partitioned data in parallel
- merge processed data (*)

Making big data smaller (sampling)

- approximate computing

when computation cost is extremely high and data is not that big

- replicate data for each machine
- making computation fine-granule

[Velocity]

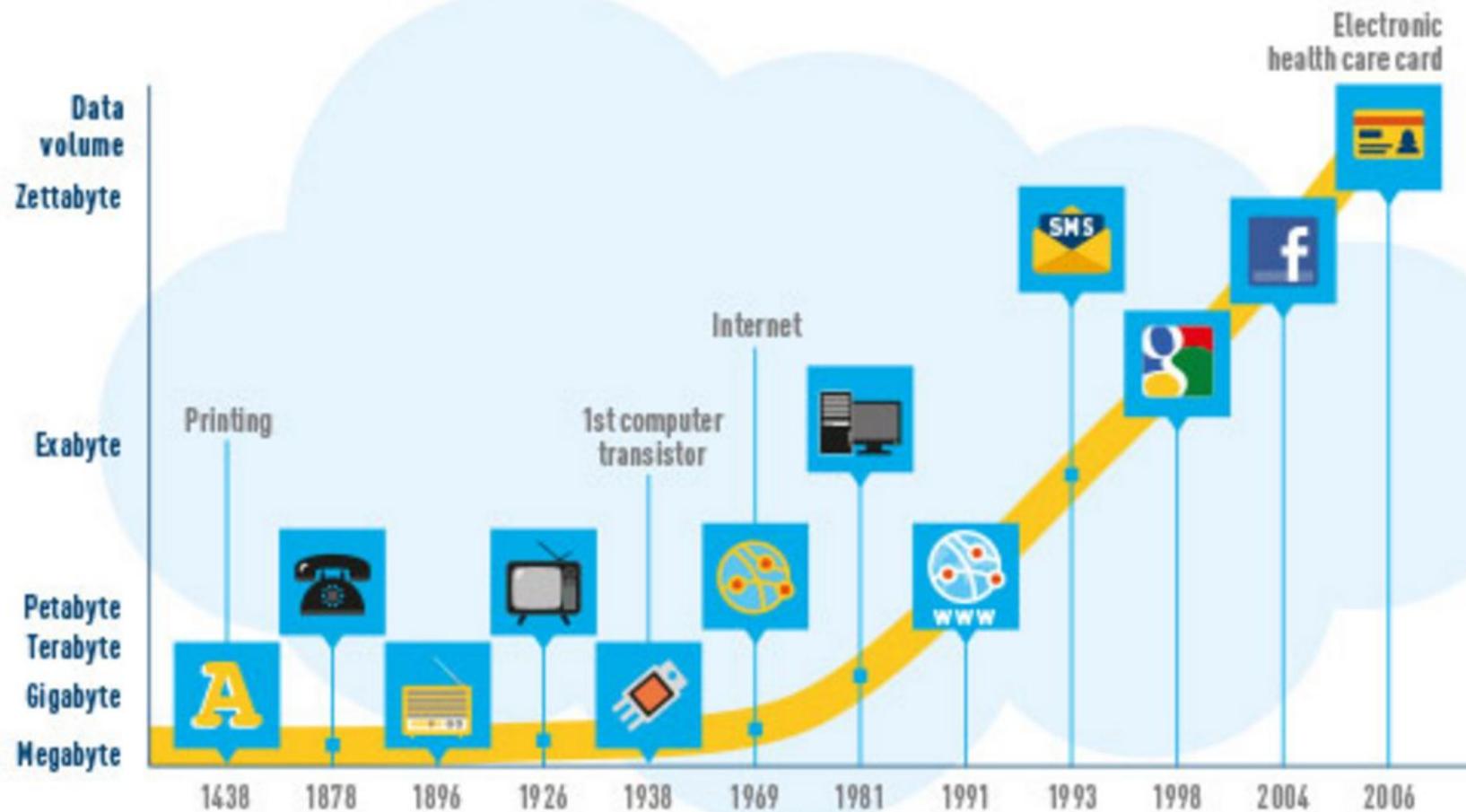
volume parallelism
 volume velocity

[Variety]

AI technology (knowledge graph, large language model)

Volume

: 18~24
가 , 가)



Source: Federal Association for Information Technology, Telecommunication and New Media (BITKOM). "Big Data im Praxiseinsatz - Szenarien, Beispiele, Effekte."

What if data does not fit in memory?

- Make algorithms scalable!
- Exploit all hardware resources at the same time
 - Macro parallelism (Scale-out)
 - Cloud computing
 - Infiniband
 - Micro parallelism (Scale-up)
 - Multi-core (Multi-core CPUs and GPUs)
 - NUMA
 - Parallel I/Os

Latency Numbers
scalable

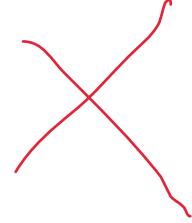
sw

latency
,

sw
hw
가

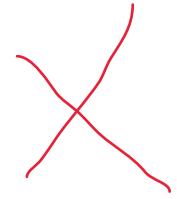
CPU
non uniform memory access:
s/w numa architecture

Scalable matrix factorization



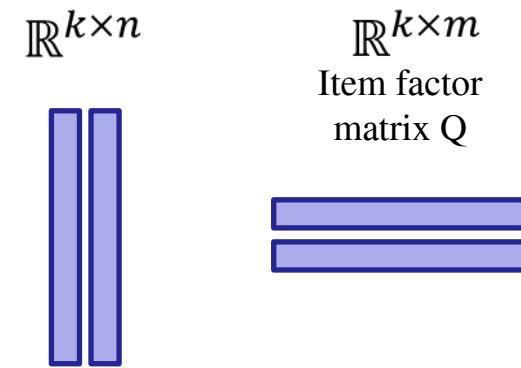
	<i>Avatar</i>	<i>The Matrix</i>	<i>Up</i>
<i>Alice</i>	?	4	2
<i>Bob</i>	3	2	?
<i>Charlie</i>	5	?	3

Matrix Factorization (MF)

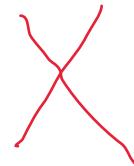


	<i>Avatar</i>	<i>The Matrix</i>	<i>Up</i>
<i>Alice</i>	?	4	2
<i>Bob</i>	3	2	?
<i>Charlie</i>	5	?	3

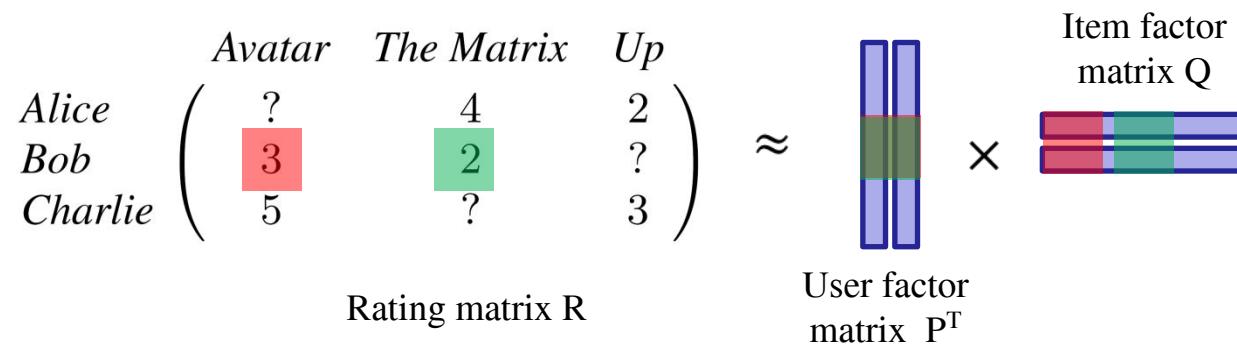
Rating matrix R
 $\mathbb{R}^{n \times m}$



