

# SUPPLEMENTARY MATERIAL

## Propagating Visual Designs to Numerous Plots and Dashboards

Saiful Khan, University of Oxford, UK  
 Phong H. Nguyen, Redsift Ltd., UK  
 Alfie Abdul-Rahman, King's College London, UK  
 Benjamin Bach, University of Edinburgh, UK  
 Min Chen, University of Oxford, UK  
 Euan Freeman, University of Glasgow, UK  
 Gagatay Turkyay, University of Warwick, UK

### A ALGORITHMS FOR COMPUTING SIMILARITY MATRICES

In this section we will describe the process of computing similarity matrices  $S_{rd}$  and  $S_{dd}$  defined in the paper and shown in Fig. 8. Matrix  $S_{rd}$  represents the pairwise similarity between the reference data streams and discovered data streams, and matrix  $S_{dd}$  represents pairwise similarities between each discovered data stream.

As described in the paper, each data stream (an OntoData instance) has several attributes: i.e., description, keywords, data type and API endpoint. Each attribute is a feature and we denote the number of features as  $f$ . The reference data streams (an ordered list) is denoted as a matrix  $R \in \mathbb{R}^{k \times f}$ , where  $k$  is the number of reference data streams. Discovered data streams are denoted as a matrix  $D \in \mathbb{R}^{n \times f}$ , where  $n$  is the number of discovered data streams.

#### A.1 Computing Matrix $S_{rd}$

The ordering algorithm will use  $S_{rd}$  to sort the data streams within each group. To compute  $S_{rd}$ , we first create a similarity matrix for each data stream feature and discovered data streams, then aggregate the four matrices.

**Data type similarity.** We derive a feature vector  $r^{(t)} \in \mathbb{R}^k$ , where  $r^{(t)}$  corresponds to the data type column of  $R$ ;  $r^{(t)} = [r_1^{(t)}, r_2^{(t)}, \dots, r_k^{(t)}]$ . We derive another feature vector  $d^{(t)} \in \mathbb{R}^n$ , where  $d^{(t)}$  corresponds to the data type column of  $D$ ;  $d^{(t)} = [d_1^{(t)}, d_2^{(t)}, \dots, d_n^{(t)}]$ .

A function  $\phi$  computes the pairwise similarity matrix between the feature vectors  $r^{(t)}$  and  $d^{(t)}$  and is defined in Equation 1.

$$\phi(r^{(t)}, d^{(t)}) \in \mathbb{R}^{k \times n} = \begin{cases} 1, & \text{if data types are similar} \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

**Keyword similarity.** Similar to the data type vectors, we derive keyword feature vectors  $r^{(w)}$  and  $d^{(w)}$  which corresponds to the keyword column of matrix  $R$  and  $D$  respectively. The keywords attribute of a data stream contains a subset of a set of all keywords used to define data streams in the system. Therefore, Jaccard [17, 35] similarity measurement function is used to compute the pairwise similarity matrix. A function  $\psi$  computes the size of the intersection divided by the size of the union of two keywords sets, defined in Equation 2.

$$\psi(r^{(w)}, d^{(w)}) \in \mathbb{R}^{k \times n} = \frac{(r^{(w)} \cup d^{(w)})}{(r^{(w)} \cap d^{(w)})} \quad (2)$$

**Description similarity.** A description field is free-form text and can be represented as a collection of words or terms. Therefore, term frequency (tf) and inverse document frequency (idf) similarity measurement algorithms [27] will be suitable here.

We derive feature vectors  $r^{(d)}$  and  $d^{(d)}$  which corresponds to the description column of  $R$  and  $D$  respectively.

Next, we derive a matrix  $U$ , where each vector  $u_i \in U$  is a tf-idf vector [27] of  $i$ -th element of  $r^{(d)}$ . Similarly, we derive another matrix  $V$ , where each vector  $v_j \in V$  is a tf-idf vector of  $j$ -th element  $d^{(d)}$ .

