



---

# On the Impact of sameAs on Schema Matching

---

**JOE RAAD** ● **ERMAN ACAR** ● **STEFAN SCHLOBACH**

Journées RoD — July 6, 2020

(work published in KCAP 2019)




Knowledge Representation & Reasoning Group

# More and more Linked Open Data...

← → ↻ [krr.triply.cc/krr/lod-a-lot/](https://krr.triply.cc/krr/lod-a-lot/)

krr / lod-a-lot Search ...

- 🏠 lod-a-lot
- 📄 Browser
- 📊 Table
- 🔗 Graphs 1
- ☁ Services 0
- 📎 Assets 0
- ★ Saved Queries 0



## lod-a-lot 🌐

by [Knowledge Representation & Reasoning \(KRR\)](#)

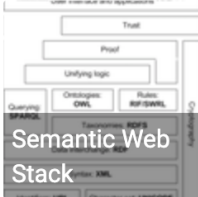
Created 8 months ago  
28.362.198.927 statements **(crawled from ~650K datasets in 2015)**

LOD-a-lot is the graph merge of the RDF graphs that were part of the LOD Laundromat. LOD-a-lot was created by [Fernández et al. 2017](#).

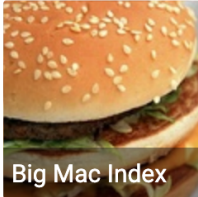
License  
[CC0 1.0](#)

Graphs  
[default](#) 28.362.198.927

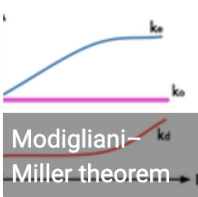
Example resources




Semantic Web Stack




Big Mac Index



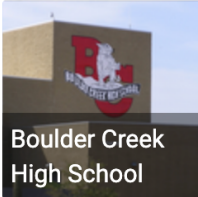
Modigliani-Miller theorem




Doppelstern



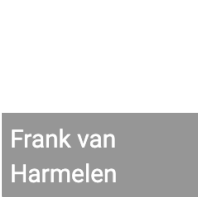
Étoile binaire



Boulder Creek High School



Winchester School of Art



Frank van Harmelen

[CommonKADS : 3rd KADS Meeting : Siemens AG Munich, March 8 - 9,1993 ; organized with GI special interest group 1.5.1 Knowledge Engineering](#)

# ...More and more (overlapping) Schemas

lov.linkeddata.es/dataset/lov/terms?q=Person&type=class&page=1

VOCABS TERMS AGENTS SPARQL/DUMP

TERMS Person

Results	URI	Score
513 results	<b>foaf:Person</b> (foaf) 2,320,027 occurrences in 72 LOD datasets <a href="http://xmlns.com/foaf/0.1/Person">http://xmlns.com/foaf/0.1/Person</a> <b>rdfs:comment</b> A person. <b>rdfs:label</b> Person <b>localName</b> Person	0.740
	<b>npg:Person</b> (npg) n/a (use in LOD) <a href="http://ns.nature.com/terms/Person">http://ns.nature.com/terms/Person</a> <b>skos:definition</b> The :Person class represents a single person entity. @en <b>skos:prefLabel</b> Person @en <b>localName</b> Person	0.556
	<b>akt:Person</b> (akt) 3,183,315 occurrences in 23 LOD datasets <a href="http://www.aktors.org/ontology/portal#Person">http://www.aktors.org/ontology/portal#Person</a> <b>localName</b> Person	0.521
	<b>bbccore:Person</b> (bbccore) n/a (use in LOD) <a href="http://www.bbc.co.uk/ontologies/coreconcepts/Person">http://www.bbc.co.uk/ontologies/coreconcepts/Person</a> <b>rdfs:label</b> Person @en-gb <b>localName</b> Person	0.512

**Type**

- vocabulary >
- property/class
- class (513)**
- agent >

**Tag**

- People (103)
- General & Upper (73)
- Society (62)
- Catalogs (33)
- Academy (24)
- Vocabularies (24)

# Schema Matching is Inevitable

- It is not possible (neither desired) to have a unique schema covering all domains
- In order to exploit this wealth of available knowledge and enhance knowledge-based systems (e.g., search engines, virtual assistants, etc.), we need to match these overlapping schemas
- Schema Matching: finding relationships between entities of different schemas
  - equivalence relations
  - subsumption
  - disjointness
  - ....

