

Muriel Foulonneau

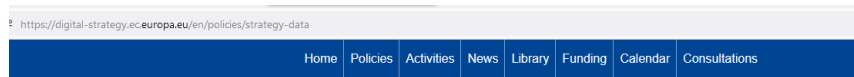
▶ Data quality as a risk

Stories of data reuse



THE EUROPEAN DATA CONFERENCE ON REFERENCE DATA AND SEMANTICS

The promise of the data economy



[Home](#) > [Policies](#) > A European Strategy for data

A European Strategy for data

The strategy for data focuses on putting people first in developing technology, and defending and promoting European values and rights in the digital world.

Data is an essential resource for economic growth, competitiveness, innovation, job creation and societal progress in general.

The [European strategy for data](#) aims at creating a single market for data that will ensure Europe's global competitiveness and data sovereignty. Common European data spaces will ensure that more data becomes available for use in the economy and society, while keeping the companies and individuals who generate the data in control.

Data driven applications will benefit citizens and businesses in many ways. They can:

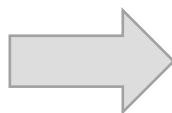
- improve health care
- create safer and cleaner transport systems
- generate new products and services
- reduce the costs of public services
- improve sustainability and energy efficiency

The Commission has proposed a [Regulation on European data governance](#) as part of its data strategy. This new Regulation will play a vital role in ensuring the EU's leadership in the global data economy.

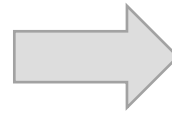
On 23 February 2022, the Commission proposed a Regulation on harmonised rules on fair access to and use of data (Data Act). [The Data Act](#) is a key pillar of the European strategy for data. Its main objective is to make Europe a leader in the data economy by harnessing the potential of the ever-increasing amount of industrial data, in order to benefit the European economy and society.

To further ensure the EU's leadership in the global data economy the European strategy for data intends to:

- adopt legislative measures on data governance, access and reuse. For example, for



**Data creates
value**



**Sharing data
creates more
value**

Reusing data to support innovation

Open data market size



- €184.45 billion open data market size in 2019
- €199.51 - €334.20 billion open data market size forecast for 2025

Open data employment

- 1.09 million open data employees in 2019
- 1.12 - 1.97 million open data employees forecast for 2025



Open data potential per sector



- 15.7% growth expected from high impact and high potential sectors

• High impact:



• High potential:



For details on calculations and assumptions see corresponding sections.



Efficiency gains

- Saving lives, e.g. 54 - 202 thousand lives saved by faster emergency response
- Saving time, e.g. 27 million hours saved in public transport
- Saving the environment, e.g. 5.8 Mtoe* saved by reducing household energy consumption
- Improving language services with open data, e.g. by increasing machine translation



Cost savings

- Saving healthcare costs, e.g. €312 - €400 thousand due to faster first aid by bystanders
- Saving labour costs, e.g. €13.7 - €20 billion by reducing time spent in traffic
- Saving costs on energy bills, e.g. €79.6 billion due to more solar energy production
- Saving public sector costs, e.g. €1.1 billion by lower translation costs



Open data in organisations

- 49% of data used by surveyed organisations is open data and 77% of organisations plan to use more data
- 46% of organisations' revenues are impacted by open data and 73% of organisations expect the impact to increase
- 70% of surveyed organisations create data internally, of which 58% publish some of it as open data



* Million tonnes of oil equivalent

For details on calculations and assumptions see corresponding sections.



What is reusable data?

- Data reusers raise data quality as a major obstacle
 - Ex. Etalab survey 2020:
 - Issues on freshness
 - Insufficient or inaccurate documentation
 - Issues on dataset uniqueness

FAIR Principles

GO FAIR is committed to making data and services **findable, accessible, interoperable** and **reusable** (FAIR).



Findable: Metadata and data should be easy to find for both humans and computers.



Accessible: The exact conditions under which the data is accessible should be provided in such a way that humans and machines can understand them.



Interoperable: The (meta)data should be based on standardized vocabularies, ontologies, thesauri etc. so that it integrates with existing applications or workflows.

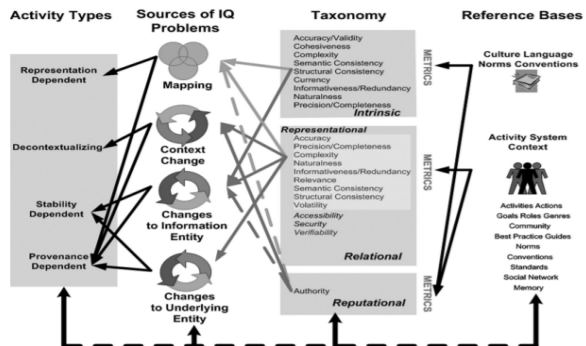


Reusable: Metadata and data should be well-described so that they can be replicated and/or combined in different research settings.

How can we measure data quality?

- Multiple frameworks and dimensions

Metadata, data and information quality



Stivilia, B., Gasser, L., Twidale, M. B., & Smith, L. C. (2007). A framework for information quality assessment. *Journal of the American society for information science and technology*, 58(12), 1720-1733

IMF

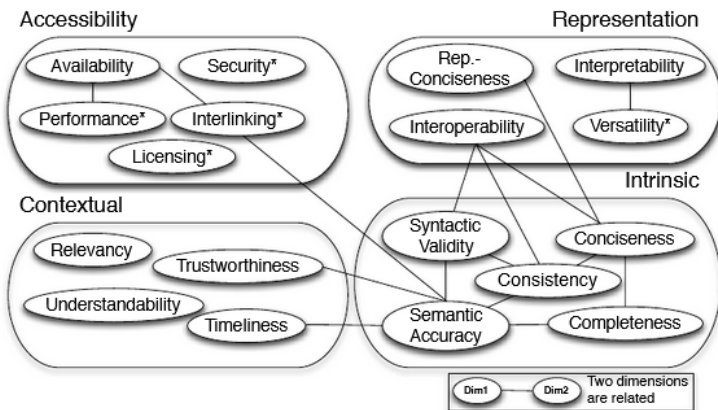
TABLE 1 The Six Dimensions of Data Quality

Relevance	The <i>relevance</i> of statistical information reflects the degree to which it meets the real needs of clients. It is concerned with whether the available information sheds light on the issues of most importance to users. Assessing relevance is a subjective matter dependent upon the varying needs of users. The NSO's challenge is to weigh and balance the conflicting needs of different users to produce a program that goes as far as possible in satisfying the most important needs and users within given resource constraints.
Accuracy	The <i>accuracy</i> of statistical information is the degree to which information correctly describes the phenomena it was to measure. It is usually characterized in terms of error in estimates and is traditionally decomposed into bias (systematic variance) and random error components. It may also be described as the degree to which the information is free from error. The major sources of error that potentially cause inaccuracies are coverage, sampling, nonresponse, and response.
Timeliness	The <i>timeliness</i> of statistical information refers to the date of the reference point (or the end of the reference period) to which the information pertains, and the date on which the information is available. It is typically involved in a trade-off against the <i>timeliness</i> of information will influence its relevance.
Accessibility	The <i>accessibility</i> of statistical information refers to the ease with which the information can be obtained from the NSO. This includes the ease of existence of information can be ascertained, as well as the form or medium through which the information can be obtained. The cost of the information may also be an aspect of accessibility.
Interpretability	The <i>interpretability</i> of statistical information reflects the degree to which the information is understandable and can be successfully brought together with other statistical information within a broad analytic framework and over time. The use of concepts, classifications and target populations promotes the use of common methodology across surveys. It does not necessarily imply full numerical consistency.
Coherence	The <i>coherence</i> of statistical information reflects the degree to which the information is brought together with other statistical information within a broad analytic framework and over time. The use of concepts, classifications and target populations promotes the use of common methodology across surveys. It does not necessarily imply full numerical consistency.

European data portal/JoinUp

- Accuracy:** is the data correctly representing the real-world entity or event?
- Consistency:** Is the data not containing contradictions?
- Availability:** Can the data be accessed now and over time?
- Completeness:** Does the data include all data items representing the entity or event?
- Conformance:** Is the data following accepted standards?
- Credibility:** Is the data based on trustworthy sources?
- Processability:** Is the data machine-readable?
- Relevance :** Does the data include an appropriate amount of data?
- Timeliness:** Is the data representing the actual situation and is it published soon enough?