SGD for MF



- SGD gradually updates factor matrices by repeating following two steps
 - Randomly pick a rating data
 - Update corresponding vectors in factor matrices



Velocity



The New York Stock Exchange captures

1 TB OF TRADE INFORMATION

during each trading session



By 2016, it is projected there will be

18.9 BILLION NETWORK CONNECTIONS

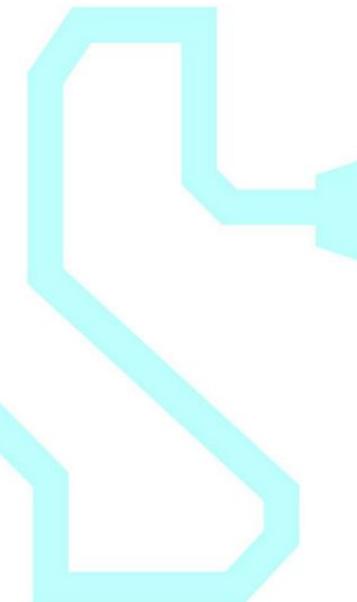
– almost 2.5 connections per person on earth

Velocity

ANALYSIS OF STREAMING DATA

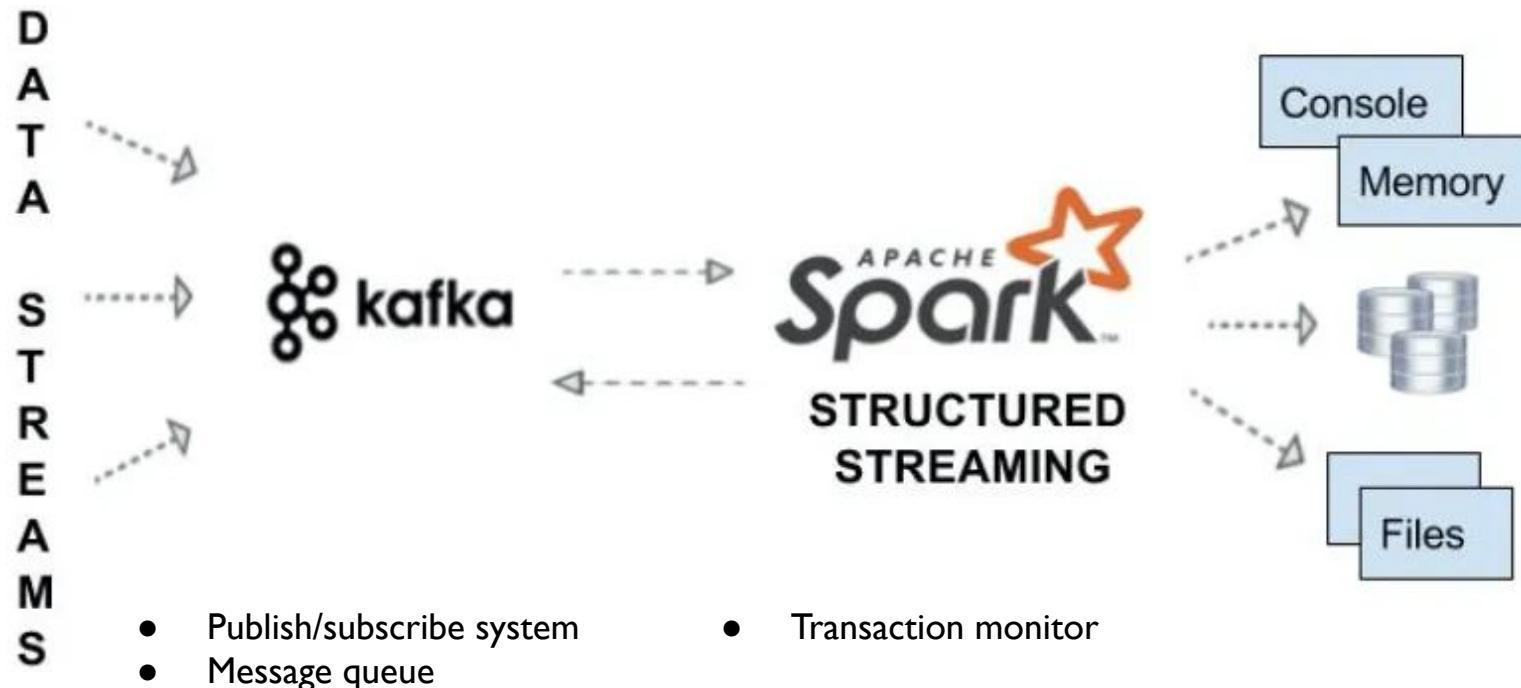


Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS

stateless vs. stateful



, , DBMS()

가

publish/subscribe system

Variety

- Excel sheets, relational data, csv files Structured data
- XML, JSON, graph Semi-structured data
- Text documents, web documents, images, audios, videos Unstructured data

schema(
table)가 data
가 schema 가

Structured data

(= tables)

	A	B	C	D	E	F	G	H	
1		Software Engg	SE II	Calculus	SE III	RDBMS	Oracle	Data Structure	Total
2	Jack	18	80	68	62	65	71	89	453
3	Billy	36	52	52	76	66	10	27	319
4	McFaden	33	75	51	34	79	52	68	392
5	Steven Shwimmer	40	10	12	66	63	28	87	306
6	Ruby Jason	56	49	76	49	81	26	10	347
7	Mark Dyne	67	44	25	67	6	30	86	325
8	Philip namdaf	86	30	62	27	66	48	64	383
9	Erik Bawn	44	5	20	42	9	32	94	246
10	Ricky ben	13	17	51	64	80	70	47	342
11	Miecky	52	16	42	13	61	73	78	335

