

Getting the Most Out of Wikidata

Markus Krötzsch
Knowledge-Based Systems, TU Dresden

Reporting on joint work with
Adrian Bielefeldt, Fredo Erxleben, Julius Gonsior,
Larry Gonzalez, Michael Günther, Stas Malyshev,
Julian Mendez, Veronica Thost, and Denny Vrandečić

and supported by the Wikimedia Foundation

Wiki Workshop 2018



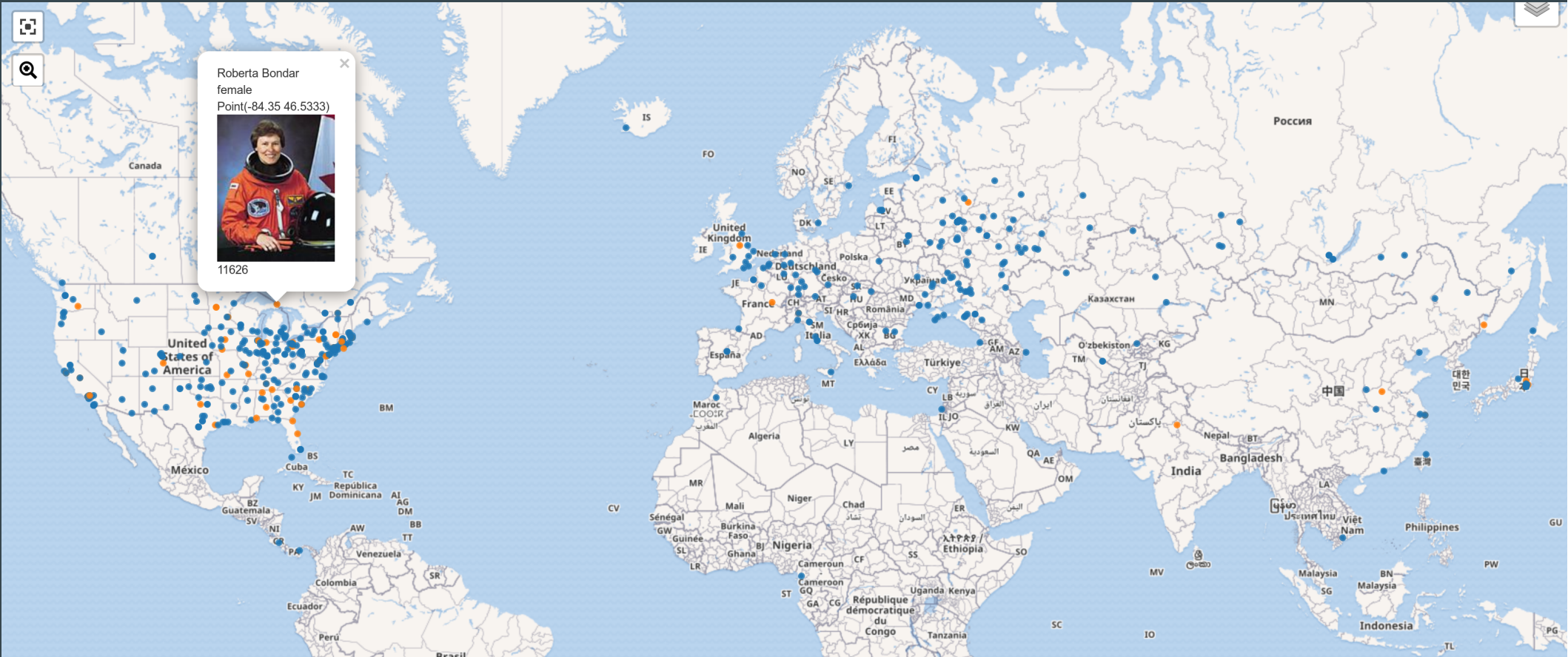
“What are the world’s largest cities with a female mayor?”

“What are the world’s largest cities with a female mayor?”

cityLabel	mayorLabel	population
Tokyo	Yuriko Koike	13742906
Hong Kong	Carrie Lam	7336585
Baghdad	Zekra Alwach	6960000
Surabaya	Tri Rismaharini	4975000
Yokohama	Fumiko Hayashi	3733234
Madrid	Manuela Carmena	3182981
Rome	Virginia Raggi	2873494
Kaohsiung City	Chen Chu	2777384
Antananarivo	Lalao Ravalomanana	2610000
Paris	Anne Hidalgo	2206488

“Where are people born who travel to space?”

(Colour-coded by gender)



“Which days of the week do disasters occur on?”

Date	Mon	Tue	Wed	Thu	Fri	Sat	Sun
1	25	33	22	18	26	28	23
2	24	26	23	23	22	32	12
3	24	27	21	31	23	28	38
4	24	25	33	25	26	26	24
5	37	23	32	18	19	17	29
6	25	28	32	20	24	33	22
7	18	22	25	16	22	18	17
8	32	28	19	25	22	23	19
9	20	25	29	29	27	21	23
10	20	20	19	14	25	25	29
11	30	34	28	23	22	20	20
12	41	33	27	30	20	20	23
13	35	26	29	21	25	24	25
14	24	23	27	23	22	28	17
15	15	22	22	24	19	22	15

“Which 19th century paintings show the moon?”



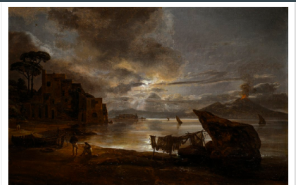
commons:J.Bernstae, Moon Pjag
Q: Moonlight Marine
Q: Moon



commons:Maufalvevskunaucaupacua1899.jpg
Q: Moonrise at Twilight
Q: Moon



commons:Johan Christian Clausen Dahl - Nysten ved Laurvig i Norge i midskyn - Thorvalds...
Q: The coast at Laurvig, Norway
Q: Moon



commons:Johan Christian Clausen Dahl - Bugten ved Napoli i midskyn (1821).jpg
Q: The Bay of Naples by moonlight
Q: Moon



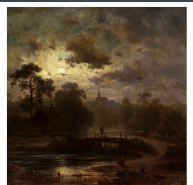
commons:Dresden in Moonlight by Johan Christian Dahl, Bergen Kunstmuseum.JPG
Q: Dresden by Moonlight
Q: Moon



commons:Johan Christian Dahl - Den Enrum spjæn i midskyn.jpg
Q: Enrum Lake in Moonlight
Q: Moon



commons:Dahl, Der Kopenhagener Hælen im Mondschein, 1831.jpg
Q: DGH18472
Q: Moon



commons:Diprd Landscape by moonlight.jpg
Q: D2899910
Q: Moon



commons:Scamertowski Morning star.jpg
Q: Morning star
Q: Moon



commons:Lampri Carthusian monastery.jpg
Q: D2730883
Q: Moon



commons:Max Tarnhauer.jpg
Q: D2309000
Q: Moon



commons:Louis Reyny Mignot Marsh in Ecuador.jpg
Q: Moonlight over a Marsh in Ecuador
Q: Moon



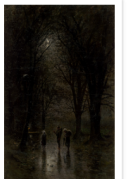
commons:Johan Christian Dahl - View of Dresden by Moonlight - Google Art Project (JWHK-NudHf7MG).jpg
Q: View of Dresden by Moonlight
Q: Moon



commons:Johan Christian Dahl - Dresden by Moonlight - Google Art Project.jpg
Q: Dresden by Moonlight
Q: Moon



commons:Elin Danielson-Gambogi Kuulamo i 1900.jpg
Q: Moonlight
Q: Moon



commons:Ladelaar Midydras...
Q: Night Travelers at a Cross
Q: Moon



commons:Night Travelers at a...
Q: Night Travelers at a Cross
Q: Moon



commons:Martin Johnson Heade - Two Owls at Sunset.jpg
Q: Two Owls at Sunset
Q: Moon



commons:Martin Johnson Heade - Sailing by Moonlight.jpg
Q: Sailing by Moonlight
Q: Moon



commons:Philo Judin, Rhode Island by Martin Johnson Heade, 1867-68.JPG
Q: Philo Judin, Rhode Island
Q: Moon



commons:Midwinter Moonlight by Regie Francois Gignoux, before 1880, oil on board...
Q: MS-Winter Moonlight
Q: Moon



commons:Stanford Robinson Gilford - Crêpuscule sur le mont Hunter.jpg
Q: Hunter Mountain, Twilight
Q: Moon



commons:Samuel Colman - The Rock of Salvation - Google Art Pr...
Q: Moonlight
Q: Moon



commons:John William Casler Mt...
Q: Moonrise on the Coast
Q: Moon



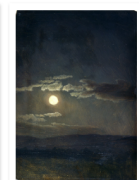
commons:Albert Pivarski Ryder - The Lowest Boat (c. 188...
Q: The Lowest Boat
Q: Moon



commons:Johan Christian Clausen Dahl - Ved den napolittanske golf, Midskyn - Staten...
Q: The Gulf of Naples, Moonlight
Q: Moon



commons:Johan Christian Clausen Dahl - Ved den napolittanske golf, Midskyn - Staten...
Q: The Gulf of Naples, Moonlight
Q: Moon



commons:Albert Bierstadt - Cloud...
Q: Cloud Study, Moonlight
Q: Moon



commons:John Martin - The Eve of the Deluge - WDA14146.jpg
Q: The Eve of the Deluge
Q: Moon



commons:Thomas Chambers - Storm-Tossed Frigate.jpg
Q: Storm-Tossed Frigate
Q: Moon



commons:Richards William Trost Moonlight On Mount Lafayette New Hampshire.jpg
Q: Moonlight on Mount Lafayette, New Hampshire
Q: Moon



commons:Moonrise by George Inness 1887.jpg
Q: Moonrise
Q: Moon



commons:“Moonlight” by George Inness, 1893.JPG
Q: Moonlight
Q: Moon



WIKIDATA

“The free knowledge base that anyone can edit”

Wikimania05/Paper-MK2

[< Wikimania05](#)

This page is part of the [Proceedings of Wikimania 2005](#), Frankfurt, Germany.

0 MISSING 1 Submitted 2 Editing 3 Author review 4 Final edit 5 DONE

Wikipedia and the Semantic Web - The Missing Links [\[edit\]](#)

- **Author(s):** Markus Krötzsch & Denny Vrandečić & Max Völkel
- **License:** CC-NC-SA 2.0 (for further license models, please contact the authors)
- **Slides:** collected but not uploaded yet
- **Video:** [16:44](#) (talk given by Denny Vrandečić)
- **Note:** Presentation, paper also at [Wikipedia and the Semantic Web](#) (PDF, 164K)