Linked data



A study in 2020 by DAMA NL found 127 quality dimensions

Nr	Dimensions of data quality	Classification	Source
1.	Ease of operation		
2.	Reproducibility		
3.	Granularity		-
4.	Retention period		
5.	Accessibility		CDDQ 2019
6.	Accuracy		CDDQ 2019
7.	Completeness		CDDQ 2019
8.	Consistency		CDDQ 2019
9.	Currency		CDDQ 2019
10.	Integrity		CDDQ 2019
11.	Lineage		CDDQ 2019
12.	Precision		CDDQ 2019
13.	Representation		CDDQ 2019
14.	Timeliness		CDDQ 2019
15.	Validity		CDDQ 2019
16.	Coverage		Daas 2010
17.	Likability		Daas 2010
18.	Accuracy		DAMA 2017
19.	Completeness		DAMA 2017
20.	Consistency		DAMA 2017
21.	Currency (of data)	Timeliness	DAMA 2017
22.	Integrity (of coherence)		DAMA 2017
23.	Latency	Timeliness	DAMA 2017
24.	Reasonability		DAMA 2017
25.	Timeliness		DAMA 2017
26.	Uniqueness		DAMA 2017
27.	Validity		DAMA 2017
28.	Volatility	Timeliness	DAMA 2017
29.	Accuracy	Core dimension	DAMA-UK 2013
30.	Completeness	Core dimension	DAMA-UK 2013
31.	Confidence		DAMA-UK 2013
32.	Consistency	Core dimension	DAMA-UK 2013
33.	Flexibility		DAMA-UK 2013
34.	Timeliness	Core dimension	DAMA-UK 2013
35.	Uniqueness	Core dimension	DAMA-UK 2013
36.	Usability		DAMA-UK 2013

37.	Validity	Core dimension	DAMA-UK 2013
38.	Value		DAMA-UK 2013
39.	Accessibility	Pragmatic	English 1999
40.	Accuracy		English 1999
41.	Accuracy to a surrogate source	Inherent	English 1999
42.	Completeness (of values)	Inherent	English 1999
43.	Concurrency (of redundant or distributed data)	Inherent	English 1999
44.	Contextual clarity	Pragmatic	English 1999
45.	Database integrity		English 1999
46.	Definition conformance (see metadata conformance)	Inherent	English 1999
47.	Entity integrity		English 1999
48.	Equivalence	Inherent	English 1999
49.	Fact completeness	Pragmatic	English 1999
50.	Flexibility		English 1999
51.	Non-duplicates	Inherent	English 1999
52.	Precision	Inherent	English 1999
53.	Stability		English 1999
54.	Timeliness	Pragmatic	English 1999
55.	Usability	Pragmatic	English 1999
56.	Validity	Inherent	English 1999
57.	Accessibility		Eurostat 2015
58.	Accuracy		Eurostat 2015
59.	Clarity		Eurostat 2015
60.	Coherence		Eurostat 2015
61.	Comparability		Eurostat 2015
62.	Confidentiality		Eurostat 2015
63.	Consistency		Eurostat 2015
64.	Punctuality		Eurostat 2015
65.	Relevance		Eurostat 2015
66.	Reliability		Eurostat 2015
67.	Timeliness		Eurostat 2015
68.	Accessibility	Inherent/System dependant	ISO 25012
69.	Accuracy	Inherent	ISO 25012
70.	Availability	System dependant	ISO 25012
71.	Completeness	Inherent	ISO 25012
72.	Compliance	Inherent/System dependant	ISO 25012
73.	Confidentiality	Inherent/System dependant	ISO 25012
74.	Consistency	Inherent	ISO 25012
75.	Credibility	Inherent	ISO 25012
76.	Currentness	Inherent	ISO 25012
77.	Efficiency	Inherent/System dependant	ISO 25012
78.	Portability	System dependant	ISO 25012
79.	Precision	Inherent/System dependant	ISO 25012
80.	Recoverability	System dependant	ISO 25012

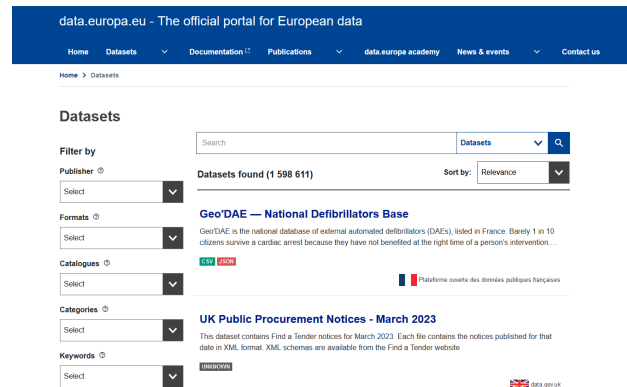
81.	Traceability	Inherent/System dependant	ISO 25012
82.	Understandability	Inherent/System dependant	ISO 25012
83.	Ability to represent null values	Representation	Redman 1996
84.	Accordance with format (of the physical instances)	Representation	Redman 1996
85.	Accuracy	Data values	Redman 1996
86.	Appropriateness	Representation	Redman 1996
87.	Clarity	Content	Redman 1996
88.	Completeness	Data values	Redman 1996
89.	Consistency	Data values	Redman 1996
90.	Currency	Data values	Redman 1996
91.	Efficient use (of storage)	Representation	Redman 1996
92.	Flexibility	Reaction to change	Redman 1996
93.	Format flexibility	Representation	Redman 1996
94.	Format precision	Representation	Redman 1996
95.	Granularity (of attributes)	Level of detail	Redman 1996
96.	Homogeneity	Composition	Redman 1996
97.	Identify-ability	Composition	Redman 1996
98.	Interpretability	Representation	Redman 1996
99.	Level of detail	Level of detail	Redman 1996
100.	Naturalness	Composition	Redman 1996
101.	Obtainability	Content	Redman 1996
102.	Portability	Representation	Redman 1996
103.	Precision (of attribute domains)	Level of detail	Redman 1996
104.	Redundancy (minimum necessary)	Composition	Redman 1996
105.	Relevance	Content	Redman 1996
106.	Robustness	Reaction to change	Redman 1996
107.	Scope	Scope	Redman 1996
108.	Semantic consistency (of the components of the model)	View consistency	Redman 1996
109.	Structural consistency (of attributes across entity types)	View consistency	Redman 1996
110.	Access security	Accessibility	Wang 1996
111.	Accessibility	Accessibility	Wang 1996
112.	Accuracy	Intrinsic	Wang 1996
113.	Appropriateness (of amount of data)	Contextual	Wang 1996
114.	Believability	Intrinsic	Wang 1996
115.	Completeness	Contextual	Wang 1996
116.	Conciseness (of representation)	Representational	Wang 1996
117.	Cost-effectiveness		Wang 1996
118.	Ease of understanding	Representational	Wang 1996
119.	Interpretability	Representational	Wang 1996
120.	Objectivity	Intrinsic	Wang 1996
121.	Relevancy	Contextual	Wang 1996
122.	Representational consistency	Representational	Wang 1996
123.	Reputation	Intrinsic	Wang 1996
124.	Timeliness	Contextual	Wang 1996
125.	Traceability		Wang 1996
126.	Value-added	Contextual	Wang 1996
127.	Variety		Wang 1996

A few examples of data reuse

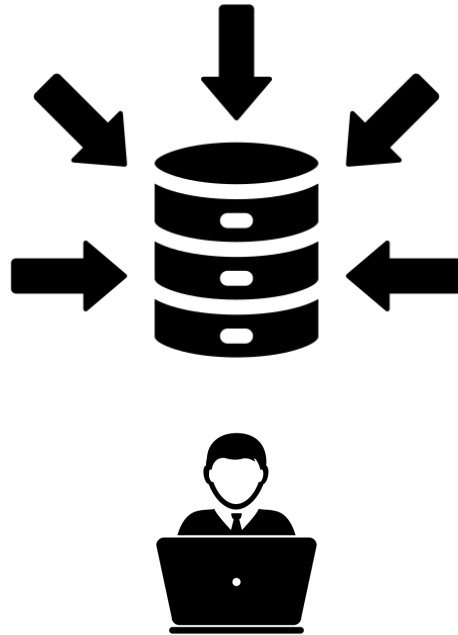
and the data quality issues they raised

A distributed digital library

Metadata aggregation



Agregate and share: Low cost, low touch



Decontextualized data



As Colonel of the Rough
Riders, 1898.

https://hollis.harvard.edu/primo-explore/fulldisplay?docid=HVD_VIAolvgroup12088&context=L&vid=HVD2&search_scope=everything&tab=everything&lang=en_US

Wendler, Robin. The Eye of the Beholder: Challenges of Image Description and Access at Harvard. In Hillmann, Diane I. and Westbrooks, Elaine L., eds., Metadata in Practice. American Library Association, Chicago, IL, 2004, 51-69.



- By using the collection description we could create a full match on 17% of multi-term searches

Foulonneau, M., Cole, T. W., Habing, T. G., & Shreeves, S. L. (2005). Using collection descriptions to enhance an aggregation of harvested item-level metadata. Proceedings of the 5th ACM/IEEE-CS Joint Conference on Digital Libraries, 2005. JCDL'05, (pp. 32–41).

