# Towards Multi-Facet Snippets for Dataset Search

**Xiaxia Wang**[1], Gong Cheng[1], Evgeny Kharlamov[2]

[1] National Key Laboratory for Novel Software Technology, Nanjing University, China
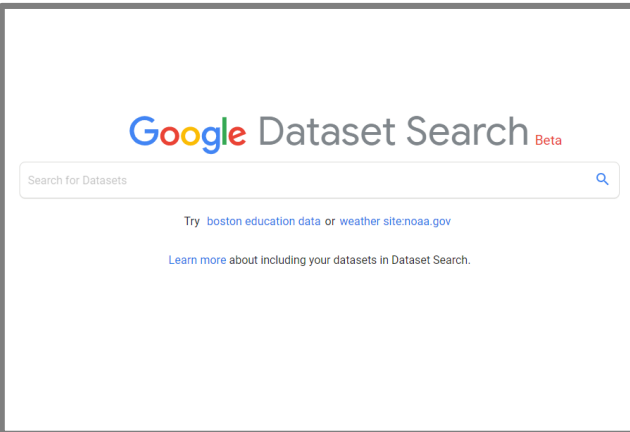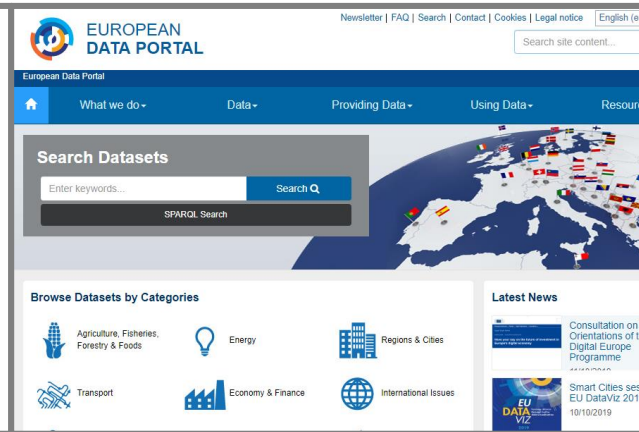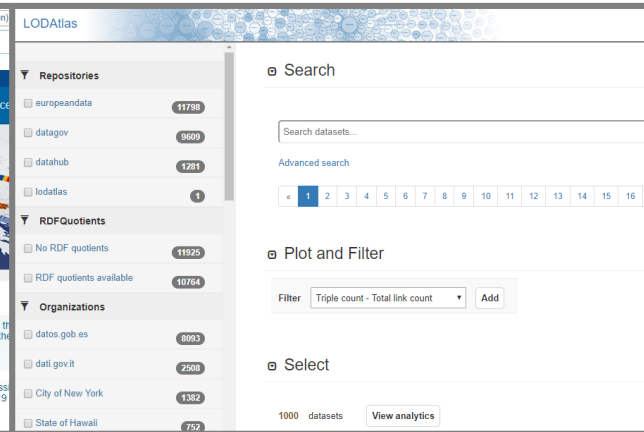[2] Bosch Center for Artificial Intelligence, Renningen, Germany