A function  $\omega$  computes similarity by measuring Cosine similarity [27] between  $U$  and  $V$ ; defined in Equation 3.

$$\omega(r^{(d)}, d^{(d)}) \in \mathbb{R}^{k \times n} = \frac{UV}{\|U\| \|V\|} \quad (3)$$

Given any two row vectors  $u \in U$  and  $v \in V$  the Cosine similarity between the vectors can be computed by Equation 4.

$$\frac{u \cdot v}{\|u\| \|v\|} = \frac{\sum_{x=1}^q u_x v_x}{\sqrt{\sum_{x=1}^q u_x^2} \sqrt{\sum_{x=1}^q v_x^2}} \quad (4)$$

where  $u_x$  and  $v_x$  are components of vector  $u$  and  $v$  respectively; and  $q$  is the number of components (all possible words from the description fields).

**API endpoint similarity.** We derive two feature vectors  $r^{(a)}$  and  $d^{(a)}$ , which correspond to the API endpoint attribute columns of  $R$  and  $D$  respectively. RESTful API endpoints contain textual features such as a data stream server address, API route, and URL encoded parameters. An API endpoint can provide information about a data stream, such as its data product, component, source, and type (described in Section 4). We tokenize the endpoint fields to extract their terms (words) and use similar functions used for description field similarity measurement,  $\omega(r^{(a)}, d^{(a)}) \in \mathbb{R}^{k \times n}$ , to compute a similarity matrix.

**Aggregated similarity.** We aggregate the four similarity matrices computed above. An aggregation function  $\gamma$ , defined in Equation 5, computes the aggregated matrix,  $S_{rd} \in \mathbb{R}^{k \times n}$ . This function takes a weighting average of the input matrices.

$$S_{rd} = \gamma(R, D) = [\alpha \psi(r^{(w)}, d^{(w)}) + \beta \omega(r^{(d)}, d^{(d)}) + \theta \omega(r^{(a)}, d^{(a)})] \odot \phi(r^{(t)}, d^{(t)}) \quad (5)$$

where  $\alpha$ ,  $\beta$ , and  $\theta$  are scalar constants that define the relative weights of keywords, description, and endpoint fields in the similarity measurement, where  $\alpha + \beta + \theta = 1$ . The pairwise similarity between any data type field is 0 or 1; therefore, for the keywords type we use the element-wise product (or Hadamard product),  $\odot$ , in the aggregation function.

#### A.2 Computing Matrix $S_{dd}$

The matrix  $S_{dd} \in \mathbb{R}^{n \times n}$  computes pairwise similarities between each discovered data streams. We use this matrix for creating uniform groups of similar data streams.

Computation of the matrix  $S_{dd}$  is almost similar to the computation steps applied for deriving the matrix  $S_{rd}$  in previous Section A.1. For each feature, e.g., keyword, description, and API endpoint of the matrix  $D$ , we derive three similarity matrices:  $\psi(d^{(w)}, d^{(w)})$ ,  $\omega(d^{(d)}, d^{(d)})$ , and  $\omega(d^{(a)}, d^{(a)})$  (using Equation 2 and 3). Finally, the aggregation function,  $\lambda(D, D)$  aggregates the three matrices, defined in Equation 6. A group can consist of data stream of different data types; therefore, we excluded the data type feature from the computation of  $S_{dd}$ .

$$S_{dd} = \lambda(D, D) = \alpha \psi(d^{(w)}, d^{(w)}) + \beta \omega(d^{(d)}, d^{(d)}) + \theta \omega(d^{(a)}, d^{(a)}) \quad (6)$$

### B ALGORITHMS FOR GROUPING AND ORDERING

Algorithm 1 outlines the brute-force grouping approach described in Section 7.2. Algorithm 2 outlines the graph spectral grouping approach described in Section 7.2. Algorithm 3 outlines the procedure for sorting and ranking groups.

