

QUANT-Question Answering Benchmark Curator

Ria Hari Gusmita, Rricha Jalota, Daniel Vollmers, Jan Reineke, Axel-Cyrille Ngonga Ngomo, and Ricardo Usbeck



September 10, 2019

- 1 *Motivation*
- 2 *Approach*
- 3 *Evaluation*
- 4 *QALD-specific Analysis*
- 5 *Conclusion & Future Work*

Motivation

Drawback in evaluating Question Answering systems over knowledge bases

- Mainly based on benchmark datasets (benchmarks)
- Challenge in maintaining high-quality and benchmarks

QALD

LC-QuAD

Free917

WebQuestions

Motivation

Challenge in maintaining high-quality and benchmarks

- Change of the underlying knowledge base

DBpedia 2016-04	DBpedia 2016-10
http://dbpedia.org/resource/Surfing	http://dbpedia.org/resource/Surfer
http://dbpedia.org/ontology/seatingCapacity	http://dbpedia.org/property/capacity
http://dbpedia.org/property/portrayer	http://dbpedia.org/ontology/portrayer
http://dbpedia.org/property/establishedDate	http://dbpedia.org/ontology/foundingDate

Motivation

Challenge in maintaining high-quality and benchmarks

- Metadata annotation errors

Question What is the revenue of IBM?

SPARQL

```
PREFIX res: <http://dbpedia.org/resource/>
PREFIX onto: <http://dbpedia.org/ontology/>

SELECT ?number
WHERE
{ res:IBM onto:revenue ?number }
```

Answer from File

[1.0363E11]

Answer from Current Endpoint

[8.1741E10]

Endpoint

Answer Type

Out of Scope

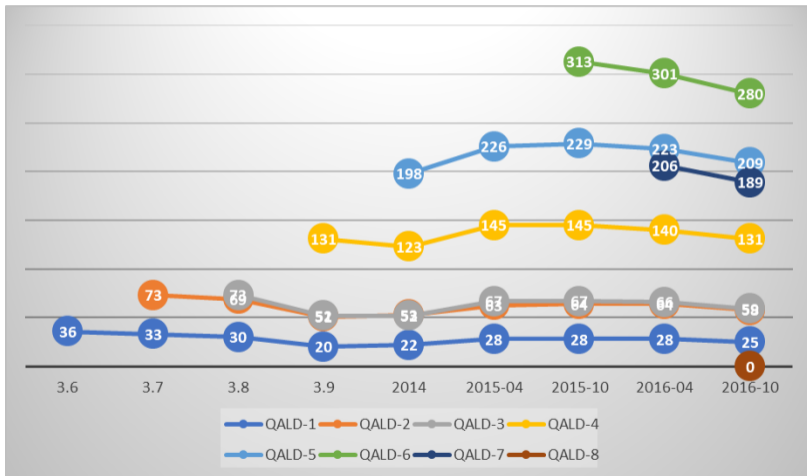
Aggregation

Onlydbo

Hybrid

Motivation

Degradation QALD benchmarks against various versions of DBpedia



- QUANT, a framework for the intelligent creation and curation of QA benchmarks

Definition

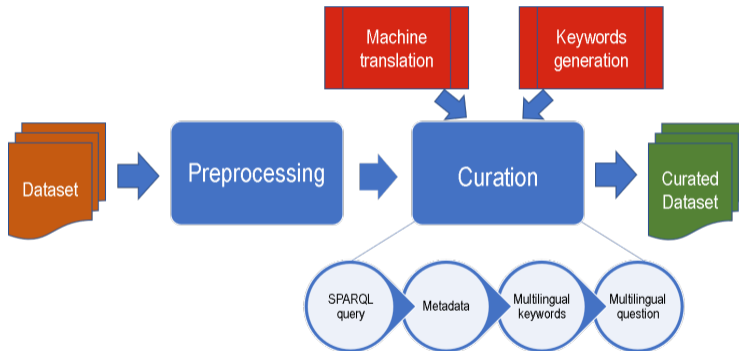
- Given B , D , and Q as benchmark, dataset, and questions respectively
- S represents QUANT's suggestions
 - i^{th} version of a QA benchmark B_i as a pair (D_i, Q_i)
 - Given a query $q_{ij} \in Q_i$ with zero results on D_k with $k > i$
 - $S : q_{ij} \rightarrow q'_{ij}$
- QUANT aims
 - to ensure that queries from B_i can be reused for B_k
 - to speed up the curation process as compared to the existing one