Pseudo-duplicates

```
<dc:xmlns="http://www.openarchives.org/OAI/2.0/oai_dc">
<identifier>http://www.usndl.umich.edu/KC-KDDMG-X:
BBN7481-UND-0002-UND-001-UND-
00000470BBN7481_0002_001_00000470</identifier>
<publisher>Knight's American Mechanical
Dictionary</publisher>
<format>img</format>
<right>These pages may be freely searched and displayed.
Permission must be received for subsequent distribution in
print or electronically. Please go to
http://www.usndl.umich.edu/ for more
information.</right>
<title>Knight's American mechanical dictionary</title>
<type>image</type>
<subject>Industrial arts, Dictionaries, Mechanical
engineering, Technology, Inventions</subject>
<subject>4; Knight's American mechanical dictionary :
being a description of tools, instruments, machines,
processes, and engineering : history of inventions : general
technological vocabulary : and digest of mechanical;
appliance in science and the arts; 3 v., 75 leaves of plates;
illustrated with upwards of five thousand
engravings.</subject>
<language>UND</language>
<creator>Knight, Edward Henry</creator>
<source>/source>
</dc>
```



```
<dc:xmlns="http://www.openarchives.org/OAI/2.0/oai_dc">
<identifier>http://www.usndl.umich.edu/KC-KDDMG-X:
BBN7481-UND-0003-UND-001-UND-
00000035BBN7481_0003_001_00000035</identifier>
<publisher>Knight's American Mechanical
Dictionary</publisher>
<format>img</format>
<right>These pages may be freely searched and displayed.
Permission must be received for subsequent distribution in
print or electronically. Please go to
http://www.usndl.umich.edu/ for more
information.</right>
<title>Knight's American mechanical dictionary</title>
<type>image</type>
<subject>Industrial arts, Dictionaries, Mechanical
engineering, Technology, Inventions</subject>
<subject>3; Knight's American mechanical dictionary :
being a description of tools, instruments, machines,
processes, and engineering : history of inventions : general
technological vocabulary : and digest of mechanical;
appliance in science and the arts; 3 v., 75 leaves of plates;
illustrated with upwards of five thousand
engravings.</subject>
<language>UND</language>
<creator>Knight, Edward Henry</creator>
<source>/source>
</dc>
```



<https://www.data.gouv.fr>

- For 2.8% of collections, every time a query matched an item of the collection, it matched all of them

ENDORSE

Foulonneau, M. (2007). Information redundancy across metadata collections. Information processing & management, 43(3), 740-751.

Standardization: Data does not always conform to expectations

DC Date

September 29–October 28, 51 AD; 1970
second half of IXth century AD; 1978
Rebuilt 1984
Possibly Vth/VIth century AD; 1935
Planted 1985
n/a
n.d.
Mid IIInd century AD; 1973
Jul-51
circa 900 AD
ca. 701 BC
Begun 14th century
184-?
1839
18–?
August 23, 2000
between 1827 and 183
VIIIth/IXth century AD ? (TC); 1965
Vth-VIth century AD (McNamee); IVth
century AD (Cribiore); 1982

XVIII Dynasty
Winter 2003
era of redevelopment
various
2002-00
1980, refurbished 1997
China: Neolithic Period (5000 BCE-ca 1600
BCE)?
1969/1968
21. Nouemb. Anno. 1564.
And finisshed on the euen of thanunciacion
of our said blissid Lady falling on the
wednesday the xxiii daye of Marche. in the
xix yeer of Kyng Edwarde the fourthe
[1479]]
19193
xxxx Oct xx
Various
1938-05-38
1963 to 1953
[not after 1579]
163[5?]

Dimension (width x height)	Descriptive name
48x48	very small
64x64	Small
96x96	Medium
128x128	Large
144x144	extra large
160x160	super large
192x192	ridiculous large



Jens Finke's reference sizes for
thumbnails



Data completeness and usability for specific user tasks

Dublin Core element	% of repositories using element at least once	No. of records containing element	Total times element used	% of total records containing element	Average times used per record	Average element length (in characters)	Mode	Frequency in %
Title	100.0	124,304	133,108	80.3	1.1	39.9	1	75.8
Creator	87.5	78,402	84,829	50.7	1.1	21.5	0	49.3
Subject	93.8	112,875	304,661	72.9	2.7	110.4	2	37.1
Description	81.3	73,298	153,088	47.4	2.1	104.1	0	52.6
Publisher	75.0	94,791	114,305	61.2	1.2	38.5	1	50.9
Contributor	62.5	10,158	16,813	6.6	1.7	47.0	0	93.4
Date	81.3	66,514	77,175	43.0	1.2	10.9	0	57.0
Type	81.3	118,419	124,853	76.5	1.1	6.6	1	72.5
Format	56.3	107,381	111,647	69.4	1.0	8.3	1	66.6
Identifier	100.0	154,113	205,719	99.6	1.3	84.4	1	71.5
Source	50.0	23,012	29,537	14.9	1.3	68.3	0	85.1
Language	75.0	85,201	85,397	55.0	1.0	3.3	1	54.9
Relation	43.8	48,356	80,629	31.2	1.7	95.6	0	68.8
Coverage	37.5	9,136	12,103	5.9	1.3	21.0	0	94.1
Rights	62.5	63,435	68,228	41.0	1.1	151.7	0	59.0



← → ↺ 🔒 https://www.europeana.eu/en/search?page=1&view=grid&query=emile zola

≡ europeana HOME COLLECTIONS STORIES FOR PROFESSIONALS LOG IN / JOIN 🔍

6,343 RESULTS FOR emile zola

Filter results

THEME
Select a theme

TYPE OF MEDIA
Select types of media

Text (5,855)
Image (471)
Video (9)
Sound (8)

LANGUAGE

A föld

Mouret abbé vétke

Stvilja, B., Gasser, L., & Twidale, M. B. (2007). Metadata quality problems in federated collections. In Challenges of Managing Information Quality in Service Organizations (pp. 154-186). IGI Global.

Assess and show impact based on user tasks

User tasks

- Find
- Identify
- Select
- Obtain
- Explore

(FRSAD and FRBR)

The screenshot displays a web interface for searching datasets. On the left, there are filter sections: 'Filter by' with a 'Publisher' dropdown (set to 'Select'), 'Formats' with a 'Select' dropdown and a list of formats (CSV: 257 709, WMS: 209 923, WFS: 125 711, JSON: 106 077, HTML: 80 528), and 'Catalogues' with a 'Select' dropdown. The main area shows 'Datasets found (1 598 611)' with a 'Sort by: Relevance' dropdown. Three dataset results are visible: 'Geo'DAE — National Defibrillators Base' (with CSV and JSON format tags), 'UK Public Procurement Notices - March 2023' (with an UNKNOWN format tag), and 'Charging stations FROTH Network'. Each result includes a brief description and a flag icon (France, UK, and an unknown flag respectively).

Datasets

Filter by

Publisher ⓘ

Select

Formats ⓘ

Select

Filter

<input type="checkbox"/> CSV	257 709
<input type="checkbox"/> WMS	209 923
<input type="checkbox"/> WFS	125 711
<input type="checkbox"/> JSON	106 077
<input type="checkbox"/> HTML	80 528

Catalogues ⓘ

Select

Search

Datasets

Datasets found (1 598 611)

Sort by: Relevance

Geo'DAE — National Defibrillators Base

Geo'DAE is the national database of external automated defibrillators (DAEs), listed in France. Barely 1 in 10 citizens survive a cardiac arrest because they have not benefited at the right time of a person's intervention....

CSV JSON

🇫🇷 Plateforme ouverte des données publiques françaises

UK Public Procurement Notices - March 2023

This dataset contains Find a Tender notices for March 2023. Each file contains the notices published for that date in XML format. XML schemas are available from the Find a Tender website

UNKNOWN

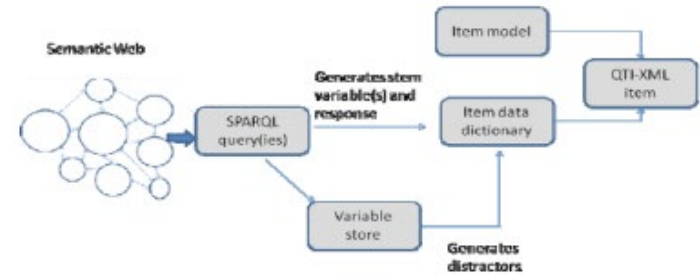
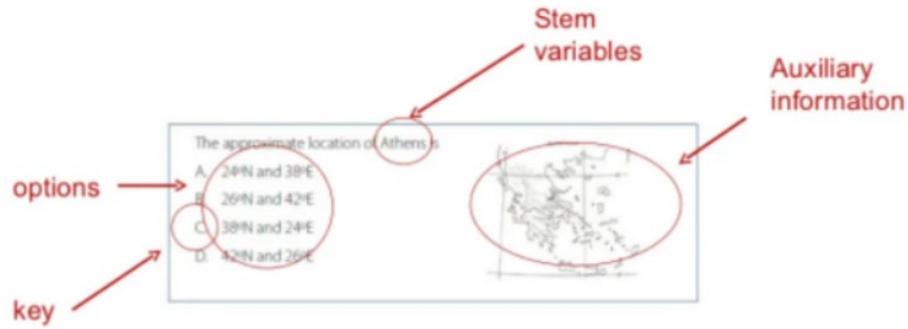
🇬🇧 data.gov.uk

Charging stations FROTH Network

Froth is a publicly available network of EVRIs deployed and operated by IZIVIA on behalf of small and medium-sized businesses and communities. Charging points are accessible with an IZIVIA pass, a...