Excel file

Region,City,Vendor,Product ID,Product,Unit,Cases Sold,Total Sales
North GA,Blue Ridge,Mountain Fruit,0100,Oranges,CASE,6168,61680
North GA,Blue Ridge,Mountain Fruit,0200,Apples,CASE,6079,85106
North GA,Blue Ridge,Mountain Fruit,0300,Kiwi,CASE,6058,66638
North GA,Blue Ridge,Mountain Fruit,0400,Bananas,CASE,6868,75548
North GA,Blue Ridge,Mountain Fruit,0500,Mixed Berries,CASE,1996,29940
North GA,Atlanta,Bob's Fruit,0100,Oranges,CASE,7818,93816
North GA,Atlanta,Bob's Fruit,0200,Apples,CASE,1565,21910
North GA,Atlanta,Bob's Fruit,0300,Kiwi,CASE,9967,99670
North GA,Atlanta,Bob's Fruit,0400,Bananas,CASE,9842,98420
North GA,Atlanta,Bob's Fruit,0500,Mixed Berries,CASE,8993,89930
North GA,Atlanta,Fruitju,0100,Oranges,CASE,4933,54263
North GA,Atlanta,Fruitju,0200,Apples,CASE,7704,107856
North GA,Atlanta,Fruitju,0300,Kiwi,CASE,5519,71747
North GA,Atlanta,Fruitju,0400,Bananas,CASE,8442,126630
North GA,Atlanta,Fruitju,0500,Mixed Berries,CASE,889,11557
North GA,Atlanta,Orange U Glad,0100,Oranges,CASE,6551,72061
North GA,Atlanta,Orange U Glad,0200,Apples,CASE,2605,31260
North GA,Atlanta,Orange U Glad,0300,Kiwi,CASE,3317,43121
North GA,Atlanta,Orange U Glad,0400,Bananas,CASE,7411,81521
North GA,Atlanta,Orange U Glad,0500,Mixed Berries,CASE,6227,93405
North GA,Blue Ridge,Mountain Fruit,0100,Oranges,CASE,6415,89810
North GA,Blue Ridge,Mountain Fruit,0200,Apples,CASE,6426,83538
North GA,Blue Ridge,Mountain Fruit,0300,Kiwi,CASE,8035,112490
North GA,Blue Ridge,Mountain Fruit,0400,Bananas,CASE,5075,60900
North GA,Blue Ridge,Mountain Fruit,0500,Mixed Berries,CASE,3064,36768
North GA,Clarkesville,Fruit Direct,0100,Oranges,CASE,686,9604
North GA,Clarkesville,Fruit Direct,0200,Apples,CASE,8203,82030
North GA,Clarkesville,Fruit Direct,0300,Kiwi,CASE,3920,58800
North GA,Clarkesville,Fruit Direct,0400,Bananas,CASE,8262,107406
North GA,Clarkesville,Fruit Direct,0500,Mixed Berries,CASE,4251,51012
Mid GA,Macon,Middle Georgia Fruit,0100,Oranges,CASE,5469,71097
Mid GA,Macon,Middle Georgia Fruit,0200,Apples,CASE,1126,15764

CSV file

Students Table

Student	ID *
John Smith	084
Jane Bloggs	100
John Smith	182
Mark Antony	219

Participants Table

ID *	Activity *
084	Tennis
084	Swimming
100	Squash
100	Swimming
182	Tennis
219	Golf
219	Swimming
219	Squash

Activities Table

Activity *	Cost
Golf	\$47
Sailing	\$50
Squash	\$40
Swimming	\$15
Tennis	\$36

Relational database

Semi-structured data

(table) schema 가 schema가

```
<? XML VERSION = "1.0" STANDALONE = "yes" ?>
<STAR-MOVIE-DATA>
  <STAR><NAME>Carrie Fisher</NAME>
    <ADDRESS><STREET>123 Maple St.</STREET>
      <CITY>Hollywood</CITY></ADDRESS>
    <ADDRESS><STREET>5 Locust Ln.</STREET>
      <CITY>Malibu</CITY></ADDRESS>
  </STAR>
  <STAR><NAME>Mark Hamill</NAME>
    <STREET>456 Oak Rd.</STREET><CITY>Brentwood</CITY>
  </STAR>
  <MOVIE><TITLE>Star Wars</TITLE><YEAR>1977</YEAR>
  </MOVIE>
</STAR-MOVIE-DATA>
```

XML

{ "users": [

{

"firstName": "Ray",
"lastName": "Villalobos",
"joined": {

"month": "January",
"day": 12,
"year": 2012

}

},

{

"firstName": "John",
"lastName": "Jones",
"joined": {

"month": "April",
"day": 28,
"year": 2010

}

}

}]}

JSON

(Mostly) tree structured data

CSV, JSON, XML

CSV

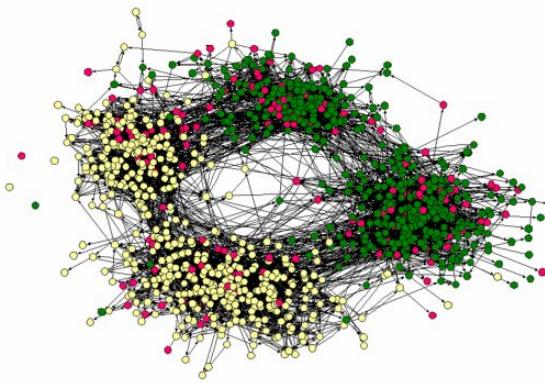
A	B	C	D
1 ID	Gender	City	Monthly
2 ID000002C	Female	Delhi	20000
3 ID000004E	Male	Mumbai	35000
4 ID000007H	Male	Panchkula	22500
5 ID000008I	Male	Saharsa	35000
6 ID000009J	Male	Bengaluru	100000
7 ID000010K	Male	Bengaluru	45000
8 ID000011L	Female	Sindhudur	70000
9 ID000012M	Male	Bengaluru	20000
10 ID000013N	Male	Kochi	75000
11 ID000014C	Female	Mumbai	30000
12 ID000016C	Male	Mumbai	25000
13 ID000018S	Female	Surat	25000
14 ID000019T	Female	Pune	24000
15 ID000021V	Male	Bhubanes	27000
16 ID000022V	Female	Howrah	28000

JSON

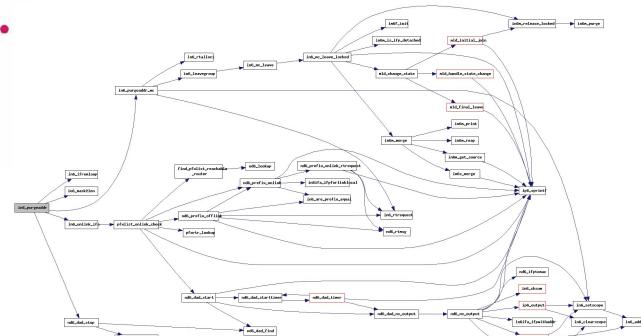
```
1 "Employee": [  
2     {  
3         "id": "1",  
4         "Name": "Ankit",  
5         "Sal": "1000",  
6     },  
7     {  
8         "id": "2",  
9         "Name": "Faizy".  
10    }  
11 ]
```

```
<?xml version="1.0"?>  
  
<contact-info>  
  
<name>Ankit</name>  
  
<company>Analytics Vidhya</company>  
  
<phone>+9187654321</phone>  
  
</contact-info>
```