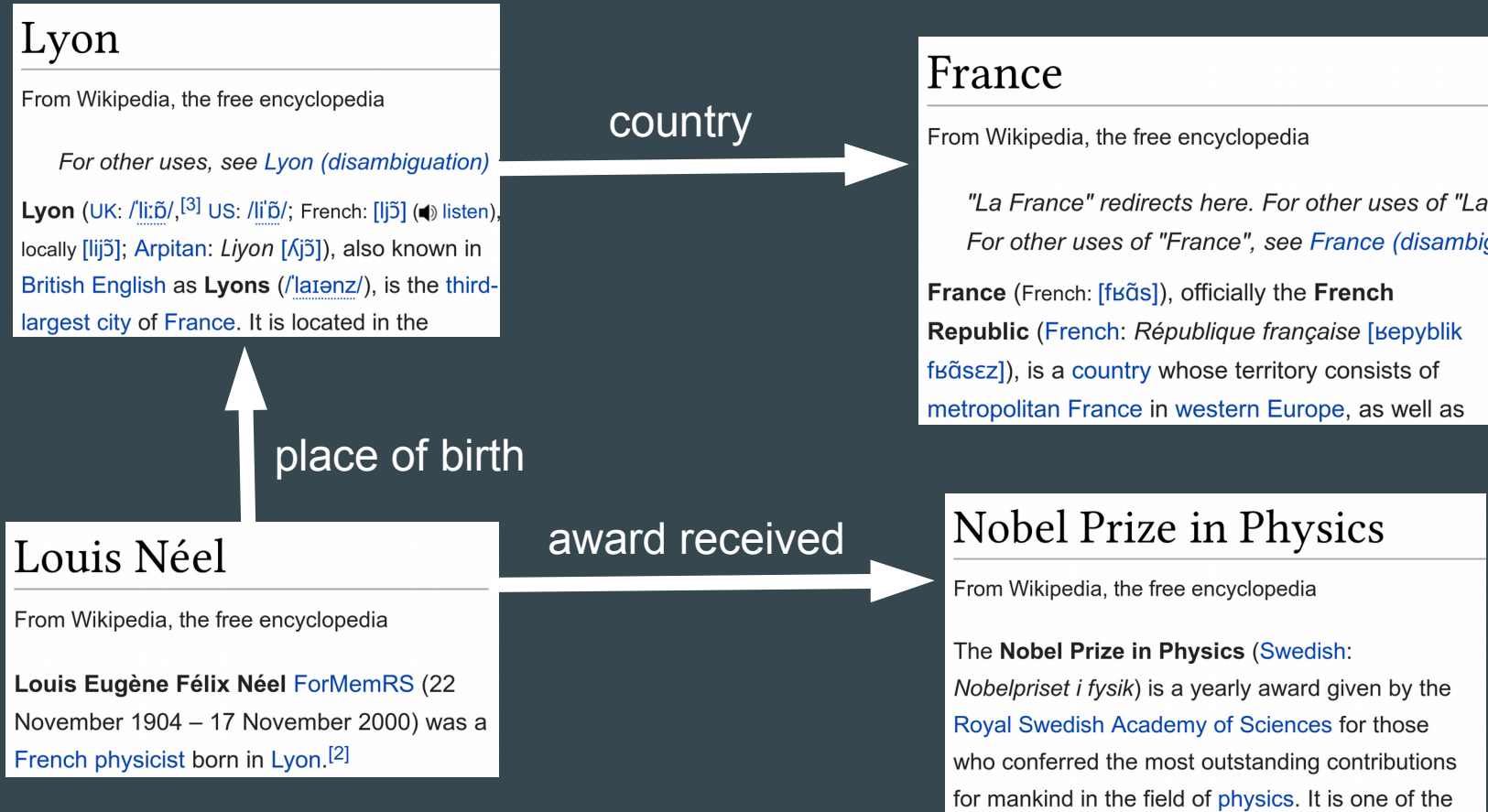
About the author: The authors are research associates at the *Institute of Applied Informatics and Formal Description Methods* (AIFB), [University of Karlsruhe](#), Germany, where they are members of the AIFB [Research Group Knowledge Management](#), an interdisciplinary team of computer scientists, mathematicians, and industrial engineers that is one of the world-wide leading institutions in the Semantic Web research community. Other relevant research topics include Semantic Web, ontologies, data and text mining, logic-based knowledge representation, peer-to-peer, and Web services.

Being enthusiastic users and contributors of various Wikis, the authors also have special interest in making emerging semantic technologies available within Wikis, where computer-assisted organization and processing of knowledge plays an important role.

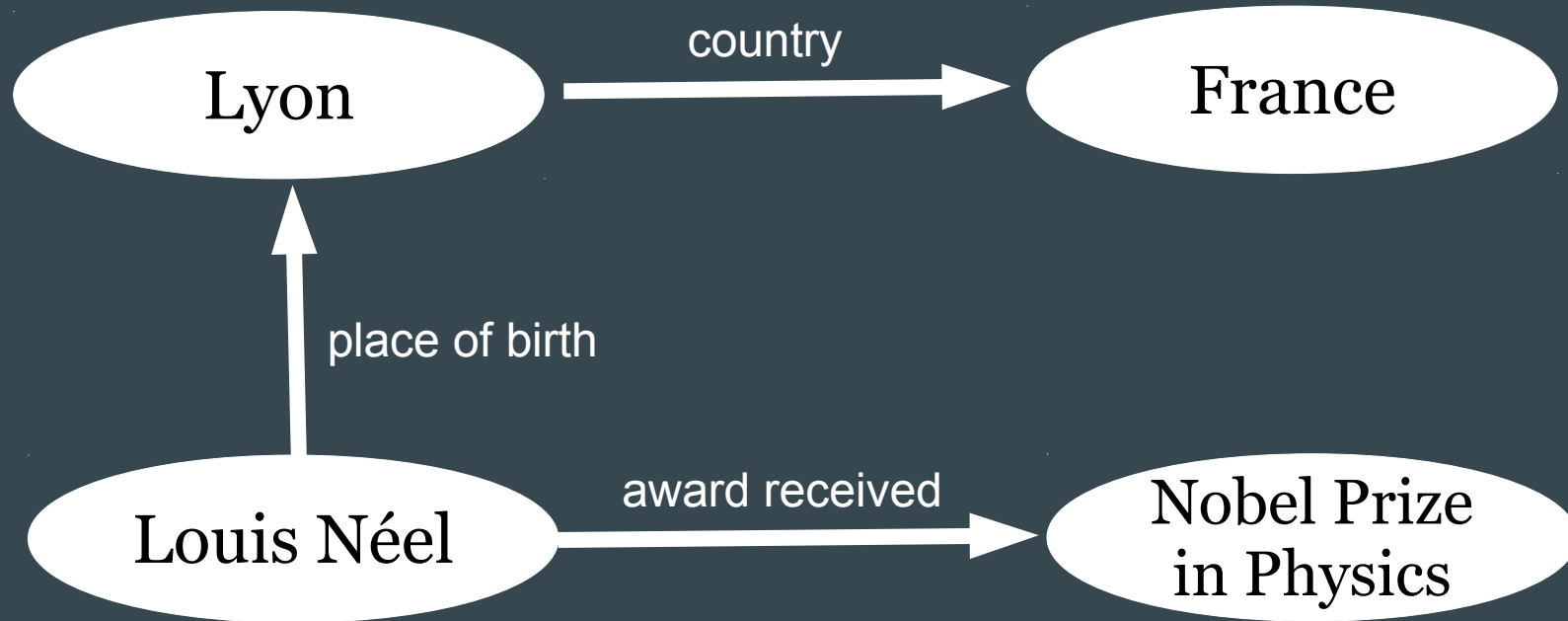
Contents [\[hide\]](#)

- 1 [Wikipedia and the Semantic Web - The Missing Links](#)
 - 1.1 [Introduction](#)
 - 1.2 [A jump start introduction to semantic technologies](#)
 - 1.3 [Design](#)
 - 1.4 [Usability aspects](#)
 - 1.5 [Implementation, performance and scalability](#)
 - 1.6 [Additional features](#)
 - 1.7 [Implementation plan](#)
 - 1.8 [Applications](#)
 - 1.9 [Related approaches](#)
 - 1.10 [Summary and conclusion](#)
 - 1.11 [Acknowledgements](#)
 - 1.12 [Bibliography](#)

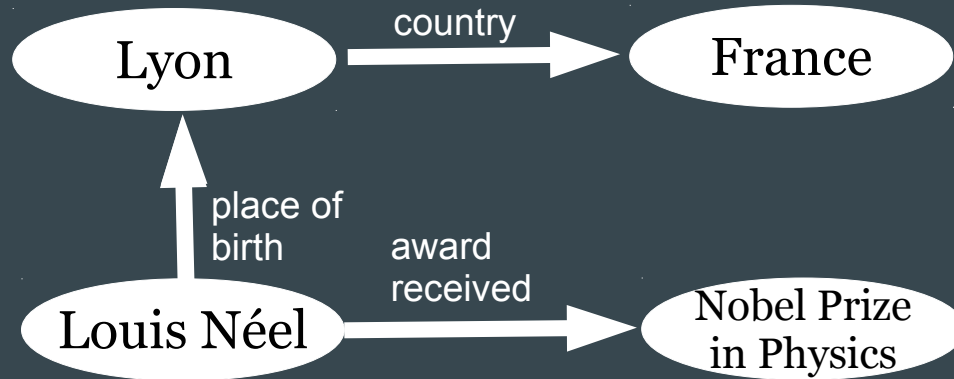
A Simple Idea (2005): “Let’s annotate Wikipedia links!”



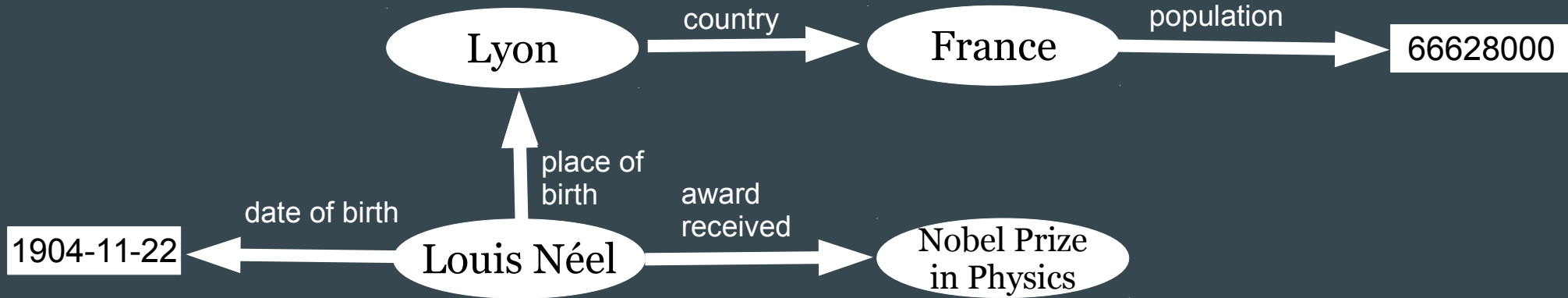
Semantic MediaWiki (2005): From Links to Graphs