---

**Algorithm 1:** A brute-force algorithm for grouping data streams

---

**Input:** A matrix  $S_{dd} \in \mathbb{R}^{n \times n}$ ; size of each group  $k$   
**Output:** Groups  $G \in \mathbb{R}^{k \times m}$

```

 $n \leftarrow \|S_{dd}\|$ 
/* Number  $m$  groups to construct */
 $m = \lfloor \frac{n}{k} \rfloor$ 
 $G \leftarrow \text{Null}$ 
 $p \leftarrow 1$ 
for  $i \leftarrow 1$  to  $n$  do
   $\text{visited}[i] \leftarrow \text{False}$ 
end
for  $i \leftarrow 1$  to  $n$  do
  if  $\text{visited}[i] == \text{True}$  then
    continue
  else
    /* Indices of  $k$  max items of  $i$ -th row of  $S_{dd}$  */
     $\text{max\_k\_idx} \leftarrow \text{MaxIndices}(S_{dd}[i], k)$ 
     $G[p] \leftarrow \text{max\_k\_idx}$ 
     $p \leftarrow p + 1$ 
    /* Mark the grouped items as visited */
    for  $j \leftarrow 1$  to  $k$  do
       $\text{idx} \leftarrow \text{max\_k\_idx}[j]$ 
       $\text{visited}[\text{idx}] \leftarrow \text{true}$ 
    end
  end
end
end

```

---



---

**Algorithm 2:** Spectral graph partitioning algorithm for grouping data streams

---

**Input:** A similarity matrix  $S_{dd} \in \mathbb{R}^{n \times n}$ ; size of each group  $k$   
**Output:** Groups  $G \in \mathbb{R}^{k \times m}$

- Number  $m$  groups to construct,  $m = \lfloor \frac{n}{k} \rfloor$ .
- Compute diagonal matrix  $D$  of  $S_{dd}$ .
- Compute Laplasian matrix,  $L = D - S_{dd}$ .
- Given  $m$  number of groups to create; compute the first  $m$  eigenvectors  $v_1, v_2 \dots v_m$  of  $L$ , such that  $Lv = \lambda Dv$ .
- Let  $V \in \mathbb{R}^{n \times m}$  be the matrix containing the vectors  $v_1, v_2, \dots v_m$  as columns.
- For  $i = 1, 2, \dots, n$ , let  $d_i \in \mathbb{R}^m$  be the vector corresponding to the  $i$ -th row of  $V$ .
- Finally, use K-means algorithm to cluster  $d_1, d_2, \dots, d_n$  into  $m$  groups  $g_1, g_2, \dots, g_m \in G$  and  $G \in \mathbb{R}^{k \times m}$ .

---



---

**Algorithm 3:** Algorithm for sorting and ranking groups

---

**Input:** Groups  $G \in \mathbb{R}^{k \times m}$ ; matrix  $S_{rd} \in \mathbb{R}^{k \times n}$   
**Output:** Sorted and ranked  $G' \in \mathbb{R}^{k \times m}$

```

 $m \leftarrow \|G\|$ 
 $k \leftarrow \|S_{rd}\|$ 
/* Priority queue for queuing based on ranking */
 $G' \leftarrow \text{Null}$ 
/* For each group */
for  $i \leftarrow 1$  to  $m$  do
   $\text{group} \leftarrow G[i]$ 
   $\text{group\_sorted} \leftarrow \text{Null}$ 
   $\text{group\_score} \leftarrow 0$ 
  for  $j \leftarrow 1$  to  $k$  do
     $\text{col\_idx} \leftarrow \text{group}[j]$ 
     $\text{col\_vec} \leftarrow S_{rd}[:, \text{col\_idx}]$ 
    /* The index of max element in  $\text{col\_vec}$  */
     $\text{row\_idx} \leftarrow \text{MaxIndices}(\text{col\_vec}, 1)$ 
     $\text{similarity} \leftarrow \text{col\_vec}[\text{row\_idx}]$ 
     $\text{group\_score} \leftarrow \text{group\_score} + \text{similarity}$ 
     $\text{group\_sorted}[\text{row\_idx}] \leftarrow \text{col\_idx}$ 
  end
  /* Add the sorted group to priority queue */
   $\text{PriorityQAdd}(G', \text{group\_sorted}, \text{group\_score})$ 
end