# Schema Matching over the years

- Active area of research from several communities, including the Semantic Web
- Ontology Alignment Evaluation Initiative (OAEI) ongoing for 15 years
- [Euzenat and Shvaiko, 2013] reviews ~100 schema-matching systems



50%

rely mostly on schema-level information  
(i.e. **schema-based approaches**)

25%

rely on schema + instance-level information  
(i.e. **mixed approaches**)

25%

rely mostly on instance-level information  
(i.e. **instance-based approaches**)

# Instance-based Schema Matching

- All instance-based schema-matching approaches share two essential ideas:

## 1. The semantics of a concept is better determined by its members rather by its annotations

**Concepts** refer to sets that possibly have named instances as *members*

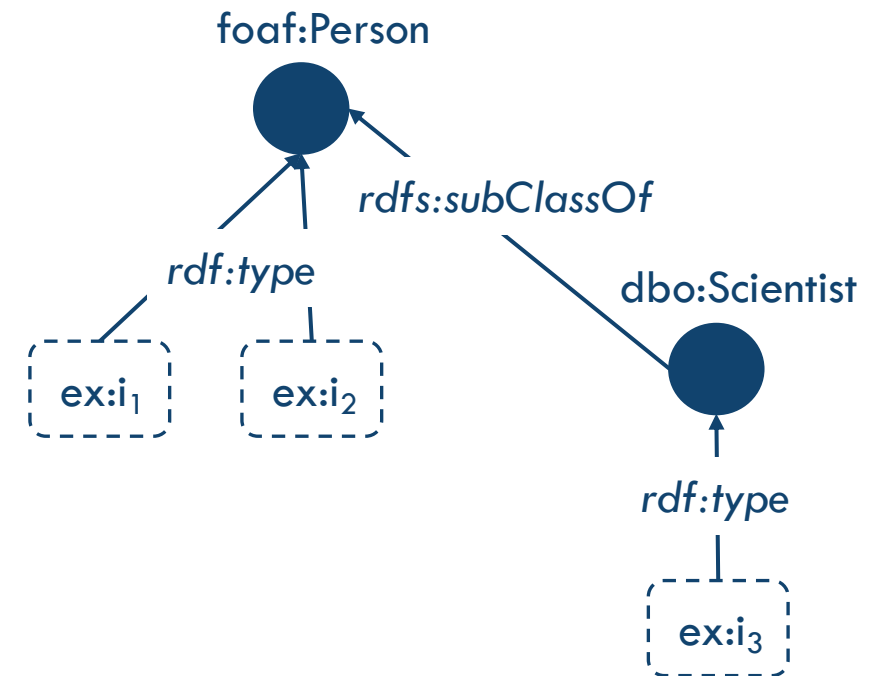
- $\text{ext}(\mathbf{C})$  refer to the set of instances which are **explicitly** stated as members of  $\mathbf{C}$

$$\text{ext}(\text{foaf:Person}) = \{\text{ex:i}_1, \text{ex:i}_2\}$$

- $\text{ext}_{\sqsubseteq}(\mathbf{C})$  refer to the set of instances which are **explicitly or implicitly** stated as members of  $\mathbf{C}$

(i.e. either explicit members or derived through concept subsumption)

$$\text{ext}_{\sqsubseteq}(\text{foaf:Person}) = \{\text{ex:i}_1, \text{ex:i}_2, \text{ex:i}_3\}$$



# Instance-based Schema Matching

- All instance-based schema-matching approaches share two essential ideas:
  2. **The more significant the overlap between two concepts' members is, the more related these concepts are**
- Multiple techniques to measure the overlap between concepts' members
  - Formal concept analysis techniques
  - Machine learning
  - Jaccard index
  - ....

# Instance-based Schema Matching using Jaccard Index

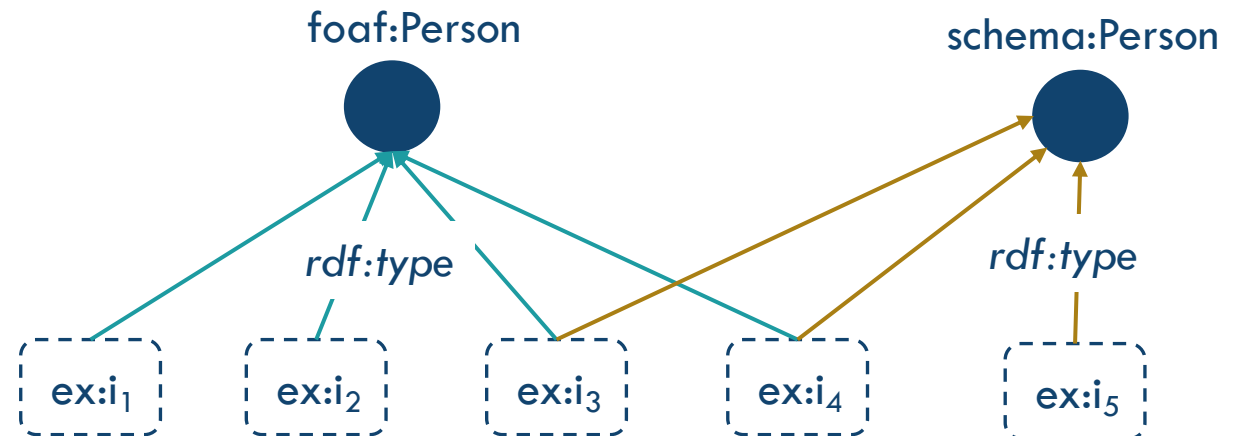
- The Jaccard index is a commonly used score to measure the similarity between two sets
- The higher the similarity of two sets is, the greater the Jaccard index