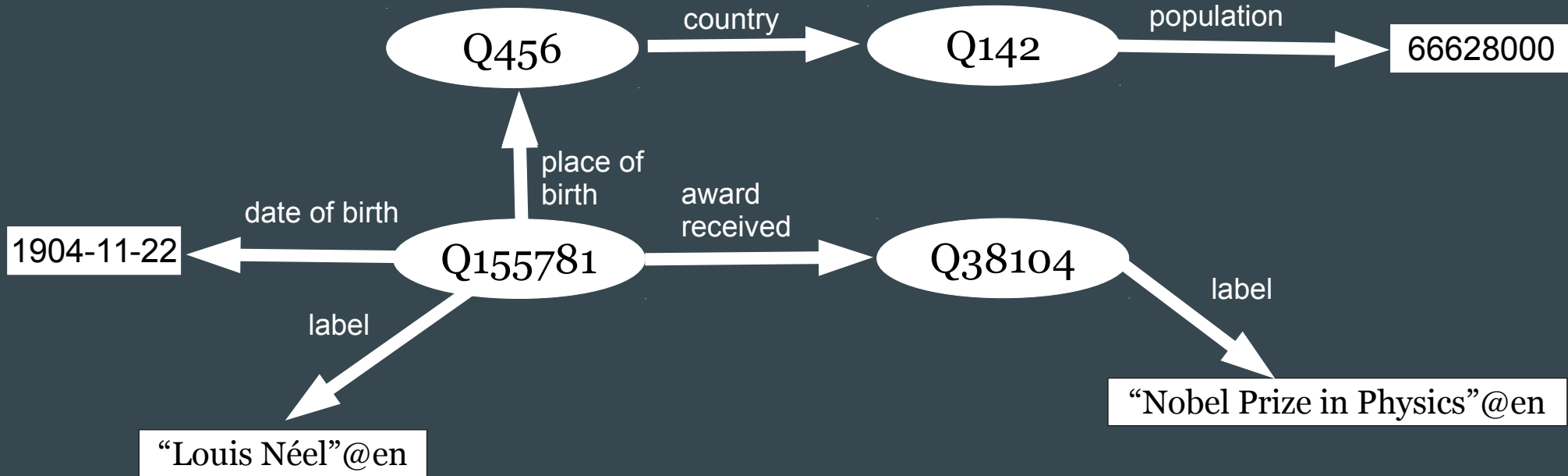
Links are not Enough: Adding Datatypes



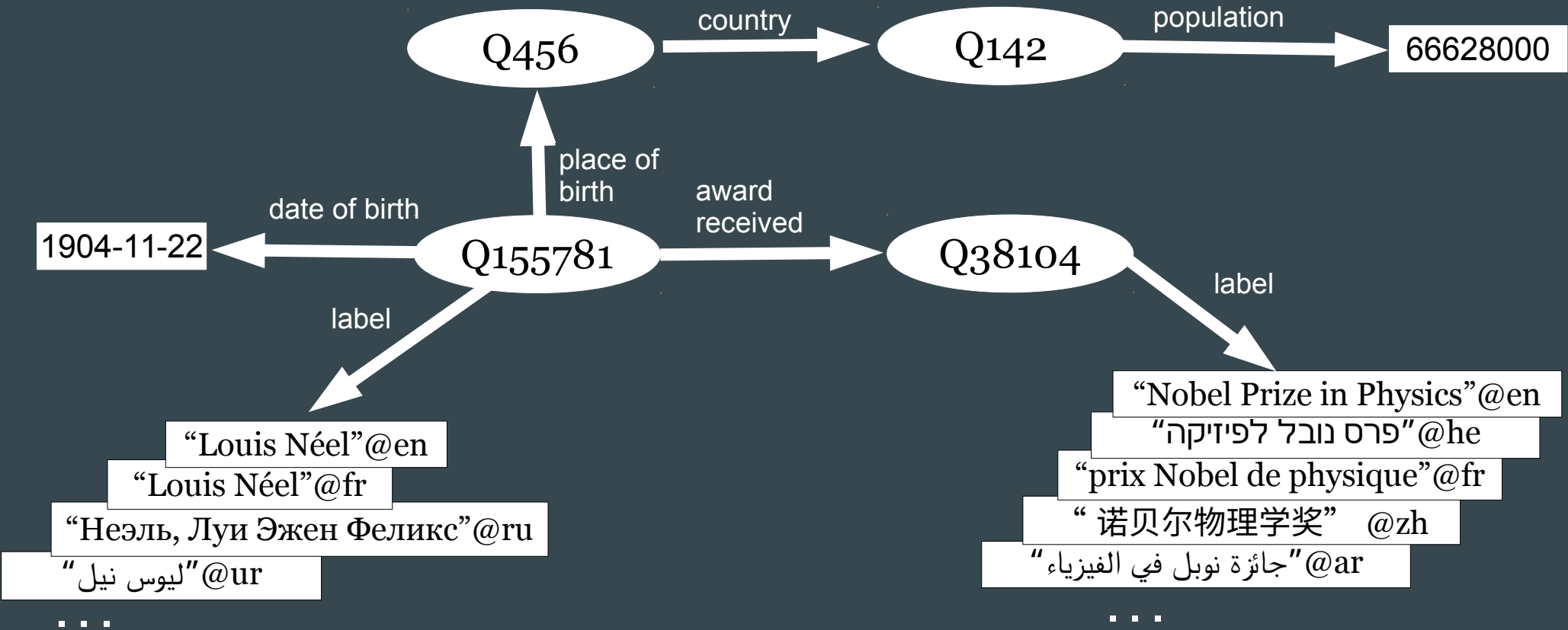
Links are not Enough: Adding Datatypes



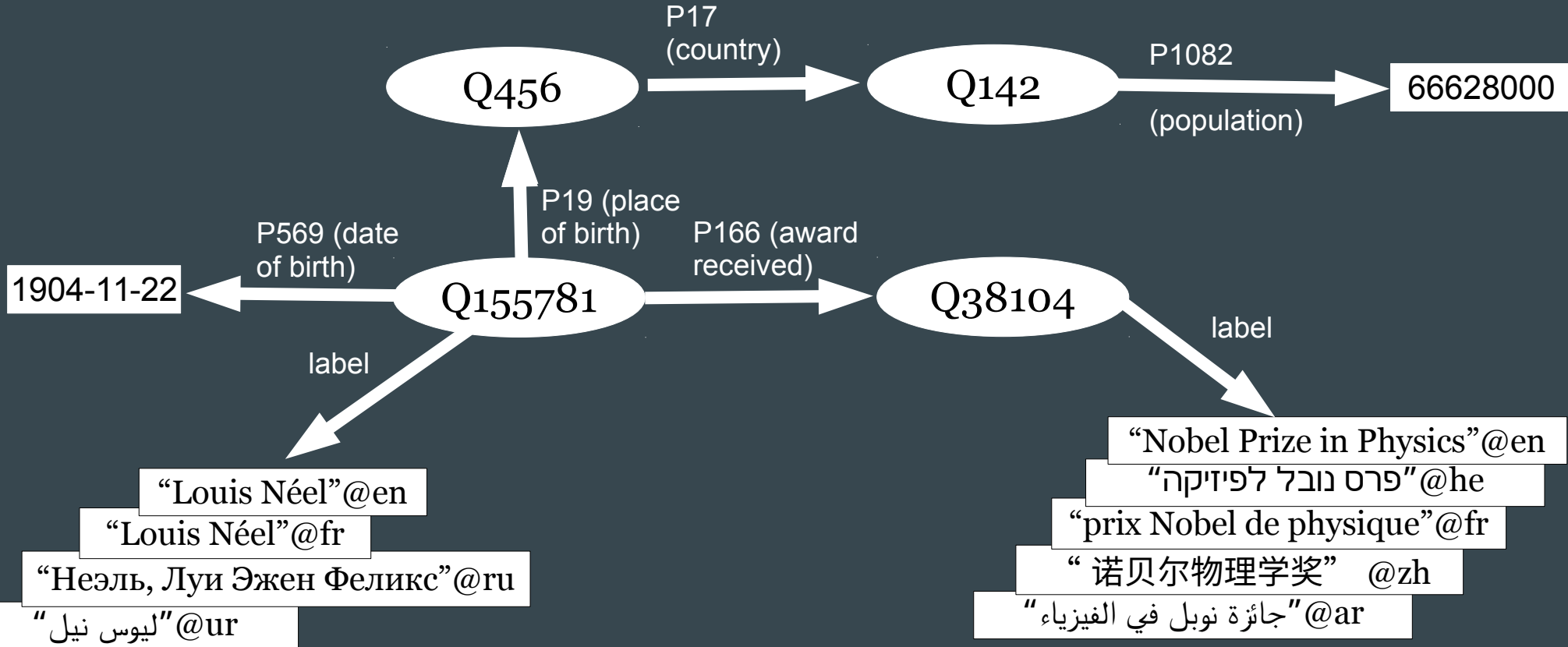
Wikidata: One Graph for Many Languages



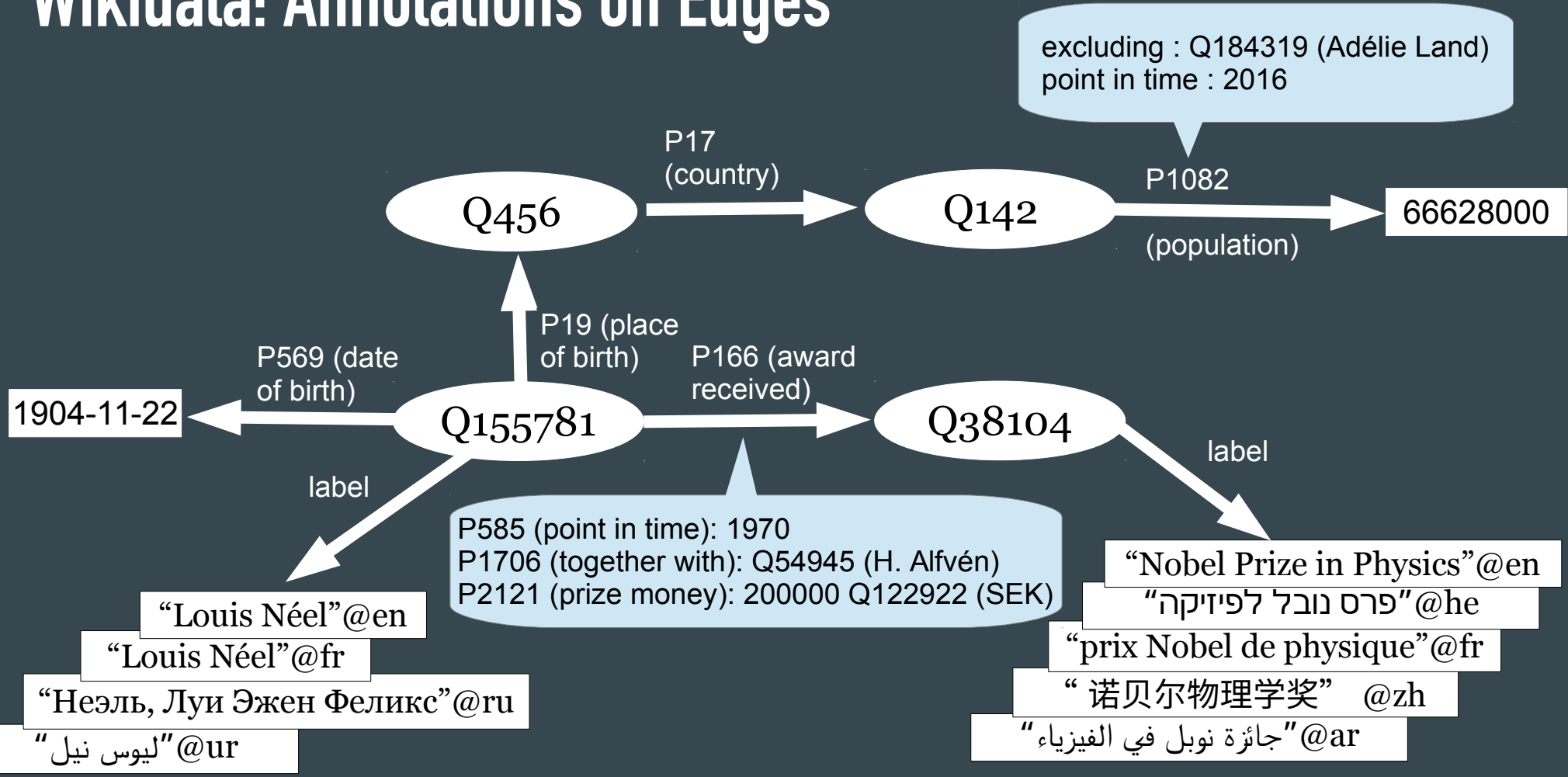
Wikidata: One Graph for Many Languages



Wikidata: One Graph for Many Languages



Wikidata: Annotations on Edges



A Not-So-Simple Idea (2012): Wikidata

Louis Néel (Q155781)

French physicist

Louis Neel | Louis Eugène Felix Néel

award received



Nobel Prize in Physics

edit

point in time

1970

together with

[Hannes Alfvén](#)

prize money

200,000 Swedish krona

▼ 2 references

copy

reference URL

http://www.nobelprize.org/nobel_prizes/physics/laureates

Wikidata in April 2018

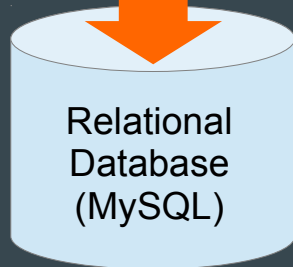
- ◆ >400M statements on >45M entities
- ◆ >60M links to Wikipedia articles
- ◆ >200M labels and aliases
- ◆ >1,200M disambiguating descriptions
- ◆ >200K registered contributors