Educational test items generation

Assessment item generation



Foulonneau, M., Ras, E. (2013). Assessment Item Generation, the way forward. International Computer Assisted Assessment (CAA) Conference, Southampton, UK
Adapted from Gierl, M. J., & Haladyna, T. M. (Eds.). (2012). Automatic item generation: Theory and practice. Routledge. New York

Foulonneau, M. (2012). Generating educational assessment items from linked open data: The case of DBpedia. In Extended Semantic Web Conference (pp. 16-27). Springer, Berlin, Heidelberg.

Risks on datasets interdependencies

```
<?xml version="1.0" encoding="UTF-8" standalone="no"
- <assessmentItem xmlns="http://www.imsglobal.org/
  instance" adaptive="false" identifier="choice" timeDe
  xsi:schemaLocation="http://www.imsglobal.org/x:
- <responseDeclaration baseType="identifiant" cardinalit
  - <correctResponse>
    <value>Option2</value>
  </correctResponse>
</responseDeclaration>
- <outcomeDeclaration baseType="integer" cardinality:
  - <defaultValue>
    <value>0</value>
  </defaultValue>
</outcomeDeclaration>
- <itemBody>
  - <p>
    
  </p>
  - <choiceInteraction maxChoices="1" responseIdentifier="RESPONSE" shuffle="false">
    <prompt>Which country is represented by this flag ? </prompt>
    <simpleChoice identifier="Option0">Bulgaria</simpleChoice>
    <simpleChoice identifier="Option1">Azerbaijan</simpleChoice>
    <simpleChoice identifier="Option2">Luxembourg</simpleChoice>
  </choiceInteraction>
</itemBody>
<responseProcessing template="http://www.imsglobal.org/question/qti_v2p0/rptemplates/match_correct" />
</assessmentItem>
```

About: [Luxembourg](#)

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

Luxembourg est une ville du comté de Kewaunee dans le Wisconsin. Sa population était de 2 515 habitants en 2010.

Property	Value
dbo:PopulatedPlace/area	<ul style="list-style-type: none">2584.8081341153282586.4
dbo:PopulatedPlace/populationDensity	<ul style="list-style-type: none">242.0232.7423811693864

- 6 out of 30 missing links

Foulonneau, M. (2012). Generating educational assessment items from linked open data: The case of DBpedia. In Extended Semantic Web Conference (pp. 16-27). Springer, Berlin, Heidelberg.

Lexical consistency

Who succeeded to ?

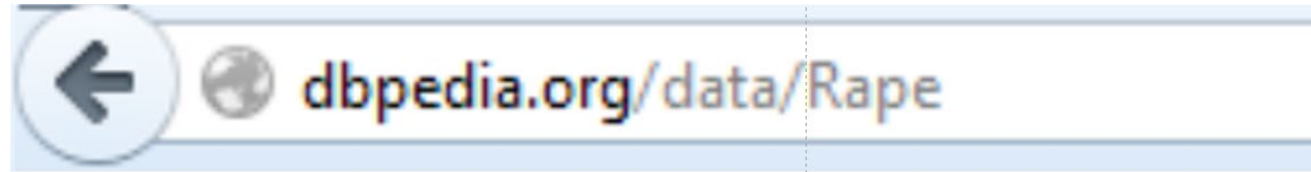
- Charles VII the Victorious
- Charles 09 Of France
- Louis VII

Risks on data accuracy



```
<dbpprop:mother rdf:resource="http://dbpedia.org/resource/Marie_de'_Medici"/>
<dbpprop:name xml:lang="en">XIII, Louis</dbpprop:name>
<dbpprop:name xml:lang="en">The Three Musketeers</dbpprop:name>
<dbpprop:name xml:lang="en">Louis XIII</dbpprop:name>
```

Data accuracy and the challenge of semantic data modelling



```
- <rdf:Description rdf:about="http://dbpedia.org/resource/Rape">  
  <rdf:type rdf:resource="http://umbel.org/umbel/rc/AilmentCondition"/>  
  <rdf:type rdf:resource="http://dbpedia.org/ontology/Disease"/>
```

Quality indicators and thresholds to assess usability

Based on quality issues in source data, what is the risk generating a “good” item?

Graph	Which part of the statements are accurate?
NELL	74%
YAGO	95%

Brumes et pluies

Ô fins d'automne, hivers, printemps trempés de boue,

Endormeuses saïZons ! je vous aime et vous loue

D'envelopper ainsi n on cœur et mon cerveau

D'un linceul vapoureux et d'un vague tombeau.

A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. Hruschka Jr, and T. M. Mitchell, “Toward an architecture for never-ending language learning,” in AAAI, 2010

F.-c. M. Suchanek, G. Kasneci, and G. Weikum, “YAGO: A core of semantic knowledge,” in WWW, pp. 697–706, 2007

A question-answering system

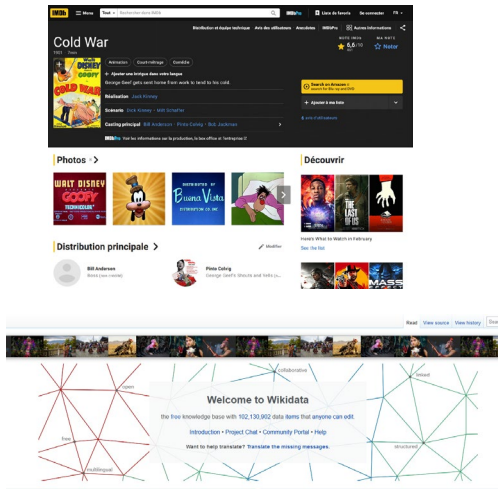
Answering factual questions



What is the
capital of
Mozambique?

Multiple sources – multiple risks

Closed
structured
data



Open
structured
data

Unstructured
data

be both more accurate and more resilient.

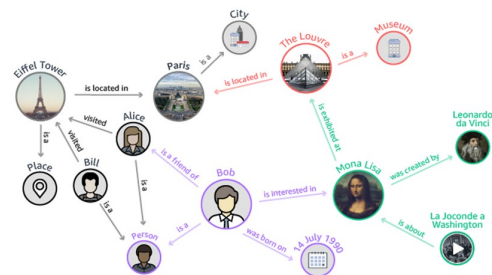
At the first FEVER workshop, [see annotated](#) the performance of 23 teams that participated in the first challenge. The top four finishers allowed us to create versions of their systems that we could host online, so that participants in the second FEVER challenge could attack them at will.

Since the first workshop, however, another 39 teams have submitted fact verification systems trained on FEVER data, pushing the top FEVER score from 64% [up to 70%](#). Three of those teams also submitted hostable versions of their systems. *Revisions that fail to resolve all instances for this second challenge to occur.*

Original REFUTED Instance:
Bullitt is a movie directed by Phillip D'Antoni
Adversarial REFUTED Instance:
There is a movie directed by Phillip D'Antoni called Bullitt.
Adversarial SUPPORTED Instance:
Bullitt is not a movie directed by Phillip D'Antoni
Evidence:
Bullitt is a 1968 American action thriller film directed by Peter Yates and produced by Philip D'Antoni

Two examples of adversarial assertions designed to confound a system trained on an assertion (the original refuted instance) in the original FEVER data set, together with supporting evidence drawn from Wikipedia.

<https://www.amazon.science/blog/the-fever-data-set-what-doesnt-kill-it-will-make-it-stronger>



<https://aws.amazon.com/fr/neptune/knowledge-graphs-on-aws/>

- How do I know what's true?

ENDORSE

How many sources agree?

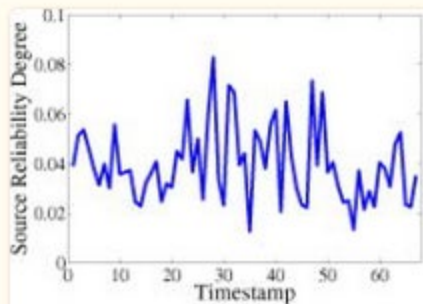
https://obamawhitehouse.archives.gov/sites/default/files/rss_viewer/birth-certificate-long-form.pdf