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$\text{ext}(\text{foaf:Person}) = \{\text{ex:i}_1, \text{ex:i}_2, \text{ex:i}_3, \text{ex:i}_4\}$

$\text{ext}(\text{schema:Person}) = \{\text{ex:i}_3, \text{ex:i}_4, \text{ex:i}_5\}$

$$J(\text{ext}(\text{foaf:Person}), \text{ext}(\text{schema:Person})) = \frac{2}{5} = 0.4$$





# Instance-based Schema Matching using Jaccard Index

With more than 558 million explicitly asserted owl:sameAs [Beek et al., ESWC 2018]

(or 35 billion after transitive closure), the reality in the Web of Data looks more like this:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Scenario 1 where J increases

$\text{ext}(\text{foaf:Person}) = \{\text{ex:i}_1, \text{ex:i}_2, \text{ex:i}_3, \text{ex:i}_4\}$

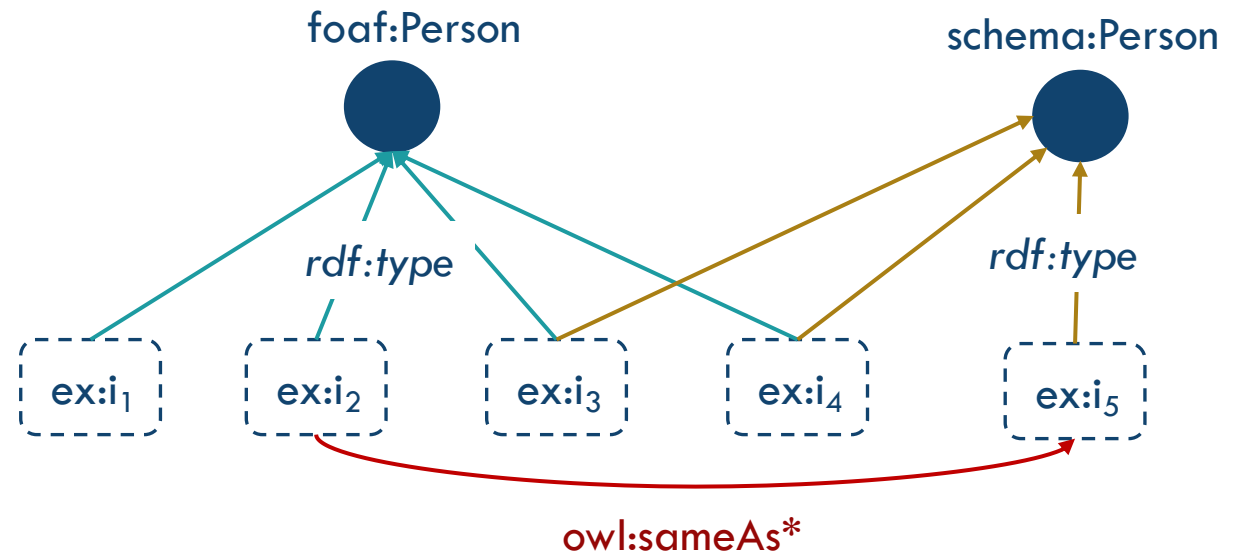
$\text{ext}(\text{schema:Person}) = \{\text{ex:i}_3, \text{ex:i}_4, \text{ex:i}_5\}$

$J(\text{ext}(\text{foaf:Person}), \text{ext}(\text{schema:Person})) = \frac{2}{5} = 0.4$

$\text{ext}^{\sim}(\text{foaf:Person}) = \{\text{ex:i}_1, \text{eq}^{\{2,5\}}, \text{ex:i}_3, \text{ex:i}_4\}$

$\text{ext}^{\sim}(\text{schema:Person}) = \{\text{ex:i}_3, \text{ex:i}_4, \text{eq}^{\{2,5\}}\}$

$J(\text{ext}^{\sim}(\text{foaf:Person}), \text{ext}^{\sim}(\text{schema:Person})) = \frac{3}{4} = 0.75$



# Instance-based Schema Matching using Jaccard Index

Or possibly like this:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Scenario 2 where J decreases

$\text{ext}(\text{foaf:Person}) = \{\text{ex:i}_1, \text{ex:i}_2, \text{ex:i}_3, \text{ex:i}_4\}$