What QUANT supports

- ① Creation of SPARQL queries
- ② The validity of benchmark metadata
- ③ Spelling and grammatical correctness of questions

Approach

Architecture



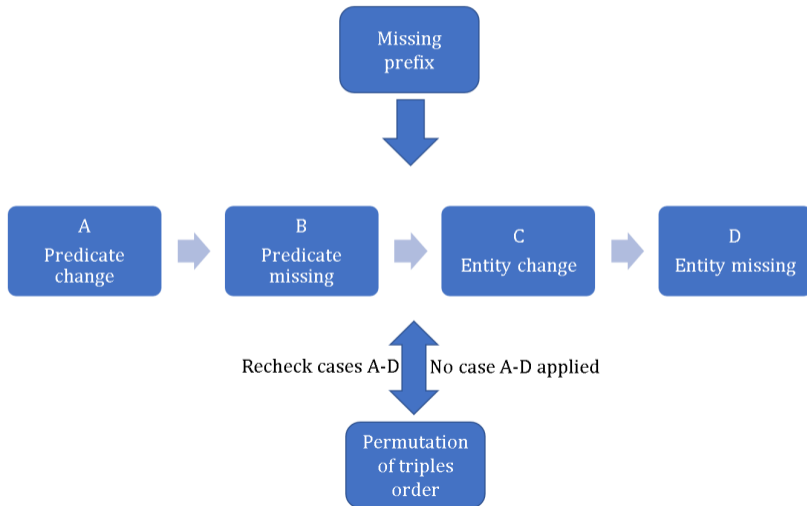
Approach

Smart suggestions

- 1 SPARQL suggestion
- 2 Metadata suggestion
- 3 Multilingual Questions and Keywords Suggestion

Smart suggestion

1. How SPARQL suggestion module works



1. SPARQL suggestion

Missing prefix

- The original SPARQL query

```
SELECT ?s
WHERE {
    res:New_Delhi dbo:country ?s .
}
```

1. SPARQL suggestion

Missing prefix

- The original SPARQL query

```
SELECT ?s
WHERE {
    res:New_Delhi dbo:country ?s .
}
```

- The suggested SPARQL query

```
PREFIX dbo: <http://dbpedia.org/ontology/>
PREFIX res: <http://dbpedia.org/resource/>
SELECT ?s
WHERE {
    res:New_Delhi dbo:country ?s .
}
```

1. SPARQL suggestion

Predicate change

- The original SPARQL query

```
SELECT ?date
WHERE {
    ?website rdf:type onto:Software .
    ?website onto:releaseDate ?date .
    ?website rdfs:label "DBpedia" .
}
```

1. SPARQL suggestion

Predicate change

- The suggested SPARQL query

```
SELECT ?date
WHERE {
    ?website rdf:type onto:Software .
    ?website rdfs:label "DBpedia" .
    ?website dbp:latestReleaseDate ?date .
}
```

1. SPARQL suggestion

Predicate missing

- The original SPARQL query

```
SELECT ?uri
WHERE {
    ?subject rdfs:label "Tom□Hanks".
    ?subject foaf:homepage ?uri
}
```


1. SPARQL suggestion

Predicate missing

- The original SPARQL query

```
SELECT ?uri
WHERE {
    ?subject rdfs:label "Tom□Hanks".
    ?subject foaf:homepage ?uri
}
```

- The suggested SPARQL query The predicate foaf:homepage is missing in ?subject foaf:homepage ?uri

1. SPARQL suggestion

Entity change

- The original SPARQL query

```
SELECT ?uri WHERE  
{ ?uri rdf:type yago:CapitalsInEurope }
```

