

ECE_5DA04_TP

Big Graph Databases

Garima Gaur, Ioana Manolescu,

Madhulika Mohanty

INRIA & Ecole Polytechnique

firstname.lastname@inria.fr

DataAI Master

Institut Polytechnique de Paris

Course goal

1. Discuss the main **characteristics (dimensions)** of Big Database systems, focused around graphs.
2. Present the main **concepts** underlying such systems.
3. Motivate and present **architectural choices** made to scale better.

Course organization

- **Instructors:**

Ioana Manolescu (Inria and Ecole Polytechnique)

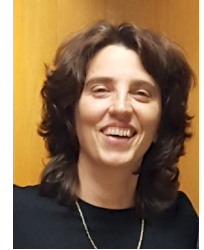
ioana.manolescu@inria.fr

Madhulika Mohanty (Inria and Ecole Polytechnique)

madhulika.mohanty@inria.fr

Garima Gaur (Inria and Ecole Polytechnique)

garima.gaur@inria.fr



- **Evaluation:** final exam (50%) + lab work (50%)

- All course material in Moodle (Self enrol):

<https://moodle.r2.enst.fr/moodle/course/view.php?id=162>



What characterizes a Big Database system?

- Functionality provided
 - What kind of data can I put in? **Data model**
 - How can I get data out of it? **Query language or API**
 - How does it handle concurrent access? **Concurrency control**
 - How long does a given operation take? **Performance**
- Implementation (internals)
 - How does it cope with scale?
 - for reads? **Data and work distribution**
 - also for writes? **Distributed concurrency control**

Topics to be covered

- Recall/crash course on **relational database management systems (RDBMS)**
 - Data model, query language, performance
- **Graph Data Model**
 - Different data models
 - Query Languages
 - Structured, Semi-structured, Unstructured
- **Distributed databases**
- **Reasoning on knowledge**
- **Heterogeneous data integration systems**
 - Local-as-view, global-as-view
 - Mediator systems
 - Dataspaces
 - Data lakes
- Massive parallelism seen in another course (Hadoop, Spark)

What about NoSQL?

1. NoSQL is mostly about **distributed concurrency control**
 - **Weaker guarantees** than a centralized RDBMS...
 - In exchange for **better performance at scale**
2. NoSQL systems also typically have **other data models and languages** than relational/SQL.
 - Key-values (Redis), JSON store (MongoDB), **graphs (Neo4J, Virtuoso)**

Today's course plan

1. Motivation: Big Data

- Characteristics
- Applications

2. From databases to architectures for Big Graph Data management

- Database management system: quick recall (or crash course)
- What needs to change to handle Big Data?

MOTIVATION: BIG DATA

Defining Big Data: the V's

- **Volume**
 - **Scale**
- **Velocity**
 - Speed of producing and consuming the data
- **Variety**
 - Very different sources and data types
- **Veracity**
 - Is the data correct / certain / true?

Where does the data volume come from?

- Human-produced data

- **Web content**: Web pages, blogs, social networks, tweets...



- **X**: 7 Terabytes (1 tera = 10^{12}) per day



- **Facebook**: 10 Terabytes per day



Where does the data volume come from? (1)

- Human-produced data

- **Web content**: Web pages, blogs, social networks, tweets...
- **Twitter**: 7 Terabytes (1 tera = 10^{12}) per day
- **Facebook**: 10 Terabytes per day



- Machine-produced data

- Log data from all kind of servers
- Real world devices: banks, telecom, energy, weather, transportation, shipment...
- Sensors, including on highways
- and trains

Gazpar (GRDF)



Linky (EDF)



Where does the data volume come from? (2)

- E.g. french railway system: surveillance trains for the normal and high-speed lines (TGV)
- TGV specially equipped for measuring while circulating at 320 km/h:
 - rail geometry
 - train/rail interaction
 - rail signalization and communication devices
 - electric power availability etc.
 - 150 sensors, 20 cameras



Gordon Bell, Microsoft, 2009

"It's like having a **multimedia transcript** of your life.

By about 2020 [...] **our entire life histories will be online and searchable.**

Location-aware smartphones and inexpensive digital memory storage in the "cloud" of the Internet make the transition possible and inevitable.

No one will have to fret about storing the details of their lives in their heads anymore. We'll have computers for that.

And this revolution will *"change what it means to be human"*

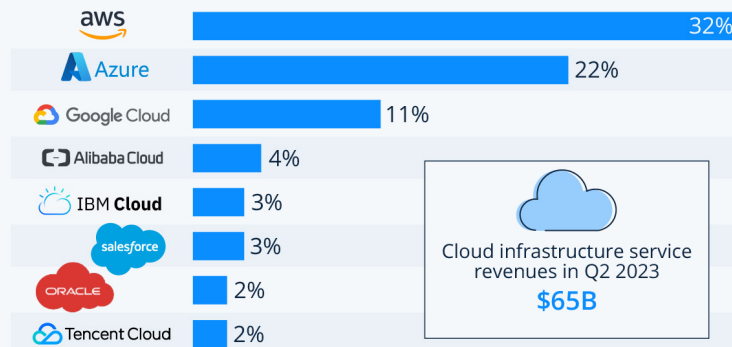


Huge data volumes lead to distributed storage

The typical architecture for large-scale distributed storage and computing is cloud-based

Amazon Maintains Lead in the Cloud Market

Worldwide market share of leading cloud infrastructure service providers in Q2 2023*



* Includes platform as a service (PaaS) and infrastructure as a service (IaaS) as well as hosted private cloud services
Source: Synergy Research Group



statista

Microsoft, Amazon and Alphabet's collective cloud capex is expected to grow
\$bn



Forecasts for 2023, 2024 and 2025. Excludes Amazon's retail investments.
Source: Bank of America Global Research
© FT

Defining Big Data: the V's

- Volume
 - Scale
- Velocity
 - Speed of producing and consuming the data
- Variety
 - Very different sources and data types
- Veracity
 - Is the data correct / certain / true?

Big Data velocity

- How much data is produced e.g. per second
- Data enters a **pipeline** consisting of storage and/or processing
 - **Store-then-process**: for off-line data analysis. Storage by itself is a challenge sometimes, e.g. data links to/from clouds are rather slow
 - **Process-then-store**: for data whose interest is maximized upon arrival (real-time processing)
 - **Process-then-discard**: sensor/monitoring (if nothing happens)

Sample high-throughput data streams

- French IT company runs a data center of **2000** servers
- **5000** energy efficiency indicators (temperature, electricity consumption etc.) are measured every **20** seconds x **50** Kb per measure result = 170 Terabytes / year
- Unable to store all data → sample (measure more rarely)
- May miss important things when they happen



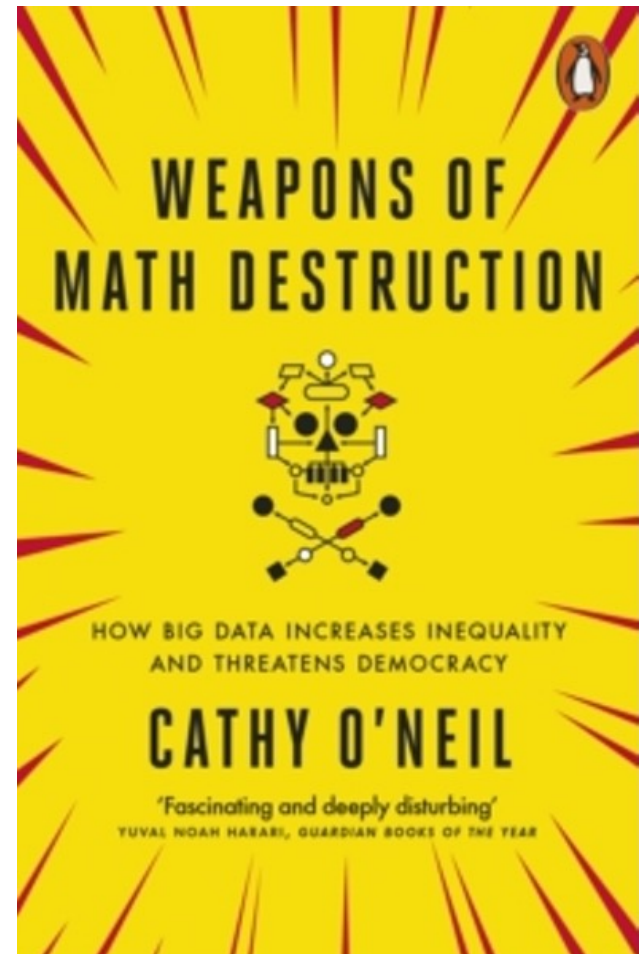
The importance of current/recent data

- Real-time applications work only / mostly with **the latest data**
 - Embedded control mechanisms based on sensor data, e.g., "*this railway wagon component is breaking*" (now!)
 - Intrusion or malfunctioning detection...
- Keeping **humans engaged**
 - Customer relationship management *while the client is on the phone* with the customer service
(see: [“Ordering pizza in the future” video](#))
 - Recommending places where your friends are hanging out *now*



Important aspect of Big Data: ethics

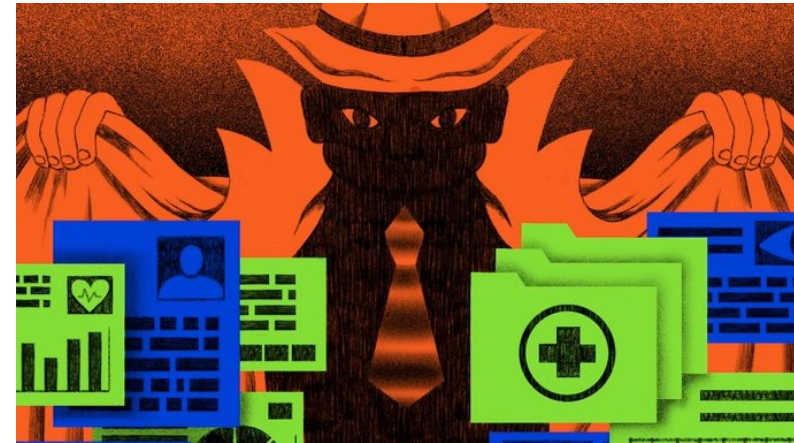
- Current technologies for gathering and processing data raise **risks** wrt personal **freedoms** and **rights**
 - Discriminations, privacy violations, consumer rights, infringement on personal freedom (cf. geo-tagging), political manipulation (cf. focused FB ads)
 - Democracy and liberty may be at risk



Data ethics problems with real consequences

Scenario 1: Anne carries a connected watch recording her movements in the city

- This makes it easy to know where her husband and children often go, even if they did not agree to sharing this info.



<https://t.co/SyLsVTsFTv?amp=1>

Scenario 2: Carol shares her DNA information with a DNA analysis company

- Then a health insurance company buys it to learn that Carol's parents and children share a gene variant associated with an expensive-treatment illness

Defining Big Data: the V's

- Volume
 - Scale
- Velocity
 - Speed of producing and consuming the data
- Variety
 - Different sources, data formats, data types
- Veracity
 - Is the data correct / certain / true?

Big data heterogeneity (variety)

- Each new data type has added up on the old ones
 - Enterprise data typically has **high per-byte value** (\$/byte)
 - Hard to explain that "we will not need this database in the future"
 - In many areas, **legal obligation** to keep old data (e.g. railway sensors, telecom, commercial...)
- Data model & data management system soup
- hierarchical, relational, object-oriented, XML, RDF, JSON, key-value pairs...

Sample relational database



Clients

NumClient	Nom	Adresse	Ville	Age
1	Julie	1 rue Dugommier	Paris	22
...

Comptes

NumCompte	Type	Découvert	NumClient
12345	Courant	1000	1
...