WIKIDATA

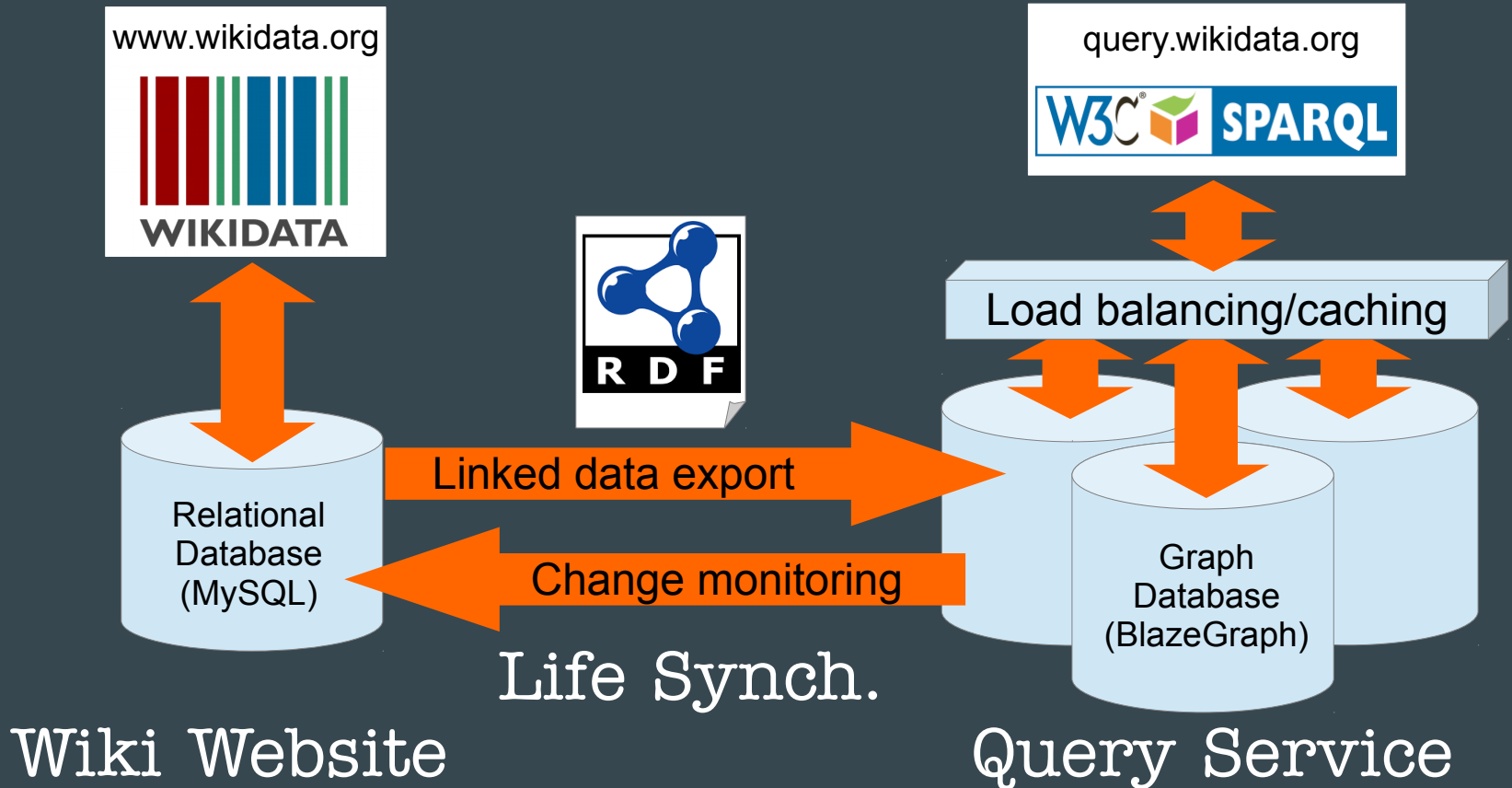
“How to query Wikidata?”

The Wikidata Query Service

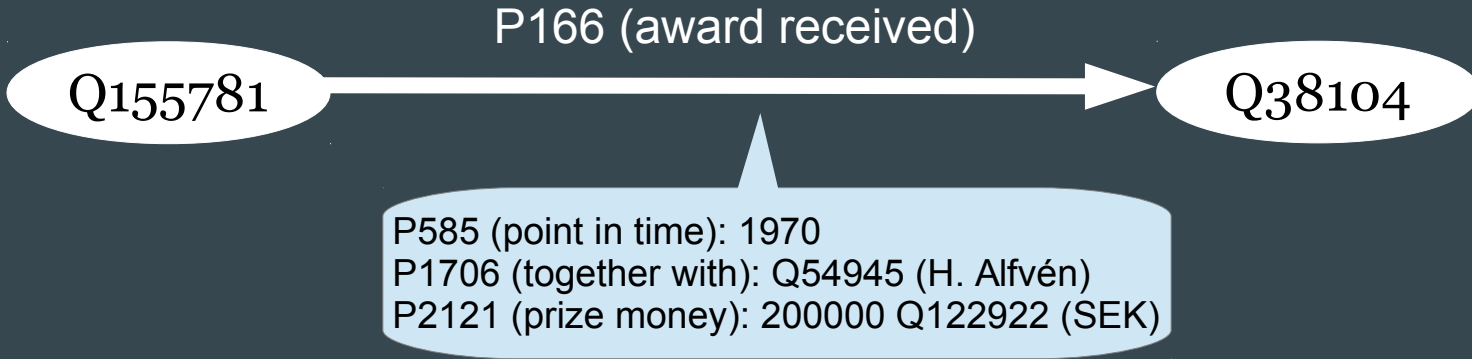


Wiki Website

The Wikidata Query Service



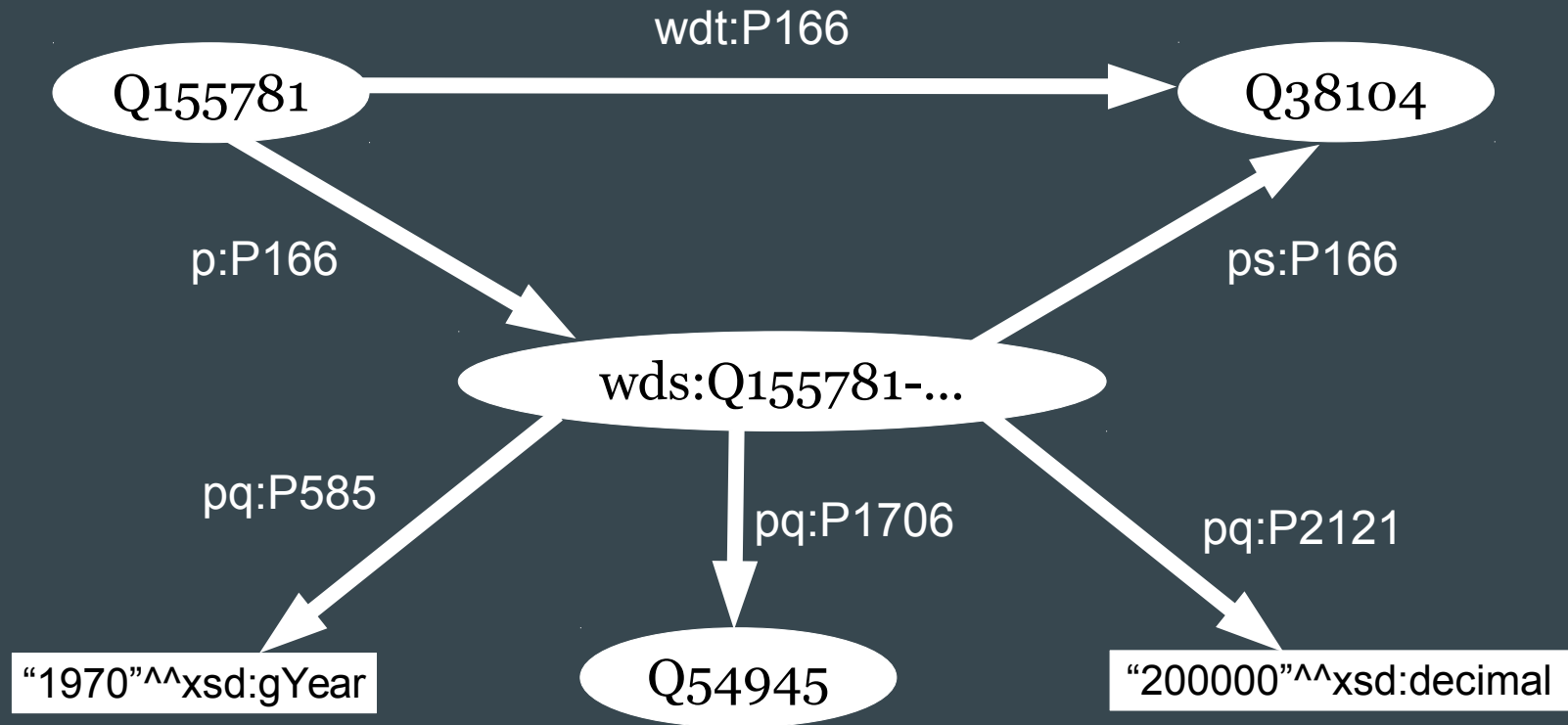
From Wikidata (rich graphs) to RDF (plain graphs)



From Wikidata (rich graphs) to RDF (plain graphs)



From Wikidata (rich graphs) to RDF (plain graphs)



From Wikidata (rich graphs) to RDF (plain graphs)

- ◆ Statements get own objects in graph
- ◆ Some simple statements also stored directly
- ◆ Each Wikidata property becomes many RDF properties
- ◆ Complex values get own objects too (not shown)

Wikidata RDF Exports

- ◆ Weekly full dumps
 - ◆ Currently 4.9 billion triples (32 GBit Turtle compressed)
 - ◆ At <https://dumps.wikimedia.org/wikidatawiki/entities/>
- ◆ Linked Data Exports
 - ◆ Live data in many formats
 - ◆ E.g., <http://www.wikidata.org/wiki/Special:EntityData/Q42.nt>

Wikidata SPARQL Query Service

- ◆ Official query service since mid 2015
 - ◆ User interface at <https://query.wikidata.org/>
- ◆ All the data (4.9B triples), live (latency<60s)
- ◆ No limits (well, almost):
 - ◆ 60sec timeout
 - ◆ No limit on result size (!)
 - ◆ No limit on query numbers per IP
 - ◆ Clients might be paused after too many parallel requests

A simple SPARQL query

Wikidata Query Examples Help Tools English

Query Helper

+ Filter

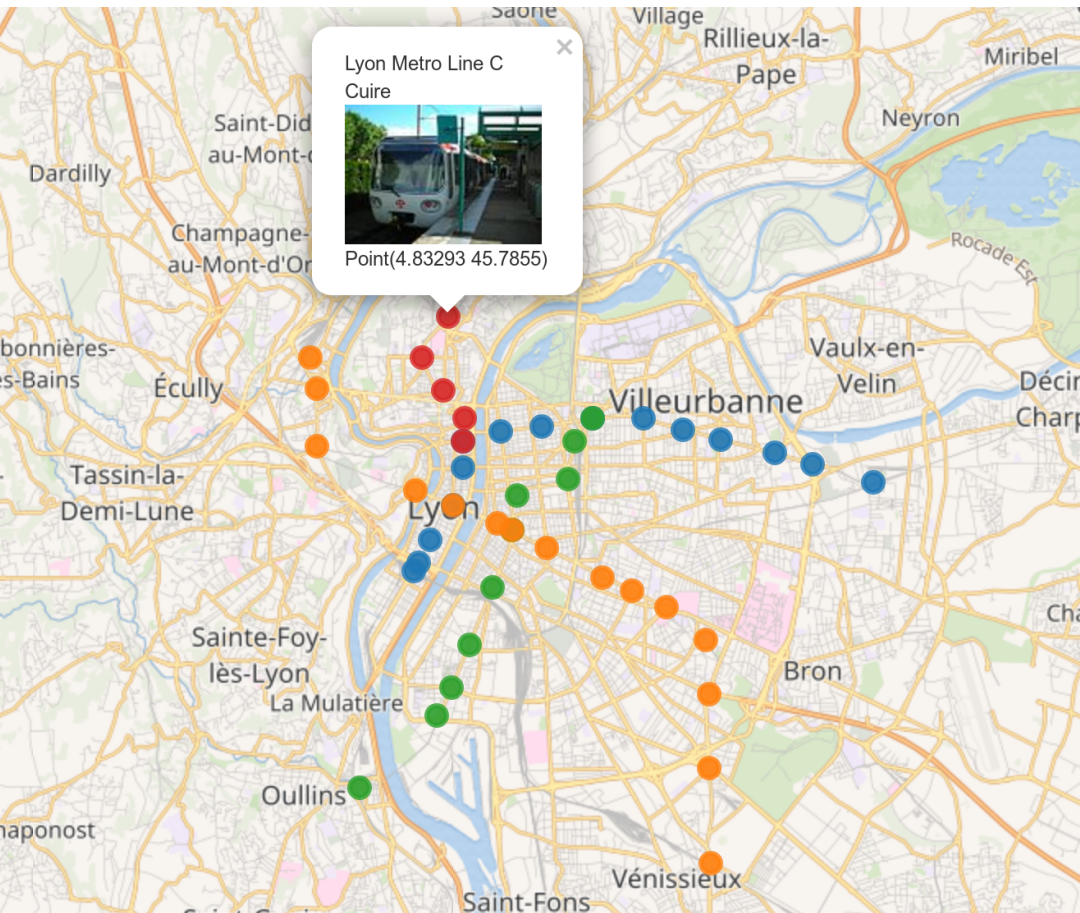
+ Show

-
-
-

Limit

```
1 #defaultView:Map{"layer": "?lineLabel"}
2 SELECT ?stationLabel ?lineLabel ?coord ?image
3 WHERE {
4     ?line wdt:P361 wd:Q1552 .
5     ?station wdt:P81 ?line;
6             wdt:P625 ?coord .
7     OPTIONAL {?station wdt:P18 ?image}
8     SERVICE wikibase:label {
9         bd:serviceParam wikibase:language "en"
10    }
11 }
```