Zoom automatique

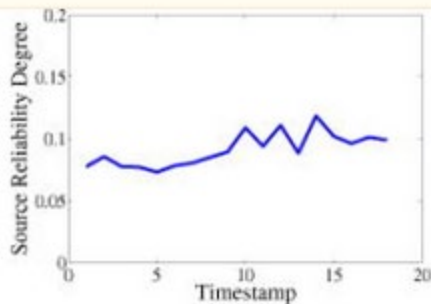
STATE OF HAWAII		CERTIFICATE OF LIVE BIRTH		DEPARTMENT OF HEALTH	
		FILE NUMBER 151		61 10641	
1a. Child's First Name (Type or print)	1b. Middle Name	1c. Last Name			
BARACK	HUSSEIN	OBAMA, II			
2. Sex	3. This Birth	4. If Twin or Triplet, Was Child Born	5a. Birth Date	5b. Month	5c. Day
Male	Single <input checked="" type="checkbox"/> Twin <input type="checkbox"/> Triplet <input type="checkbox"/>	1st <input type="checkbox"/> 2nd <input type="checkbox"/> 3rd <input type="checkbox"/>	August	4	1961
6a. Place of Birth: City, Town or Rural Location	6b. Island		6c. Hour		
Honolulu	Oahu		7:24 P.M.		
6d. Name of Hospital or Institution (If not in hospital or institution, give street address)			6e. Is Place of Birth Inside City or Town Limits?		
Kapiolani Maternity & Gynecological Hospital			Yes <input checked="" type="checkbox"/> No <input type="checkbox"/>		
7a. Usual Residence of Mother: City, Town or Rural Location	7b. Island	7c. County and State or Foreign Country			
Honolulu	Oahu	Honolulu, Hawaii			
10. Street Address	11. In Residence Inside City or Town Limits?		12. Is Residence on a Farm or Plantation?		
6085 Kalanianaʻole Highway	Yes <input checked="" type="checkbox"/> No <input type="checkbox"/>		Yes <input type="checkbox"/> No <input checked="" type="checkbox"/>		
13. Mother's Mailing Address	14. Race of Mother		15. Date Last Worked		
	Caucasian		17. Date Last Worked		
8. Full Name of Father	9. Race of Father	10. Age of Father			
BARACK HUSSEIN OBAMA	African	25			
11. Birthplace (State, Town or Foreign Country)	12. Usual Occupation	13. Kind of Business or Industry			
Kenya, East Africa	Student	University			
14. Full Maiden Name of Mother	15. Race of Mother		16. Date of Signature		
STANLEY ANN DONHAM	Caucasian		8-7-61		
17. Age of Mother	18. Birthplace (State, Town or Foreign Country)	19. Type of Occupation Outside Home During Pregnancy		20. Date of Signature	
18	Wichita, Kansas	None		8-8-61	
I certify that the above stated information is true and correct to the best of my knowledge.		18a. Signature of Father or Other Informant		18b. Date of Signature	
		Stanley Ann Obama		8-7-61	
I hereby certify that this child was born alive on the date and hour stated above.		19a. Signature of Accoucheur		19b. Date of Signature	
		W. A. Simola		8-8-61	
20. Date Accepted by Local Reg.		21. Signature of Local Registrar		22. Date Accepted by Reg. General	
AUG - 8 1961		W. A. Simola		AUG - 8 1961	
23. Evidence for Delayed Filing or Alteration					

https://obamawhitehouse.archives.gov/sites/default/files/rss_viewer/birth-certificate-long-form.pdf

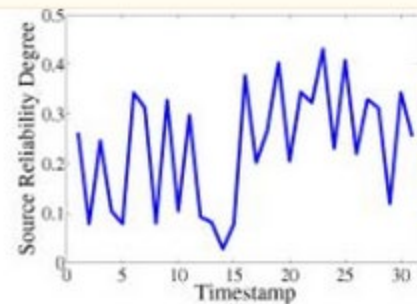
Data source quality is a moving target



(a) Weather



(b) Stock



(c) Flight

Li, Yaliang, et al. "On the discovery of evolving truth." Proceedings of the 21th acm sigkdd international conference on knowledge discovery and data mining. 2015.

Consistency with prior knowledge

- Rules:
 - a person cannot manage a company if they are dead
 - a horse cannot fly

- Data:

Table 7: Top 20 male professions in FB3M relative to female using ComplEx embeddings

Profession	Score	C_{male}	$C_{fem.}$
/m/0513qg	0.186	160	8
detective	0.163	27	2
trumpeter	0.161	346	6
gangster	0.146	45	0
private investigator	0.142	18	4
assn. football manager	0.132	587	5
Trombonist	0.131	196	1
session musician	0.130	184	7
sailor	0.119	429	23
bodyguard	0.117	33	2
bandleader	0.115	533	32
assn. football player	0.115	13321	227
samurai	0.114	26	0
music director	0.114	643	29
mastering engineer	0.111	33	1
clergy	0.107	78	4
baseball umpire	0.107	88	0
rabbi	0.105	180	5
Mafioso	0.103	60	0
statistician	0.103	205	3

Table 8: Top 20 female professions in FB3M relative to male using ComplEx embeddings

Profession	Score	$C_{fem.}$	C_{male}
gravure idol	0.210	62	0
fitness professional	0.184	24	12
Nude Glamour Model	0.177	511	1
showgirl	0.171	41	0
nun	0.167	41	0
socialite	0.164	81	11
art model	0.157	22	2
Key Hair Stylist	0.157	43	11
jewellery designer	0.154	39	9
fashion model	0.153	508	32
nurse	0.152	185	20
supermodel	0.151	95	9
Memoirist	0.148	30	35
Adult model	0.147	24	1
pin-up girl	0.146	55	0
dialect coach	0.143	14	8
Prostitute	0.140	63	0
flight attendant	0.137	34	3
ballet dancer	0.135	237	104
Cheerleader	0.133	20	1



https://www.europeana.eu/en/item/03919/public_mistral_joconde_fr_ACTION_CHERCHER_FIELD_1_REF_VALUE_1_50030_026886

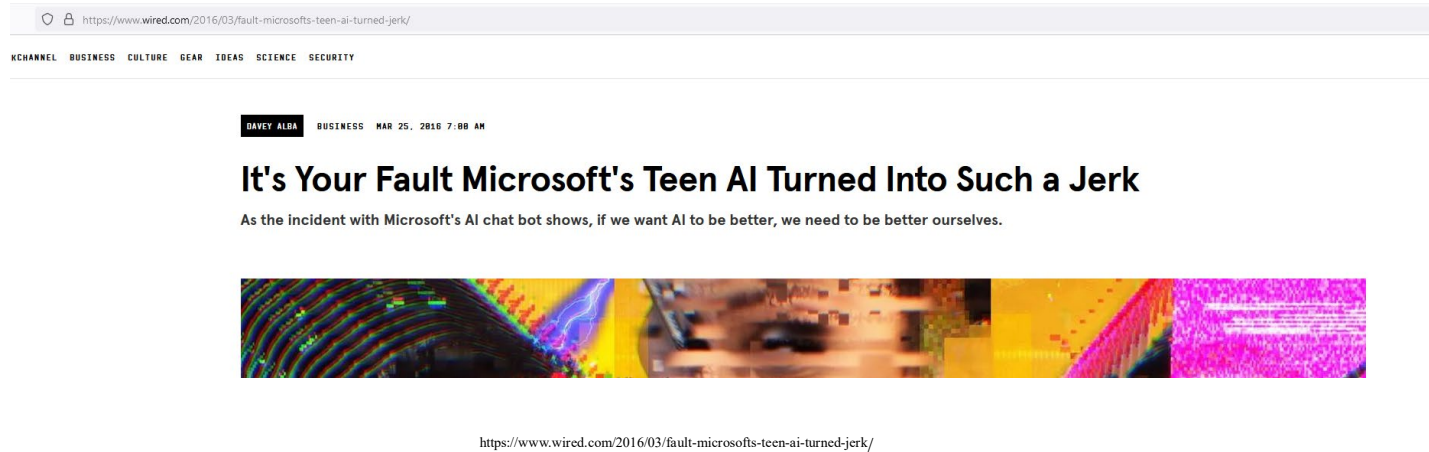
Fisher, J., Palfrey, D., Christodoulopoulos, C., & Mittal, A. (2019). Measuring social bias in knowledge graph embeddings. arXiv preprint arXiv:1912.02761.

Labelling, annotating and fixing data manually

- Humans are not perfect



The risk of getting it wrong



What's the truth?



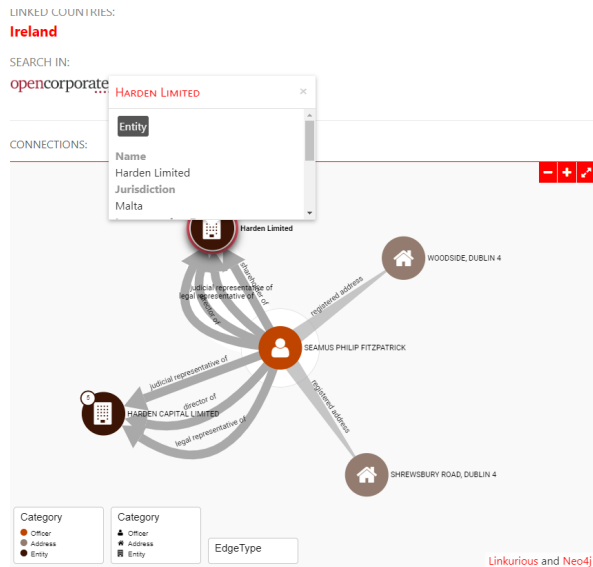
Who is the president of North Korea?

