

Clustering Semantic Predicates in the Open Research Knowledge Graph

Omar Arab Oghli^[0000-0002-9092-9096], Jennifer D’Souza^[0000-0002-6616-9509],
and Sören Auer^[0000-0002-0698-2864]

TIB Leibniz Information Centre for Science and Technology, Hannover, Germany
{omar.araboghli, jennifer.dsouza, auer}@tib.eu

Abstract. When semantically describing knowledge graphs (KGs), users have to make a critical choice of a vocabulary (i.e. predicates and resources). The success of KG building is determined by the convergence of shared vocabularies so that meaning can be established. The typical lifecycle for a new KG construction can be defined as follows: nascent phases of graph construction experience terminology divergence, while later phases of graph construction experience terminology convergence and reuse. In this paper, we describe our approach tailoring two AI-based clustering algorithms for recommending predicates (in RDF statements) about resources in the Open Research Knowledge Graph (ORKG) <https://orkg.org/>. Such a service to recommend existing predicates to semantify new incoming data of scholarly publications is of paramount importance for fostering terminology convergence in the ORKG.

Our experiments show very promising results: a high precision with relatively high recall in linear runtime performance. Furthermore, this work offers novel insights into the predicate groups that automatically accrue loosely as generic semantification patterns for semantification of scholarly knowledge spanning 44 research fields.

Keywords: Content-based recommender systems · Open research knowledge graph · Artificial Intelligence · Clustering algorithms.

1 Introduction

Traditional, discourse-based scholarly communication in “pseudo-digitized” PDF format is being now increasingly transformed to a completely new representation leveraging semantified digital-born formats e.g. within the Open Research Knowledge Graph (ORKG) [7] among other initiatives [3,8,11,19,26,35,50]. This “digital-first” scholarly information representation is based on a fundamentally new information organization paradigm that creates and uses *structured, fine-grained scholarly content*. Specifically, in the ORKG, scholarly communication is based on a large, interconnected knowledge graph (KG) of *fine-grained scholarly content*. Such an information organization paradigm facilitates the evolution of scholarly communication from documents for humans to read towards human *and* machine-readable knowledge with the aim of alleviating human reading cognitive tie-ups. To this end, the ORKG-based scholarly communication comprises

a crucial machine-actionable unit of scholarly content in the form of *human and machine-readable comparisons* of semantified *scholarly contributions* [44]. These comparisons are meant to be used by researchers to quickly get familiar with existing work in a specific research domain. For example, determining the reproduction number estimate R0 of the Sars-Cov-2 virus from a number of studies in various regions across the world <https://orkg.org/comparison/R44930>. The semantically represented scholarly contribution comparisons in ORKG are especially necessary in our era of the deluge of peer-reviewed publications [29] and preprints [18] to help researchers stay on top of the fast-paced scientific progress. It concretely helps scientists to still keep an oversight over scientific progress by freeing unnecessary human cognitive tie-ups involved when searching for key information buried in large volumes of text.

The ORKG *machine-readable comparisons* depend on the availability of a knowledge base of *machine-actionable, semantified scholarly contributions*. The scholarly contributions are a unit of information defined in the context of the ORKG that describe the addressed problem and comprise the utilized materials, employed methods and yielded results in a scholarly article – a model which subsumes LEADERBOARDS [27,31]. A large community of researchers has recently been growing around the crowdsourced curation of *scholarly contributions* in the ORKG (e.g., <https://orkg.org/paper/R163747>).¹ To describe the scholarly contributions, RDF statements are used as structured semantic units that are machine-actionable as a result. A core semantic construct of these contribution-centric statements are the *predicates* or *properties* used to describe the contribution of an article. While the *subject* and *object* are content-based, *predicates* can generically span contributions across articles. E.g., *task name*, *dataset name*, *metric*, and *score* are a group of four predicates used to semantically describe the leaderboard contribution across AI articles [31] in the Computer Science domain; the predicates *basic reproduction number*, *confidence interval (95%)*, *location*, and *time period* are used to describe Covid-19 reproductive number estimates in epidemiology articles [43].