1. SPARQL suggestion

Entity change

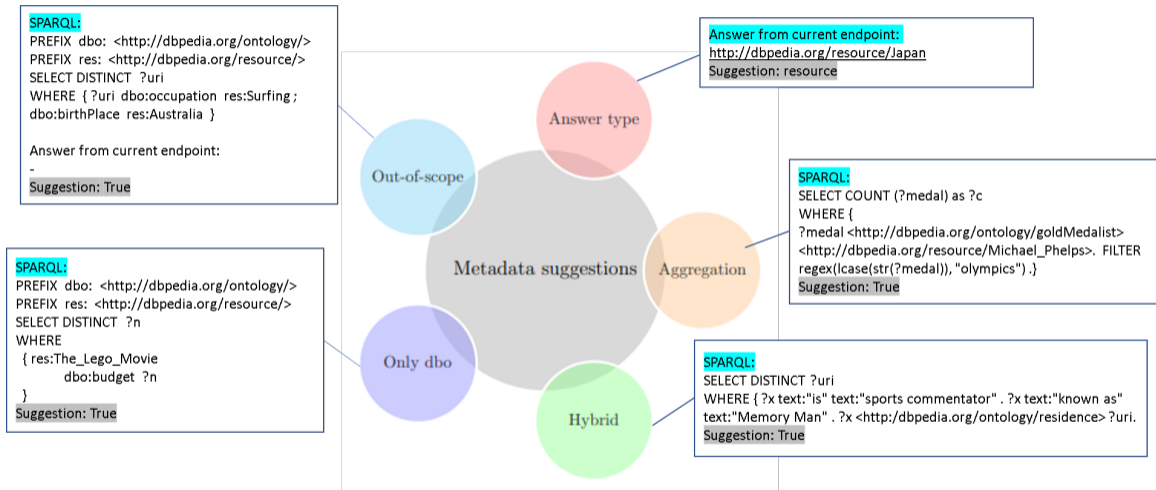
- The original SPARQL query

```
SELECT ?uri WHERE  
{ ?uri rdf:type yago:CapitalsInEurope }
```

- The suggested SPARQL query

```
SELECT ?uri WHERE  
{ ?uri rdf:type yago:WikicatCapitalsInEurope }
```

2. Metadata suggestion



3. Multilingual questions and keywords suggestion

Question with missing keywords and translations

```
"id": "2",
"datasetVersion": null,
"answerType": "",
"aggregation": false,
"hybrid": false,
"onlydbo": false,
"sparqlQuery": "PREFIX yago: <http://dbpedia.org/class/yago/>\nPREFIX rdf: <http://dbpedia.org/property/>\nSELECT ?uri ?string\nWHERE \n{\n\t?uri rdf:type ya\n\t?string. FILTER (lang(?string) = 'en') }\n} ORDER BY ASC(?density) LIMIT 1",
"pseudoSparqlQuery": null,
"outOfScope": null,
"languageToQuestion": {
  "en": "Which state of the United States of America has the highest density?"
},
"languageToKeyword": {},
"goldenAnswer": [
  "http://dbpedia.org/resource/New\_Jersey"
]
```

3. Multilingual questions and keywords suggestion

- Generated keywords: state, united, states, america, highest, density
- Utilizing Trans Shell tool→Generated keywords translations suggestion

```
"languageToKeyword" : {  
  "de" : [ "Zustand ", "vereinigt ", "Zustände ", "Amerika ", "höchste ", "Dichte " ],  
  "ru" : [ "государство ", "единый ", "состояния ", "Америка ", "наибольший ", "плотность " ],  
  "pt" : [ "Estado ", "Unidos ", "estados ", "América ", "maior ", "densidade " ],  
  "en" : [ "state ", "united ", "states ", "america ", "highest ", "density " ],  
  "hi_IN" : [ "राज्य ", "संयुक्त ", "राज्यो ", "अमेरिका ", "उच्चतम ", "घनत्व " ],  
  "it" : [ "stato ", "unito ", "stati ", "America ", "massimo ", "densità " ],  
  "fa" : [ "تلاخ ", "دحتم ", "اه تلاب ", "اکیرما ", "نیرنالاب ", "بلاغج " ],  
  "fr" : [ "Etat ", "uni ", "États ", "Amérique ", "le plus élevé ", "densité " ],  
  "ro" : [ "stat ", "Unit ", "statele ", "America ", "cel mai inalt ", "densitate " ],  
  "es" : [ "estado ", "unido ", "estados ", "America ", "más alto ", "densidad " ],  
  "nl" : [ "staat ", "verenigd ", "staten ", "Amerika ", "hoogst ", "dichtheid " ]  
},
```