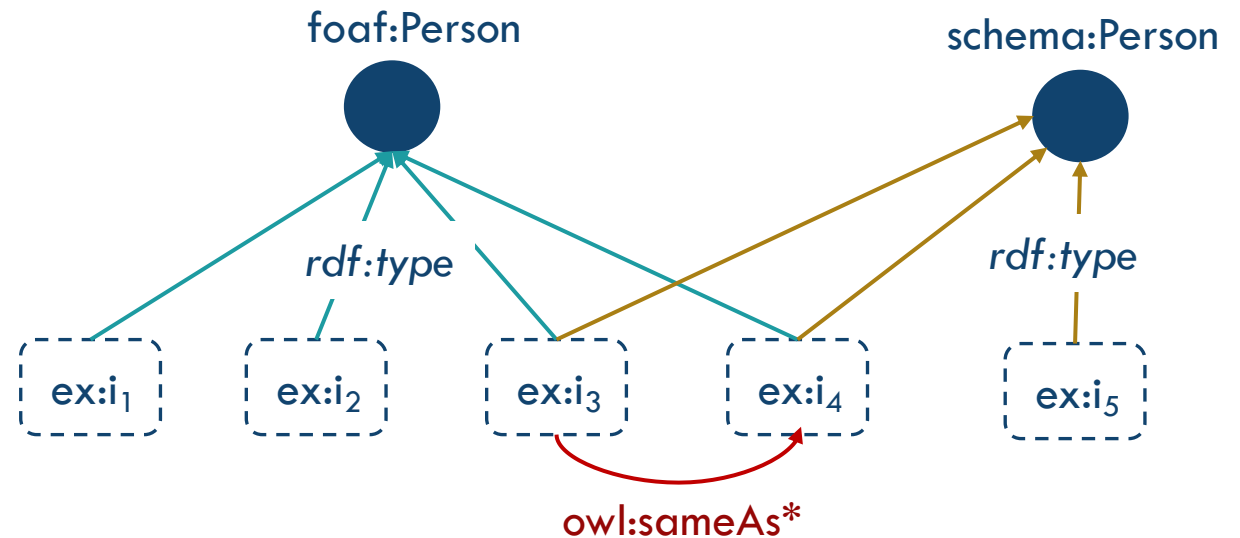
$\text{ext}(\text{schema:Person}) = \{\text{ex:i}_3, \text{ex:i}_4, \text{ex:i}_5\}$

$$J(\text{ext}(\text{foaf:Person}), \text{ext}(\text{schema:Person})) = \frac{2}{5} = 0.4$$

$\text{ext}^{\sim}(\text{foaf:Person}) = \{\text{ex:i}_1, \text{ex:i}_2, \text{eq}^{\{3,4\}}\}$

$\text{ext}^{\sim}(\text{schema:Person}) = \{\text{eq}^{\{3,4\}}, \text{ex:i}_5\}$

$$J(\text{ext}^{\sim}(\text{foaf:Person}), \text{ext}^{\sim}(\text{schema:Person})) = \frac{1}{4} = 0.25$$



# Research Question

# Research Question

**Does the inclusion of instance-level interlinks (i.e. owl:sameAs) positively impact instance-based schema alignments ?**

- a.** Does the inclusion of owl:sameAs increase the Jaccard Index of equivalent concepts?
- b.** Does the inclusion of owl:sameAs increase the Jaccard Index of non-equivalent concepts?

# | Why should we care?

- Provides empirical evidence for schema-matching designers on whether exploiting a large external collection of instance-level interlinks (e.g. from the LOD Cloud) is beneficial for improving the accuracy of schema-matching techniques
- Shows the risks/benefits of using owl:sameAs after a number of studies suggesting that a large\* number of the existing owl:sameAs links in the Web are actually erroneous
  - \* **20%** of evaluated owl:sameAs are erroneous [Halpin et al., ISWC 2010]
  - \* **3%** of evaluated owl:sameAs are erroneous [Hogan et al., JWS 2012]
  - \* **4%** of evaluated owl:sameAs are erroneous [Raad et al., ISWC 2018]

# Dataset Description

# Dataset

← → ↻ [krr.triply.cc/krr/lod-a-lot/](https://krr.triply.cc/krr/lod-a-lot/)

krr / lod-a-lot

Search ...

lod-a-lot

- Browser
- Table
- Graphs 1
- Services 0
- Assets 0
- Saved Queries 0



## lod-a-lot

by [Knowledge Representation & Reasoning \(KRR\)](#)

Created 8 months ago

28.362.198.927 statements **(crawled from ~650K datasets in 2015)**

LOD-a-lot is the graph merge of the RDF graphs that were part of the LOD Laundromat. LOD-a-lot was created by [Fernández et al. 2017](#).

License

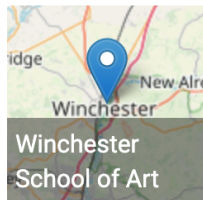
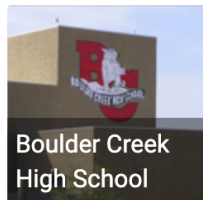
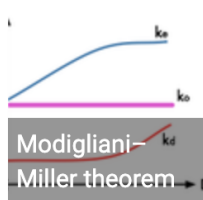
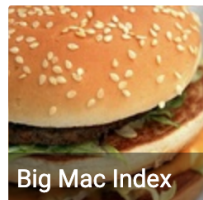
[CC0 1.0](#)

Graphs

[default](#)