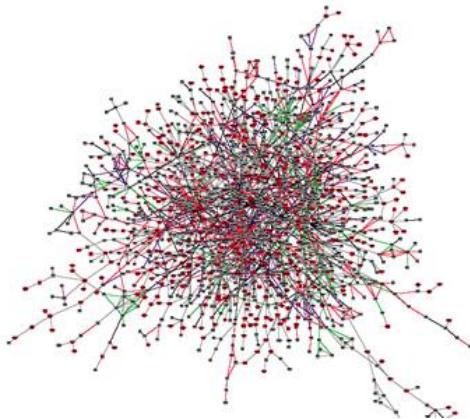
Graph data



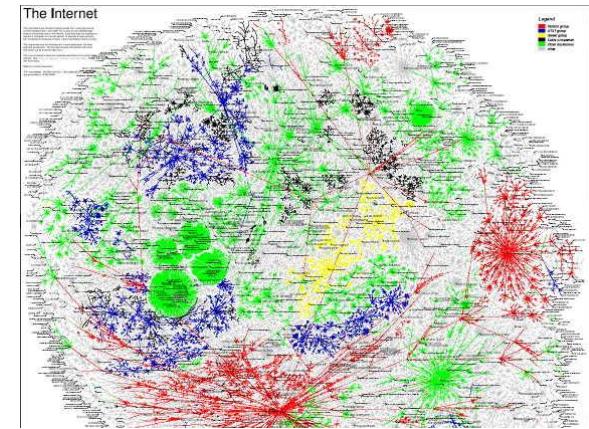
Friendship network



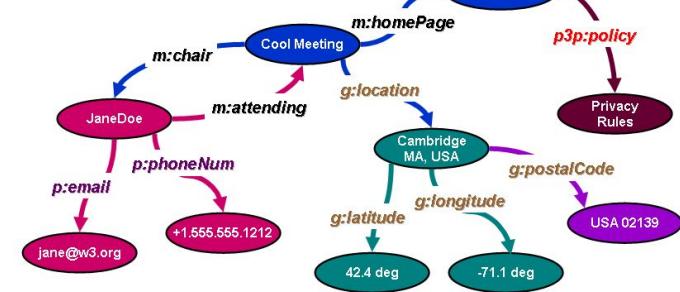
Call Graph



Protein interactions



Internet map



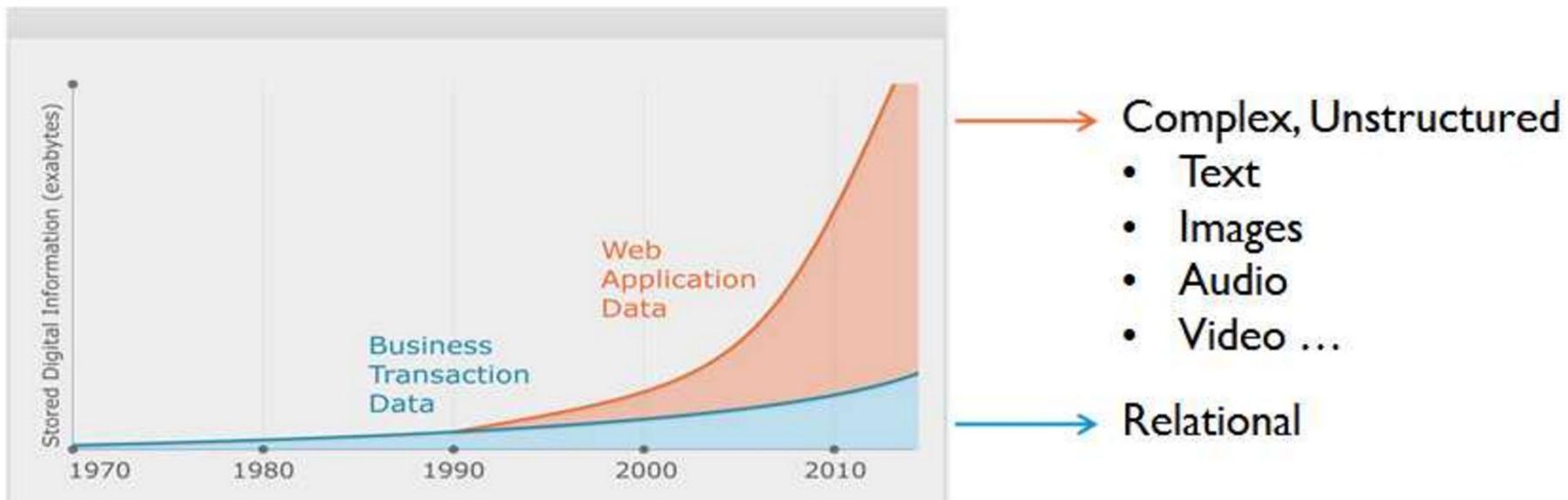
Semantic web

Everything can be easily represented as graphs

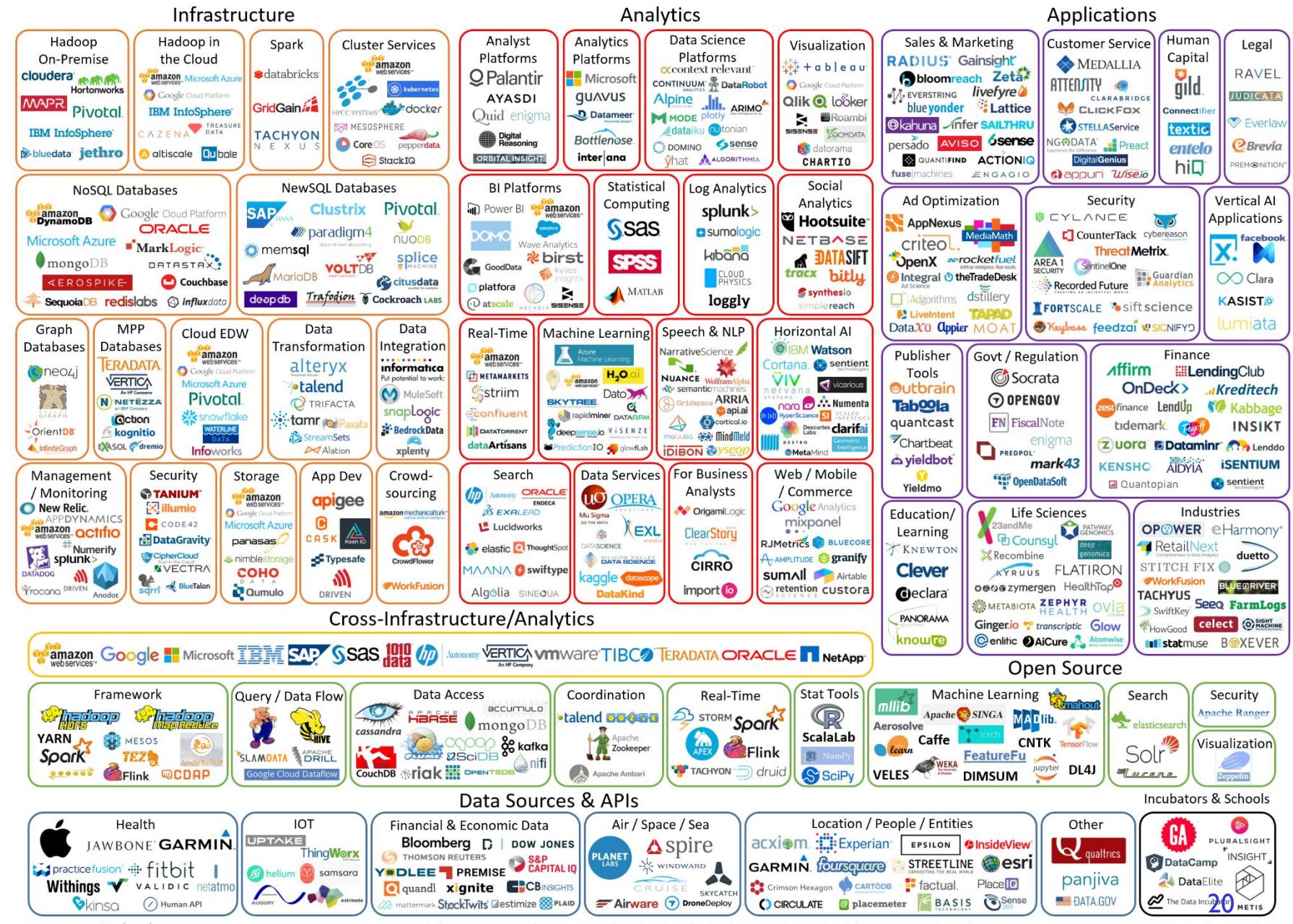
- Right!
- All structured and semi-structured data can be viewed as graphs