3. Multilingual questions and keywords suggestion

Suggested Question Translations

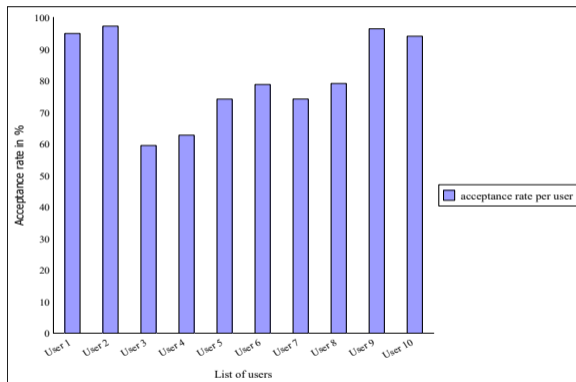
```
"languageToQuestion": {  
  "de": "Welche Zustand von das Vereinigt Zustände von Amerika hat das höchste Dichte? ",  
  "ru": "Который государство из объединенный состояния из Америка имеет наибольший плотность? ",  
  "pt": "Qual Estado do a Unidos Estados do América tem a maior densidade? ",  
  "en": "Which state of the United States of America has the highest density? ",  
  "hi_IN": "कौन कौन से राज्य का यूनाइटेड राज्य अमेरिका का अमेरिका है उच्चतम घनत्व? ",  
  "it": "Quale stato di il Unito stati di America ha il massimo densità? ",  
  "fa": "مادک تلاج زا نیا دنیانوی اه تلایا زا اکیرما یاراد نیا نیرتلاب ؟مکارت ",  
  "fr": "Lequel Etat de la Uni États de Amérique a la le plus élevé densité? ",  
  "ro": "Care stat de Unit statele de America are cel mai inalt densitate? ",  
  "es": "Cual estado de el Unido Estados de America tiene el más alto ¿densidad? ",  
  "nl": "Welke staat van de Verenigd Staten van Amerika heeft de hoogst dichtheid? "  
},
```

- Three goals of the evaluation:
 - ① QUANT vs manual curation
 - Graduate students curated 50 questions using QUANT and another 50-question manually
 - 23 minutes vs 278 minutes
 - ② Effectiveness of smart suggestions
 - 10 expert users got involved in creating a new joint benchmark, called QALD-9, with 653 questions
 - ③ QUANT's capability to provide a high-quality benchmark dataset
 - The inter-rater agreement between each two users amounts up to 0.83 on average

Group	Inter-rater Agreement
1st Two-Users	0.97
2nd Two-Users	0.72
3rd Two-Users	0.88
4th Two-Users	0.77
5th Two-Users	0.96
Average	0.83

Evaluation

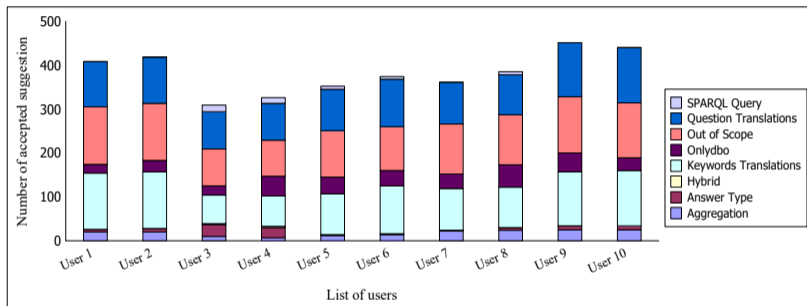
Users acceptance rate in %



- QUANT provided 2380 suggestions and user acceptance rate on average is 81%
- The top 4 acceptance-rate are for QALD-7 and QALD-8

Evaluation

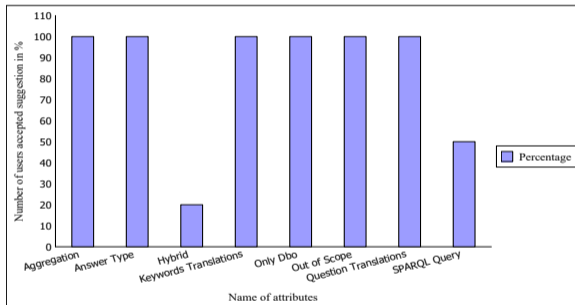
Number of accepted suggestions from all users



- Most users accepted suggestion for out-of-scope metadata
- Keyword and question translation suggestions yielded the second and third highest acceptance rates.

