

Towards More Usable Dataset Search: From Query Characterization to Snippet Generation

Jinchi Chen
National Key Laboratory for Novel
Software Technology, Nanjing
University, China
jcchen@smail.nju.edu.cn

Xi Xia Wang
National Key Laboratory for Novel
Software Technology, Nanjing
University, China
xxwang@smail.nju.edu.cn

Gong Cheng
National Key Laboratory for Novel
Software Technology, Nanjing
University, China
gcheng@nju.edu.cn

Evgeny Kharlamov
Bosch Center for AI
Bosch GmbH, Germany
University of Oslo, Norway
evgeny.kharlamov@de.bosch.com
evgeny.kharlamov@ifi.uio.no

Yuzhong Qu
National Key Laboratory for Novel
Software Technology, Nanjing
University, China
yzqu@nju.edu.cn

ABSTRACT

Reusing published datasets on the Web is of great interest to researchers and developers. Their data needs may be met by submitting queries to a dataset search engine to retrieve relevant datasets. In this ongoing work towards developing a more usable dataset search engine, we characterize real data needs by annotating the semantics of 1,947 queries using a novel fine-grained scheme, to provide implications for enhancing dataset search. Based on the findings, we present a query-centered framework for dataset search, and explore the implementation of snippet generation and evaluate it with a preliminary user study.

CCS CONCEPTS

• **Information systems** → **Query intent; Summarization.**

KEYWORDS

dataset search; data needs; query annotation; snippet generation

1 INTRODUCTION

Reusing existing datasets helps in improving productivity and reducing cost for application developers. In order to support convenient search of datasets that match a developer’s *data needs*, Google Dataset Search and other *dataset search engines* have recently emerged. However, there is much room for improving their usability. Existing efforts mostly focus on metadata management [8], result filtering [6], and dataset browsing [9]. At the same time little attention has been given to the queries in dataset search, which may differ considerably from general Web search queries.

In this ongoing work towards developing a more usable dataset search engine, we first analyze real data needs collected from user-submitted posts on a variety of websites. We reveal the complex constitution of a typical dataset search query, where various elements of both metadata and data content can be mentioned. This brings us to the idea that as a specialization of information retrieval [6], dataset search poses unique challenges.

Based on our empirical findings, in Figure 1 we present a new *query-centered framework* for dataset search that has four main components. (i) *Query Processing* component understands how a data

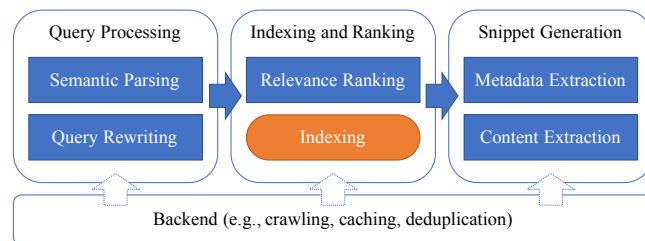


Figure 1: Proposed framework for dataset search.

need is formulated as an input query and pre-processes the query in order to achieve best query results. The query is interpreted by semantic parsing, and will be rewritten (e.g., relaxed) when needed (e.g., to avoid empty results). (ii) *Indexing and Ranking* component determines which datasets are relevant to the query. We compute relevance and construct indexes for fast computation. (iii) *Snippet Generation* component explains how a retrieved dataset is relevant to the query and reflects the underlying data need. A query-biased snippet is extracted from its metadata and/or its content. (iv) All these components are supported by the *backend module*, which crawls, caches, and cleans datasets.

Some components in the framework have attracted research interests, such as dataset ranking [10]. Others are still rather ad hoc; they lack rigorous solutions or their performance has not been rigorously evaluated. In this paper, we report our progress in two aspects: *characterizing data needs* (Section 2) and *generating dataset snippets* (Section 3). Our contribution is summarized as follows.