Unstructured data (\approx dark data)

가



Big Data Landscape 2016 (Version 3.0)



Infrastructure for Big data is very important

- But, we skip this part completely
 - Students who are interested in building systems love this part
- The rest of the talk focuses on handling variety of data

Text classification

Text classification

- Pop vs Heavy Metal
 - Given verse from lyrics, recognize genre

"I'm a rolling thunder, a pouring rain
I'm comin' on like a hurricane
My lightning's flashing across the sky
You're only young but you're gonna die
I won't take no prisoners, won't spare no lives
Nobody's putting up a fight
I got my bell, I'm gonna take you to hell
I'm gonna get you, Satan get you"

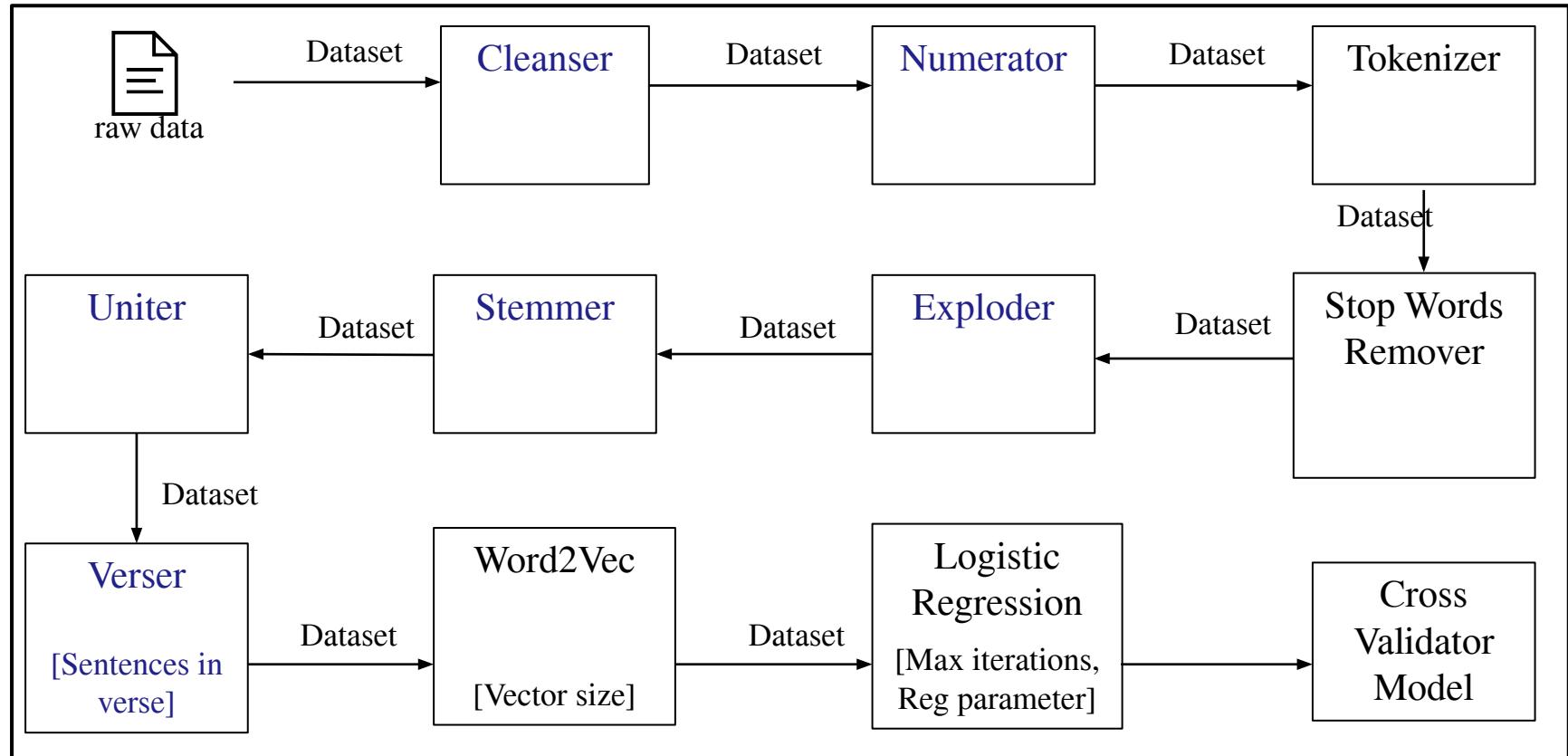
- Collect raw data set of lyrics (~65k sentences)
 - [Pop] Abba, Ace of base, Backstreet Boys, Britney Spears, Christina Aguilera, Madonna, etc.
 - [Heavy Metal] Black Sabbath, In Flames, Iron Maiden, Metallica, Moonspell, Nightwish, Sentences, etc.
- Create training set, i.e. label (0|1) + features
- Train logistic regression

Wait a minute!

- ML takes only numeric vectors as input!
- How about text documents?

Word2Vec

ML pipeline



Cleanser

I'm a rolling thunder, a pouring rain
I'm comin' on like a hurricane
My lightning's flashing across the sky
You're only young but you're gonna die
I won't take no prisoners, won't spare no lives
Nobody's putting up a fight
I got my bell, I'm gonna take you to hell
I'm gonna get you, Satan get you

Im a rolling thunder a pouring rain
Im comin on like a hurricane
My lightnings flashing across the sky
Youre only young but youre gonna die
I wont take no prisoners wont spare no lives
Nobodys putting up a fight
I got my bell Im gonna take you to hell
Im gonna get you Satan get you

Stemmer

im rolling thunder pouring rain
im comin like hurricane
lightnings flashing across sky
you're young you're gonna die
wont take prisoners wont spare lives
nobodys putting fight
got bell im gonna take hell
im gonna get satan get

im **roll** thunder **pour** rain
im comin like **hurrican**
lightn flash across sky
your young your gonna die
wont take **prison** wont spare **live**
nobodi put fight
got bell im gonna take hell
im gonna get satan get

Word2Vec & Logistic regression

```
[0.036463763926011056,  
 -0.013076733228398295,  
 0.044362547532774695,  
 0.03816963326281462,  
 .....  
 -0.013962931134021625,  
 0.049275818325650804,  
 -0.058982484615766086]
```

Probability:

```
[0.9212126972383768,  
 0.07878730276162313]
```

P(Heavy Metal)

P(Pop)

Prediction:

0.0

How to handle dark data?