Transaction

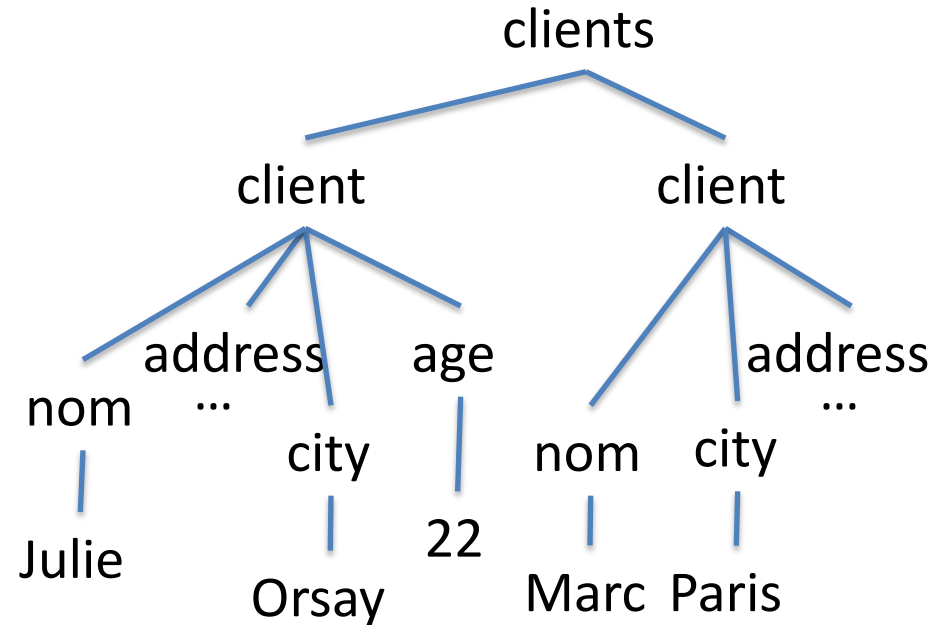
NumCompte	Montant	Date	Info
12345	-40,00	5/10/11	Retrait
12345	+23,45	6/10/11	Remb. MAIF
12345	-300,00	7/10/11	Chaussures

XML: extensible markup language

W3C, 2008

clients.xml:

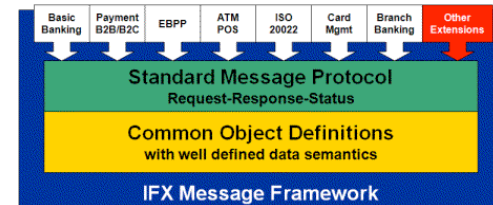
```
<clients>
<client><nom>Julie</nom>
<address>1, rue Dugommier</address>
<city>Paris</city><age>22</age>
</client>
<client><nom>Marc</nom>...
</client>
</clients>
```



Flexible
Platform-independent
Separate content from presentation
Schema possible (not compulsory)

XML applications

- Main language for the Web: **XHTML**, **XML Schema**, **SVG**, **RSS**, ...
- Web Services: **SOAP**, **WSDL**
- **MathML** (mathematical markup language)
- **CML** (chemical markup language)
- **SMILE** (synchronized multimedia integration language)
- Financial Exchange (**IFX**)
- The Text Encoding Initiative (**TEI**)



JavaScript Object Notation (JSON)

Human-readable XML

1. Object = set of (attribute, value) pairs
2. Array = list of values.
3. Value = string | number | true | false | null | object | Array

```
{"menu": {  
  "header": "SVG Viewer",  
  "items": [  
    {"id": "Open"},  
    {"id": "OpenNew", "label": "Open New"},  
    {"id": "ZoomIn", "label": "Zoom In"},  
    {"id": "ZoomOut", "label": "Zoom Out"},  
    {"id": "OriginalView", "label": "Original View"},  
    null,  
    {"id": "Quality"},  
    {"id": "Pause"},  
    {"id": "Mute"},  
    {"id": "Help"},  
    {"id": "About", "label": "About SVG Viewer..."}  
  ]  
}}
```

JavaScript Object Notation

- Among the most popular data interchange formats today
- There exist JSON notations for other data formats, e.g., RDF

```
{
  "http://example.org/about" : {
    "http://purl.org/dc/terms/creator" : [ { "value" : "_:anna",
                                             "type" : "bnode" } ] ,
    "_:anna" : {
      "http://xmlns.com/foaf/0.1/name" : [ { "value" : "Anna",
                                              "type" : "literal" } ]
    }
  }
}
```

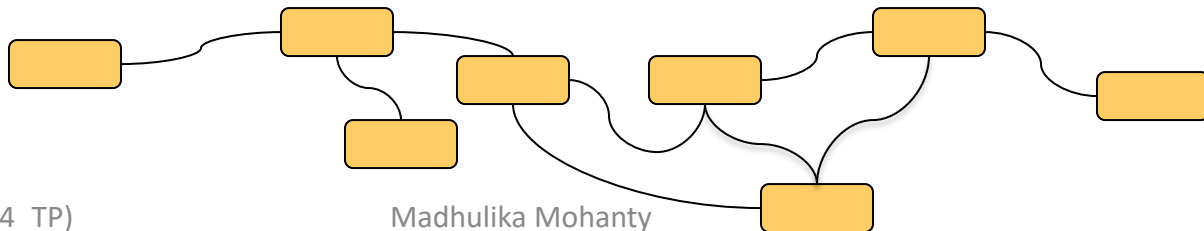
Critique of tree data models (XML, JSON)

- Each information ends up in only one place
- OK for "classification" applications, structured text
- Fundamentally restrictive for **data = real world!**

Tim Berners-Lee, WWW proposal, CERN, 1998:

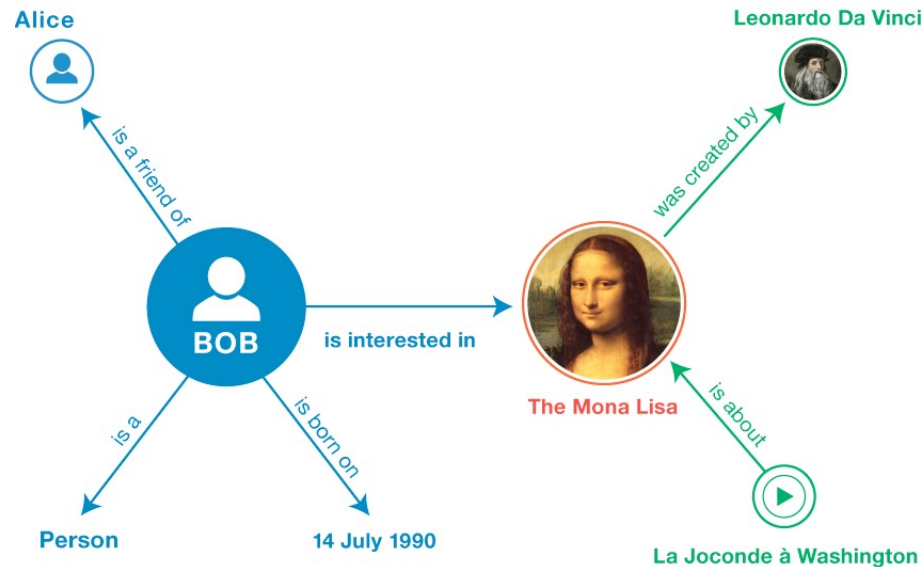
*"Many systems are organised hierarchically. A tree has the practical advantage of giving every node a unique name. However, **it does not allow the system to model the real world.**"*

*(On newsgroups): "Typically, a discussion under one newsgroup will develop into a different topic, at which point **it ought to be in a different part of the tree.**"*



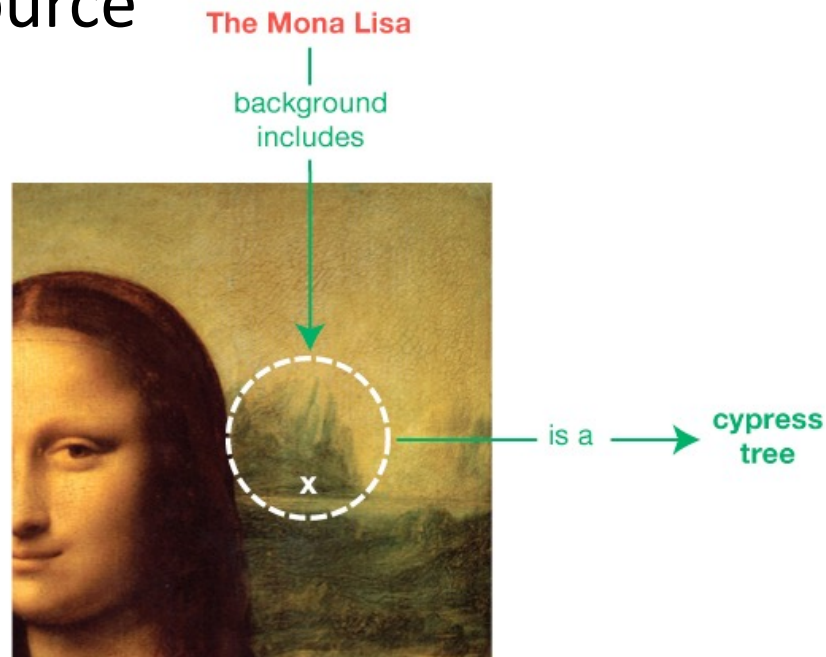
Graph data format for the Web: RDF

- Resource Description Format, W3C, 2003
- **Resources** have **properties** with **values**.
- **URIs** (Universal Resource Identifiers) identify resources
- Resources, properties, or values may be specified by an URI.
- Values may be constants



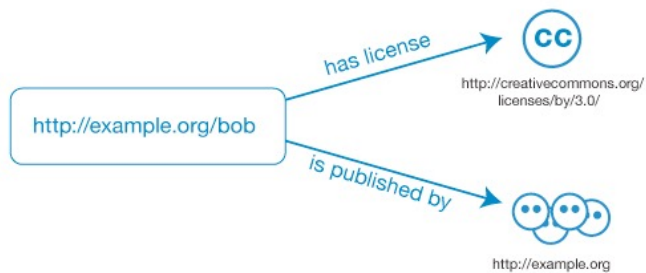
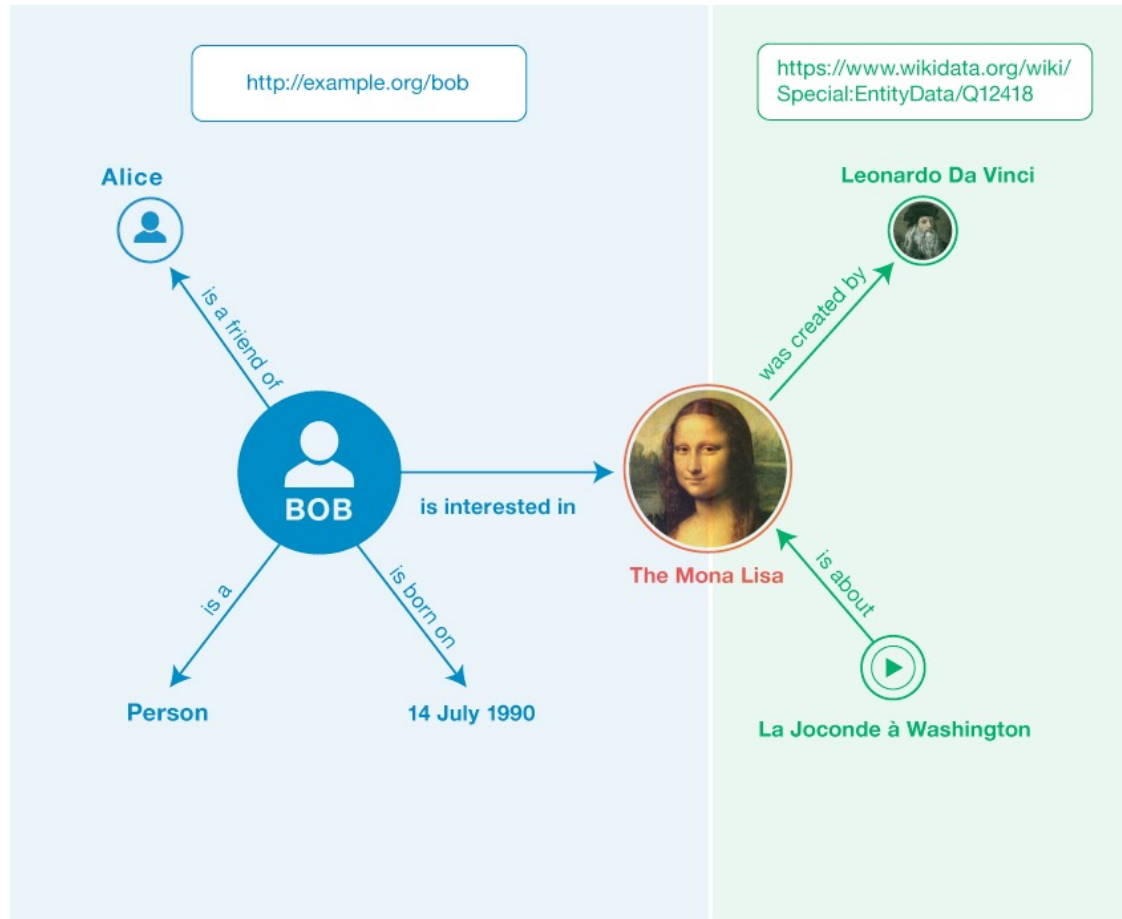
RDF feature: blank nodes