Evaluation

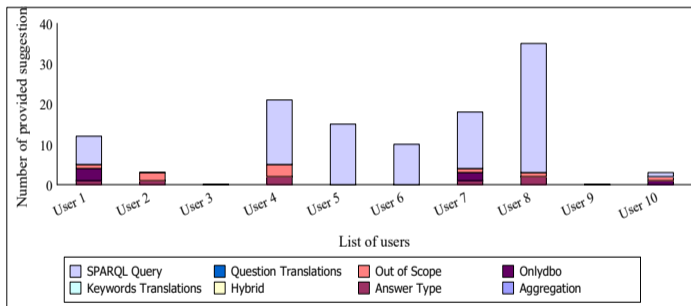
Number of users who accepted QUANT's suggestions for each question's attribute.



- 83.75% of the users accepted QUANT's smart suggestions on average
- Hybrid and SPARQL suggestions were only accepted by 2 and 5 users respectively.

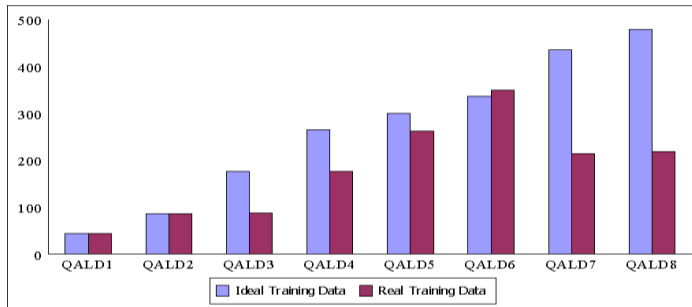
Evaluation

Number of suggestions provided by users



- Answer type, onlydbo, out-of-scope, and SPARQL query metadata were attributes whose value redefined by users

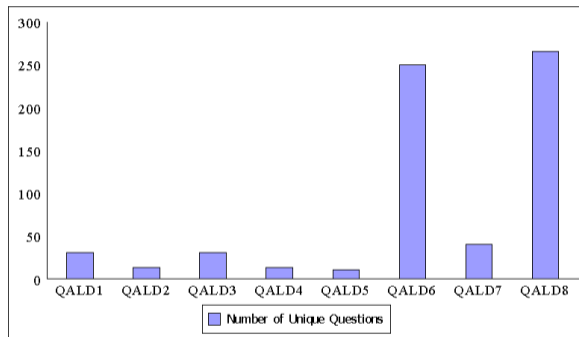
QALD-specific Analysis



There are 1924 questions where 1442 questions are training data and 482 questions are test data

QALD-specific Analysis

- Duplication removal resulted 655 unique questions
- Removing 2 semantically similar questions produced 653 questions
- Using QUANT with 10 expert users, we got 558 total benchmark questions → increase QALD-8 size by 110.6%
- The new benchmark formed QALD-9 dataset



Distribution of unique questions in all QALD versions

- QUANT's evaluation highlights the need for better datasets and their maintenance
- QUANT speeds up the curation process by up to 91%.
- Smart suggestions motivate users to engage in more attribute corrections than if there were no hints



- There is a need to invest more time into SPARQL suggestions as only 5 users accepted them
- We plan to support more file formats based on our internal library



Thank you for your attention!

Ria Hari Gusmita

ria.hari.gusmita@uni-paderborn.de

<https://github.com/dice-group/QUANT>

DICE Group at Paderborn University

[https:](https://dice-research.org/team/profiles/gusmita/)

[//dice-research.org/team/profiles/gusmita/](https://dice-research.org/team/profiles/gusmita/)