A simple SPARQL query



English

```
1 #defaultView:Map{"layer":"?lineLabel"}
2 SELECT ?stationLabel ?lineLabel ?coord ?image
3 WHERE {
4   ?line wdt:P361 wd:Q1552 .
5   ?station wdt:P81 ?line;
6           wdt:P625 ?coord .
7   OPTIONAL {?station wdt:P18 ?image}
8   SERVICE wikibase:label {
9     bd:serviceParam wikibase:language "en"
10  }
11 }
```

A not-so-simple SPARQL query

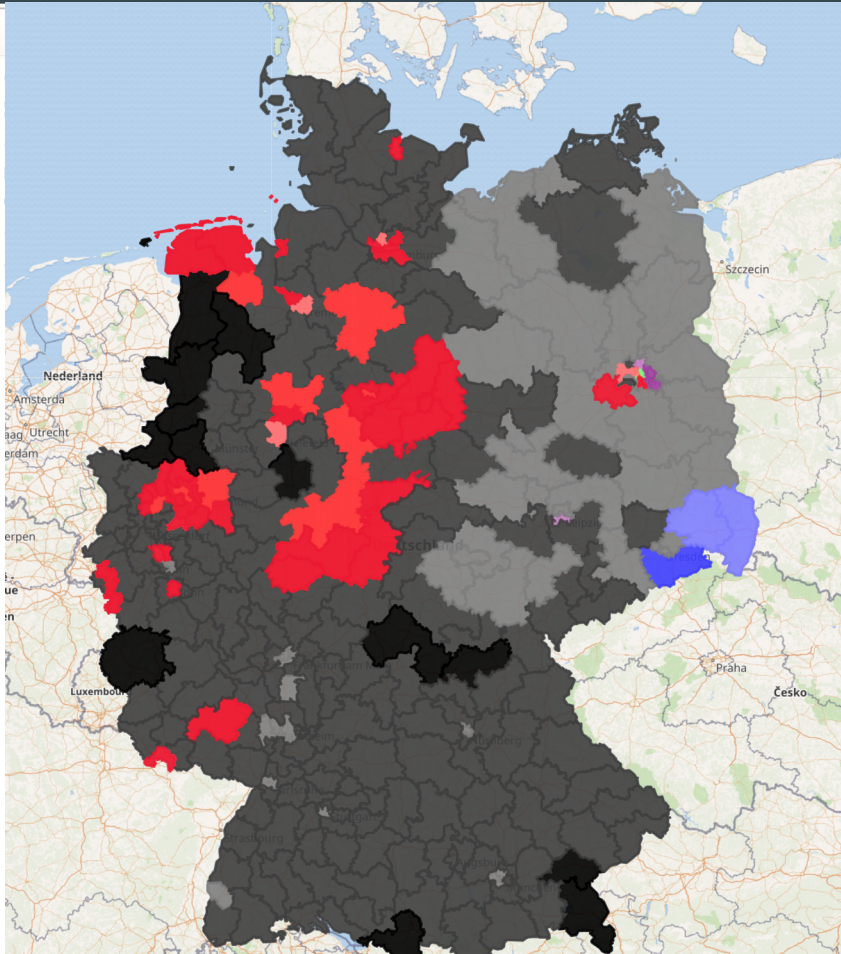
Wikidata Query Examples Help Tools English

Query Helper x

district	successful candidate	._b5
._b5	successful candidate	mdb
._b5	parliamentary term	19th German Bundestag
._b5	represents	party
._b5	votes received	._b4
._b4	http://wikiba.se/ontology#quantityAmount	votesPercentage
._b4	http://wikiba.se/ontology#quantityUnit	percentage
district	catalog code	._b6
._b6	catalog code	districtNumberString
._b6	catalog	list of constituencies for the election to the German Bundestag 2017
mdb	position held	._b7
._b7	position held	member of the German Bundestag
._b7	parliamentary term	19th German Bundestag
._b7	electoral district	district
._b7	parliamentary group	party

```
1 #defaultView:Map
2 # constituencies for the election to the German Bundestag 2017, with winning candidate and party
3 SELECT ?district ?districtLabel ?districtNumber ?mdb ?mdbLabel ?party ?partyLabelCONF (?partyLabel AS ?layer) ?votesPercentage ?rgb ?shape :
4 # find districts with shape
5 ?district wdt:P3896 ?shape;
6 # successful candidate for 19th German Bundestag with party and % votes
7 p:P991 [
8 ps:P991 ?mdb;
9 pq:P2937 wd:Q30579723;
10 pq:P1268 ?party;
11 pqv:P1111 [ wikibase:quantityAmount ?votesPercentage; wikibase:quantityUnit wd:Q11229 ]
12 ];
13 # district number in 2017 Bundestag constituencies
14 p:P528 [
15 ps:P528 ?districtNumberString;
16 pq:P972 wd:Q26971257
17 ].
18 # turn string district number into integer
19 BIND(xsd:integer(?districtNumberString) AS ?districtNumber)
20 # sanity check
21 OPTIONAL {
22 ?mdb p:P39 [
23 ps:P39 wd:Q1939555;
24 pq:P2937 wd:Q30579723;
25 pq:P768 ?district;
26 pq:P4100 ?party
27 ].
28 BIND(true AS ?sanityCheckMdb)
29 }
30 # find original color of party
31 ?party wdt:P462/?wdt:P465 ?rgbOriginal.
32 # fade color depending on % votes, knowing that the original colors are only composed of FF, 80, 00: shift 80 to A0 or C0, and 00 to 40 or
33 # (using separate calls to replace R, G, and B components so that the replacements are aligned to them)
34 BIND(IF(?votesPercentage >= (100/2),
35 ?rgbOriginal,
36 IF(?votesPercentage >= (100/3),
37 REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(?rgbOriginal, "80(..)", "A0$1$2"), "(..)80(..)", "$1A0$2"), "(..)(..)80", "
38 REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(?rgbOriginal, "80(..)", "C0$1$2"), "(..)80(..)", "$1C0$2"), "(..)(..)80", "
39 )
40 ) AS ?rgb)
41 SERVICE wikibase:label {
42 bd:serviceParam wikibase:language "[AUTO_LANGUAGE],en".
43 ?district rdfs:label ?districtLabel.
44 ?party rdfs:label ?partyLabel.
45 ?mdb rdfs:label ?mdbLabel.
46 }
47 }
48 ORDER BY ?districtNumber
```

A not-so-simple SPARQL query



```
1 #defaultView:Map
2 # constituencies for the election to the German Bundestag 2017, with winning candidate and party
3 SELECT ?district ?districtLabel ?districtNumber ?mdb ?mdbLabel ?party ?partyLabelCONF (?partyLabel AS ?layer) ?votesPercentage ?rgb ?shape :
4 # find districts with shape
5 ?district wdt:P3896 ?shape;
6 # successful candidate for 19th German Bundestag with party and % votes
7 p:P991 [
8   ps:P991 ?mdb;
9   pq:P2937 wd:Q30579723;
10  pq:P1268 ?party;
11  pqv:P1111 [ wikibase:quantityAmount ?votesPercentage; wikibase:quantityUnit wd:Q11229 ]
12 ];
13 # district number in 2017 Bundestag constituencies
14 p:P528 [
15   ps:P528 ?districtNumberString;
16   pq:P972 wd:Q26971257
17 ];
18 # turn string district number into integer
19 BIND(xsd:integer(?districtNumberString) AS ?districtNumber)
20 # sanity check
21 OPTIONAL {
22   ?mdb p:P39 [
23     ps:P39 wd:Q1939555;
24     pq:P2937 wd:Q30579723;
25     pq:P768 ?district;
26     pq:P4100 ?party
27   ].
28   BIND(true AS ?sanityCheckMdb)
29 }
30 # find original color of party
31 ?party wdt:P462/?wdt:P465 ?rgbOriginal.
32 # fade color depending on % votes, knowing that the original colors are only composed of FF, 80, 00: shift 80 to A0 or C0, and 00 to 40 or
33 # (using separate calls to replace R, G, and B components so that the replacements are aligned to them)
34 BIND(IF(?votesPercentage >= (100/2),
35   ?rgbOriginal,
36   IF(?votesPercentage >= (100/3),
37     REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(?rgbOriginal, "80(..)", "A0$1$2"), "(..)80(..)", "$1A0$2"), "(..)(..)80", "
38     REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(?rgbOriginal, "80(..)", "C0$1$2"), "(..)80(..)", "$1C0$2"), "(..)(..)80", "
39   )
40 ) AS ?rgb)
41 SERVICE wikibase:label {
42   bd:serviceParam wikibase:language "[AUTO_LANGUAGE],en".
43   ?district rdfs:label ?districtLabel.
44   ?party rdfs:label ?partyLabel.
45   ?mdb rdfs:label ?mdbLabel.
46 }
47 }
48 ORDER BY ?districtNumber
```


An advanced SPARQL query

Wikidata Query Examples Help Tools English

Query Helper

film	instance of	Volver a Empezar	
	any		
	subclass of		
headOfGovernment	instance of	human	
headOfGovernment	position held	_:b2	
+ Filter	_:b2	position held	position
	_:b2	start time	startTime
position	subclass of	head of government	
http://www.bigdata.com/queryHints#Prior	http://www.bigdata.com/queryHints#runLast	"false"^^http://www.w3.org/2001/XMLSchema#boolean	
film	publication date	publicationDate	
film	cast member	headOfGovernmentStatement	

```
1 # films starring more than one future head of government
2 SELECT ?film ?filmLabel ?filmDescription (COUNT(DISTINCT ?headOfGovernmentLabel)
3 ?film wdt:P31/wdt:P279* wd:Q11424;
4 wdt:P577 ?publicationDate;
5 p:P161 ?headOfGovernmentStatement.
6 ?headOfGovernmentStatement ps:P161 ?headOfGovernment.
7 OPTIONAL { ?headOfGovernmentStatement pq:P453 ?character. ?character rdfs:label
8 ?headOfGovernment wdt:P31 wd:Q5;
9 p:P39 [
10 ps:P39 ?position;
11 pq:P580 ?startTime
12 ]}.
13 ?position wdt:P279+ wd:Q2285706.
14 FILTER(?startTime > ?publicationDate) # *future* head of government
15 FILTER NOT EXISTS {
16 ?headOfGovernment p:P39 [
17 ps:P39 ?otherPosition;
18 pq:P580 ?otherStartTime
19 ]}.
20 ?otherPosition wdt:P279+ wd:Q2285706.
21 FILTER(?otherStartTime < ?publicationDate) # not already a head of government
22 }
23 SERVICE wikibase:label {
24 bd:serviceParam wikibase:language "[AUTO_LANGUAGE],en".
25 ?film rdfs:label ?filmLabel;
26 schema:description ?filmDescription.
27 ?headOfGovernment rdfs:label ?headOfGovernmentLabel.
28 ?position rdfs:label ?positionLabel.
29 } hint:Prior hint:runLast false.
30 BIND(IF(BOUND(?characterLabel), CONCAT(?characterLabel, " / " @en, ?positionLabel
31 })
32 GROUP BY ?film ?filmLabel ?filmDescription
33 HAVING(?count > 1)
```

You expect normal people to use SPARQL?!