28.362.198.927

Example resources



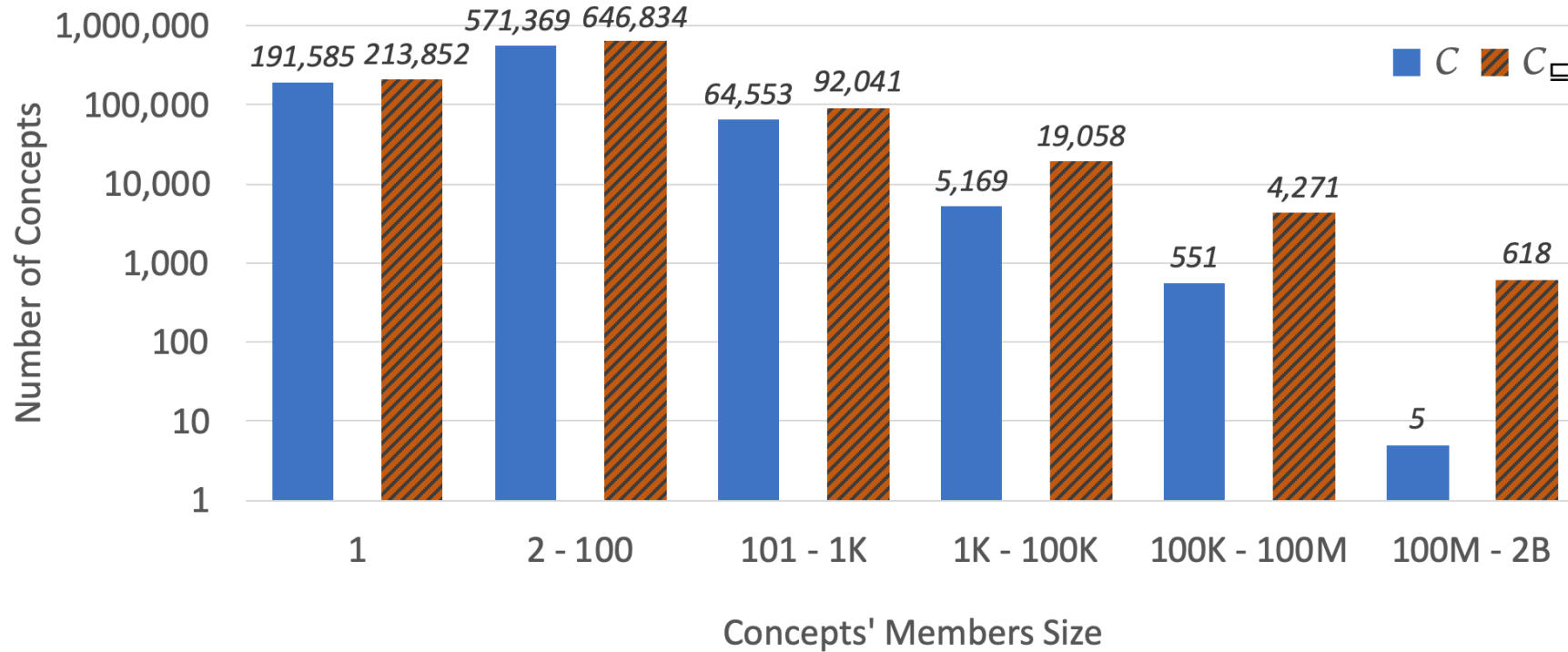
[CommonKADS : 3rd KADS Meeting : Siemens AG Munich, March 8 - 9,1993](#) ; organized with GI special interest group 1.5.1 Knowledge Engineering

# Dataset

# triples	28,362,198,927
# rdf:type statements	3,321,354,308
# rdfs:subClassOf statements	4,461,717
# owl:equivalentClass statements	1,051,979
# explicit owl:sameAs statements	558,943,116
# implicit owl:sameAs statements	35,201,120,188
# equivalence classes (after closure of owl:sameAs)	48,999,148
# concepts with at least one explicit member $ C $	833,232
# concepts with at least one explicit or implicit member $ C_{\sqcup} $	976,674



# Size distribution of the Concepts' members



- **23%** of the concepts have one explicit member
- **92%** of the concepts have  $\leq 100$  explicit members
- **618 concepts** have more than 100M explicit or implicit members
- **5 concepts** have more than 100M explicit members

# Concepts with more than $>100M$ explicit members

Concept	Cardinality	%
<a href="http://purl.org/linked-data/cube#Observation">http://purl.org/linked-data/cube#Observation</a>	1,306,389,396	39.3
<a href="http://data-gov.tw.rpi.edu/2009/data-gov-twc.rdf#DataEntry">http://data-gov.tw.rpi.edu/2009/data-gov-twc.rdf#DataEntry</a>	304,878,654	9.2
<a href="http://geovocab.org/geometry#Geometry">http://geovocab.org/geometry#Geometry</a>	167,808,111	5
<a href="http://knoesis.wright.edu/ssw/ont/sensorobservation.owl#MeasureData">http://knoesis.wright.edu/ssw/ont/sensorobservation.owl#MeasureData</a>	144,044,989	4.3
<a href="http://xmlns.com/foaf/0.1/Person">http://xmlns.com/foaf/0.1/Person</a>	132,919,327	4
<b>Total</b>	<b>2,056,040,477</b>	<b>61.9</b>