```

---

manager formulates a search query and discover a possible list of data streams. From the discovered data streams, a function (a) creates possible groups of data streams and orders each group (as described in Section 7 and Section A). For discovered data streams, another function scans the ontology to retrieve (b) all possible pages or bindings visualizing the groups of (a). From the (a) and (b) list of possible groups, both data stream groups and page groups are created. We then use similar grouping, ordering, and ranking algorithms (described in Section 7 and Section A).

## C PROPAGATING DASHBOARDS WITH LINKS

In the paper, we described the process for propagating a visualization function with multiple data streams. Propagating a dashboard is more complicated, because of the need to match data streams (OntoData instances) and web page links (OntoPage instances). If incorrectly matched, the visual designs in the dashboard would not be linked to the correct webpage.

Fig. 9 shows the process used to propagate a dashboard with data streams and links. After the infrastructure manager selects a VIS function for propagation, its reference data streams and links are retrieved from the ontology. The metadata of the data streams and links are forwarded to the UI. While, the process of propagation for data streams is described in Section 7, the process of propagating links involves additional steps.

Links are web pages (i.e., OntoPage instances) and their attributes include a VIS function and data streams (described in Section 5.1). Following the process described earlier in Section 6.1, the infrastructure

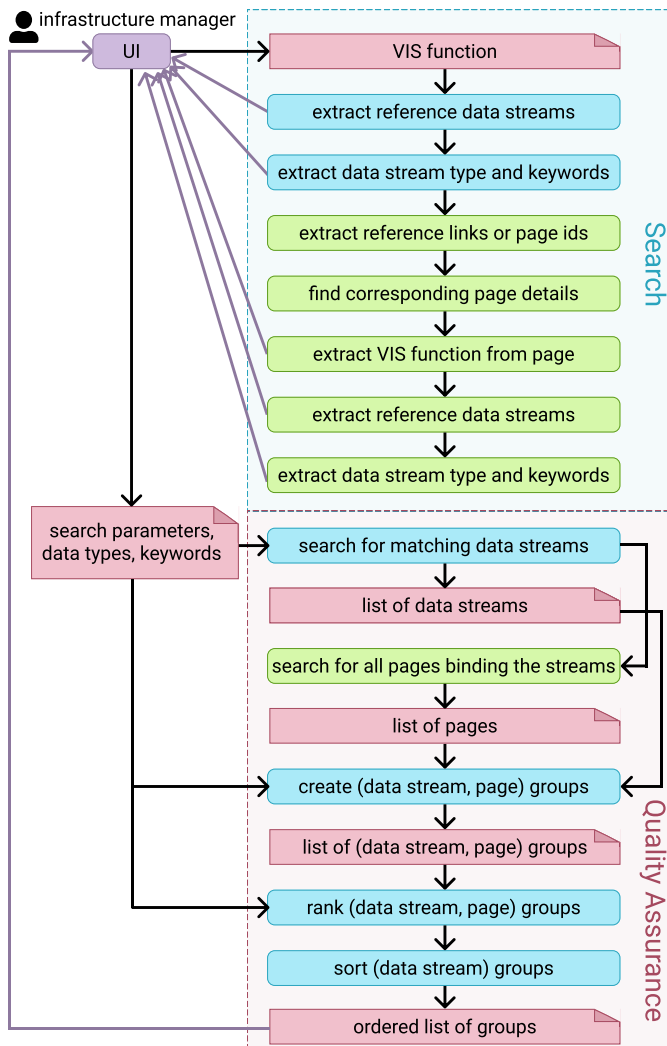


Fig. 9. A flowchart illustrating a search and ranking workflow for finding relevant links and propagating the links to a dashboard.