Which Personal Assistant is the Smartest?

Here are the results of our research:

Personal Assistant	% Questions Answered	100% Complete & Correct
The Google Assistant on Google Home	68.1%	90.6%
Cortana (MS)	56.5%	81.9%
Siri	21.7%	62.2%
Alexa on the Amazon Echo	20.7%	87.0%

<https://www.stonetemple.com/digital-personal-assistants-test>

Do you know the secret of success?

Job posting information extraction

unstructured data structured data

>

Title: Web Development Engineer

Location: Beaverton, Oregon

2017.04.21

This individual is responsible for design and implementation of the web-interfacing components of the AccessBase server, and general back-end development duties.

A successful candidate should have experience that includes:

One or more of: Solaris, Linux, IBM AIX, plus Windows/NT

Programming in C/C++, Java

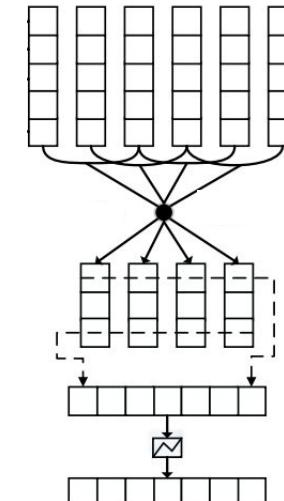
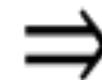
Database access and integration: Oracle, ODBC

CGI and scripting: one or more of Javascript, VBScript, Perl,

PHP, ASP

Exposure to the following is a plus: JDBC, Flash/ Shockwave,
FrontPage and/or Cold Fusion.

A BSCS and 2+ years experience (or equivalent) is required.



정형화



Title
Web
Development
Engineer

Date
04-21-
2017

Location
Beaverton,
Oregon

Skills
Windows/NT
Programming
in C/C++, Java

...

...

Demo

- Is being developed in our lab

The screenshot shows a web browser window with the following details:

- Address Bar:** 141.223.199.44:8983/solr/techproducts/browse
- Solr Logo:** On the left side of the search bar.
- Search Fields:** Keyword (with magnifying glass icon) and Attribute (with magnifying glass icon).
- Filter Options:** Company, Country, State, City, Education, Sort.
- Job Listings:**
 - Manuscripts Processing Archivist** LA, Baton Rouge
 - Louisiana State University
 - Required Education: Master's degree
 - Required Experience: 1 year
 - Posting Date: 20170629
 - Digital Library Projects Coordinator** MO, St. Louis
 - Federal Reserve Bank
 - Required Education: Master of Library Science or Information Science
 - Required Experience: 5 years related experience2 years supervisory experience
 - Posting Date: 20170629
 - Reference & Outreach Librarian** NJ, Wayne
 - William Paterson University
 - Required Education: A second graduate degree
 - Required Experience: 1
 - Posting Date: 20170629
 - Electronic Resources Librarian** ID, Pocatello
 - Idaho Division of Human Resources
 - Required Education: Master's degree
- Video Recording Application:** oCam (1, 24, 1365, 703)
 - Buttons: 메뉴, 화면 녹화, 게임 녹화, 소리 녹음.
 - Controls: 중지 (Stop), 일시중지 (Pause), 캠처 (Capture).
 - Display: 00:00:00, 0bytes / 697.8GB.
 - Text: 원도우 7 이상에서 마이크 녹음 시 유용한 팁.

Structured data in search results

Recognizing Set-Oriented Queries

Association Queries

Behind the search results

- **Knowledge graph**

knowledge a set of facts
a fact: a triple <subject, predicate, object>

- Large-scale extraction of facts from Web text
- Incorporate it into an IR-style search engine

Resource Description Framework (RDF) data

<Bob> <is a> <person>.

<Bob> <is a friend of> <Alice>.

<Bob> <is born on> <the 4th of July 1990>.

<Bob> <is interested in> <the Mona Lisa>.

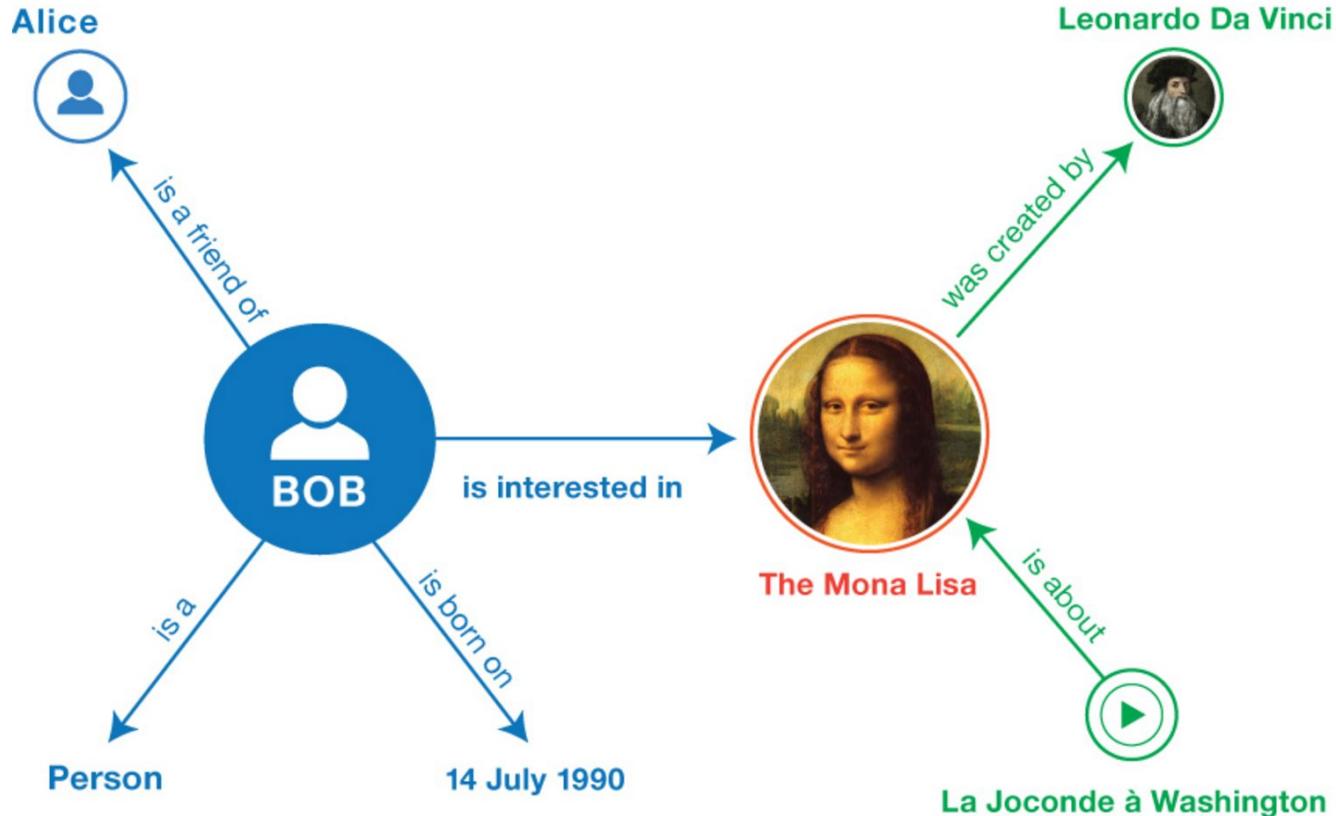
<the Mona Lisa> <was created by> <Leonardo da Vinci>.