- Unnamed resource



- « Labeled null »

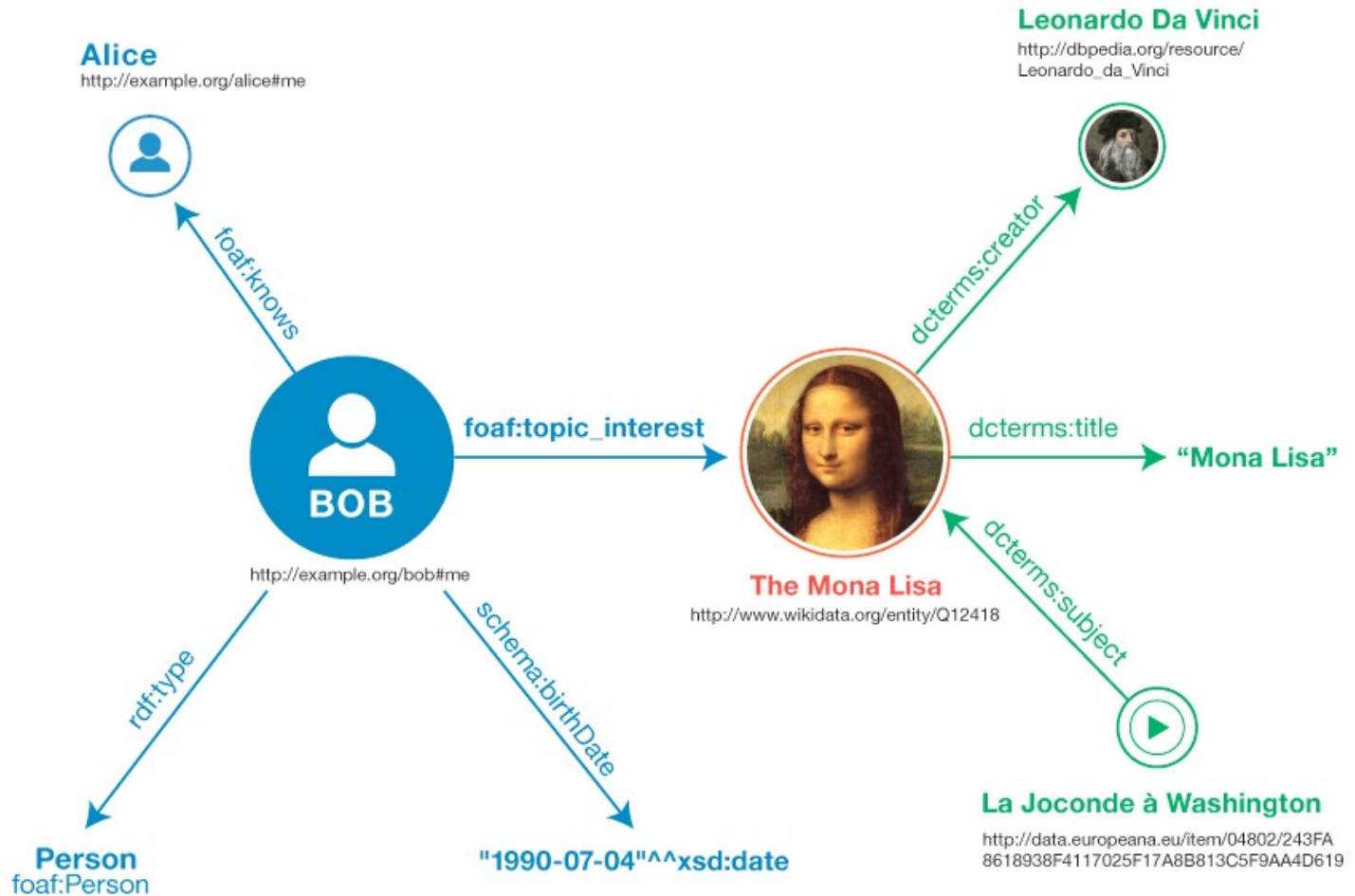
RDF graphs



RDF Schema constructs

Construct	Syntactic form	Description
Class (a class)	C rdf:type rdfs:Class	C (a resource) is an RDF class
Property (a class)	P rdf:type rdf:Property	P (a resource) is an RDF property
type (a property)	I rdf:type C	I (a resource) is an instance of C (a class)
subClassOf (a property)	C1 rdfs:subClassOf C2	C1 (a class) is a subclass of C2 (a class)
subPropertyOf (a property)	P1 rdfs:subPropertyOf P2	P1 (a property) is a sub-property of P2 (a property)
domain (a property)	P rdfs:domain C	domain of P (a property) is C (a class)
range (a property)	P rdfs:range C	range of P (a property) is C (a class)

Typed RDF graph

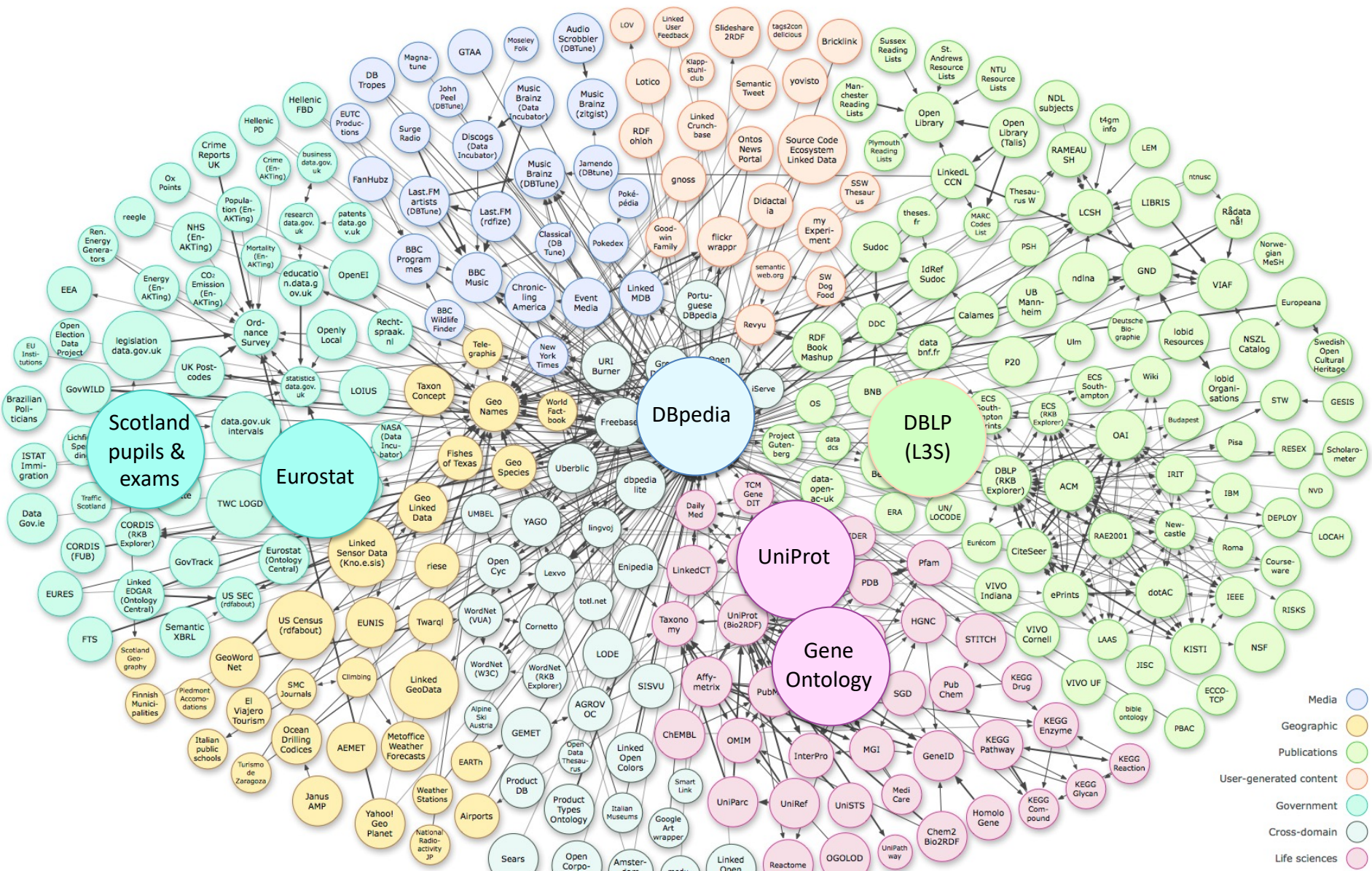


RDF reasoning

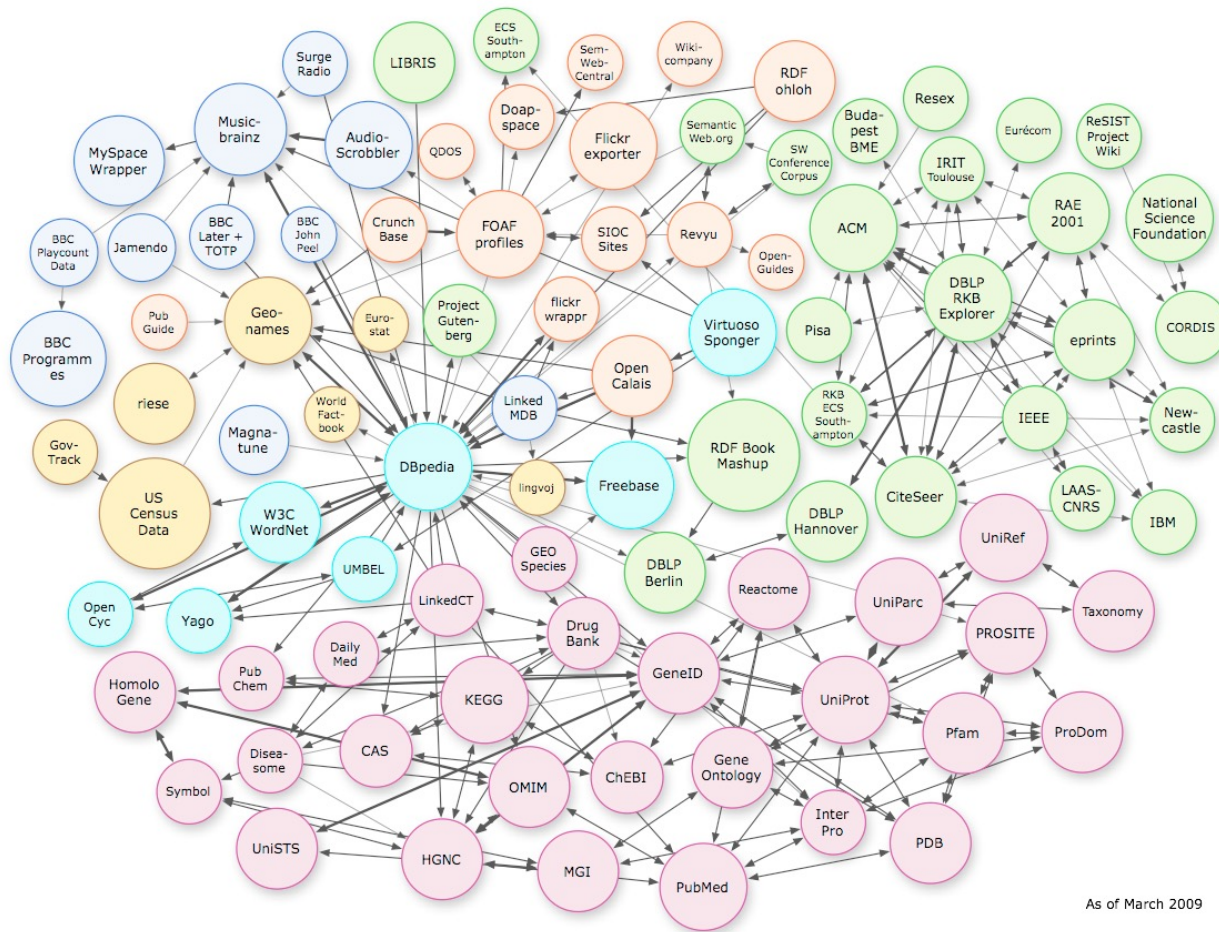
- RDF allow expressing **data *and* knowledge**
- Example:
 - if X *teaches a class*
 - then X *is a person*, *is an instructor*, *belongs to the school giving the class*, and *works for the university which includes the school*
- Reasoning exploits knowledge to infer *implicit data*