# Dataset search systems: Conveniently find relevant datasets.
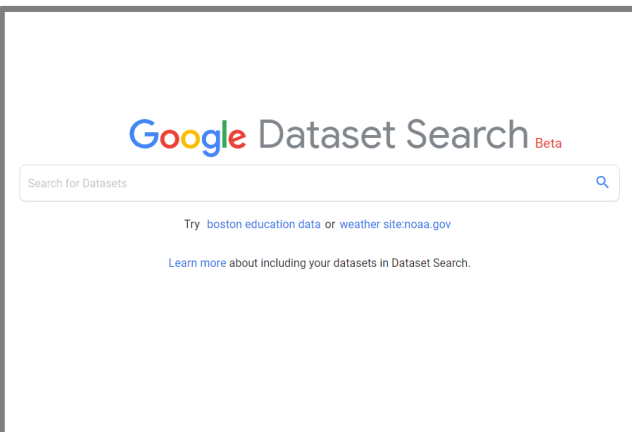


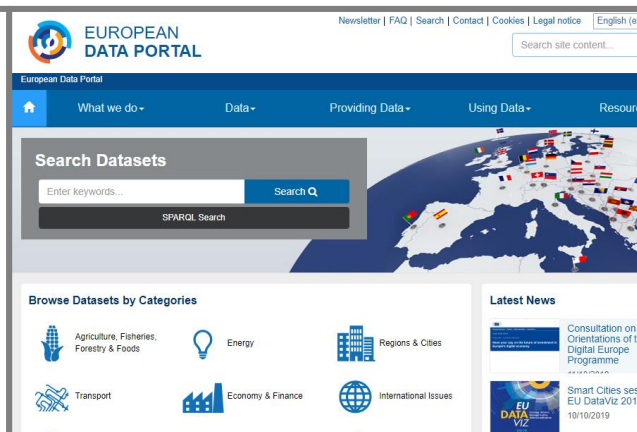*Google Dataset Search*

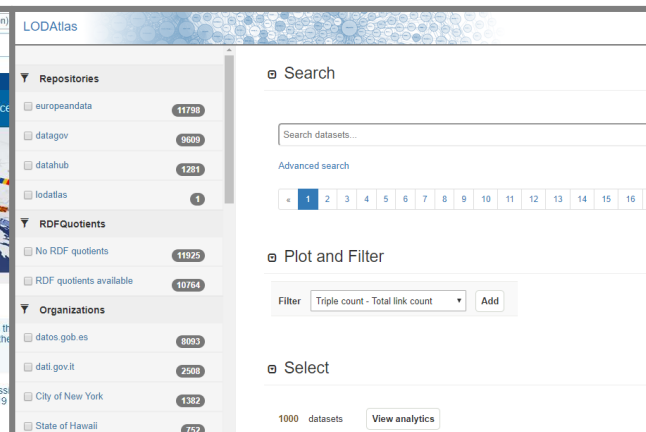*European Data portal*

*LODAtlas*

# Dataset search systems: Conveniently find relevant datasets.



*Google Dataset Search*          *European Data portal*          *LODAtlas*

# Existing systems provide only metadata for users.

# Metadata:

- No detailed information of dataset content
- <span style="color:red">Limited utility</span> for users to judge the relevance.



*Google Dataset Search system*

# Dataset Snippet: Complementary to metadata

- a size-limited subset of triples
- exemplify dataset content
- illustrate the relevance to the query

# Dataset Snippet: Complementary to metadata

- a size-limited subset of triples

- exemplify dataset content

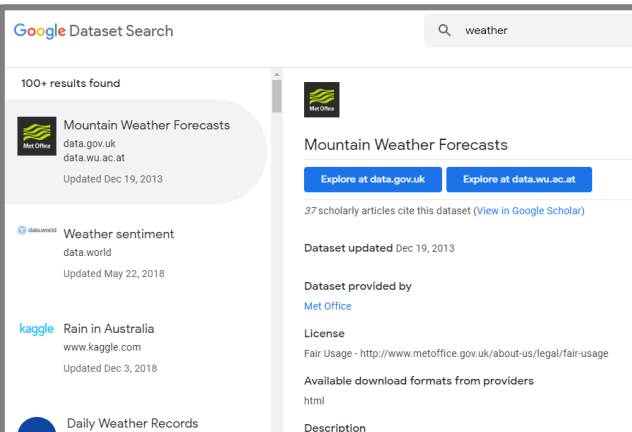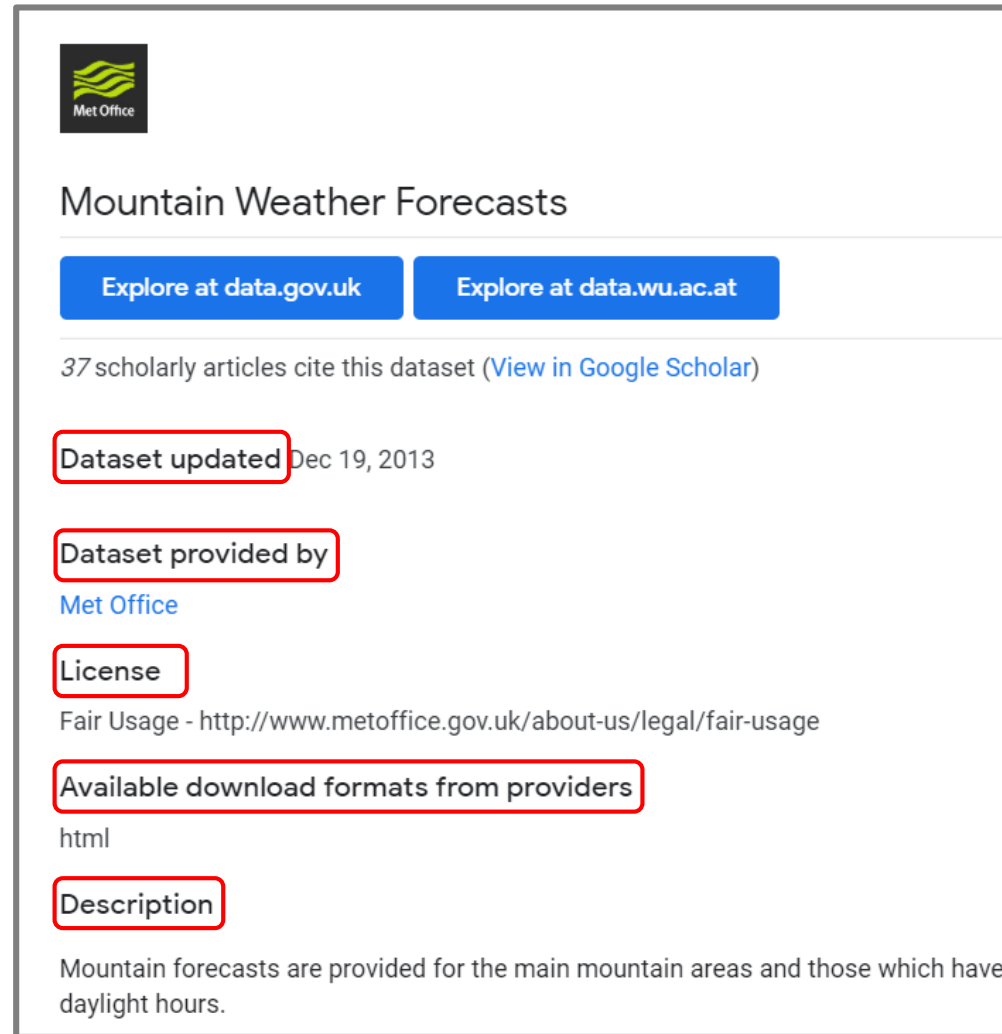- illustrate the relevance to the query