<the video 'La Joconde à Washington'> <is about>
<the Mona Lisa>

Resource Description Framework (RDF) data (cont.)

- represent a knowledge
- A knowledge is a set of facts, each of which is represented by a triple (subject, predicate, object)

RDF data is a kind of graph

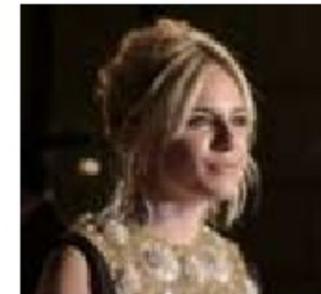


Freebase is a database of entities

One entity per thing in the world

Stable, long-lived identifiers

Inclusive policy



[/en/sienna_miller](#)



[/en/sony_dsc_s750](#)

Practical data

Focus on available data

People, places, products, etc.

Data to build apps

Names, images, descriptions

Dates, measurements and relationships

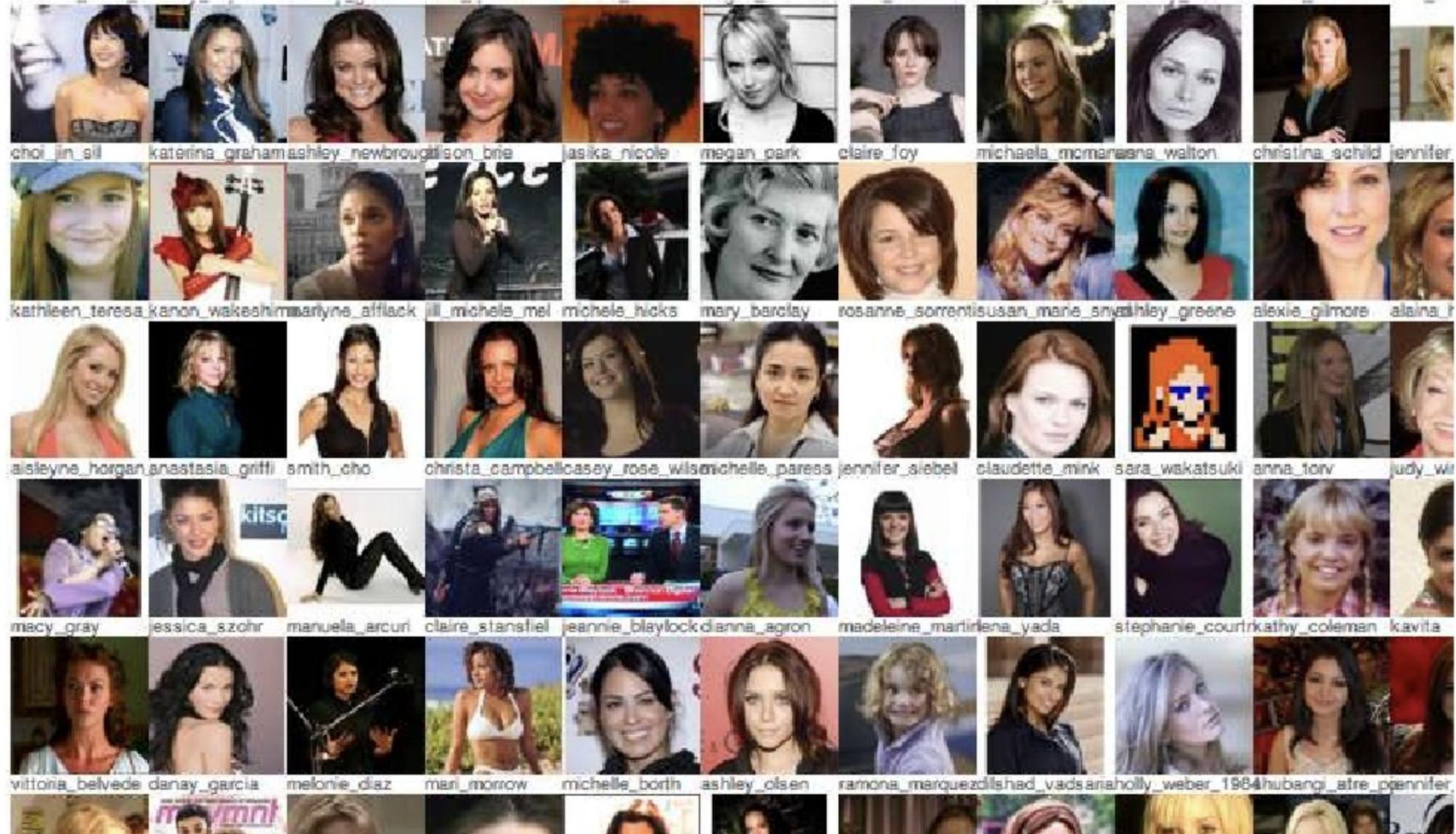
Frost/Nixon (2008)



Director:	Ron Howard
Released:	Oct 15, 2008
Runtime:	2 hr. 2 min.
Genres:	Drama
Box office:	\$18,593,156

[/en/frost_nixon_2008](#)

Actresses (37,079)