 CC BY-SA 3.0

https://commons.wikimedia.org/wiki/File:Mansudae_Grand_Monument_08.JPG

Anti-money laundering

Finding connections

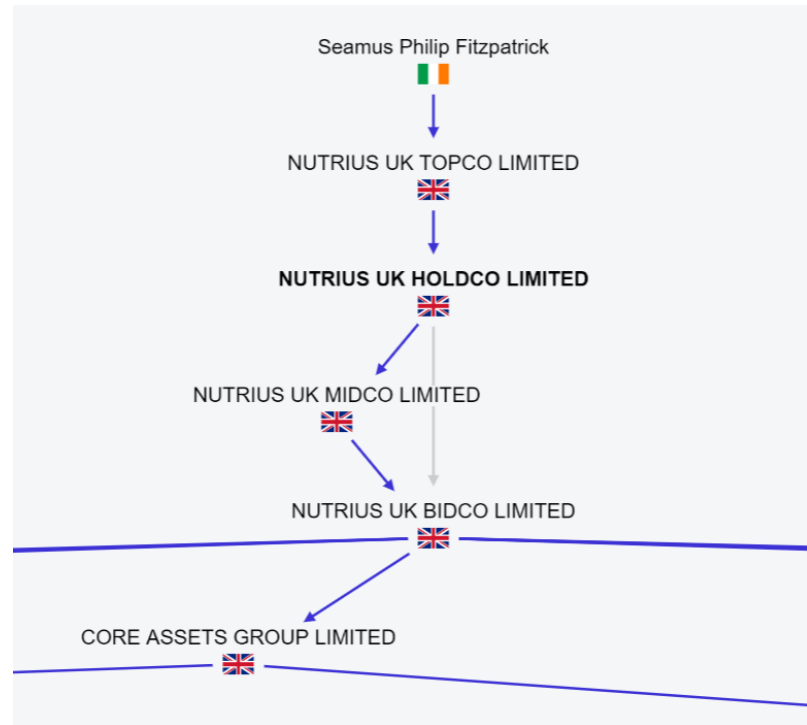
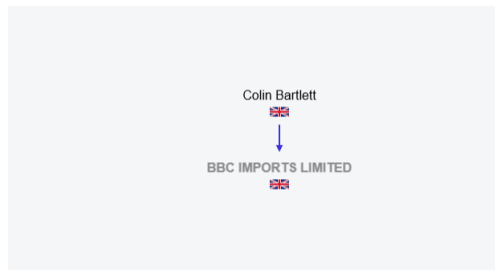


Entity (2)

	Role	From	To	Incorporation	Jurisdiction	Status	Data From
HARDEN CAPITAL LIMITED	Judicial representative	-	-	03-MAY-2013	Malta	-	Paradise Papers
HARDEN CAPITAL LIMITED	Director	-	-	03-MAY-2013	Malta	-	Paradise Papers
HARDEN CAPITAL LIMITED	Legal representative	-	-	03-MAY-2013	Malta	-	Paradise Papers
Harden Limited	Shareholder	-	-	12-JAN-2009	Malta	-	Paradise Papers

<https://offshoreleaks.icij.org/nodes/56101617>





Finding company structure



Entity resolution to establish connections

Is this the same person?

Is this the same company?

DURAND Jean-Michel	JM2D	12 RUE VICTOR MASSÉ 75009 PARIS	504 631 227 00015 RCS PARIS Siège social		kbis 
Durand Jean-Michel	POLE POSITION	60 RUE MONSIEUR LE PRINCE 75006 PARIS	791 407 091 00010 RCS PARIS Siège social		kbis 

<https://data.inpi.fr/>

Some issues due to international data collection

- Transliteration
- Multiple names
- Uneven variability of patronyms
- Geographic locations names change
- Geopolitical changes

Yevgeny

🌐 4 languages ▾

[Article](#) [Talk](#)

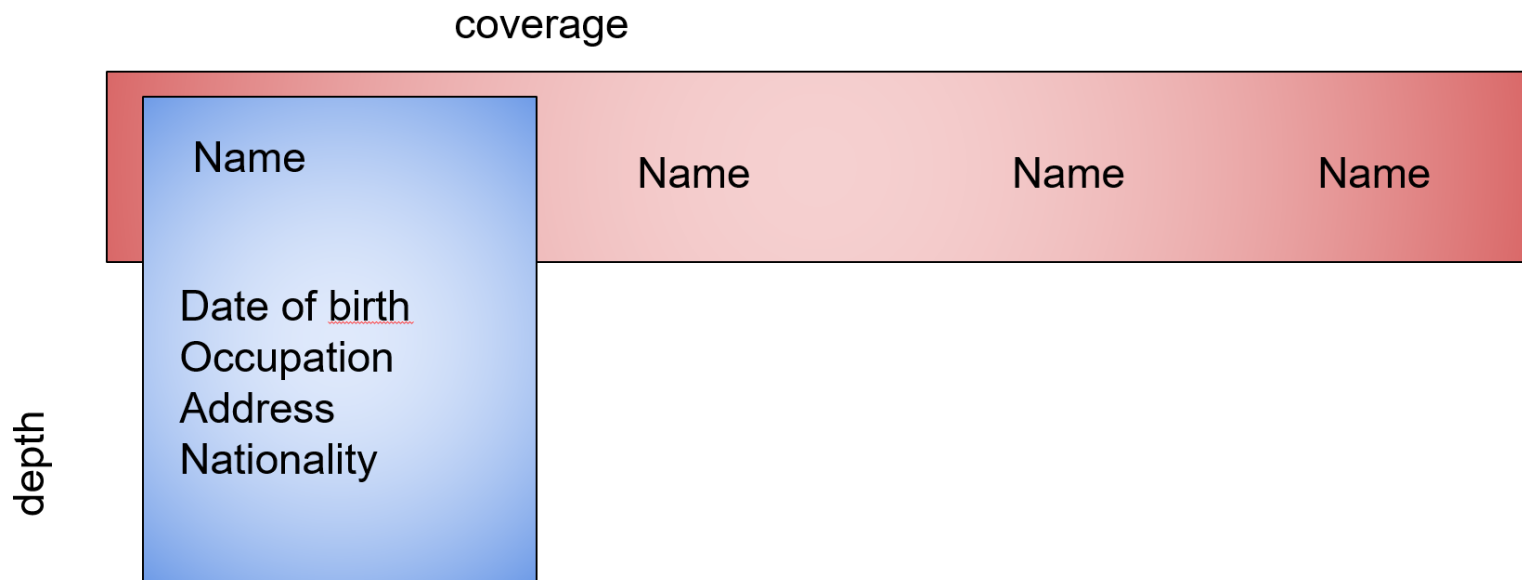
[Read](#) [Edit](#) [View history](#)

From Wikipedia, the free encyclopedia

Yevgeni, **Yevgeny**, **Yevgenii** or **Yevgeniy** (*Russian*: Евге́ний), also transliterated as **Evgeni**, **Evgeny**, **Evgenii**, **Evgeniy** or **Evgenij** is the Russian form of the masculine given name [Eugene](#). People with the name include:

<https://en.wikipedia.org/wiki/Yevgeny>

Completeness issues



Sparse data makes it difficult to answer:
'is this the same person?'

Assess a risk on data freshness

“Experts say 2 percent of records in a customer file become obsolete in one month because customers die, divorce, marry, and move.”

Eckerson, W. W. (2002). Data quality and the bottom line: Achieving business success through a commitment to high quality data. The Data Warehousing Institute, 1-36.

<https://data.inpi.fr/entreprises/450516737?q=wimi#450516737>

- The risk of data obsolescence

ENDORSE

Identité

Dénomination

SCI WIMI

SIREN (siège)

450 516 737

N° de gestion

2003D00300

Début d'activité

12/09/2003

Durée de la personne morale

99 ans

Date de clôture

31 Décembre

Forme juridique

Société civile

Activité principale

Acquisition, administration, gestion par location ou autrement de tous immeubles et plus spécialement, acquisition sur la commune de Vic-Fezensac d'une grange avec terrain attenant, le tout sis chemin de l'Abattoir.

Capital social

1 600.00 €

Adresse du siège

chemin de l'Abattoir 32190 Vic-Fezensac FRANCE

Département du siège

32

Représentants

Pour plus d'informations sur les représentants, veuillez vous connecter

Nom, Prénom(s)

LAMBERT Remi, Robert, Paul

(Gérant, Associé indéfiniment responsable)

Date de naissance (mm/aaaa)

06/1966

Nom, Prénom(s)

COCHONNEAU William, Louis, Marcel

(Associé indéfiniment responsable)

Date de naissance (mm/aaaa)

04/1956

Établissements

Type d'établissement

Siège et principal

Début d'activité

12/09/2003

Origine du fonds

Création

Type d'exploitation

Exploitation directe

Activité

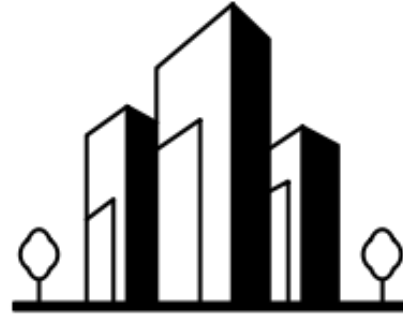
Acquisition, administration, gestion par location ou autrement de tous immeubles et plus spécialement, acquisition sur la commune de Vic-Fezensac d'une grange avec terrain attenant, le tout sis chemin de l'Abattoir.