Predicates are a core construct for semantically describing contributions in ORKG. To base the ORKG on meaningfully described semantic scholarly contributions, certain, specific groups of predicates that can capture key contribution aspects of the scholarly articles are essential. Each such group then becomes a *contribution-centric* predicate group. Further, the group varies in applicability from being applicable to only a specific scholarly contribution or generalizing across a group of contributions from different papers. In this respect, the ORKG follows an agile, iterative Wiki-style collaboration approach giving curators the autonomy to coin new properties easily, but aims in the long-term trajectory to be coherent in terms of vocabulary for both predicates and resources. Note that contributions can only be compared based on standard predicates terminology for the *machine-readable ORKG comparisons*. Further, the typical lifecycle of a new KG construction must also be accounted which starts with nascent

¹The related construct to ORKG contributions, of LEADERBOARDS in AI <https://paperswithcode.com/> has also garnered large-scale crowdsourcing interest.

phases of graph construction experiencing terminology divergence, while later phases of graph construction aim at terminology convergence and reuse. In this background setting of building the ORKG, the overarching research question investigated in this paper is: *How to ensure that individuals, free to use arbitrary terminology, converge towards shared vocabularies for contribution-centric semantic predicates?*

Allowing users to make arbitrary statements is important, since it ensures that the expression of the diverse discoveries in Science are not being lost or unrepresented due to restricted semantic vocabularies. However, some authoring considerations need to be made. Without further considerations, the authoring freedom of contributions in the ORKG would result in statements with different vocabularies, defying the purpose of the need to semantify contributions. A terminology policy could be enforced but that would highly restrict users. Instead, a suggestion mechanism, recommending terminology based on the dataset, would help converge terminology without forcing users, as demonstrated in collaborative tagging [34,37]. In collaborative data entry, participants construct a dataset by continuously and independently adding further statements to existing data. Each curation participant faces the question: Which vocabulary elements to use? To ensure convergence, the answer is: use the most relevant and frequently occurring vocabulary elements. Finding the most frequent vocabulary elements is straightforward: one can simply count the occurrences. We therefore focus on finding the relevant vocabulary elements. Science comprises very heterogeneous contributions. Finding the vocabulary that is relevant for one contribution therefore means: finding similar contributions and reuse their vocabulary.

To this end, this work describes our implementation of an unsupervised AI service based on clustering similar papers and recommending contribution-centric predicate groups from the existing ORKG contributions. *Similar scholarly contributions should be semantified with a homogeneous contribution-centric semantic predicate groups.* This is our intuition behind adopting clustering since the method aims to group the data points having similar features, where data points in different groups should have highly offbeat features. We chose hierarchical (Agglomerative¹) and non-hierarchical (K-means²) clustering strategies. We avoid computationally intensive methods (e.g., Affinity) or methods that can handle only small cluster sizes (e.g., Spectral clustering).

In summary, the contributions of our work are:

1. a formalization of the application of homogeneous related groups of predicates to semantically describe scholarly contributions;
2. the evaluation of two contrasting flavors of clustering objectives (hierarchical and non-hierarchical) to semantify contributions based on contribution-centric predicate groups. Since the task itself is formalized for the first time in this work, the application of an AI approach is correspondingly novel;
3. detailed empirical evaluations of four machine learning model variants resulting from testing two different embedding representations; and

¹<https://scikit-learn.org/stable/modules/clustering.html#hierarchical-clustering>

²<https://scikit-learn.org/stable/modules/clustering.html#k-means>

4. the demonstration of a predicates recommendation service for the ORKG scholarly knowledge digitalization platform. Its objectives are two-fold: *i*) expedite adding a new contribution to the graph, and *ii*) semantify the contributions with a shared vocabulary. The recommender service takes as input a paper’s title and abstract and in turn recommends a group of semantically related predicates based on earlier similar semantified papers if found by the clustering method, otherwise an empty set of predicates is returned. Such a system is described for the first time.