These **5 concepts** with more than 100M explicit members are the objects of 62% of the total `rdf:type` statements in the LOD-a-lot

# Experiments

# Research Question

**Does the inclusion of instance-level interlinks (i.e. owl:sameAs) positively impact instance-based schema alignments ?**

- a. Does the inclusion of owl:sameAs increase the Jaccard Index of equivalent concepts?**
- b. Does the inclusion of owl:sameAs increase the Jaccard Index of non-equivalent concepts?**

# Dataset

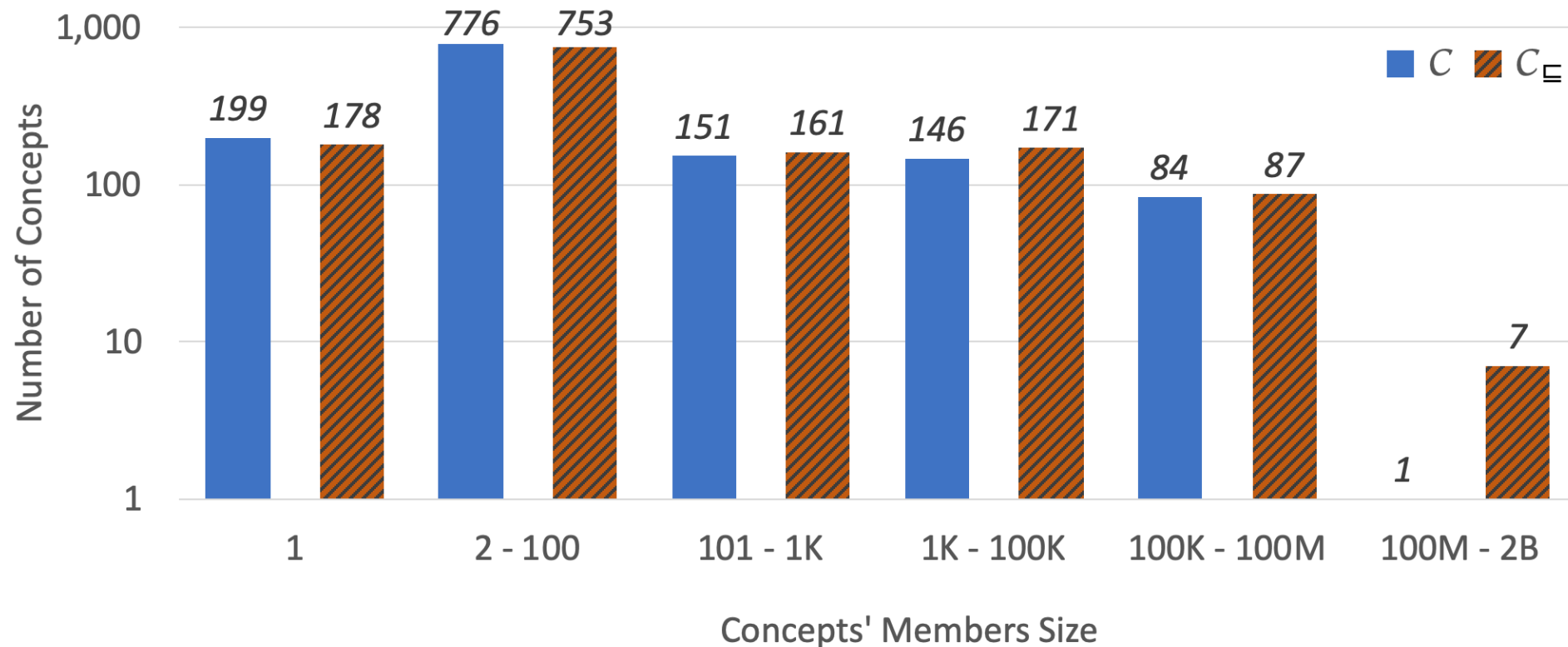
# triples	28,362,198,927
# rdf:type statements	3,321,354,308
# rdfs:subClassOf statements	4,461,717
# owl:equivalentClass statements	1,051,979
# explicit owl:sameAs statements	558,943,116
# implicit owl:sameAs statements	35,201,120,188
# equivalence classes (after closure of owl:sameAs)	48,999,148
# concepts with at least one explicit member $ C $	833,232
# concepts with at least one explicit or implicit member $ C_{\sqcup} $	976,674

# a. Does the inclusion of owl:sameAs increase the Jaccard Index of equivalent concepts?

- 1,051,979 owl:equivalentClass statements in the LOD-a-lot
  - Hypothesis: all these existing statements are correct alignments
  - Only 972 owl:equivalentClass statements where both concepts have explicit members
    - 208 reflexive alignments ( $C1, owl:equivalentClass, C1$ )
    - 22 duplicate symmetric alignments ( $C1, owl:equivalentClass, C2$ ) and ( $C2, owl:equivalentClass, C1$ )
    - **742 alignments between 1,357 distinct concepts (i.e. gold standard)**