Linked Open Data (LOD) cloud

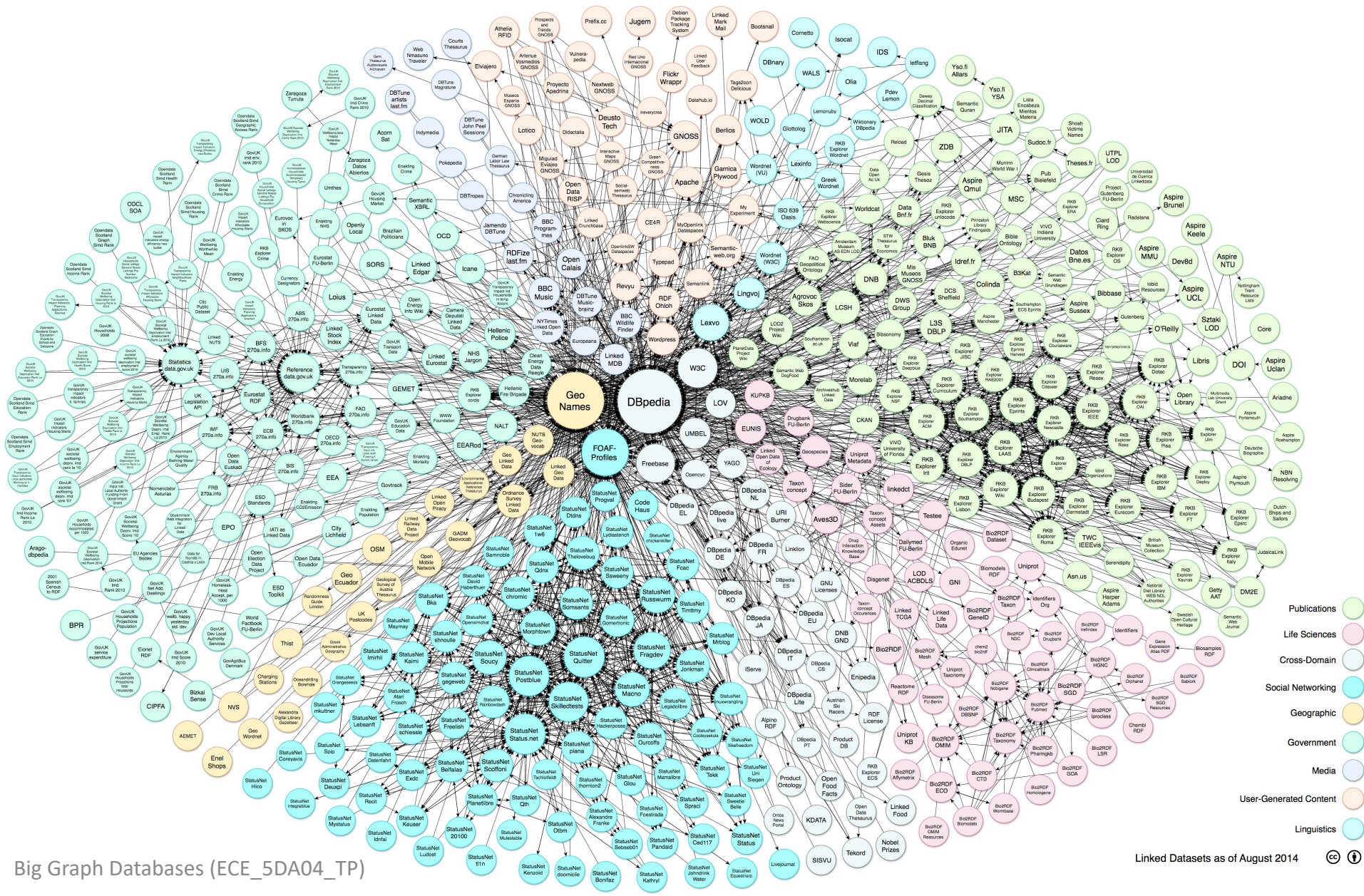
Linked Open Government Data project (logd.rw.rpi.edu): 10^{10} triples.



LOD cloud 2009

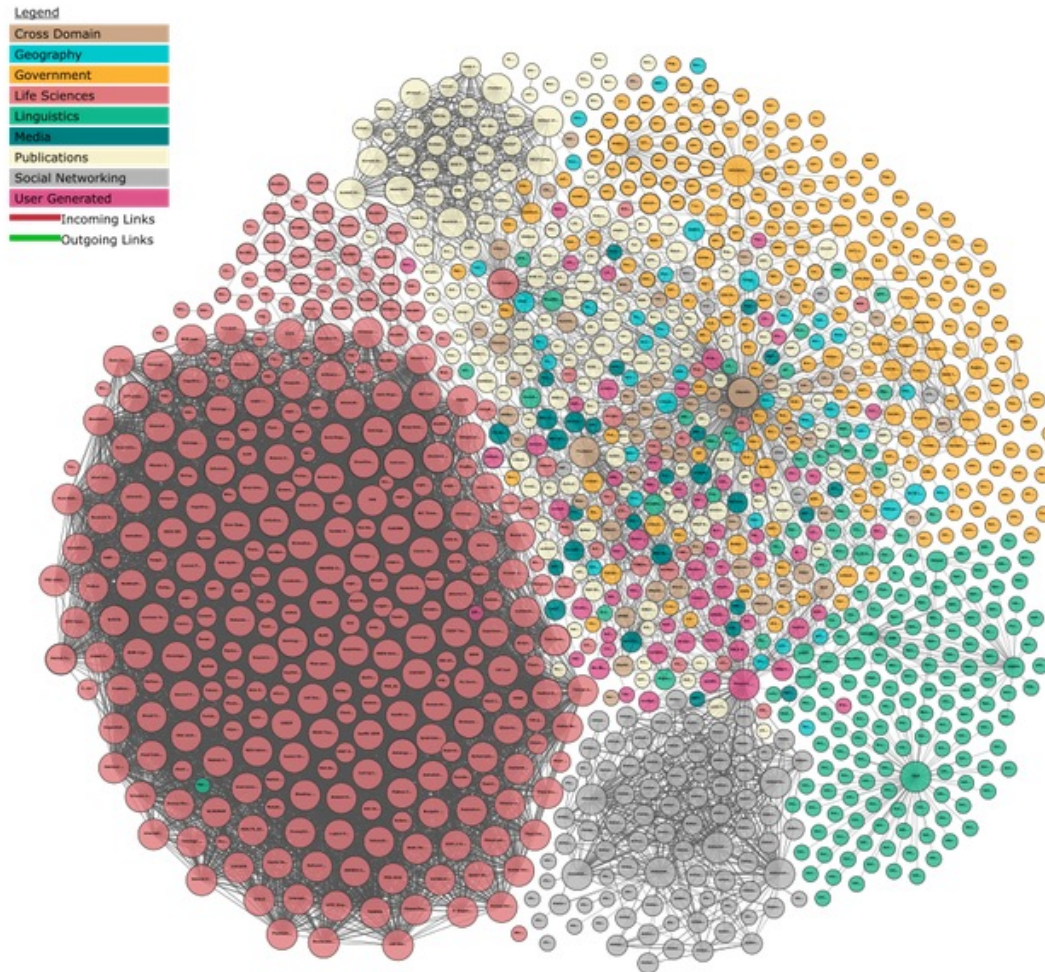


LOD cloud 2014



Publications
Life Sciences
Cross-Domain
Social Networking
Geographic
Government
Media
User-Generated Content
Linguistics

LOD cloud 2017



Linking Open Data cloud diagram 2017, by Andrejs Abele, John P. McCrae, Paul Buitelaar, Anja Jentzsch and Richard Cyganiak.
<http://lod-cloud.net/>

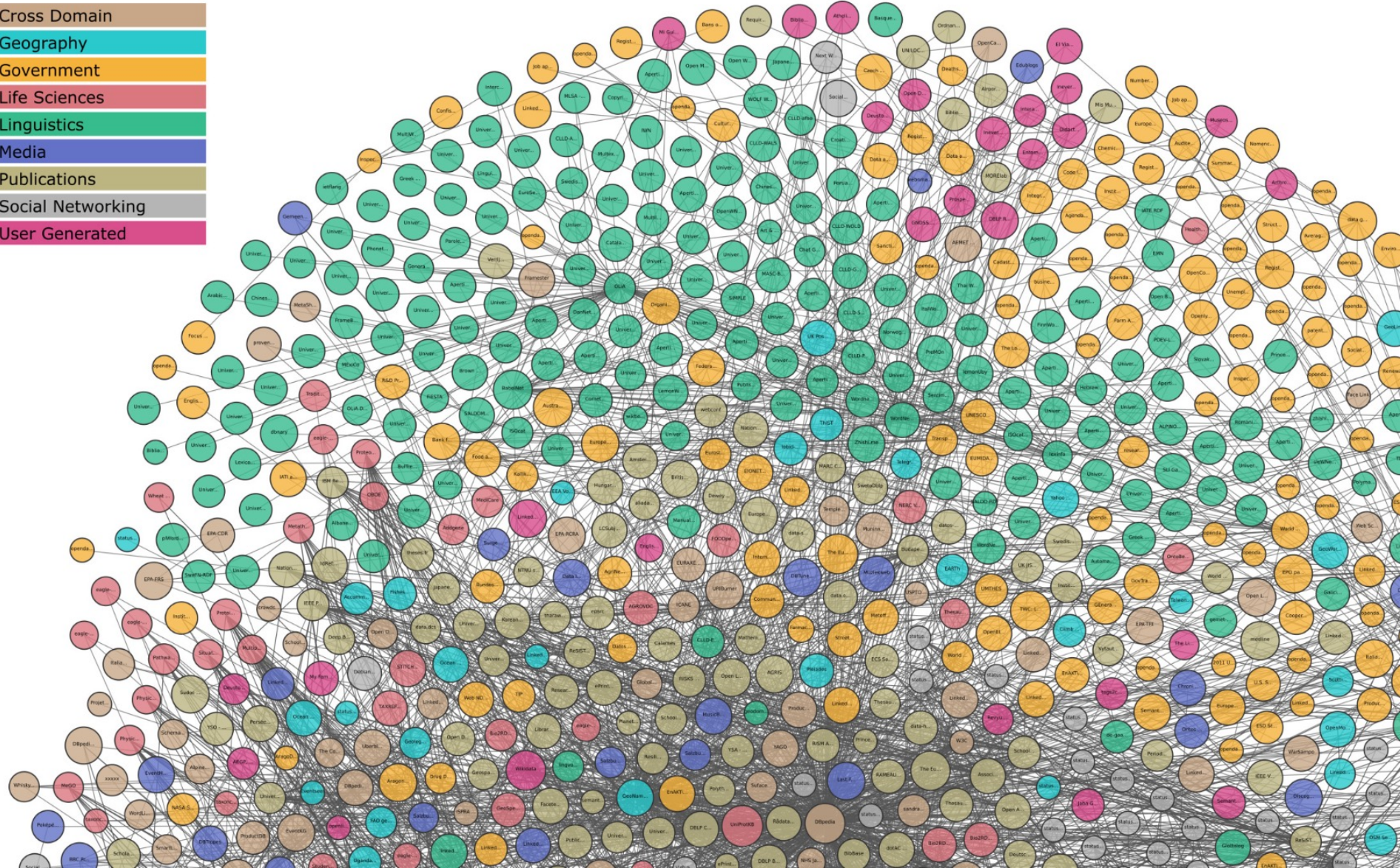
LOD cloud 2021

<http://lod-cloud.net/>

→ ↻ 🏠 🔒 <https://lod-cloud.net/#about>

Legend

- Cross Domain
- Geography
- Government
- Life Sciences
- Linguistics
- Media
- Publications
- Social Networking
- User Generated



Open vs. linked data

1. Linked Data:

"recommended **best practice** for exposing, sharing, and connecting pieces of data, information, and knowledge on the Semantic Web using **URIs** and **RDF**"

- (Tim Berners-Lee) vision for the Web

2. Open Data:

"**idea** that certain data should be **freely available to everyone to use and republish as they wish**, without restrictions from copyright, patents or other mechanisms of control"

- In principle, orthogonal to the Linked aspect
- In practice, Linked is a technical means toward Open

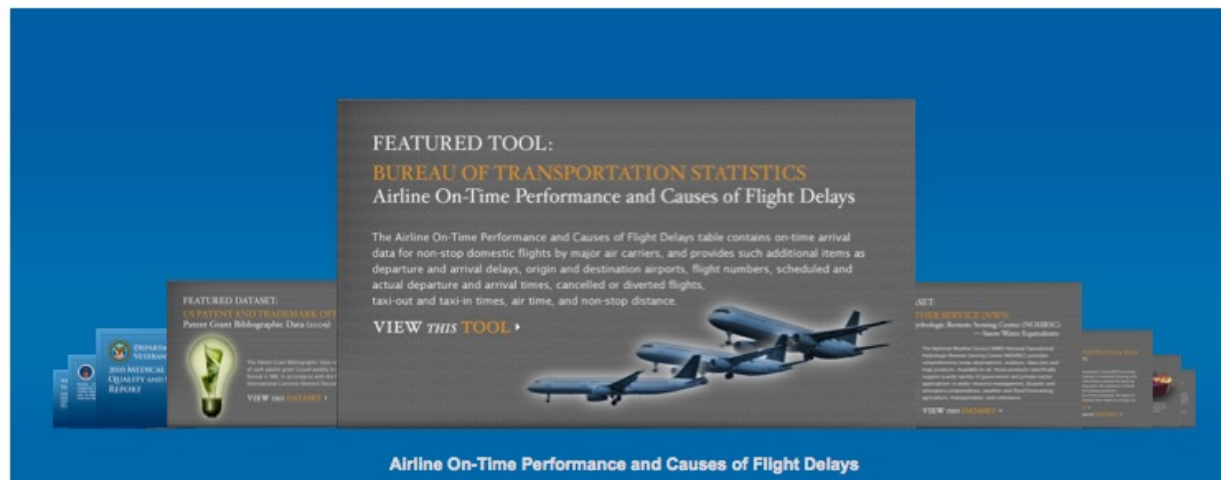
Open Data: data.gov (US)



Below is a gallery of datasets and tools that have been highlighted on the Data.gov home page. Click on "View More" to learn more about the data and link to the data itself.