- ◆ If they want ... it's really not that difficult
 - ◆ Extensive online documentation
 - ◆ Over 300 example queries
 - ◆ Tutorials and workshops at community events
- ◆ But SPARQL is often hidden from users
 - ◆ Embedded results on Web pages (incl. Wikipedia)
 - ◆ Mobile apps and online apps
 - ◆ Crowdsourcing platforms

Wikidata:Request a query

Shortcut: [WD:RAQ](#)

This is a page where [SPARQL 1.1 Query Language \(Q32146616\)](#) queries can be requested. Please provide feedback if a query is written for you.

For sample queries, see [Examples](#). Property talk pages include also summary queries for these.

For help writing your own queries, or other questions *about* queries, see [Wikidata talk:SPARQL query service/queries](#).

Help resources about [Wikidata Query Service \(Q20950365\)](#) and SPARQL: [Wikidata:SPARQL query service/Wikidata Query Help](#) and [Category:SPARQL](#).

Contents [\[hide\]](#)

- [1 Slide show with images](#)
- [2 Retrieve property if available](#)
- [3 Surname lookup](#)
- [4 What's in Wikipedia lists?](#)
- [5 Properties missing a label or description in a given language](#)
- [6 P: Properties for a set of items](#)
- [7 About population](#)
- [8 SPARQL for Q5 externalid statistics](#)
- [9 Who held what position in the year 420 ?](#)



Fishing in the [Wikidata river](#) ✉
requires both an idea where to look for fish and a suitable fishing method. If you have the former, this page can help you find the latter.

Current Usage

- ◆ SPARQL is widely used
 - ◆ >100M requests per month (3.8M per day) in 2018
- ◆ Excellent availability and performance
 - 50% of queries answered in <40ms (95% in <440ms; 99% in <40s)
 - Less than 0.05% of queries time out
 - Service has never been down so far
- ◆ All software/customisations free & open source
 - See <https://github.com/wikimedia/wikidata-query-rdf>



WIKIDATA

“What can we learn from
all these SPARQL queries?”

SPARQL Queries Are Interesting

- ◆ Which data is actually asked for?
- ◆ Which SPARQL features are most important?
- ◆ Who is using SPARQL through which tools?

We have analysed complete Wikidata SPARQL query logs (Wikimedia Research Collaboration)

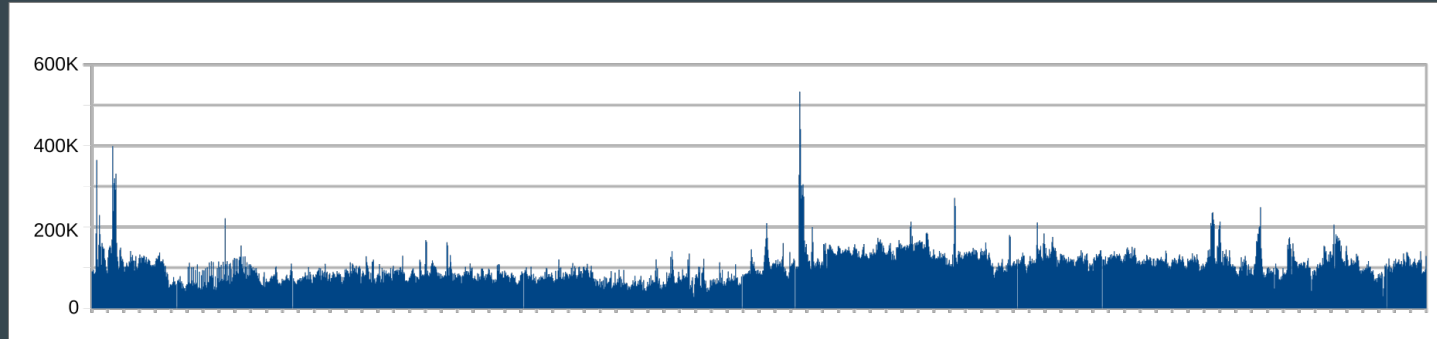
[Bielefeldt et al., “Linked Data on The Web” @ WWW 2018]

Analysing SPARQL logs: The Bot Problem

Analysing SPARQL logs: The Bot Problem

- Query traffic is **ruled** by a few bots

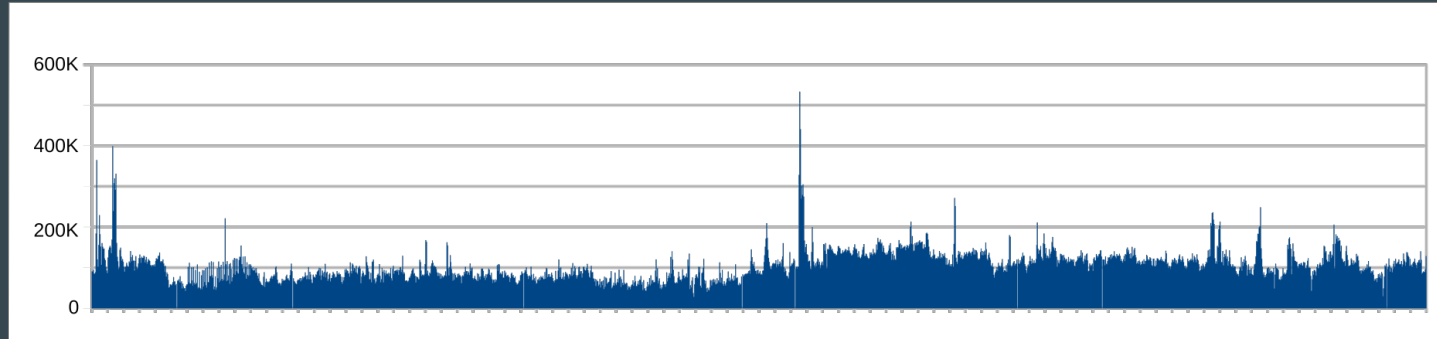
Fig.: Wikidata SPARQL traffic Jun-Sep 2017



Analysing SPARQL logs: The Bot Problem

- Query traffic is **ruled** by a few bots

Fig.: Wikidata SPARQL traffic Jun-Sep 2017

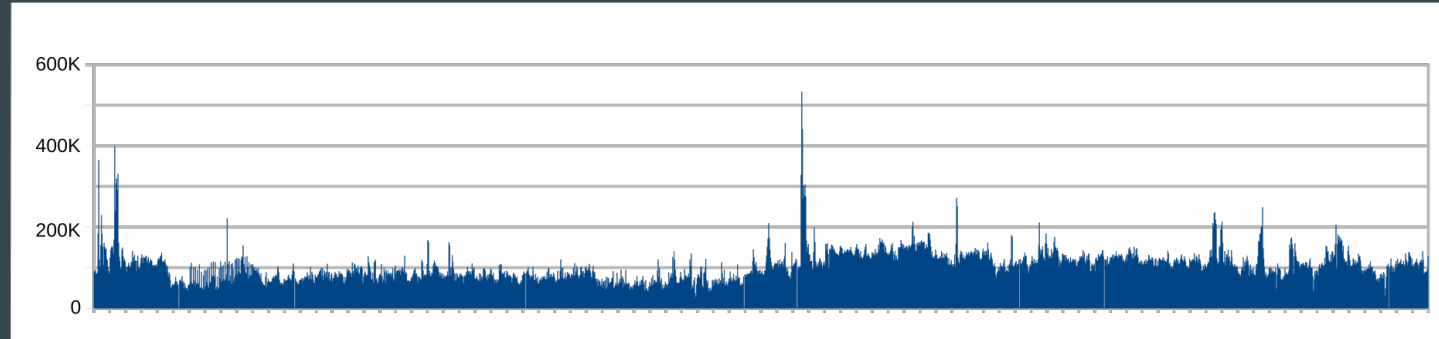


- 41% of all Wikidata query traffic from June – September 2017 caused by one super-power user (Magnus Manske)

Analysing SPARQL logs: The Bot Problem

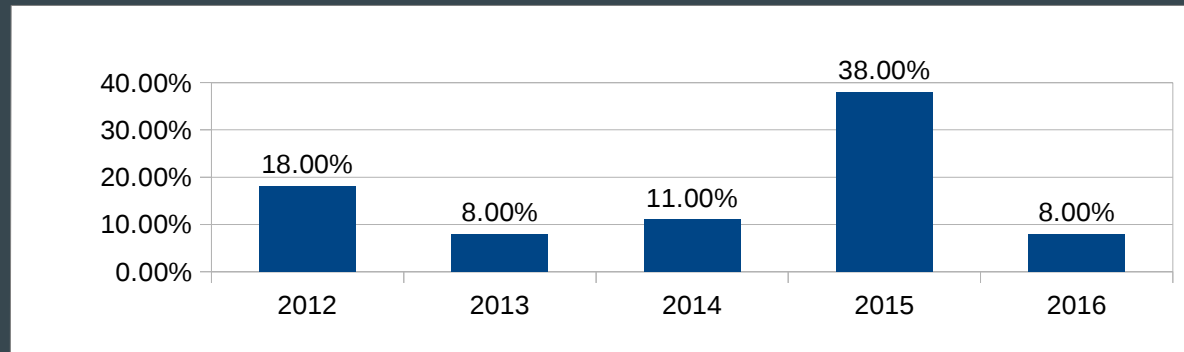
- Query traffic is **ruled** by a few bots

Fig.: Wikidata SPARQL traffic Jun-Sep 2017



- 41% of all Wikidata query traffic from June – September 2017 caused by one super-power user (Magnus Manske)
- The effect does **not** average out, and it affects other sites too

Fig.: Usage of DISTINCT on DBpedia [Bonifati et al. 2017]



Analysing SPARQL logs: The Bot Problem

- Query traffic is **ruled** by a few bots

Fig.: Wikidata SPARQL traffic Jun-Sep 2017



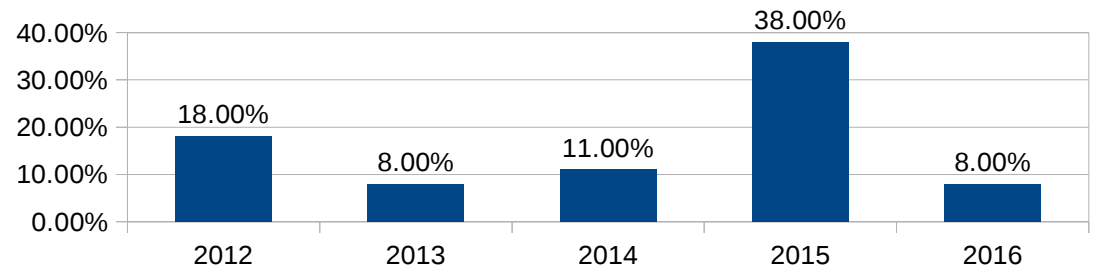
No trends!
No predictability!
No insights!

- 41%

June – September 2017
(Magnus Manske)

sites too

Fig.: Usage of DISTINCT on DBpedia [Bonifati et al. 2017]



Are SPARQL queries interesting after all?

- ◆ Observation: Robotic traffic dominates
 - ◆ May not represent any real interest
 - ◆ Governed by very few sources
 - ◆ Random changes – not uniform on any observed scale

Are SPARQL queries interesting after all?