Adresse

chemin de l'Abattoir 32190 Vic-Fezensac FRANCE

Exploration of the UK company registry

UK Companies House



Exploration of the UK company registry



[Sign in / register](#)

Search for a company or officer

[Advanced company search](#)

MINISTRY OF FINANCE OF SAUDI ARABIA LTD
Company number **10625481**

[Follow this company](#) [File for this company](#)

Overview [Filing history](#) [People](#) [More](#)

Registered office address
27 Old Gloucester Street, London, United Kingdom, WC1N 3AX

Company status
Active

Company type
Private limited Company

Incorporated on
17 February 2017

Some financial data on the company

MINISTRY OF FINANCE OF SAUDI ARABIA LTD

Registered Number 10625481

Balance Sheet as at 28 February 2022

	2022	2021
	£	£
Called up share capital not paid	1000000000	1000000000
Net assets	1000000000	1000000000
Issued share capital		
10000000 Ordinary Shares of £ 100 each	1000000000	1000000000
Total Shareholder funds	1000000000	1000000000

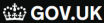
STATEMENTS

Financial statements for the financial year ending 28 February 2022

<https://find-and-update.company-information.service.gov.uk/company/10625481/filing-history>

Business associates



 **GOV.UK** Find and update company information

Companies House does not verify the accuracy of the information filed

[Sign in / Register](#)

[Advanced company search](#)

MINISTRY OF FINANCE OF SAUDI ARABIA LTD

Company number **10625481**

[Follow this company](#) [File for this company](#)

Overview

Filing history

People

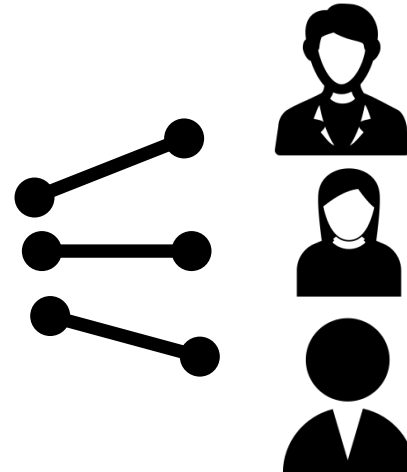
More

Registered office address
27 Old Gloucester Street, London, United Kingdom, WC1N 3AX

Company status
Active

Company type
Private limited Company

Incorporated on
17 February 2017



Indirect connections to other companies and business associates

BANQUE DU CANADA LTD

Company number **12342074**

[Follow this company](#)

[File for this company](#)

BANQUE CANTONALE DE GENÈVE LTD

Company number **09931624**

[Follow this company](#)

[File for this company](#)

Filter officers

☐

Current officers

[Overview](#)

[Filing history](#)

[People](#)

[More](#)

Registered office address

27 Old Gloucester Street, London, WC1N 3AX

Company status

Active

Company type

Private limited Company

Incorp

29 Dec

[Advanced company search](#)

MINISTRY OF FINANCE OF THE STATE OF QATAR LTD

Company number **11656934**

[Follow this company](#)

[File for this company](#)

[Overview](#)

[Filing history](#)

[People](#)

[More](#)

Registered office address

27 Old Gloucester Street, London, United Kingdom, WC1N 3AX

Balance Sheet as at 28 December 2021

	2021	2020
	£	£
Called up share capital not paid	1000000000	1000000000
Net assets	1000000000	1000000000
Issued share capital		
10000000 Ordinary Shares of £ 100 each	1000000000	1000000000
Total Shareholder funds	1000000000	1000000000

ENDORSE

STATEMENTS

Data collection method

GOV.UK Find and update company information

Companies House does not verify the accuracy of the information filed

Search for a company or officer

[Advanced company search](#)

UKCH LIMITED

Company number **10115452**

Follow this company | File for this company

Companies House is dysfunctional and facilitating fraud, MPs told

Less verification for someone to set up fraudulent shell firm than to borrow a library book, risk managers say



Anti-fraud bosses at NatWest and HSBC have criticised the online register of UK-based companies. Photograph: Dominic Lipinski/PA

The anti-fraud leader at the trade body UK Finance has said the government needs to fix the “dysfunctional” Companies House because it is helping to facilitate business fraud.

Data input interfaces and honest(?) mistakes

ELLI FINANCE (UK) PLC

Company number **08094161**

[Follow this company](#)

[File for this company](#)

[Overview](#)

[Filing history](#)

[People](#)

[Charges](#)

[Insolvency](#)

[More](#)

[Officers](#)

[Persons with significant control](#)

1 active person with significant control / 0 active statements

Elli Group (Uk) Limited ACTIVE

Correspondence address

Norcliffe House, Station Road, Wilmslow, United Kingdom, SK9 1BU

ELLI GROUP (UK) LIMITED

Company number **08092763**

[Follow this company](#)

[File for this company](#)

Elli Finance (Uk) Plc CEASED

Correspondence address

C/O Alvarez & Marsal Europe Llp, Suite 3 Regency House, 91 Western Road, Brighton, United Kingdom, BN1 2NW

Notified on
6 April 2016

Ceased on
30 April 2019

Governing law

United Kingdom (England And Wales)

Legal form

Public Limited Company

Place registered

Companies House

Registration number

08094161

Incorporated in

England And Wales

Lack of data standardizations makes connections difficult

London Laundromat: Police seize £2 million profits of Italian mafia gang held in British banks



this company

- 29 Chichele Road, London, United Kingdom, NW2 3AN
- Office 2092, No.1, Fore Street, London, England, EC2Y 5EJ
- 20-22, Wenlock Road, London, England, N1 7GU
-

Overview Filing history People More

Registered office address

7 29 Chichele Road, London, NW2 3AN

Company status

Active

Company type

Private limited Company

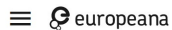
Incorporated on

14 October 2010

Aim for excellence but build services that
are resilient to imperfections

Data has multiple roles

Provide access in a new context



256 RESULTS FOR josephine baker



Sarano - award to orphans
Luce Institute



Sans Amour
Society for historical sound recordings

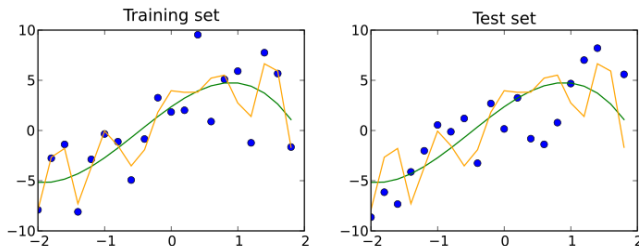


Milan - Josephine Baker at mother's day
Luce Institute



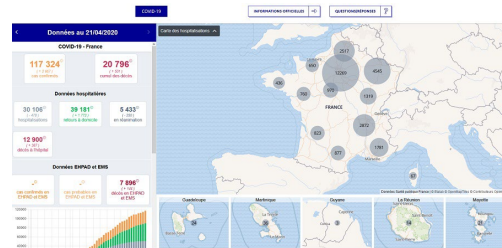
Josephine Baker at the Mother's Day
Luce Institute

Train model

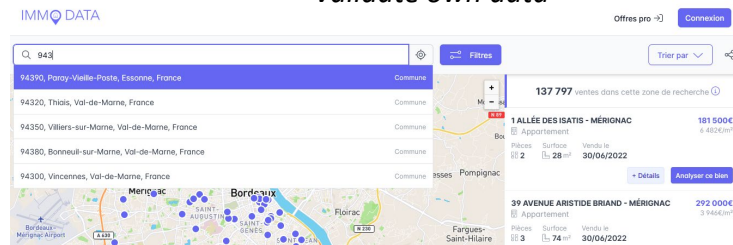


Skbekas, CC BY 3.0 <<https://creativecommons.org/licenses/by/3.0/>>, via Wikimedia Commons

Represent data



Validate own data



Validate algorithm performance

MovieLens 25M Dataset

MovieLens 25M movie ratings. Stable benchmark dataset. 25 million ratings and one million tag applications applied to 62,000 movies by 162,000 users. Includes tag genome data with 15 million relevance scores across 1,129 tags. Released 12/2019

- [README.txt](#)
- [ml-25m.zip](#) (size: 250 MB, [checksum](#))

Permalink: <https://grouplens.org/datasets/movielens/25m/>

MovieLens Tag Genome Dataset 2021

10.5 million computed tag-movie relevance scores from a pool of 1,084 tags applied to 9,734 movies. Released 12/2021. This dataset also contains input necessary to generate the tag genome using both the original process (Vig et al. 2012) and a more recent improvement (Kolkov et al. 2021)

- [genome_2021_readme.txt](#)
- [genome_2021.zip](#) (size: 1.6GB)