- We collect real data needs from diverse sources, including user-submitted posts from online communities and data requests submitted to a national open data portal. We derive 1,947 dataset search queries and semantically annotate them using a novel fine-grained scheme. Our analysis provides implications for enhancing dataset search.
- For RDF data, in contrast to query-independent illustrative dataset snippet [1], we extract a query-biased snippet from the content of a dataset by adapting a state-of-the-art algorithm for the group Steiner tree problem [7]. We conduct a

Table 1: Query Annotation Scheme and Its Distribution in Dataset Search Queries

Category	% of Queries	Example Query
Metadata	Name	3.54% <i>HUST-ASL Dataset</i>
	Domain/Topic	94.45% <i>weather dataset with solar radiance and solar energy production</i>
	Data Format	16.23% <i>jpg images for all unicode characters</i>
	Language	3.90% <i>annotated moive review dataset in German</i>
	Accessibility	7.40% <i>open source handwritten English alphabets dataset</i>
	Provenance	0.21% <i>FDA datasets about medicine name and the result has adverse events</i>
	Statistics	2.98% <i>dataset contains at least 1000 examples of opinion articles</i>
	Overall	96.05%
Content	Concept	50.59% <i>dataset about people, include gender, ethnicity, name</i>
	Geospatial	19.21% <i>judicial decisions in France</i>
	Other Entities	0.41% <i>datasets with nutrition data for many commercial food products (i.e., Lucky Charms, Monster Energy, Nutella, etc.)</i>
	Temporal	9.35% <i>2011–2013 MoT failure rates on passenger cars</i>
	Other Numbers	1.59% <i>businesses that employ over 1000 people in Yorkshire region</i>
	Overall	63.79%

user study to evaluate the usefulness of these two types of snippets, and also quantitatively analyze their features.

2 CHARACTERIZATION OF DATA NEEDS

Comparing to general Web search, the knowledge about dataset search is rather limited. Therefore, we collect and analyze real data needs to provide empirical findings and implications for enhancing dataset search. We will now discuss this in details.

2.1 Collection and Preprocessing

Collection. We collected descriptions of real data needs from four sources, including three online communities, namely Stack Overflow, Open Data Stack Exchange, and Reddit, and one national open data portal, namely data.gov.uk. For each online community, we leveraged its search function to retrieve posts with the query “looking for dataset”. Note that since Reddit covers a wide range of topics, our search was performed within its r/datasets community. From the search results, we manually identified 50 top-ranked posts where a clear data need was described, e.g.,

*I am looking for datasets that lists the location of accidents or traffic (latitude and longitude) with date and time in many countries. I found datasets for USA and UK, now looking for datasets for other countries. Any type of road accident would be great.*¹

For data.gov.uk, we reused 50 user-submitted data requests which were sampled and published by the authors of [5]. To summarize, a total of 200 descriptions of data needs² were collected.

Preprocessing. The description of a data need may be long and noisy. To improve the accuracy of analysis, we recruited ten human experts to summarize each description and remove irrelevant information. Specifically, each description was independently presented to every human expert. The expert summarized the data

need by reformulating it as a concise free-form *dataset search query* containing 1–20 words. For the 150 posts collected from online communities, we received 1,500 reformulated queries from ten experts; an example such query is:

location of accidents or traffic with date and time in many countries .

From the resulting set we removed 2 queries that were arguably meaningless. Moreover, the 50 data requests submitted to data.gov.uk had been processed in the same way by crowd workers [5], producing 449 queries. To summarize, a total of 1,947 dataset search queries³ were derived from the collected descriptions of data needs.

2.2 Analysis and Implications

Annotation Scheme. We propose a fine-grained scheme in Table 1 to annotate dataset search queries with their semantics. The scheme allows to distinguish between mentions of metadata and of data content. The latter is further divided into schema-level elements (i.e., concepts), instance-level elements (i.e., geospatial or other entities), and data values (i.e., temporal or other numbers). Using the proposed scheme, all the 1,947 queries have been annotated manually by two human experts. The results are accessible online⁴.