- ◆ Observation: Robotic traffic dominates
 - ◆ May not represent any real interest
 - ◆ Governed by very few sources
 - ◆ Random changes – not uniform on any observed scale
- ◆ Hypothesis: Organic traffic also exists
 - ◆ Representing human information need during some interaction
 - ◆ Composed of many diverse sources
 - ◆ Continuous change over months

Note: “Organic” ≠ “hand-written SPARQL” (user apps might use SPARQL to get user-requested data without users actually writing queries)

Extracting organic traffic

- ♦ Main signal: User Agents
 - ♦ Assumption: organic traffic generally from browser-like agents

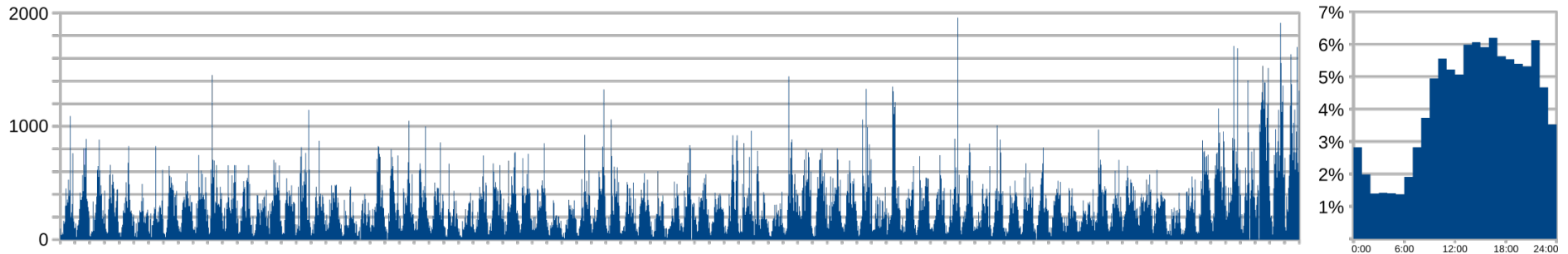
Extracting organic traffic

- ◆ Main signal: User Agents
 - ◆ Assumption: organic traffic generally from browser-like agents
- ◆ 2nd signal: query comments
 - ◆ Some browser-based tools mark queries using comments
- ◆ 3rd signal: activity spikes
 - ◆ Group queries by query pattern (following [Raghuveer, USEWOD'12])
 - ◆ Find agent-pattern pairs that spike (>2K requests/month)
 - ◆ Manually inspect these queries to decide if organic or robotic
 - About 300 further browser-based sources classified “robotic”

Results: Organic component

- ♦ Jun–Sep 2017: 658,890 queries (<0.5%)

Temporal distribution of organic queries (12 weeks / time of day)



- ♦ More triples
organic 17%: 1, 97%: ≤ 11 vs. robotic 57%: 1, 96%: ≤ 7
- ♦ More varied (vocabulary, SPARQL features)

Insights on Wikidata Usage

- ◆ Robotic traffic:
 - ◆ Mainly information integration bots (comparing database contents)
 - ◆ Potentially also selective data download (spider-like)
 - ◆ Most queries from a few dominant bots (>60% from top-three bots)
- ◆ Organic traffic:
 - ◆ Data browsers (often general-purpose)
 - ◆ Mobile apps (often topical)
 - ◆ Most queries from of unidentified “small” sources
- ◆ Reified statements in 4%–10% of queries



WIKIDATA

What's Next?

More data

- ♦ Wikidata: >45M items with >400M statements

More data

- ♦ Wikidata: >45M items with >400M statements
 - ♦ OSM: >4B nodes, >230M buildings, >10M trees
 - ♦ WDC: >9.5B entities, >38B RDF triples
- ♦ Why don't we just import everything?!

More data

- ♦ Wikidata: >45M items with >400M statements
 - ♦ OSM: >4B nodes, >230M buildings, >10M trees
 - ♦ WDC: >9.5B entities, >38B RDF triples
- ♦ Why don't we just import everything?!
 - ♦ Notability? Well, sometimes ...
 - ♦ Community support! Who will maintain this?

More data: current efforts

- ◆ Data donation guidelines
- ◆ Wikidata aligns with >2500 databases and catalogues
- ◆ Supervised data alignment with crowdsourcing (Mix'n'Match)

AcademiaNet	Database for excellent female scientists	99%
Austrian Parliament ID	Austrian Parliament's "Who's Who" database	99%
International World Games Ass	Sportspeople	24% 59%
botanist author abbreviation	standard form (official abbreviation) of a personal name for use in an :	99%
Mactutor	identifer of the person's biography in the MacTutor History of Mathe	98%
South Australian Football Hall	Australian rules football players	72%
AIBL members	Membres de l'Académie des Inscriptions et Belles Lettres (AIBL)	70%
Lotsawa House Tibetan author	Tibetan authors in the Lotsawa House library	40%
parliament.uk	UK MP or Peer's biography	95%
EPHE	identifier of a researcher on the online prosopographical dictionary of	74% 16%
Sport Australia Hall of Fame	Sportspeople linked to Australia	73% 10%
North Carolina Sports Hall of F	Sportspeople linked to North Carolina	54% 16%

More data: current efforts

Soccerdonna

Soccerdonna website female association football player db

Markus Krötzsch

Load next entry on empty search results

Casey Short

'player, born 23.08.1990 at Naperville, Illinois plays '



The screenshot shows the Soccerdonna website profile for Casey Short. The header includes the Soccerdonna logo and navigation tabs for 'STARTSEITE', 'WETTBEWERBE', '1. BUNDESLIGA', and 'EUROPA'. The profile section features a photo of Casey Short, her name '6 Casey Short', and her team 'Chicago Red Stars, NWSL (Vereinigte Staaten)'. It also lists her current national player status as 'Vereinigte Staaten U23' and her manager 'Katie Naughton'. A 'Facebook' widget shows the team's page with 2.6K likes. Below the profile is a table of performance data for the current season.

Wettbewerb	Spiele	Goals	Assists
SheBelieves Cup	2	-	-

The screenshot shows a search interface with 'Casey Short' entered in the search bar. Below the search bar is a 'Find' button. The search results list several entries related to Casey Short, including her Wikidata ID [Q16766251], her profession as a US-American association football player, and several short films she has appeared in or directed, such as 'Alaska – Die raue Eiswelt' and '1912 silent short film from United States of America'.

The screenshot shows the Wikidata page for Casey Short (Q16766251). The page includes a search bar, a star icon for bookmarks, and the description 'American association football player'. Below this, there is a section for 'Statements' which shows that Casey Short is an instance of 'human' and has 0 references.



New kinds of data

- ◆ Coming soon: **lexical data** (dictionary/thesaurus)
 - ◆ Exciting & dangerous
- ◆ Planned: **media (meta-)data** (Wikimedia Commons)
- ◆ Factual knowledge that is not in catalogues?
- ◆ Common sense?

In many cases: technical changes/extensions needed

Quality!

- ◆ Errors, spam, vandalism
- ◆ Global coherency of modelling
- ◆ Sources & alignments
- ◆ Incompleteness
- ◆ Change & data rot

Germany (Q183)

federal parliamentary republic in central-western Europe

FRG | BRD | Bundesrepublik Deutschland | Federal Republic of Germany | de | 

basic form of government



federal parliamentary republic



 edit

Potential issues



conflicts-with constraint

[Help](#) [Discuss](#)

An entity should not have a statement for [basic form of government](#) if it also has a statement for [instance of](#) with value [republic](#).

one-of constraint

[Help](#) [Discuss](#)

The value for [basic form of government](#) should be one of the following:

- [republic](#)
- [constitutional monarchy](#)
- [federal republic](#)
- [representative democracy](#)
- [parliamentary system](#)
- [soviet republic](#)

 copy

 copy



Inferring new knowledge with ontologies

[edit label](#)

Nauru (Q697)

Republic of Nauru | Pleasant Island | Naoero | nr | 

republic in Oceania

head of state

2+28 statements ▼

[Baron Waqa](#) (Nauruan politician) ★ [▶](#)
start time : 2013-06-11

[Sprent Dabwido](#) (president of Nauru) [▶](#)

[Frederick Pitcher](#) (President of Nauru) [▶](#)
start time : 2011-11-10
end time : 2011-11-15

(Proposal) ✓
Source: MARS

[Marcus Stephen](#) (Nauruan sportperson and politician) [▶](#)
start time : 2007-12-19
end time : 2011-11-10

(Proposal) ✓
Source: MARS

[Ludwig Scotty](#) (Nauruan politician, president) [▶](#)
start time : 2004-06-22
end time : 2007-12-19

(Proposal) ✓
Source: MARS

**[Marx & MK, International
Semantic Web Conf. 2017]**

<https://tools.wmflabs.org/sqid/>

[Dorothy H. ...](#) (1947-2022) [▶](#)

(Proposal) ✓

Frederick Pitcher (Q917601)

[edit label](#)

position held

President of Nauru (head of state and government in Nauru) >

start time : 2011-11-10

end time : 2011-11-15

replaces : [Marcus Stephen](#) (Nauruan sportperson and politician)

replaced by : [Sprent Dabwido](#) (president of Nauru)

Nauru (Q697)

[edit label](#)

office held by head of government President of Nauru (head of state and government in Nauru) >

A rule of inference:

```
(?headOfState.position heldP39 = ?headOffice)@?X,  
(?country.office held by head of stateP1906 = ?headOffice)@?Y  
→ (?country.head of stateP35 = ?headOfState)@{start timeP580 = ?X.start timeP580,  
end timeP582 = ?X.end timeP582}
```

[Marx et al., International Joint Conf. On Artif. Intellig. 2017]

Conclusion and Outlook

- ◆ Wikidata is a fascinating, fast-moving project
 - ◆ Large amounts of quality data & much more to come
 - ◆ Data export and analysis services for all needs
 - ◆ Innovation-friendly community
- ◆ Many unsolved questions for research
 - ◆ Quality, provenance, social aspects, performance challenges, data integration, internationalisation, ...

Literature

- Adrian Bielefeldt, Julius Gonsior, Markus Krötzsch: “Practical Linked Data Access via SPARQL: The Case of Wikidata” Proceedings of the WWW2018 Workshop on Linked Data on the Web (LDOW-18), CEUR Workshop
- Fredo Erxleben, Michael Günther, Markus Krötzsch, Julian Mendez, Denny Vrandečić: “Introducing Wikidata to the Linked Data Web” In Proceedings of the 13th International Semantic Web Conference (ISWC 2014)
- Maximilian Marx, Markus Krötzsch: “SQID: Towards Ontological Reasoning for Wikidata” In Proceedings of the ISWC 2017 Posters & Demonstrations Track, CEUR Workshop Proceedings. CEUR-WS.org
- Maximilian Marx, Markus Krötzsch, Veronika Thost: “Logic on MARS: Ontologies for generalised property graphs” Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI'17), 1188-1194, 2017