This gallery displays just a tiny fraction of the datasets available to you on Data.gov. As we continue to add datasets, tools and highlights, we encourage you to explore all the valuable resources in our raw data, tools, and geodata catalogs.

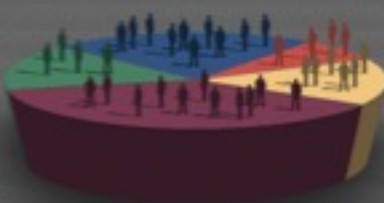
* Displaying 62 datasets and tools.



Open Data: data.gov (US)

FEATURED TOOL: US CENSUS BUREAU
DataFerrett

The DataFerrett is an online analytically oriented, self-service tool designed to deliver a wide variety of population, health, economic, geographic and housing information about the United States. It searches American Community Survey Public Use Microdata, Current Population Survey(CPS), CPS supplemental surveys, Survey of Income and Program Participation (SIPP), SIPP Topical Module surveys, Survey of Program Dynamics, the American Housing Survey, National Survey of Fishing, Hunting, and Wildlife Associated Recreation, The New York City Housing and Vacancy Survey, Local Employment Dynamics.



[VIEW THIS TOOL ▶](#)

Search our catalogs.. **SEARCH ▶**

ITEMS **GALLERY** **WHAT'S NEW** UPDATED

DATASETS AND TOOLS


highlighted on the Data.gov home page. Click on "View More" to

able to you on Data.gov. As we continue to add datasets, tools and resources in our raw data, tools, and geodata catalogs.

FEATURED TOOL:
BUREAU OF TRANSPORTATION STATISTICS
Airline On-Time Performance and Causes of Flight Delays

The Airline On-Time Performance and Causes of Flight Delays table contains on-time arrival data for non-stop domestic flights by major air carriers, and provides such additional items as departure and arrival delays, origin and destination airports, flight numbers, scheduled and actual departure and arrival times, cancelled or diverted flights, taxi-out and taxi-in times, air time, and non-stop distance.

[VIEW THIS TOOL ▶](#)



FEATURED DATASET:
EV PRESENT AND FUTURE MARK OFF
Percent Count Bibliographic Data (Novel)

[VIEW THE DATASET ▶](#)

WHAT'S NEW:
CHINA SERVICE INCOME
Inflight, Revenue, Airline, Carrier, INCIDENK2
— State, Water, Expenditure

[VIEW THE DATASET ▶](#)

INTEGRATION BLOG

Airline On-Time Performance and Causes of Flight Delays

Open Data: data.gov (US)

The screenshot displays the data.gov website interface. At the top, there is a search bar with the text "Search our catalogs.." and a "SEARCH" button. Below the search bar, there are navigation tabs for "DATASETS", "GALLERY", and "WHAT'S NEW" (with a "UPDATED" indicator). A large blue banner reads "DATASETS AND TOOLS".

On the left side, there is a section titled "FEATURED TOOL: US CENSUS BUREAU DataFerrett". The text describes it as an online analytical tool for population, health, economic, geographic, and housing information. Below this is a 3D pie chart with human figures on top.

In the center, there is a "FEATURED DATASET: ENERGY INFORMATION ADMINISTRATION (EIA) Residential Energy Consumption Survey (RECS)". It includes a lightbulb icon and a house icon. The text describes the survey's scope and provides links for "COMPLETE DATASET" and "CONSUMPTION PORTION OF DATASET".

At the bottom, there is a section for "Airline On-Time Performance and Causes of Flight Delays". It features an image of several airplanes and a "VIEW THIS TOOL" link. To the left of this section, there are smaller thumbnails for other datasets, including "CURRENT AND PAST MARK OFF Percent Count Bibliographic Data (level)".

Text on the right side of the image reads: "v home page. Click on 'View More' to v. As we continue to add datasets, tools data, tools, and geodata catalogs."

Open Data: data.gov (US)

The screenshot shows the data.gov website interface. At the top right, there is a search bar with the text "Search our catalogs.." and a "SEARCH" button. Below the search bar are navigation tabs for "DATASETS", "GALLERY", and "WHAT'S NEW" (with a "UPDATED" badge). A large blue banner reads "DATASETS AND TOOLS".

Overlaid on the left side of the screenshot are three feature callouts:

- FEATURED TOOL: US CENSUS BUREAU DataFerrett**
The DataFerrett is an online analytically oriented, self-service tool designed to deliver a wide variety of population, health, economic, geographic and housing information about the United States. It includes data from the American Community Survey, the Census of Population and Housing, the Survey of Income and Program Dynamics, the National Health and Medical Expenditure Survey, the National Longitudinal Survey of Youth, and the National Longitudinal Survey of the Youth.
- FEATURED DATASET: ENERGY INFORMATION ADMINISTRATION (EIA) Residential Energy Consumption Survey (RECS)**
The RECS provides information on the energy consumption in the residential sector. It includes data on energy use, fuels, and energy efficiency.
- FEATURED TOOL: RECREATION INFORMATION DATABASE (RIDB)**
The Recreation Information Database (RIDB) is a warehouse of information about Federal recreation sites. This web service has the ability to export the data to state tourism portals, recreation-related businesses, etc. It is also the "back end" supplying data to the Recreation.gov portal for trip planning information regarding more than 3,000 Federal recreation sites.

At the bottom of the screenshot, there is a blue banner with the text "Airline On-Time Performance and Causes of Flight Delays".

view home page. Click on "View More" to view. As we continue to add datasets, tools, data, tools, and geodata catalogs.

Open Data: data.gov (US)

The screenshot displays the data.gov website interface. At the top, there is a search bar with the text "Search our catalogs.." and a "SEARCH" button. Below the search bar, there are navigation tabs for "ITEMS", "GALLERY", and "WHAT'S NEW" (with a "UPDATED" indicator). A large blue banner reads "DATASETS AND TOOLS".

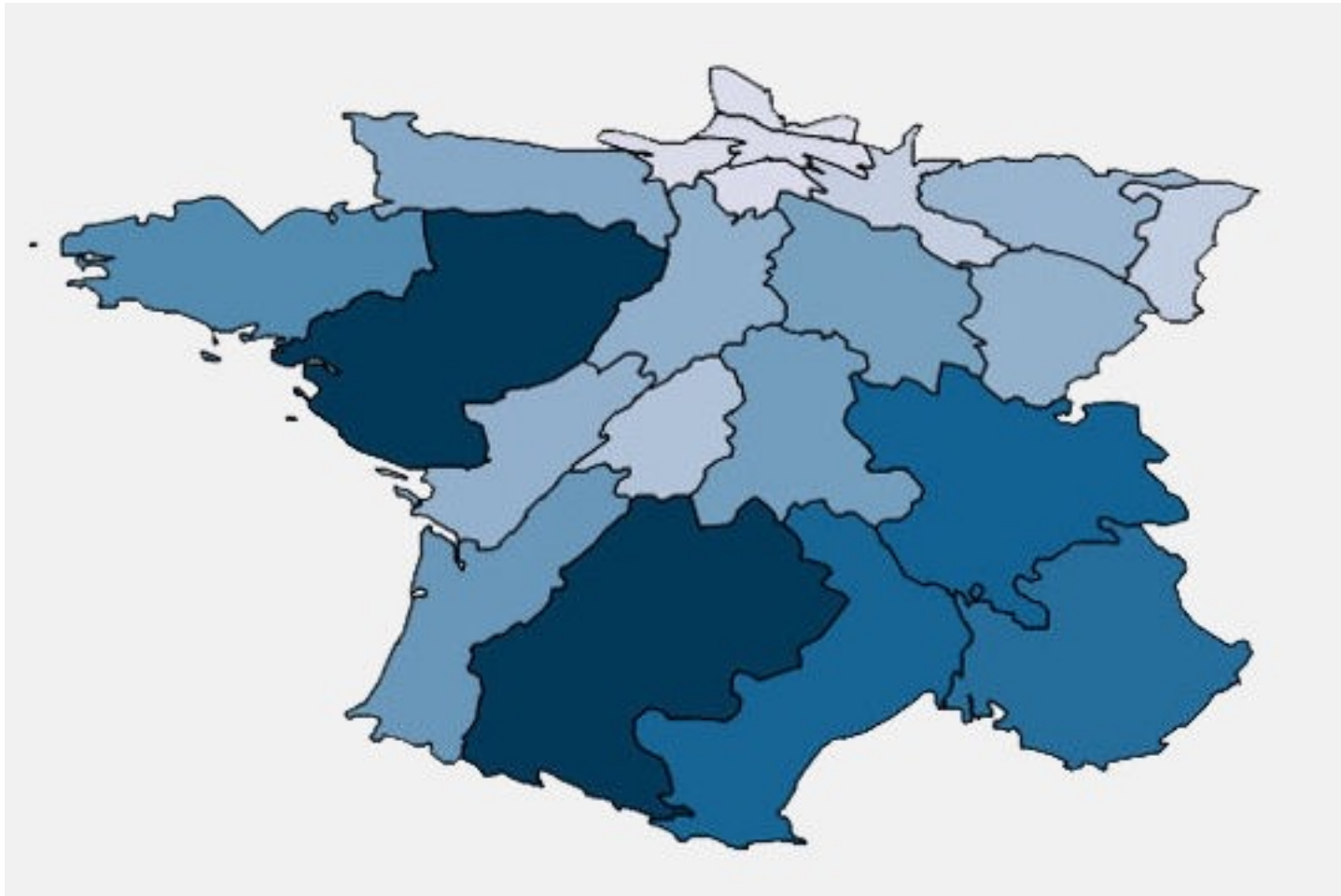
Several featured items are highlighted in dark grey boxes with orange and white text:

- FEATURED TOOL: US CENSUS BUREAU DataFerrett**
The DataFerrett is an online analytically oriented, self-service tool designed to deliver a wide variety of population, health, economic, geographic and housing information about the United States. Community Survey, Population Survey, Survey of Income and Program Dynamics, Wildlife Associates, Employment Dynamics.
- FEATURED DATASET: ENERGY INFORMATION ADMINISTRATION (EIA) Residential Energy Consumption Survey (RECS)**
The RECS provides information on energy consumption in the residential sector, including energy efficiency measures, fuels, and appliances.
- FEATURED TOOL: RECREATION INFORMATION DATABASE (RIDB)**
The Recreation Information Database (RIDB) provides information on recreational activities, including hiking, fishing, and hunting. It also includes information on recreational facilities, such as trails, camps, and picnic areas.
- FEATURED DATASET: NATIONAL WEATHER SERVICE (NWS) National Operational Hydrologic Remote Sensing Center (NOHRSC) — Snow Water Equivalents**
The National Weather Service (NWS) National Operational Hydrologic Remote Sensing Center (NOHRSC) provides comprehensive snow observations, analyses, data sets and map products. Available to all, these products specifically support a wide variety of government and private-sector applications in water resource management, disaster and emergency preparedness, weather and flood forecasting, agriculture, transportation, and commerce.

Other visible elements include a "VIEW THIS" link, a "FEATURED DATA" section with a lightbulb icon, and a "FEATURED TOOL" section with a hiker icon. A blue banner at the bottom of the featured items reads "Airline On-Time Performance and...".