Results. Among the 1,947 queries we collected, a query in average contains 9.22 words. More than half (58.60%) of the queries contain 5–11 words. Three types of queries are identified: phrases (63.33%), keywords (31.38%), and sentences (5.29%).

As summarized in Table 1, queries usually mention some meta-data (96.05%), especially the domain/topic of interest (94.45%). Data format and accessibility also occupy notable proportions (> 5%). On the other hand, most queries (63.79%) mention the data content, mainly some schema-level concepts (50.59%), followed by instance-level geospatial entities (19.21%) which appear more often in the data requests submitted to data.gov.uk (46.55%) but less often in

¹<https://opendata.stackexchange.com/questions/11146/dataset-for-road-accidents-or-traffic>

²<http://ws.nju.edu.cn/datasetsearch/query-cikm2019/dataneeds.txt>

³<http://ws.nju.edu.cn/datasetsearch/query-cikm2019/queries.txt>

⁴<http://ws.nju.edu.cn/datasetsearch/query-cikm2019/annotations.txt>

Table 2: Human-rated Usefulness of Snippets (1–5)

	Mean \pm Standard Deviation
Query-biased Snippets	1.91 \pm 1.22
Illustrative Snippets	3.04 \pm 1.23
<i>t</i> -test: $p = 0.0127$	

other queries (11.01%). Besides, temporal information is not neglectable (9.35%). Moreover, 60.61% of the queries mention both metadata and data content.

Implications. First, it would be insufficient to only take metadata into account in indexing, ranking, and result presentation like Google Dataset Search [8]. Data content should also be considered. Second, the constitution of a dataset search query is complex, requiring novel semantic parsing techniques to process, understand, and finally improve the accuracy of search. Third, data needs collected from different sources exhibit different features. The results of analyzing a single type of data needs (e.g., those submitted to national open data portals [4, 5]) may not be generalizable.

3 SNIPPET GENERATION

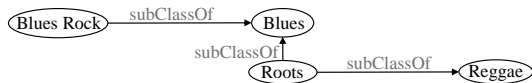
The above analysis reveals that dataset search queries mention the elements of both metadata and data content. However, current dataset search engines only present metadata about each dataset in search results pages. In order to complement this metadata and help the users to quickly identify relevant datasets, we propose to extract a snippet from data content. We now present our two such methods and report preliminary evaluation results.

3.1 Methods

We currently focus on graph data, or more precisely RDF datasets. The content of an RDF dataset is a directed graph where nodes and edges are associated with meaningful labels.

Query-biased Snippet. In the first method, a dataset search query is treated as a set of keywords after removing stop words. Each keyword is mapped to a set of nodes in the graph-structured data content. We intend to generate a query-biased snippet by extracting a subgraph that not only fully covers the query but also reflects the relationships between query keywords. We formulate the problem as finding an optimal connected subgraph that spans at least one mapped node from each query keyword. Optimality is defined by the minimization of total edge weights. We follow [3] to weight edges. This gives rise to an instance of the group Steiner tree problem (GST), and we implement a state-of-the-art algorithm [7] for solving this NP-hard problem.

Illustrative Snippet. Our second method we implement is to generate an illustrative snippet [1] which is query-independent. It is an optimal size-constrained connected subgraph extracted from data content. Optimality is defined by a linear combination of (i) covering the most frequent schema-level elements that appear in the data content, i.e., classes and properties in the RDF schema, and (ii) covering the most central instance-level elements, i.e., entities. We implement the approximation algorithm presented in [1] for solving this NP-hard problem.

**Figure 2: A query-biased snippet for the dataset *Learned Genre Ontology Intl* w.r.t. the query *blues rock reggae*.**

3.2 Preliminary Results

Design of User Study. We recruited 15 researchers/developers to participate in a user study. We crawled 311 RDF datasets from DataHub. Each participant was assigned 5 datasets and had access to their metadata and schema-level elements. For each dataset, the participant was asked to describe a data need that could be fulfilled by the dataset, by formulating a query where the keywords could be fully covered by the dataset. Then, a query-biased snippet was computed online and compared with the precomputed illustrative snippet. The size of the former was automatically determined by the algorithm, whereas the latter was constrained to have at most 20 edges (i.e., RDF triples). Both snippets were blindly visualized as two node-link diagrams. The participant rated, on a scale of 1–5, the usefulness of each snippet in supporting relevance judgment.