# a. Does the inclusion of owl:sameAs increase the Jaccard Index of equivalent concepts?

Size distribution of the Concepts' members of our Gold Standard



# a. Does the inclusion of owl:sameAs increase the Jaccard Index of equivalent concepts?

Jaccard Index distribution for the 742 alignments

Compare

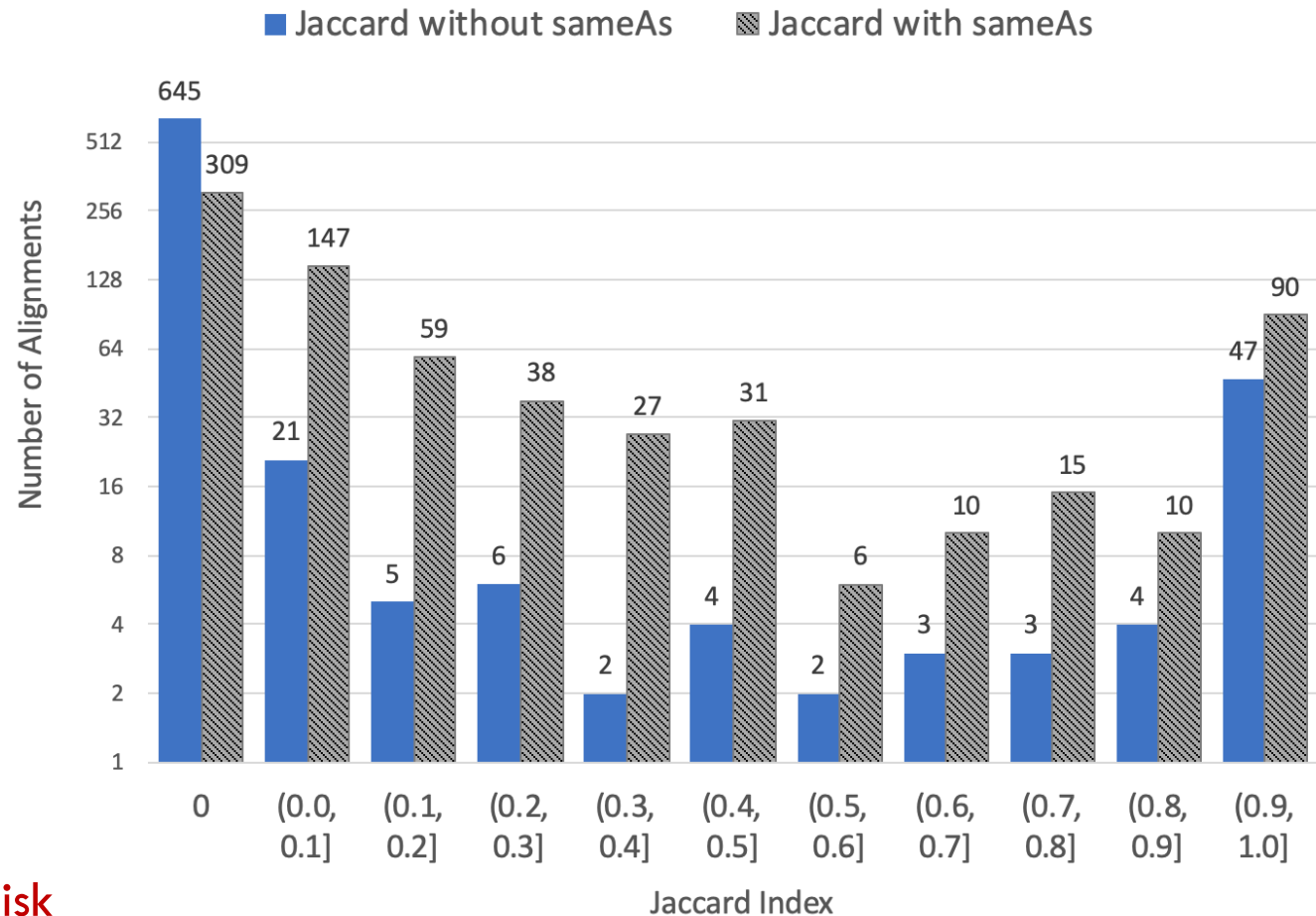
$$J(\text{ext}_{\sqsubseteq}(C1), \text{ext}_{\sqsubseteq}(C2))$$

with

$$J(\text{ext}_{\sqsubseteq\sim}(C1), \text{ext}_{\sqsubseteq\sim}(C2))$$

such that

$$(C1, \text{owl:equivalentClass}, C2)$$



Runtime: 4 hours on 64GB SSD disk



# a. Does the inclusion of owl:sameAs increase the Jaccard Index of equivalent concepts?

Jaccard Index variation for the 742 alignments

- When owl:sameAs is considered, the Jaccard index increases for 381 / 742 of the correct alignments (52%)
- Out of these 381 cases, Jaccard increases from 0 to 1 in 44 cases (6%)
- When owl:sameAs is considered, the Jaccard index decreases for 25 / 742 of the correct alignments (3%)
- Slight drop in impact when only explicit members are considered