⟨Augsburg-TYPE-City⟩
⟨Berlin-capitalOf-Germany⟩
⟨Berlin-locatedIn-Germany⟩
⟨Berlin-neighboringCity-Dresden⟩
⟨Berlin-TYPE-Capital⟩
⟨Berlin-TYPE-City⟩
⟨Germany-isPartOf-CentralEurope⟩
⟨Germany-TYPE-Country⟩
⟨Munich-locatedIn-Germany⟩
⟨Munich-TYPE-City⟩
⟨Munich-neighboringCity-Augsburg⟩

*An RDF dataset*

Query:
Munich Europe

Snippet
Generation

⟨Germany-isPartOf-CentralEurope⟩
⟨Munich-locatedIn-Germany⟩

CentralEurope
↑ isPartOf
Germany ← locatedIn ← Munich

*An example snippet*

# Dataset Snippet: Complementary to metadata

- a **size-limited subset** of triples —— Size limit $k$

- **exemplify** dataset content —— Contain central elements of Schema and Instance

- **illustrate** the relevance to the query —— Contain query keywords

⟨Augsburg-TYPE-City⟩
⟨Berlin-capitalOf-Germany⟩
⟨Berlin-locatedIn-Germany⟩
⟨Berlin-neighboringCity-Dresden⟩
⟨Berlin-TYPE-Capital⟩
⟨Berlin-TYPE-City⟩
⟨Germany-isPartOf-CentralEurope⟩
⟨Germany-TYPE-Country⟩
⟨Munich-locatedIn-Germany⟩
⟨Munich-TYPE-City⟩
⟨Munich-neighboringCity-Augsburg⟩

*An RDF dataset*

Query: Munich Europe

**Snippet Generation**

⟨Germany-isPartOf-CentralEurope⟩
⟨Munich-locatedIn-Germany⟩

CentralEurope
↑ isPartOf
Germany ← locatedIn ← Munich

*An example snippet*

# Problem Formulation: A Weighted Maximum Coverage Problem

- Input: a collection of sets
- Select at most $k$ sets to maximize the total weight of covered elements
- Consider *keywords, classes, properties, entities* as elements



$$U = Q \cup \mathrm{Cls}(T) \cup \mathrm{Prp}(T) \cup \mathrm{Ent}(T)$$

| Snippet | | WMC Problem |
|---|---|---|
| Each triple | ⟷ | A set |
| keyword  class  property  entity | ⟷ | Weighted elements |
| Size limit | ⟷ | At most $k$ sets |
| Relevant to keywords and content | ⟷ | Maximum weight |