Results. As summarized in Table 2, illustrative snippets significantly outperform query-biased snippets ($p < 0.05$). However, they are both not very satisfying, leaving room for improvement.

Post-experiment interviews partially explained the above results. For query-biased snippets, a snippet like the one illustrated in Figure 2 could cover all the query keywords, which was appreciated by the participants. However, most participants complained that such a snippet contained very limited information, due to the minimality of group Steiner trees. Indeed, in average, this kind of snippet only covered 6% of the schema-level elements after weighting⁵.

In contrast, for illustrative snippets, most participants confirmed their richness and the diversity of information they contain. Such snippets in average covered 75% of the (weighted) classes and properties that appeared in the data content. Thus, the illustrations (or exemplifications) of schema-level concepts that form illustrative snippets were beneficial for user comprehension of the dataset. On the other hand, illustrative snippets often failed to match the keywords appearing in the query. This kind of weak relevance to the query was criticized by the participants.

Therefore, we conclude that the two types of snippets showed complementary features. A promising direction for future work would be to combine their advantages.

4 RELATED WORK

Dataset Search. Google Dataset Search [8] supports keyword search and presents the metadata for each retrieved dataset. A query is processed using the same methodology as a Web search engine, but is only matched with the metadata of a dataset. Other prototype systems [6, 9] support faceted search. All these systems exploit *metadata* but ignore the *content* of a dataset. By comparison, our proposed framework of dataset search is centered around the relationship between the query and the content of a retrieved dataset. We highlight semantic query parsing and query-biased

⁵Classes and properties are weighted by their frequencies in data content.

snippet generation in order to understand and reflect the data need underlying a query, respectively.

Snippet Generation. The utility of metadata in relevance judgment is limited because the information it spans is almost orthogonal to the content of a dataset where elements are often mentioned in a query. A promising way to complement this approach is to also present a snippet generated from the data content. IlluSnip [1] generates an *extractive* snippet by selecting a small illustrative subset of data. HIEDS [2] generates an *abstractive* snippet by producing a hierarchical grouping of the data content. However, none of these static snippets can precisely explain how the dataset is relevant to the query. We are among the first who explore query-biased snippet generation for datasets. Our preliminary results suggest combining the coverage of the data content [1] with the relevance to the query [7] for future work.

Query Analysis. Currently there is no public access to the query logs of commercial dataset search engines. Research efforts are restricted to the queries submitted to national open data portals. In [5], data requests submitted to data.gov.uk are transformed into queries by crowd workers. In [4], queries are extracted from the query logs of four national open data portals. The analysis in [4, 5] mainly considers query length and query annotation. Their annotation scheme only includes geospatial, temporal, file/data type, numbers, and abbreviations. Their *lexical* annotation is carried out automatically based on a predefined keyword mapping. They show that dataset search queries differ from general Web search queries in their length, topic, and structure. By comparison, our extended annotation scheme is more fine-grained, focusing on the *semantics* of a query. Besides, the data needs we collect are more diverse, including data requests submitted to a national open data portal as well as user-submitted posts on three online communities. Therefore, our results may exhibit better generalizability.

5 CONCLUSION AND FUTURE WORK

On the way of developing a more usable dataset search engine, we summarize our findings presented in this paper as follows.

- Real data needs mention both metadata and data content.
- The constitution of a dataset search query is complex.
- Snippet generation for dataset search should combine query relevance with schema coverage.

These empirical findings and implications support our idea of designing a query-centered framework for dataset search. Extending the methods that have been implemented and preliminarily evaluated in this work, we have identified the following further steps for future development of our system.