**The inclusion of owl:sameAs does increase the overlap of two equivalent concepts in half of the cases**

Jaccard Index		0	(0, 1)	1	Total
C	<b>Total</b>	<b>655</b> (88%)	<b>73</b> (10%)	<b>14</b> (2%)	<b>742</b>
	<i>Decreases</i>	N/A	25 (34%)	0 (0%)	<b>25</b> (3%)
	<i>No variation</i>	333 (51%)	9 (12%)	14 (100%)	<b>356</b> (48%)
	<i>Increases (J &lt; 1)</i>	278 (42%)	39 (54%)	N/A	<b>317</b> (43%)
	<i>Increases (J = 1)</i>	44 (7%)	0 (0%)	N/A	<b>44</b> (6%)
C <sub>E</sub>	<b>Total</b>	<b>645</b> (87%)	<b>81</b> (11%)	<b>16</b> (2%)	<b>742</b>
	<i>Decreases</i>	N/A	25 (31%)	0 (0%)	<b>25</b> (3%)
	<i>No variation</i>	309 (48%)	11 (14%)	16 (100%)	<b>336</b> (45%)
	<i>Increases (J &lt; 1)</i>	292 (45%)	45 (55%)	N/A	<b>337</b> (46%)
	<i>Increases (J = 1)</i>	44 (7%)	0 (0%)	N/A	<b>44</b> (6%)

# Research Question

**Does the inclusion of instance-level interlinks (i.e. owl:sameAs) positively impact instance-based schema alignments ?**

- a.** Does the inclusion of owl:sameAs increase the Jaccard Index of equivalent concepts?
- b.** Does the inclusion of owl:sameAs increase the Jaccard Index of non-equivalent concepts?

# Dataset

# triples	28,362,198,927
# rdf:type statements	3,321,354,308
# rdfs:subClassOf statements	4,461,717
# owl:equivalentClass statements	1,051,979
# explicit owl:sameAs statements	558,943,116
# implicit owl:sameAs statements	35,201,120,188
# equivalence classes (after closure of owl:sameAs)	48,999,148
# concepts with explicit members $ C $	833,232
# concepts with explicit or implicit members $ C_{\sqcup} $	976,674

## b. Does the inclusion of owl:sameAs increase the Jaccard Index of non-equivalent concepts?

- 833,232 concepts with explicit members in the LOD-a-lot
  - Create one random alignment for each concept, such that each concept is paired only once
  - Hypothesis: all these random alignments are erroneous
    - **416,616 random alignments**

# b. Does the inclusion of owl:sameAs increase the Jaccard Index of non-equivalent concepts?

Jaccard Index variation for the 416,616 random alignments

- When owl:sameAs is considered, the Jaccard index increases for only 94 / 416,616 of the random alignments (0.02%)
- When owl:sameAs is considered, the Jaccard index decreases for 3 / 416,616 of the random alignments

**owl:sameAs rarely increases the overlap of two non-equivalent concepts**

Jaccard Index	0	(0, 1)	1	Total
<b>Total</b>	<b>412,828</b> (99.1%)	<b>2,808</b> (0.67%)	<b>980</b> (0.23%)	<b>416,616</b>
<i>Decreases</i>	N/A	3 (0.1%)	0 (0%)	<b>3</b> (0%)
<i>No variation</i>	412,751 (99.98%)	2,788 (99.3%)	980 (100%)	<b>416,519</b> (99.98%)
<i>Increases (J &lt; 1)</i>	77 (0.02%)	17 (0.6%)	N/A	<b>94</b> (0.02%)
<i>Increases (J = 1)</i>	0 (0%)	0 (0%)	N/A	<b>0</b> (0%)

Take away message

# Take away message

**This work provides an empirical study on the impact of including instance-level interlinks on the overlap between concepts members**

- Including instance-level interlinks can enhance the performance of instance-based schema alignments
  - Increases the overlap for 52% of the existing (i.e. correct) alignments in the LOD-a-lot
  - Increases the overlap for less than 0.3% of randomly created (i.e. erroneous) alignments
- Inference does positively impact instance-based schema alignments
  - Considering also the implicit members enhances the results on the Gold Standard by 3 pp

**Additional findings in the paper:**

- Discarding only isolated owl:sameAs links in the network can increase the quality of instance-based schema alignments (owl:sameAs links are probably not as bad as we first thought)
  - Reduces the cases where Jaccard index increases for non-equivalent concepts by 71%



---

# On the Impact of sameAs on Schema Matching

---

JOE RAAD ● ERMAN ACAR ● STEFAN SCHLOBACH



Knowledge Representation & Reasoning Group

## Thank you for your attention!

Code & Results

<https://github.com/raadjoe/impact-sameAs-schema-matching>