# Problem Formulation: A Weighted Maximum Coverage Problem

Optimize the coverage of Keywords, Classes, Properties, and Entities.



| Snippet | | WMC Problem |
|---|---|---|
| Each triple | ⟷ | A set |
| keyword class property entity | ⟷ | Weighted elements |
| Size limit | ⟷ | At most $k$ sets |
| Relevant to keywords and content | ⟷ | Maximum weight |

$$U = Q \cup \mathrm{Cls}(T) \cup \mathrm{Prp}(T) \cup \mathrm{Ent}(T)$$

# Experiment

Datasets: 311 real datasets from *Datahub*

Queries:
- Real queries from *data.gov.uk*
- Artificial queries selected from *DMOZ*

Evaluation Metrics:
- Coverage of keywords
- Coverage of Paths between keywords
- Coverage of Dataset Schema
- Coverage of Data instance

Details about our Evaluation Metrics & Baselines Tomorrow 12:00, Session 1B, FPAA level 0.

# Result

- Our approach achieved a Balance among four evaluation metrics

| | coKyw | coCnx | coSkm | coDat | Average |
|---|---|---|---|---|---|
| IlluSnip | 0.1000 | 0.0540 | 0.6820 | 0.3850 | 0.3053 |
| TA+C | 0.9590 | 0.4703 | 0.0425 | 0.0915 | 0.3908 |
| PrunedDP++ | **1** | **1** | 0.0898 | 0.2133 | 0.5758 |
| CES | 0.9006 | 0.3926 | 0.3668 | 0.2684 | 0.4821 |
| **coKSD** | 0.8352 | 0.3595 | **0.8651** | **0.4247** | **0.6211** |

- Our approach got consistent scores on different query groups

| | coKyw | coCnx | coSkm | coDat | Average |
|---|---|---|---|---|---|
| data.gov.uk | 0.7643 | 0.2882 | 0.8249 | 0.3870 | 0.5661 |
| DMOZ-1 | 0.8977 | 0.7955 | 0.8873 | 0.4726 | 0.7633 |
| DMOZ-2 | 0.8433 | 0.2444 | 0.8710 | 0.4569 | 0.6039 |
| DMOZ-3 | 0.8395 | 0.2337 | 0.8693 | 0.4145 | 0.5893 |
| DMOZ-4 | 0.7936 | 0.1877 | 0.8521 | 0.3731 | 0.5516 |

# Conclusion

- Our approach (coKSD) achieved a Balance between evaluation metrics, can be used as a better dataset snippet than existing baselines.

# Future Work

- Better snippet for dataset search
- Faster generation process

# Thanks for your time!
# Q&A

Contact: xxwang@smail.nju.edu.cn

- ## Keywords

  ### To match user's data needs as much as possible

$$\text{maximizes } q(S) = \sum_{x \in \cup_{t_i \in S} \text{cov}(t_i)} w(x), \qquad \text{subject to } |S| \le k$$

$$w(x) = \begin{cases} \alpha \cdot \dfrac{1}{|Q|}, & x \in Q \\[2mm] \beta \cdot \text{frqCls}(x), & x \in \text{Cls}(T) \\[2mm] \beta \cdot \text{frqPrp}(x), & x \in \text{Prp}(T) \\[2mm] \gamma \cdot \left( \dfrac{\log(d^+(x) + 1)}{\sum_{e \in \text{Ent}(T)} \log(d^+(e) + 1)} + \dfrac{\log(d^-(x) + 1)}{\sum_{e \in \text{Ent}(T)} \log(d^-(e) + 1)} \right), & x \in \text{Ent}(T) \end{cases}$$

Query:
Munich Europe

⟨Augsburg-TYPE-City⟩
⟨Berlin-capitalOf-Germany⟩
⟨Berlin-locatedIn-Germany⟩
⟨Berlin-neighboringCity-Dresden⟩
⟨Berlin-TYPE-Capital⟩
⟨Berlin-TYPE-City⟩
⟨Germany-isPartOf-CentralEurope⟩
⟨Germany-TYPE-Country⟩
⟨Munich-locatedIn-Germany⟩
⟨Munich-TYPE-City⟩
⟨Munich-neighboringCity-Augsburg⟩