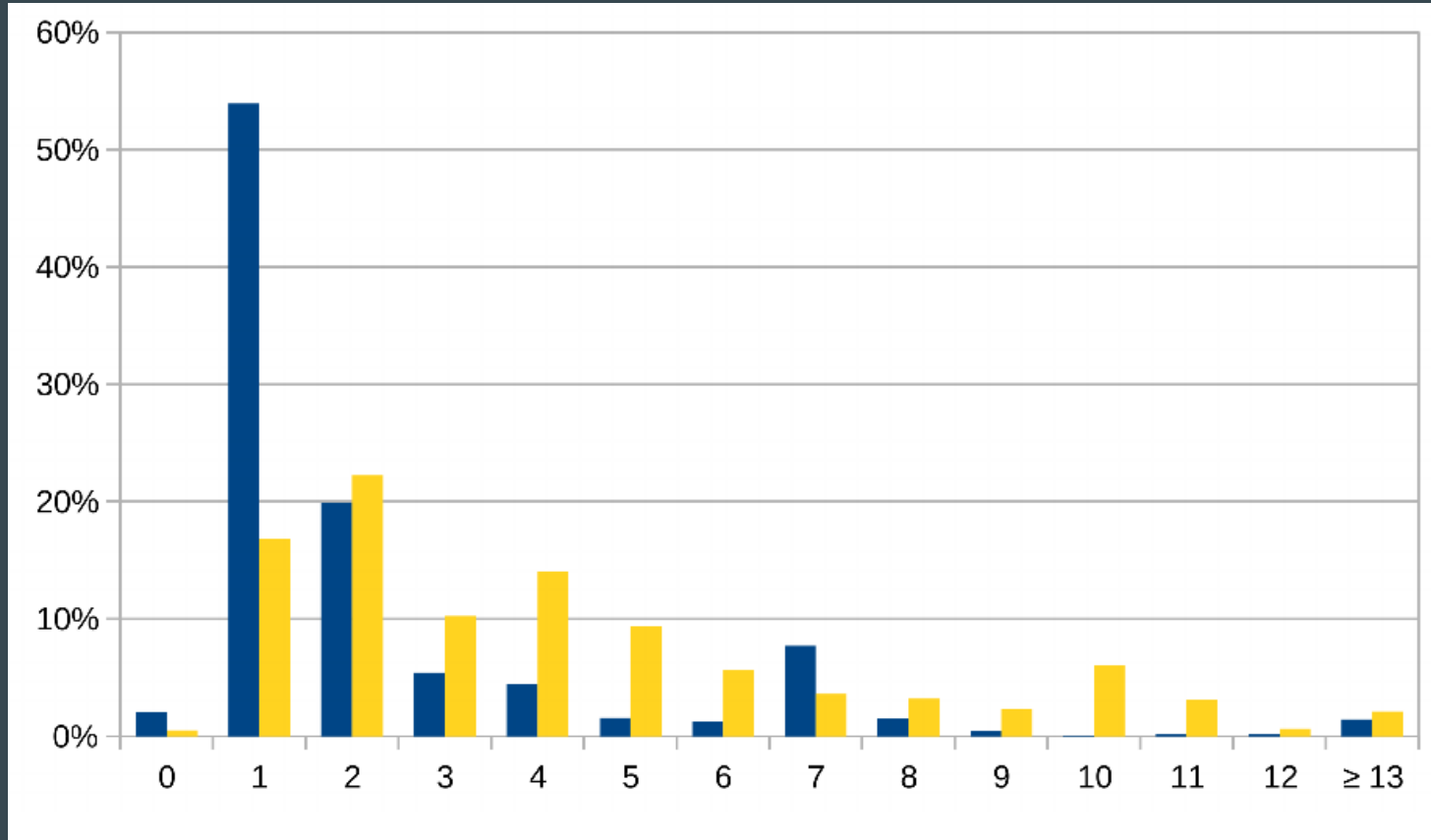
Films with future heads of government

Star in the Dust	1956 film by Charles F. Haas	2	Clint Eastwood, mayor; George Wallace, Governor of Alabama
The Two Who Stole the Moon	1962 Polish film by Jan Batory	2	Jarosław Kaczyński, Prime Minister of Poland; Lech Kaczyński, Mayor of Warsaw
Ragasiya Police 115	1968 film by B. R. Panthulu	2	M. G. Ramachandran, Chief Minister of Tamil Nadu; Jayalalithaa, Chief Minister of Tamil Nadu
Québec : Duplessis et après...	documentary	2	Bernard Landry, Premier of Quebec; René Lévesque, Premier of Quebec
Q3541438	1994 film by Claude Lanzmann	2	Ariel Sharon, Prime Minister of Israel; Ehud Barak, Prime Minister of Israel
Batman & Robin	1997 American superhero film based on the DC Comics character Batman	2	Arnold Schwarzenegger, Mr. Freeze / Governor of California; Jesse Ventura, Governor of Minnesota

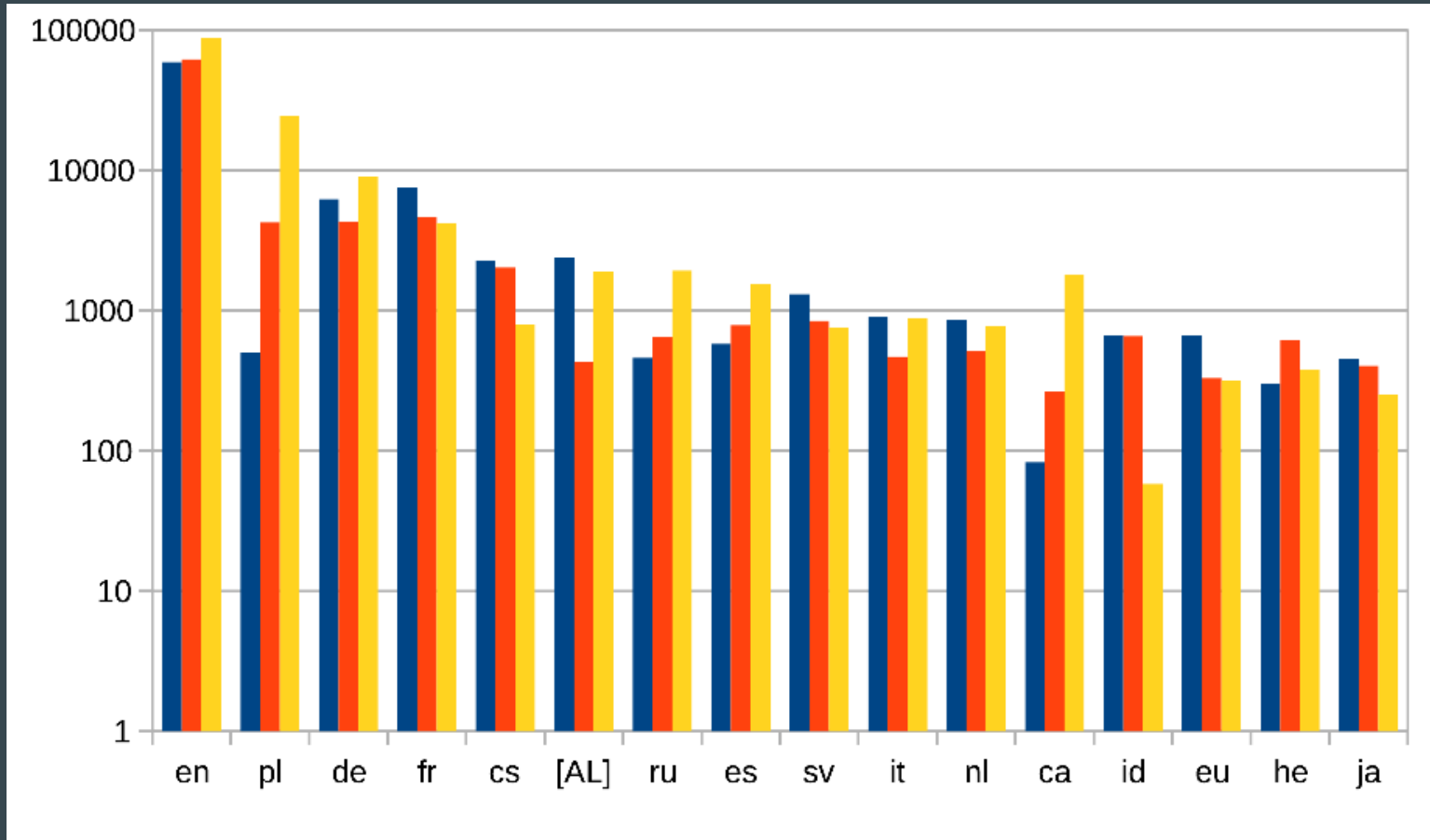
SPARQL Feature Distribution (2017/2018)

	organic						robotic					
	Jun 2017	Jul 2017	Aug 2017	Jan 2018	Feb 2018	Mar 2018	Jun 2017	Jul 2017	Aug 2017	Jan 2018	Feb 2018	Mar 2018
Limit	31.08%	39.55%	46.56%	52.31%	51.23%	36.87%	21.12%	16.86%	17.42%	20.38%	11.47%	15.17%
Distinct	26.50%	31.40%	19.05%	59.30%	60.42%	63.78%	15.84%	5.48%	4.27%	4.32%	7.54%	12.25%
Order By	17.29%	14.75%	13.22%	46.89%	46.99%	34.53%	12.97%	8.01%	6.78%	8.76%	7.68%	17.46%
Offset	0.40%	2.92%	0.37%	0.09%	0.08%	0.06%	7.73%	6.07%	6.29%	0.10%	0.07%	0.10%
Join	87.59%	87.82%	89.76%	82.50%	91.70%	87.02%	88.48%	78.53%	67.41%	73.26%	61.39%	70.19%
Optional	42.36%	46.24%	55.92%	50.90%	41.30%	41.15%	25.08%	11.63%	11.45%	12.73%	15.41%	30.71%
Filter	25.89%	29.12%	22.24%	12.59%	11.76%	11.76%	21.64%	17.92%	13.79%	14.70%	16.83%	29.02%
Path with *	15.02%	15.59%	12.88%	40.92%	32.43%	30.34%	16.43%	19.19%	14.80%	20.56%	17.26%	24.81%
Subquery	13.09%	15.30%	12.79%	6.45%	5.07%	5.39%	0.34%	0.28%	0.33%	0.09%	0.13%	0.11%
Bind	9.85%	9.23%	8.68%	4.72%	3.99%	4.15%	16.29%	12.07%	9.60%	11.94%	13.79%	24.03%
Union	5.10%	5.76%	12.62%	2.56%	2.07%	3.39%	11.26%	8.63%	7.61%	13.96%	13.05%	18.57%
Values	4.44%	3.07%	10.88%	3.29%	3.23%	3.20%	35.72%	30.74%	28.92%	29.82%	23.80%	11.90%
Not Exists	3.31%	3.37%	2.46%	1.24%	0.94%	0.69%	0.19%	0.21%	0.19%	0.27%	0.29%	0.35%
Minus	2.04%	2.91%	1.60%	0.82%	0.57%	0.71%	0.53%	0.92%	1.07%	1.46%	1.26%	1.78%
Service (lang)	44.63%	42.09%	54.78%	50.88%	41.71%	42.95%	10.40%	6.15%	4.27%	7.15%	7.91%	8.90%
Service (other)	11.49%	10.53%	10.32%	7.30%	13.14%	2.31%	4.51%	0.19%	1.16%	0.17%	0.18%	0.51%
Group By	17.12%	19.93%	13.04%	7.00%	5.40%	5.07%	0.41%	0.37%	0.48%	0.22%	0.23%	0.39%
Sample	8.85%	10.93%	4.60%	1.61%	1.63%	1.69%	0.04%	0.04%	0.06%	0.05%	0.04%	0.10%
Count	7.55%	7.60%	8.15%	5.22%	3.88%	3.73%	1.15%	4.30%	0.30%	1.52%	0.65%	0.89%
GroupConcat	1.80%	2.79%	1.17%	0.86%	0.86%	0.74%	0.06%	0.09%	0.02%	0.03%	0.02%	0.28%
Having	1.17%	1.14%	0.72%	0.65%	0.26%	0.33%	0.01%	0.01%	0.00%	0.00%	0.00%	0.04%

Triples per query: organic (blue) /robotic (yellow)



Languages of labels in organic queries



SPARQL feature co-occurrence

		organic		robotic				organic		robotic				
J	F O U P V S	I1-I3	I4-I6	I1-I3	I4-I6	J	F O U P V S	I1-I3	I4-I6	I1-I3	I4-I6			
	(none)	8.04	9.22	19.67	27.67	J	F O	2.66	1.32	2.13	1.18			
J		13.29	31.35	11.26	10.09	J	O U	3.49	0.25	0.02	0.01			
	F	1.10	0.98	1.92	1.31	J	O V	3.38	0.41	0.11	0.43			
J	F	6.68	2.39	2.61	1.68	J	O P V	1.01	0.06	0.16	0.07			
J		P	2.98	1.62	13.50	13.94	J		S	2.76	1.41	0.06	0.01	
J	F		P	2.48	0.58	0.39	0.07	J	O	S	4.78	0.62	0.00	0.01
J			V	0.39	2.01	30.42	17.47	J	F	S	3.19	2.28	0.03	0.01
		O	1.26	1.64	0.11	0.63	J	F O	S	1.02	0.13	0.00	0.00	
J	O	22.32	7.04	1.86	1.95	J	F O P	0.79	0.31	0.64	1.58			
J	O P	2.07	29.10	0.35	0.05	J		U P V	0.01	0.02	0.05	1.92		