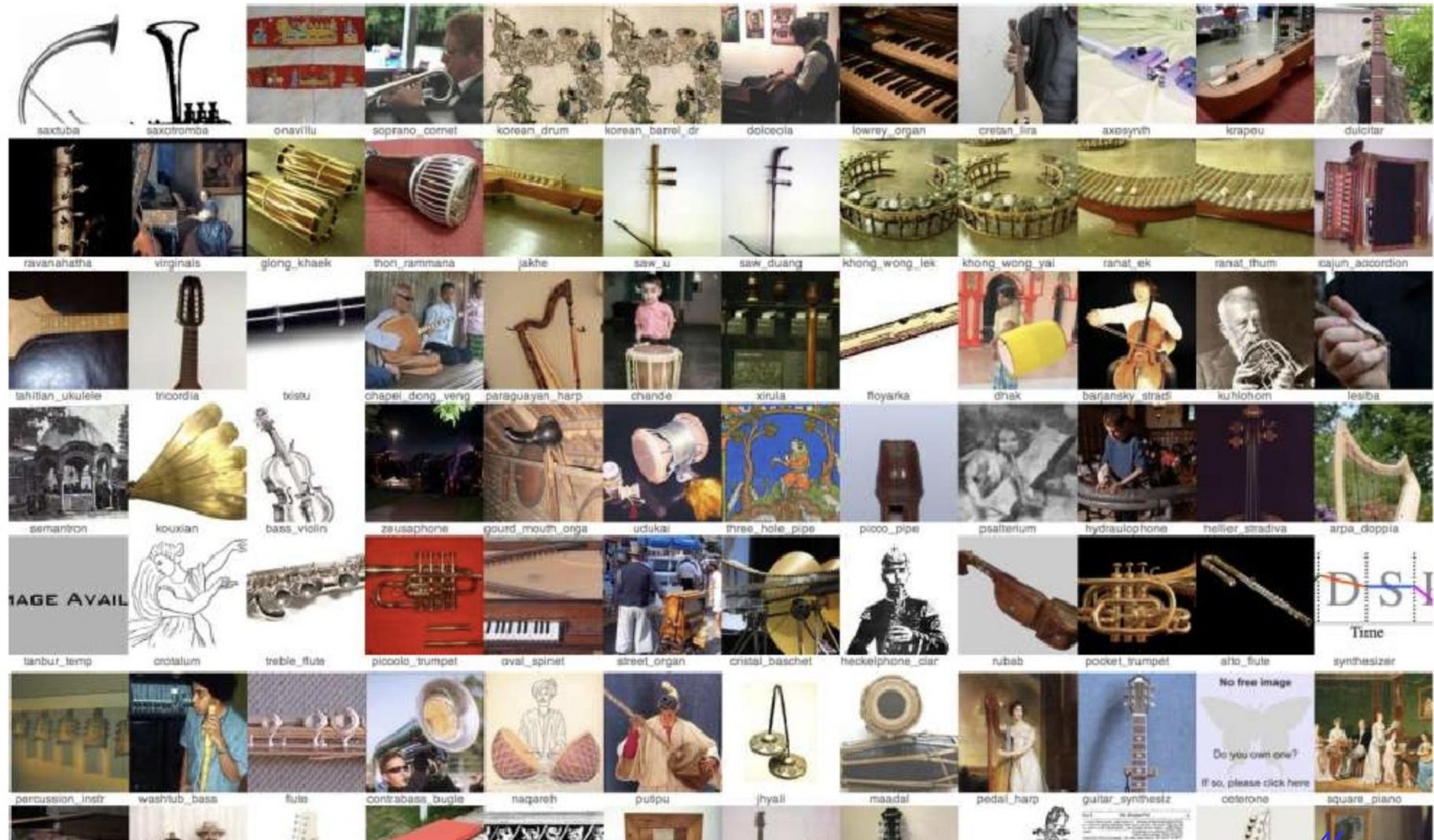
Football Players (16,568)



Cheeses (488)



Musical Instruments (1,034)



Airports (11,556)



TV Programs (33,630)



Related entities are connected, forming a graph



Current stats:

- ~10M entities
 - Celebrity
 - Movie
 - TV show
 - Book
 - Company
 - Location
 - Sports team
 - Product
 - Etc.
- ~1,800 “types”
 - Celebrity
 - Movie
 - TV show
 - Book
 - Company
 - Location
 - Sports team
 - Product
 - Etc.
- ~275M facts
- Continuous data input, cleanup, and syncing

DBpedia

- **Extracting Structured Information from Wikipedia**
- **Wikipedia**
 - **6.9 million articles**
 - **in 251 languages**
 - **monthly growth-rate: 4%**

Structured information in Wikipedia

- Wikipedia articles contain infoboxes which use a template mechanism
- Images depicting the article's topic
- Categorization of the article
- Links to external webpages
- Intra-wiki links to other articles
- Inter-language links to articles about the same topic in different languages

province, in an area
"Calgary–Edmonton

pulation of 1,392,609,

alth and wellness,
ions.^[5]

old" and "garden",
(pasture); or Gaelic

Inhabited by the
ompson spent the
was the first

NWMP detachment
ois, after NWMP

an important
Montreal in 1996.^[10]
ed as "The City of
municipal, local duties

200. Although no one
to pre

As a re
Generated from: Infobox settlement
itre of C

Calgary

City

City of Calgary

image_skyline

image_size

image_caption

image_flag

image_seal

image_shield

nickname

motto

image_map

Flag

Coat of arms

Nickname(s): Cowtown, Stampede City, more... [1][2]

Motto: *Onward*

Template

Edit

Generated from: Infobox settlement

Infobox

Cancel Transclusion Apply changes

Infobox settlement

an Infobox for human settlements (cities, towns, villages, communities) as well as other administrative districts, counties, provinces, etcetera

There might be some additional information about the "Infobox settlement" template on its page.

name
Calgary

official_name
City of Calgary

settlement_type
City

image_skyline

Calgarymontage5.jpg <!--Montage determined through past

Hide options

Extracting Infobox Data

<http://en.wikipedia.org/wiki/Calgary>

```
<http://dbpedia.org/resource/Calgary>
  dbpedia:native_name "Calgary" ;
  dbpedia:altitude "1048" ;
  dbpedia:population_city "988193" ;
  dbpedia:population_metro "1079310" ;
  mayor_name
    dbpedia:Dave_Bronconnier ;
  governing_body
    dbpedia:Calgary_City_Council ;
  ...
```

- Altogether 9,100,000 RDF triples extracted from 754,000 infoboxes

Calgary	
	
Downtown Calgary.	
Government	
- Mayor	Dave Bronconnier (Past mayors)
- Governing body	Calgary City Council
- Manager	Owen A. Tobert
Area [1]	
- City	726.50 km ² (280.5 sq mi)
- Metro	5,107.43 km ² (1,972 sq mi)
Elevation	1,048 m (3,438.3 ft)
Population (2006) [1]	
- City	988,193
- Density	1,360.2/km ² (3,522.9/sq mi)
- Metro	1,079,310
rank	3rd
- Metro rank	5th

Christian Bizer et al: DBpedia – Querying Wikipedia Like a Database (May 11, 2007)

Extracting Other Article Data

■ Short and long abstracts in 10 different languages

```
dbpedia:Calgary dbpedia:abstract "Calgary is the largest ..."@en ;
                 dbpedia:abstract "Calgary ist eine Stadt ..."@de .
```

■ Categorization information

```
dbpedia:Calgary skos:subject dbpedia:Category_Cities_in_Alberta ;
                 skos:subject dbpedia:Host_cities_Olympic_Games .
```

■ Links to the original Wikipedia articles, pictures and relevant external web pages

```
dbpedia:Calgary
        foaf:page <http://en.wikipedia.org/wiki/Calgary> ;
        dbpedia:wikidata <http://de.wikipedia.org/wiki/Calgary> ;
        foaf:depiction <http://upload.wikimedia.org/thumb/3/32> ;
        dbpedia:reference <http://www.calgary.ca> ;
        dbpedia:reference <http://www.tourismcalgary.com>.
```

The DBpedia Dataset

- **1,600,000 concepts**
 - 58,000 persons
 - 70,000 places
 - 35,000 music albums
 - 12,000 films
- **described by 91 million triples using 8,141 different properties.**
- **557,000 links to pictures**
- **1,300,000 links to relevant external web pages**

Example queries

- **People who were born in Berlin before 1900**
- **German musicians with German and English descriptions**
- **Musicians who were born in Berlin**
- **Soccer players, who are born in a country with more than 10 million inhabitants, who played as goalkeeper for a club that has a stadium with more than 30.000 seats and the club country is different from the birth country**

Current research interests

- Parallel (/Distributed/Concurrent)
 - Effectively utilizing all hardware resources
 - Graph
 - Flexible representation
 - Intelligent
 - Incorporating the new wave of deep learning/AI for dark data
- 
- COMPUTATION**