Permalink: <https://grouplens.org/datasets/movielens/tag-genome-2021>

Enrich own data



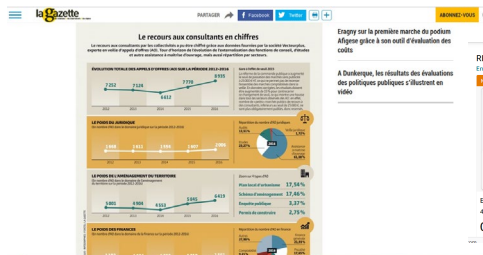
1. Choisir un fichier



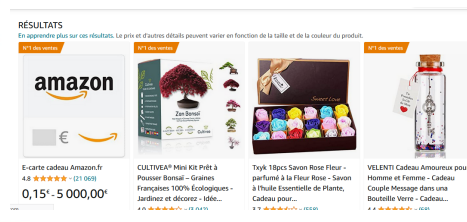
2. Aperçu du fichier et vérification de l'encodage

3. Choisir les colonnes à utiliser pour construire les adresses

Analyze data to support (political) decision making



Personalize experience



ENDORSE

Data quality as fitness for use or fitness for purpose

- A multiplicity of reuse contexts
- The data has a multiplicity of roles

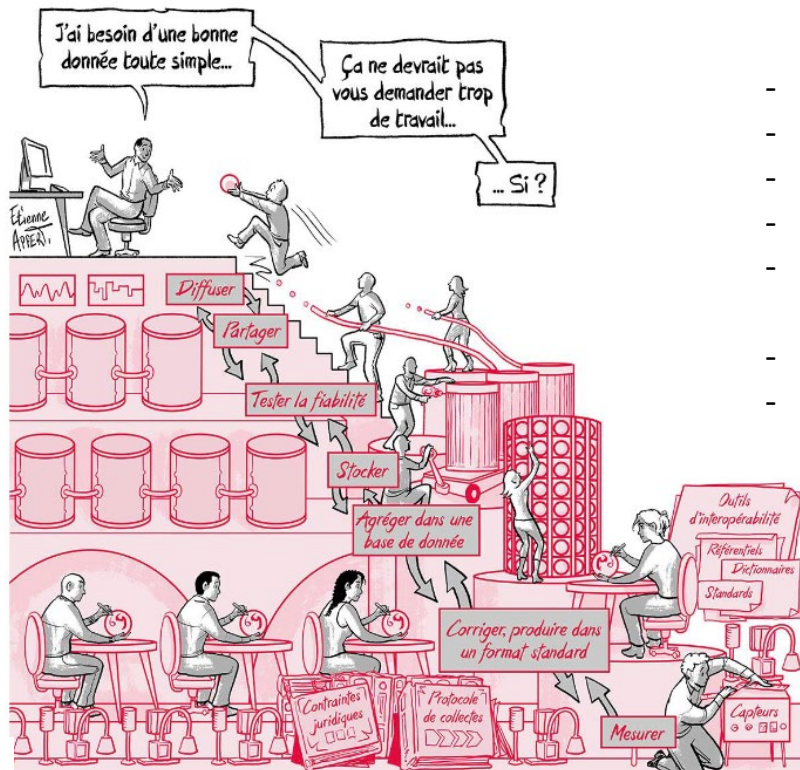
 not all are well known in open environments

Quality is a construction

- I need good data, it's very simple

... that should not require a lot of work


... or does it?



- Distribute
- Share
- Assess quality
- Store
- Aggregate in a database
- Fix, standardize
- Measure

Show impact

Most of the data is not bad in its original environment



Internet Archive

Track from Damien Rice Live at Orpheum Theatre on 2006-12-08

Disc 1: 01: Intro/Happy Birthday to Damien [01:17] 02: Delicate [05:14] 03: Woman Like A Man [04:29] 04: 9 Crimes [06:15] 05: Me, My Yoke + I [06:51] 06: Slow* > [02:19] 07: Rootless Tree [04:32] 08: Cannonball [04:56] 09: The Anim...

IN COPYRIGHT SOUND SAUVEGARDER MENTION J'AIME

Internet Archive

Track from Damien Rice Live at Orpheum Theatre on 2006-12-08

Disc 1: 01: Intro/Happy Birthday to Damien [01:17] 02: Delicate [05:14] 03: Woman Like A Man [04:29] 04: 9 Crimes [06:15] 05: Me, My Yoke + I [06:51] 06: Slow* > [02:19] 07: Rootless Tree [04:32] 08: Cannonball [04:56] 09: The Anim...

IN COPYRIGHT SOUND SAUVEGARDER MENTION J'AIME

Internet Archive

Track from Damien Rice Live at Orpheum Theatre on 2006-12-08

Disc 1: 01: Intro/Happy Birthday to Damien [01:17] 02: Delicate [05:14] 03: Woman Like A Man [04:29] 04: 9 Crimes [06:15] 05: Me, My Yoke + I [06:51] 06: Slow* > [02:19] 07: Rootless Tree [04:32] 08: Cannonball [04:56] 09: The Anim...

IN COPYRIGHT SOUND SAUVEGARDER MENTION J'AIME

HOME MAIL NEWS FINANCE SPORT CELEBRITY STYLE WEATHER MORE...

yahoo!news

Sign in Mail

News Cost of living crisis Yahoo Originals Quizzes TV + Celebrity Royals Crime Science & Tech Motoring Viral UK World ...

sky news | Sky News

Coronavirus: Boris Johnson unable to say how many people weren't traced due to 16,000 missed cases



Acknowledge excellence is a target

- But understand the risks
- To some extent it is possible to account for imperfection and work around it

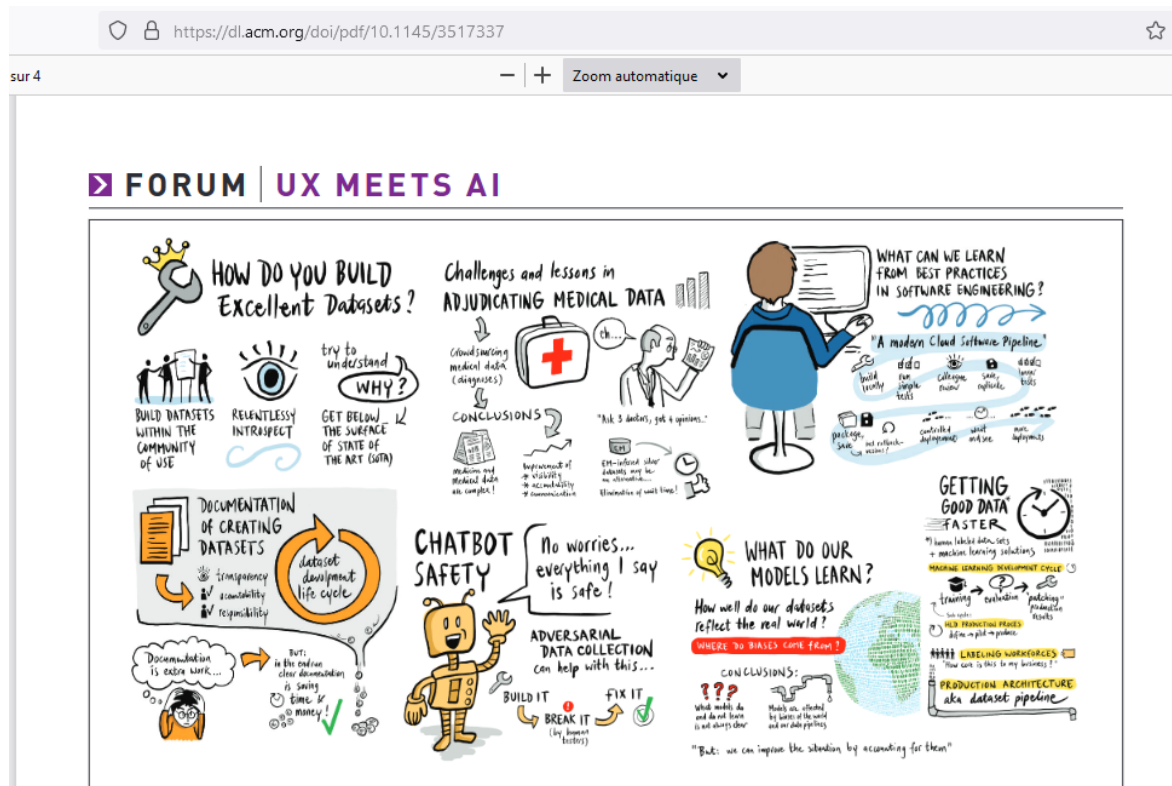


Figure 2 Visual overview of Data Excellence Workbench panels

Arovo, L., Lease, M., Paritosh, P., & Schackermann, M. (2022). Data excellence for AI: why should you care?. *Interactions*, 29(2), 66-69.

Quality, excellence and perfectionism

“Osborn employs the concept of perfectionism to describe a hyper-emphasis on exactitude and precision-in this case, quality gone awry by being taken to an extreme.”



Thomas, Sarah E. "Quality in bibliographic control." (1996).
<https://www.ideals.illinois.edu/items/7998/bitstreams/27643/stream>



Muriel Foulonneau