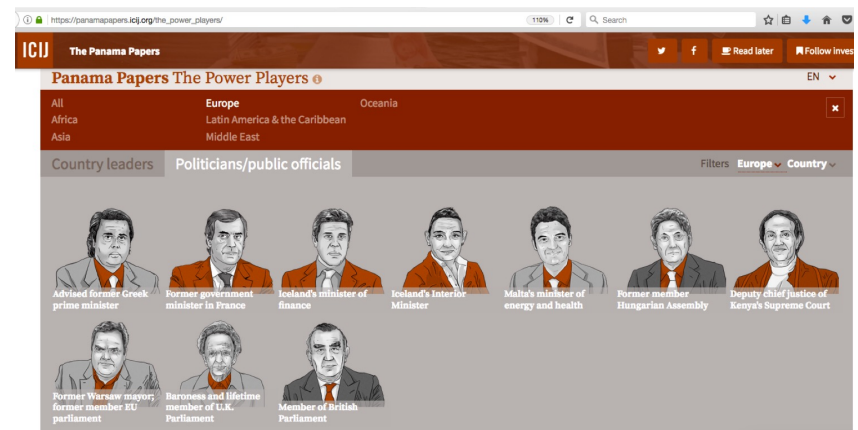
Open Data from Etalab (FR)

Organic agriculture per French region



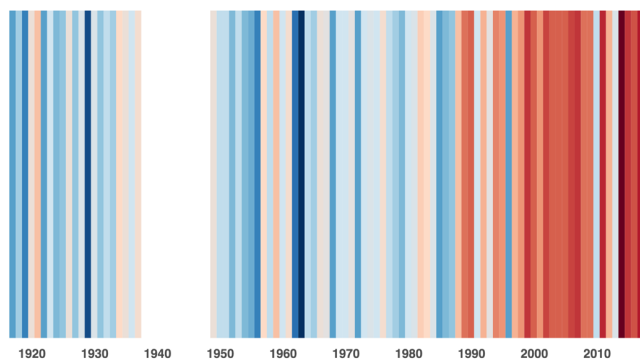
Open Data important for democracy

- Journalists and/or NGO workers play increasingly important role explaining society functioning
 - E.g., ICIJ Panama Papers investigation

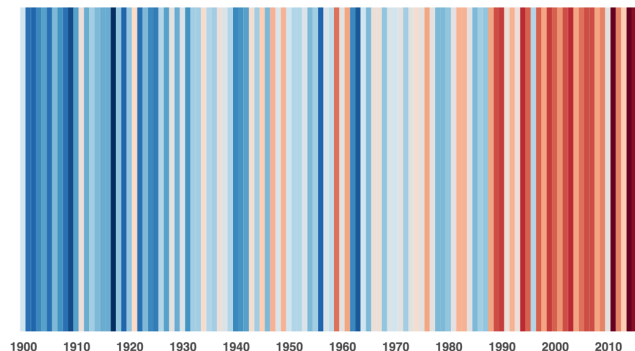


Open Data important for democracy

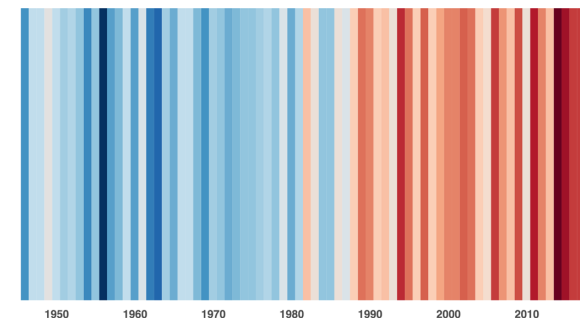
- Journalists and/or NGO workers play increasingly important role explaining society functioning
 - E.g., Météo France data on average yearly temperature, Les Décodeurs (Le Monde)
https://www.lemonde.fr/les-decodeurs/article/2021/01/06/visualisez-le-rechauffement-climatique-en-france-et-dans-votre-ville-avec-nos-barres-de-rechauffement_6065388_4355770.html



Evolution de la température par année depuis 1917
Dunkerque (59) : de 8,9 °C à 12,8 °C



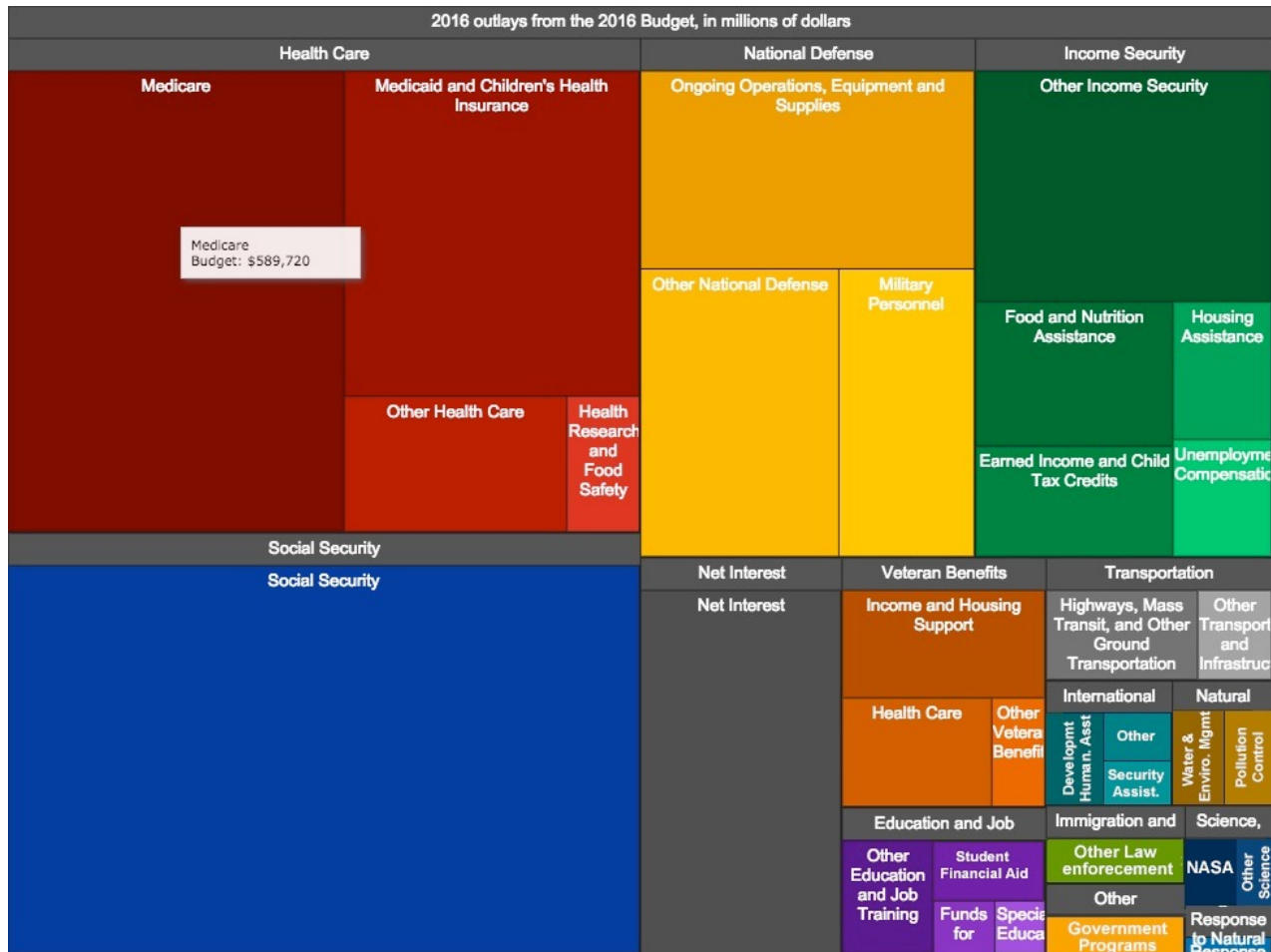
Evolution de la température par année depuis 1900
Paris (75) : de 10 °C à 13,7 °C



Evolution de la température par année depuis 1946
Montpellier (34) : de 12,8 °C à 16,5 °C



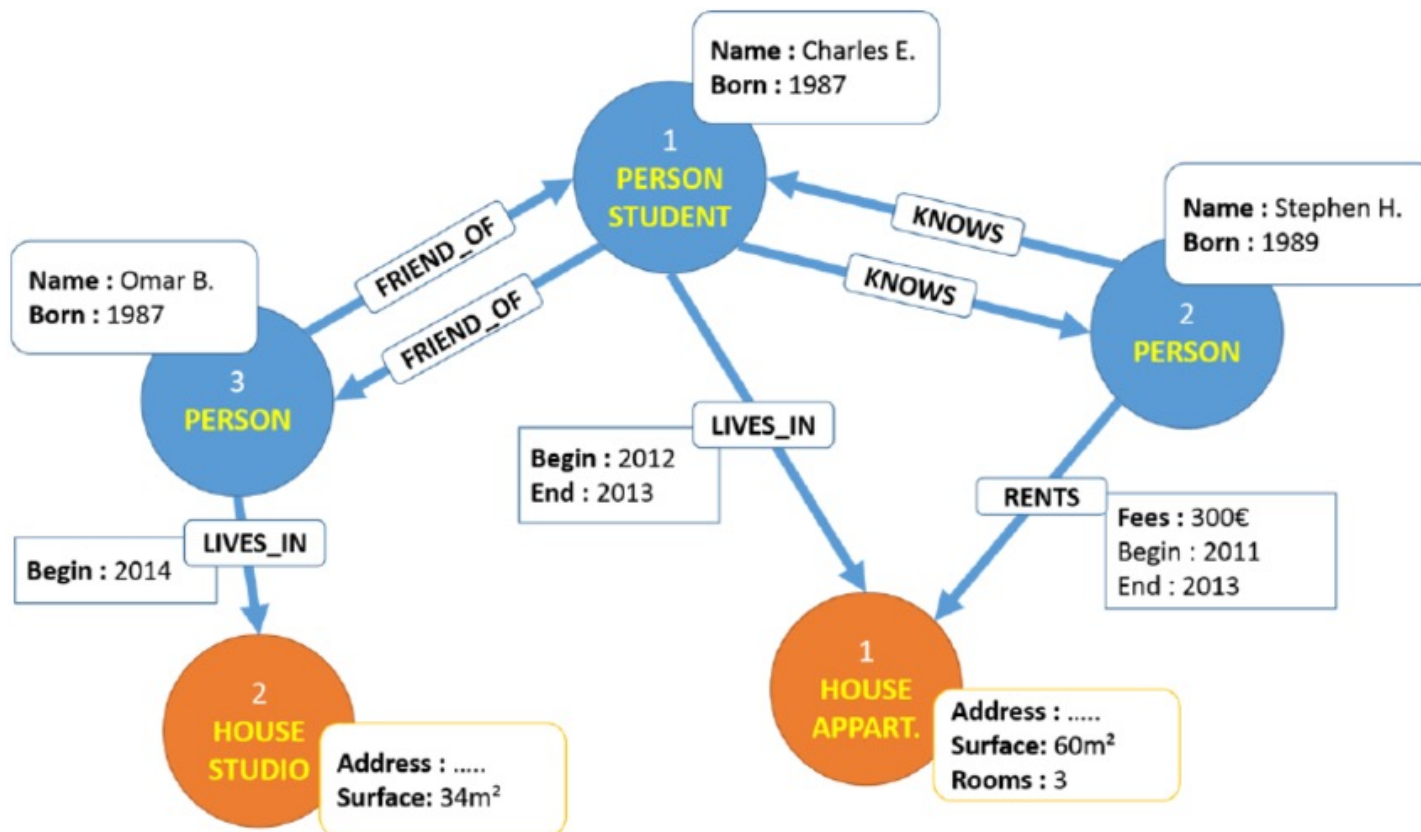
US Budget allocation, 2016



https://cdns.tblsft.com/sites/default/files/pages/5_2016_budget_visualization_safe_img.jpg