CentralEurope
isPartOf
Germany — locatedIn — Munich

- **Classes and properties**

  To exemplify central schema elements

$$\text{maximizes } q(S) = \sum_{x \in \bigcup_{t_i \in S} \text{cov}(t_i)} w(x), \qquad \text{subject to } |S| \le k$$

$$w(x) = \begin{cases} \alpha \cdot \dfrac{1}{|Q|}, & x \in Q \\ \beta \cdot \text{frqCls}(x), & x \in \text{Cls}(T) \\ \beta \cdot \text{frqPrp}(x), & x \in \text{Prp}(T) \\ \gamma \cdot \left( \dfrac{\log(d^+(x) + 1)}{\sum_{e \in \text{Ent}(T)} \log(d^+(e) + 1)} + \dfrac{\log(d^-(x) + 1)}{\sum_{e \in \text{Ent}(T)} \log(d^-(e) + 1)} \right), & x \in \text{Ent}(T) \end{cases}$$

⟨Augsburg-TYPE-**City**⟩
⟨Berlin-capitalOf-Germany⟩
⟨Berlin-**locatedIn**-Germany⟩
⟨Berlin-**neighboringCity**-Dresden⟩
⟨Berlin-TYPE-Capital⟩
⟨Berlin-TYPE-**City**⟩
⟨Germany-isPartOf-CentralEurope⟩
⟨Germany-TYPE-Country⟩
⟨Munich-**locatedIn**-Germany⟩
⟨Munich-TYPE-**City**⟩
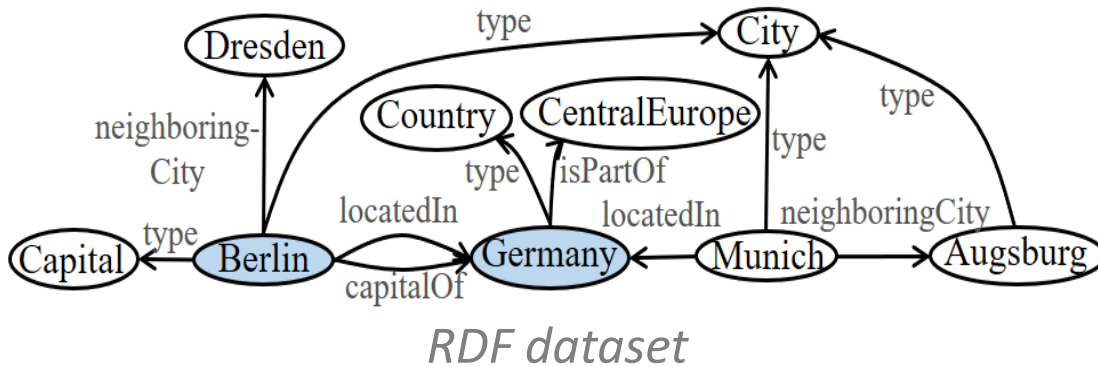⟨Munich-**neighboringCity**-Augsburg⟩

Central schema elements:

- Frequent classes & properties

  - e.g., **City**: 3 times

    **locatedIn**: 2 times

    **neighboringCity**: 2 times

15

$$\text{maximizes } q(S) = \sum_{x \in \bigcup_{t_i \in S} \text{cov}(t_i)} \text{w}(x), \qquad \text{subject to } |S| \le k$$

$$w(x) = \begin{cases} \alpha \cdot \dfrac{1}{|Q|}, & x \in Q \\[2mm] \beta \cdot \text{frqCls}(x), & x \in \text{Cls}(T) \\[2mm] \beta \cdot \text{frqPrp}(x), & x \in \text{Prp}(T) \\[2mm] \gamma \cdot \left( \dfrac{\log(d^+(x) + 1)}{\sum_{e \in \text{Ent}(T)} \log(d^+(e) + 1)} + \dfrac{\log(d^-(x) + 1)}{\sum_{e \in \text{Ent}(T)} \log(d^-(e) + 1)} \right), & x \in \text{Ent}(T) \end{cases}$$



*RDF dataset*

Central entity:

- High <span style="color:red">in-degree</span> and <span style="color:red">out-degree</span>

# Approach

/*Greedy Algorithm: **coKSD***/

Input: A dataset $T$ , a keyword query $Q$, a size bound $k$

Output: An optimum dataset snippet $S \subseteq T$

1. $S \leftarrow \emptyset$;

2. while $|S| < k$ do

3.    $t^* \leftarrow argmax_{t \in (T \backslash S)} \big( q(S \cup \{t\}) - q(S) \big)$;

4.    $S \leftarrow S \cup \{t^*\}$;  *//select an optimal triple in each step*

5. end while

6. return $S$;

Approximation ratio: $1 - \frac{1}{e}$