First of all, dataset search requires specialized query processing. The analyzed data needs have exhibited some query patterns that are captured by our annotation scheme. It inspires us to formulate automated query parsing as a sequence labeling task and we plan to solve it using supervised learning techniques. To this end, a large set of data needs or dataset search queries should be labeled as training data. This in turn may reveal new patterns and requires extending our annotation scheme, until convergence.

Query processing is tightly coupled with indexing, ranking, and snippet generation. As a query may mention both the metadata

and the content of a target dataset, it would be desirable to consider both of them in retrieval models and snippets. Metadata is relatively easy to handle, as it can be viewed as a semi-structured document and hence existing Web search methods may apply. By contrast, indexing, ranking, and summarizing data content pose new challenges related to both effectiveness and efficiency, as a dataset is much larger than a webpage.

With all these components being implemented, we plan to assemble a prototype of a new dataset search engine, and conduct user study to evaluate its end-to-end performance. We will start with RDF datasets, and progressively extend to other formats.

ACKNOWLEDGMENTS

This work was funded partially by the NSFC under Grant 61772264, and partially by the SIRIUS Centre, Norwegian Research Council project number 237898. Gong Cheng was supported by the Six Talent Peaks Program of Jiangsu Province under Grant RJFW-011.

REFERENCES

- [1] Gong Cheng, Cheng Jin, Wentao Ding, Danyun Xu, and Yuzhong Qu. 2017. Generating Illustrative Snippets for Open Data on the Web. In *WSDM 2017*. 151–159.
- [2] Gong Cheng, Cheng Jin, and Yuzhong Qu. 2016. HIEDS: A Generic and Efficient Approach to Hierarchical Dataset Summarization. In *IJCAI 2016*. 3705–3711.
- [3] Bolin Ding, Jeffrey Xu Yu, Shan Wang, Lu Qin, Xiao Zhang, and Xuemin Lin. 2007. Finding Top-k Min-Cost Connected Trees in Databases. In *ICDE 2007*. 836–845. <https://doi.org/10.1109/ICDE.2007.367929>
- [4] Emilia Kacprzak, Laura Koesten, Luis Daniel Ibáñez, Tom Blount, Jeni Tennison, and Elena Simperl. 2019. Characterising dataset search - An analysis of search logs and data requests. *J. Web Semant.* 55 (2019), 37–55. <https://doi.org/10.1016/j.websem.2018.11.003>
- [5] Emilia Kacprzak, Laura Koesten, Jeni Tennison, and Elena Simperl. 2018. Characterising Dataset Search Queries. In *WWW 2018*. 1485–1488. <https://doi.org/10.1145/3184558.3191597>
- [6] Sven R. Kunze and Sören Auer. 2013. Dataset Retrieval. In *IEEE ICSC 2013*. 1–8. <https://doi.org/10.1109/ICSC.2013.12>
- [7] Rong-Hua Li, Lu Qin, Jeffrey Xu Yu, and Rui Mao. 2016. Efficient and Progressive Group Steiner Tree Search. In *SIGMOD 2016*. 91–106. <https://doi.org/10.1145/2882903.2915217>
- [8] Natasha Noy, Matthew Burgess, and Dan Brickley. 2019. Google Dataset Search: Building a search engine for datasets in an open Web ecosystem. In *WWW 2019*. 1365–1375. <https://doi.org/10.1145/3308558.3313685>
- [9] Emmanuel Pietriga, Hande Gözükan, Caroline Appert, Marie Destandau, Sejla Cebiric, François Goasdoué, and Ioana Manolescu. 2018. Browsing Linked Data Catalogs with LODAtlas. In *ISWC 2018, Part II*. 137–153. https://doi.org/10.1007/978-3-030-00668-6_9
- [10] Nikolai Toupikov, Jürgen Umbrich, Renaud Delbru, Michael Hausenblas, and Giovanni Tummarello. 2009. DING! Dataset Ranking using Formal Descriptions. In *LDDW 2009*.