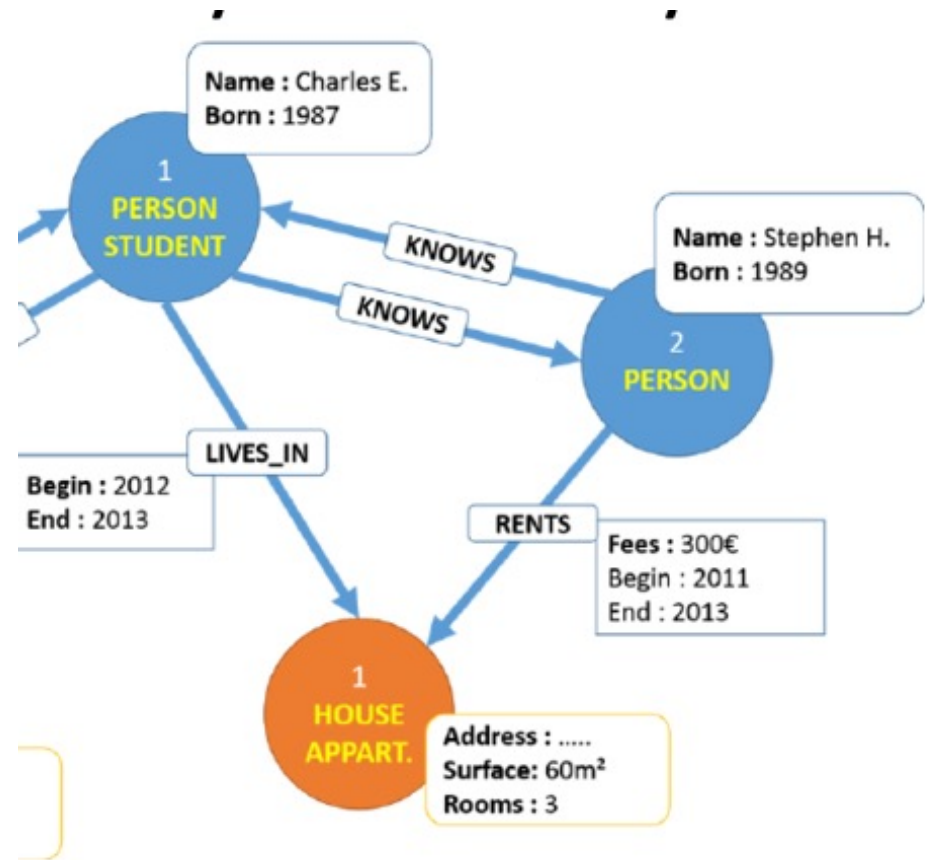
Property graphs: the other “standard” for graph data

- Initially introduced by the Neo4J system



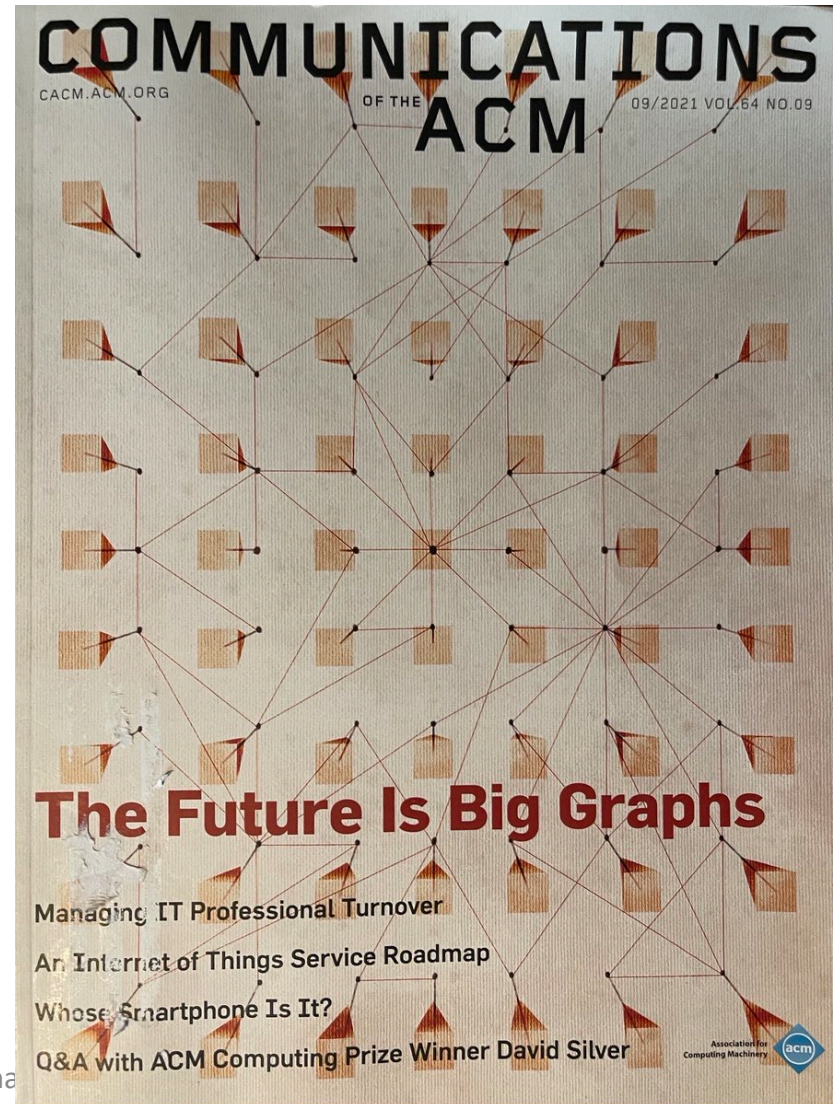
Property graphs: the other “standard” for graph data

- A property graph can be loaded from or dumped into a set of CSV files.
 - One CSV for each node type, e.g., one for PERSON, one for STUDENT, one for HOUSE, etc.
 - One CSV for each edge label, e.g., FRIENDS_OF, LIVES_IN, KNOWS etc.
- A graph can be queried using the Cypher or GQL languages (SQL for graphs)



Property graphs: the other “standard” for graph data

- Embraced by all major database vendors (Oracle, DB2, Microsoft)
- Linked Data Benchmark Council (LDBC)
- Many exciting developments



Object-oriented databases

- Studied in 1980's-2000
- Goal: adopt benefits of object-oriented programming (rich structure, encapsulation, inheritance, strong typing, etc.) to databases
- C++ , then Java classes with persistence

```
class myClass1 { ... }  
    persistent class myClass2 { ... }  
    public static void main { ... commit() }
```

- Native OO databases mostly disappeared.
- Major relational databases have *OO extensions*
 - Declare a table as containing objects of a certain type
 - Some method support
 - Little-known, not central to the SQL standard.

Key-value databases

- Simplest possible data model
 - V1: 1 key \rightarrow 1 value
 - e.g., “1” \rightarrow “Julie, Paris 12e”
 - V2: 1 key \rightarrow 1 record \rightarrow 1 property \rightarrow 1 value
 - e.g., “1” \rightarrow • \rightarrow “name” \rightarrow “Julie”
 \rightarrow “address” \rightarrow “Paris 12e”
- Records are grouped in named collections
- In a collection, records may have different property sets \rightarrow **Heterogeneity!**
- A property may be defined with unique values or multiple values, with or without duplicates
- **No query language**, only put and get.

Big data heterogeneity (variety)



Data model & data management system soup

- hierarchical, relational, object-oriented, XML, RDF, JSON, key-value pairs...

Traditionally this has been solved (time and \$ permitting) with **data migration** / **ETL** (extract-transform-load)

- Heterogeneous data and high throughput may make ETL **impractical for Big Data** →
Need to **exploit big, heterogeneous data as is**

Varied Big Data has huge value potential

- Real estate ad from Zillow (US):

www.zillow.com/homedetails/217-Newark-Ave-APT-515-Jersey-City-NJ-07302/106623623_zpid/

CORRECT HOME FACTS SAVE GET UPDATES SHARE

Public Owner

217 Newark Ave APT 515
Jersey City, NJ 07302

217 Newark Ave APT 515
Jersey City, NJ 07302

-- beds -- baths - 871 sqft Edit

Edit home facts for a more accurate Zestimate.

Thinking About Selling?
Find a local agent who can give you a professional estimate of your home value.

[Find an Agent](#)

2 Bedroom 2 Bath w/ Garage Parking at The Saffron. Sleek modern design and details that include stainless steel appliances, Caesarsstone counters, Bamboo floors and soaring ceiling heights. Master bedroom with walk in closet and master bathroom. Great storage space, W/D and dishwasher. Manhattan views from spacious roof deck lounge and sun deck or relax in common courtyard. On site fitness center. Low maint. & taxes. 2.5 blocks from the Grand Central PATH.

Home price
\$ 505,000

Down payment
\$ 101,000 20 %

Loan program
30-year fixed

Interest rate See current rates
3.304 %

Include taxes/ins.

Your payment
\$1,770

P&I
\$1,770

[Get pre-qualified](#)

Home Expenses

INTERNET, TV & PHONE



Bundle Services and Save

Bundle your monthly services and consider ways to "cut the cord" with low-cost online alternatives to cable such as streaming and on-demand video.

\$50-\$100 / month

[Learn More](#)

SECURITY



Make Your Home a Fortress

- 24/7 protection for just \$14.99/mo
- No long-term contracts, so you're free to cancel anytime

\$14.99 / month

[Try It Today](#)

Powered by **SimpliSafe**

Nearby Schools in Jersey City

SCHOOL RATING

3 out of 10

Number 4 Middle (assigned)

GRADES DISTANCE

6-8 & ungraded 0.3 mi

10 out of 10

Dr. Ronald Mc Nair Academy High (assigned)

9-12 & ungraded 0.3 mi

5 out of 10

Number 5 Elementary

PK-8 & ungraded 0.4 mi

[More schools in Jersey City](#)

Data by [GreatSchools.org](#)

Price / Tax History

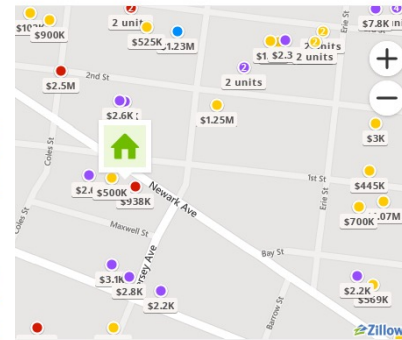
Price History Tax History

DATE	EVENT	PRICE	AGENTS
05/19/14	Sold	\$505,000 +1.2%	Stacey Bonavitacola
03/09/14	Listed for sale	\$499,000 +23.5%	
		\$404,000	

Downtown

Home values will increase 6.3% next year, compared to a 4% decline. Among Downtown homes, this home is valued 8.4% less than the average, and is valued 5.2% less per square foot.

paradise)



4 photos

OFF MARKET

750 sqft

2 bds • 2 ba • -- sqft

1 day on Zillow • 217 Newark Ave APT 513, J...

cher, centre...

Comparez et économisez jusqu'à 55% sur votre nuit d'hôtel.

trivago.fr

Nearby Similar Sales

SOLD: \$551,700
Sold on 4/4/2016
2 beds, 1.0 baths, 800 sqft
158 Wayne St APT 401A, Jersey City, NJ 07302

SOLD: \$589,000
Sold on 10/15/2015
2 beds, 1.0 baths, 987 sqft
280 Monmouth St APT B, Jersey City, NJ 07302

SOLD: \$594,000
Sold on 6/1/2016
2 beds, 1.0 baths, 1013 sqft
227 Christopher Columbus Dr APT 217B, Jer...

SOLD: \$599,000
Sold on 5/27/2016
2 beds, 1.0 baths, 978 sqft
341 Monmouth St APT 311D, Jersey City, NJ ...

SOLD: \$600,000
Sold on 10/5/2015
2 beds, 2.0 baths, 1016 sqft
287 8th St APT 4B, Jersey City, NJ 07302

[See sales similar to 217 Newark Ave APT 515](#)

ROCKY FLATS
Colorado

Did you own property near the Rocky Flats Nuclear Weapons Plant on June 7, 1989? Are you an heir of someone who did? Are you the successor of an entity that did? You could get money from a \$375 million Settlement.

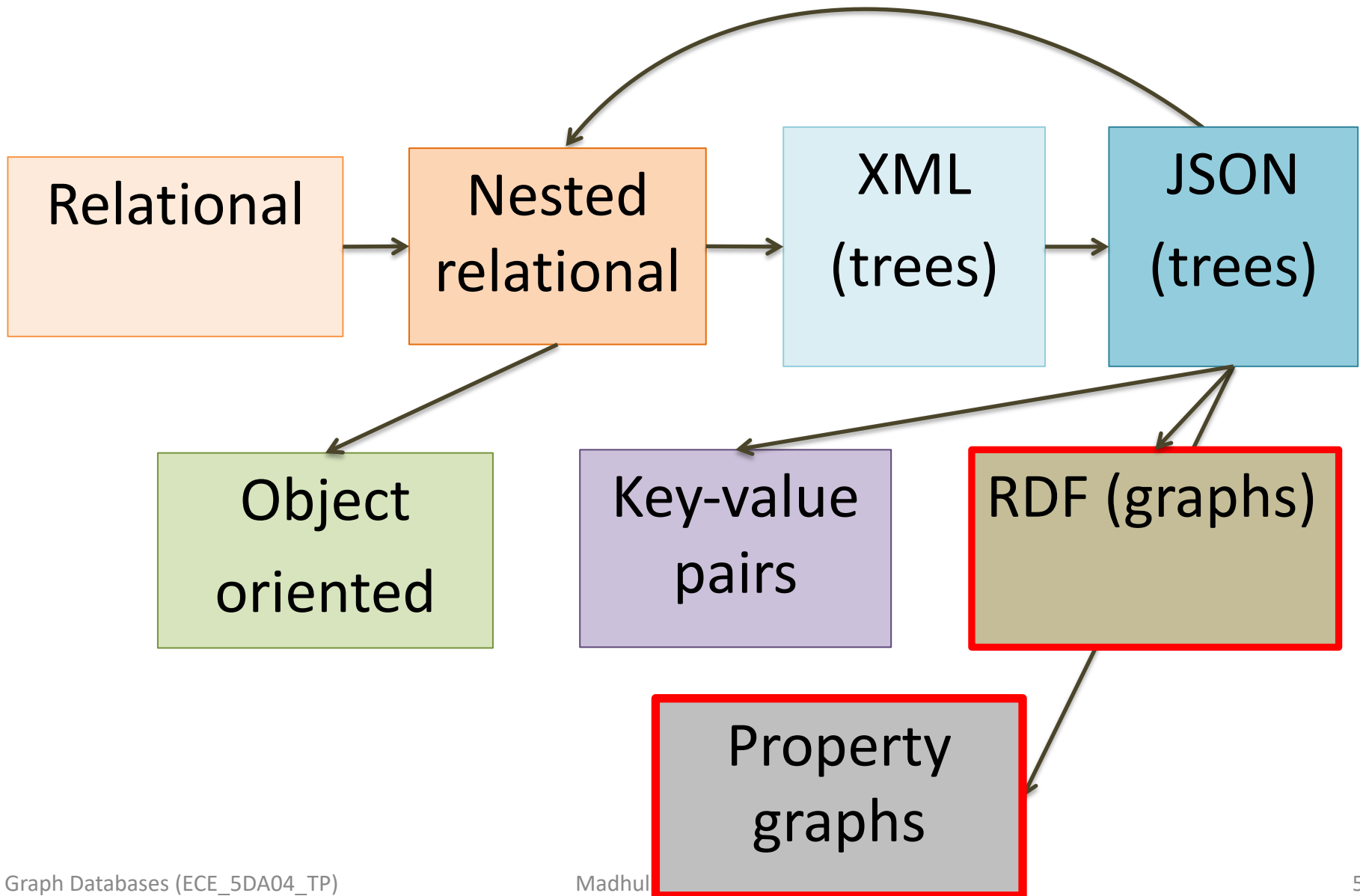
ROCKY FLATS
Colorado

Improve Your Home Value

PROJECT	ADDED VALUE	PROJECT COST
Entry Door Replacement	+\$2,337	\$3,156
Vinyl Window Replacement	+\$10,353	\$15,030
Bathroom Remodel	+\$7,297	\$18,841
Major Kitchen Remodel	+\$44,675	\$61,399

[See More Home Improvement Inspiration](#)

A brief history of data models



Defining Big Data: the V's

- Volume
 - Scale
- Velocity
 - Speed of producing and consuming the data
- Variety
 - Very different sources and data types
- Veracity
 - Is the data correct / certain / true?

Big Data veracity

- Is this **true**? (What is the **probability**?)
- Contradictory sources (1 vs. 2 clocks)
- Errors in the data
 - Humans introduce many errors
 - Sensors may have failures or erroneous readings
 - Light or heating sensors in a building
 - Wear and tear
- Tackled by **data curation / cleaning / quality** tools for regular (or at least homogeneous) data
- Recent ML methods for this

Big Data veracity

- Data reconciliation / entity extraction
- Large-scale **Knowledge Bases** such as YAGO, DBPedia, Google Knowledge Base (> Freebase)
 - Reference database of core facts

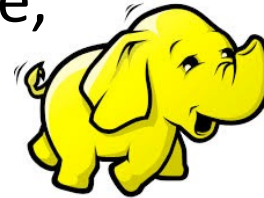
- Places (city/country etc.), people (public figures, scientists, artists etc.), events (born, died, emigrated, was created...), time

- Ontologies automatically extracted from Web and other specific sources → they may have *errors* and are *incomplete*
 - Manual curation and improvement, e.g., YAGO 4.5

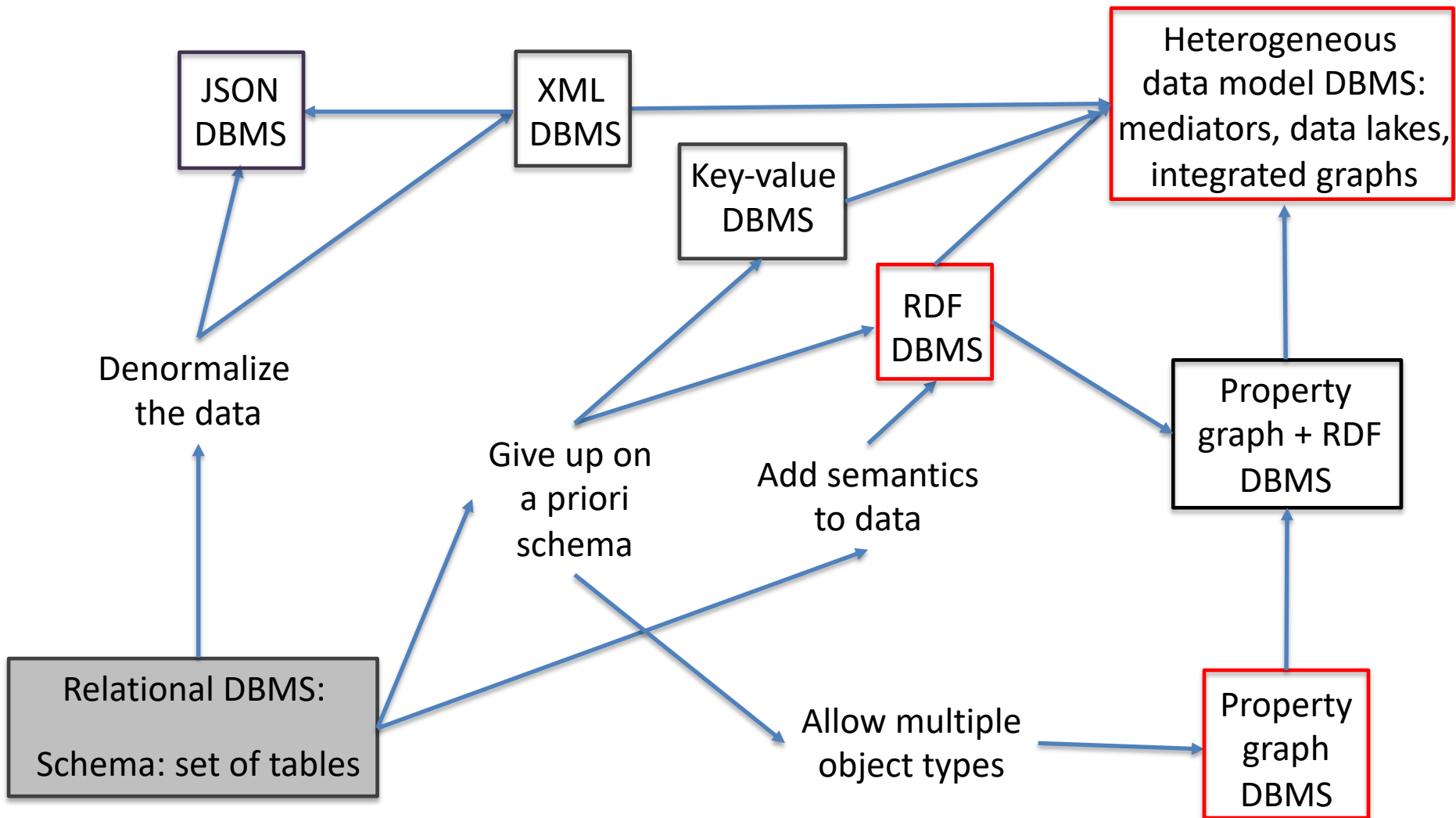


Big picture on Big Data

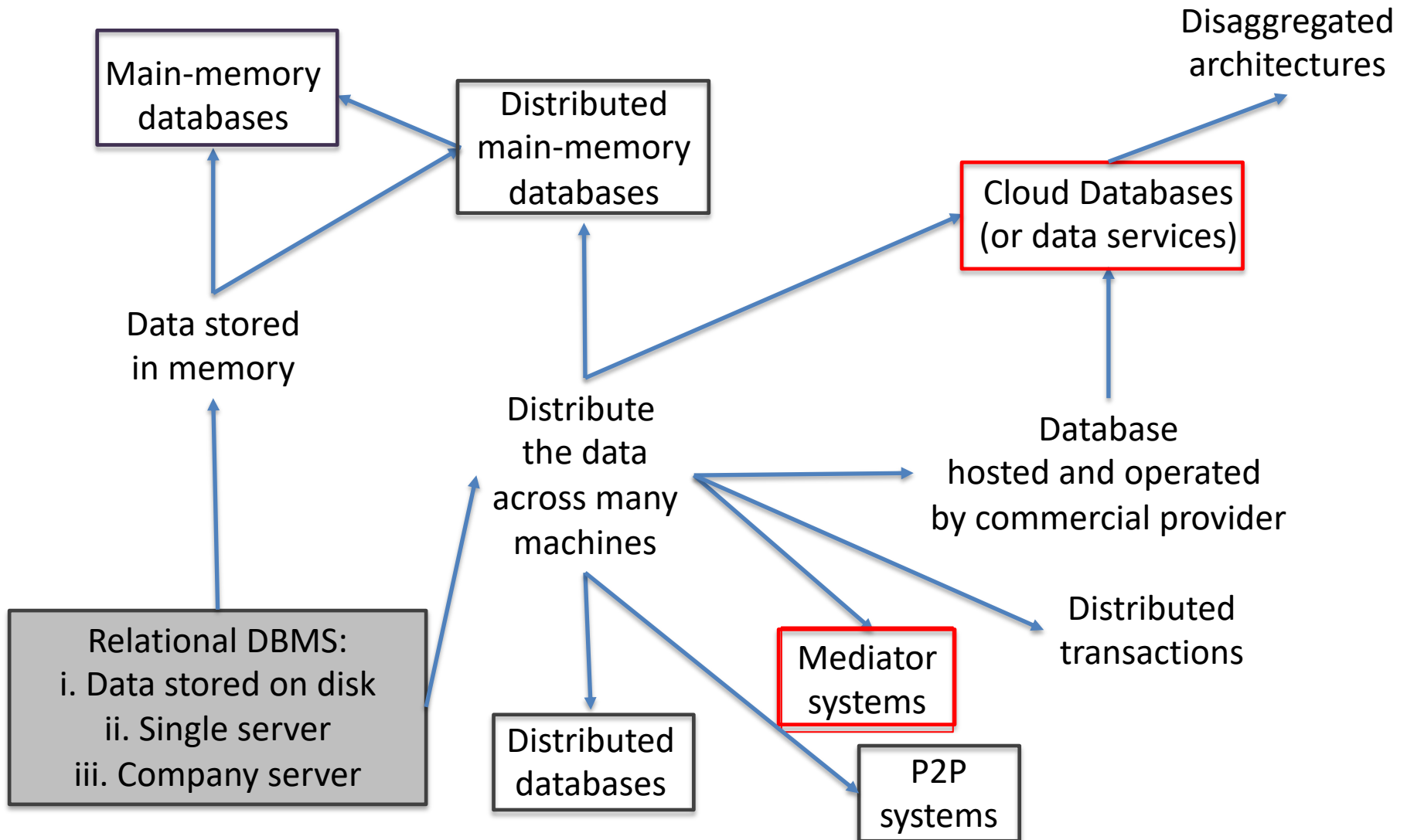
1. **Volume, velocity, variety, veracity**
2. Probably not all the data has the same **value** \$/B
 - This is why enterprise databases are preserved
 - Very large, unstructured, uncertain-value data may instead be stored in larger-scale, lower-performance systems, mostly based on Hadoop or Spark
 - Massively parallel processing
 - Iterative, machine learning
3. Deployment model changes (→ **cloud**)



From databases to Big Data



From databases to Big Data



Questions?