The remainder of the paper is organized as follows. We first define the core concepts relevant to this work but which may be new in the community in Section 2. We then offer the formalized definition for our *contributions-centric predicates grouping* task in Section 3, following which, in Section 4, we explain the custom dataset created from the ORKG RDF data dump incorporating our novel task. Next, we introduce our method for the contribution-centric predicates group recommendation service in Section 5. We then show the experimental results from our methods on our created custom task corpus in Section 6. Finally, we conclude with discussions on the possibility of further improvement and future work in Section 7.

2 Definitions

We first define the central concepts to the task attempted in this work.

Contribution. Highlights the findings of a research endeavour. An ORKG *contribution* addresses a research problem, and can be described in terms of the materials and methods used and the results achieved. Contributions in different papers addressing the same research problem can be expected to have comparable semantic descriptions at least for their key properties whose values, i.e. resources and literals, then are specific to the research endeavour.

Contribution Triple. Contributions are semantically described in a series of (*subject, predicate, object*) RDF triple statements that build the ORKG.

Contribution Predicates’ Set (CPS). Is a set of predicates in contribution triples.

Comparison. The ORKG supports downstream smart applications such as the creation of comparisons/surveys over its structured contributions. In other words, given the ORKG structured contributions, it is possible to compare the values of several such machine-actionable contributions provided their CPSs are more or less similar. Comparisons can either be generated over several contributions of a single article (e.g., comparison of an AI benchmark characteristics having similar CPSs but differing values over the different data domains annotated <https://orkg.org/comparison/R163843/>); or over contributions with similar CPSs in different articles (e.g., comparison of the Covid-19 reproductive number (R0) estimate set of studies, respectively, conducted by different research groups for different countries <https://orkg.org/comparison/R44930/>).

Contribution Template. An ORKG template is a set of predicates manually specified by a domain expert to describe contributions over a specific research problem. It helps standardize the process in ORKG of semantifying contributions with similar properties but different values. E.g., the leaderboard template <https://orkg.org/template/R107801>.

Contribution Predicates' Group (CPG). A CPG is similar to a template. Where in templates, the predicates' groups are manually created by a domain expert, in the case of CPGs, they are heuristically obtained from several similar CPSS. E.g., since in the process of obtaining the crowdsourced contributions in ORKG, it may surface that some CPSS have been used similarly and repeatedly (say, above a designated threshold of contributions) to structure new contributions, they are then designated as a CPG. CPGs are indeed potential candidate templates.

The goal of this work is to recommend CPGs discovered from CPSS for describing new contributions, thereby offering ORKG users intelligent assistance in the process of structuring the new contributions. In the next section, we offer our actual task definition.

3 Task Definition

As mentioned above, the task addressed in this paper deals with discovering CPGs from CPSS in the ORKG knowledge base (KB) to describe a new article's contribution. At a high-level, given an article and the ORKG KB of crowdsourced, structured contributions w.r.t. their CPSS, the most relevant CPG, if found, should be recommended for describing the new article contribution.

Our task formalism is as follows. The ORKG KB comprising structured contributions defined only w.r.t. predicates is $CPS = \{CPS_1, \dots, CPS_N\}$ which were used to structure contributions in the set $C = \{c_1, \dots, c_N\}$, respectively. Here, the base N represents the total number of contributions in the ORKG, CPS is the knowledge base of predicates sets, and CPS_i is the predicates set used to structure contribution c_i . Furthermore, the set of predicates in each CPS formally is, $CPS_i = \{p_{1i}, p_{2i}, \dots, p_{xi}\}$. Finally, $P = \{p_1, p_2, \dots, p_y\}$ represents the set of unique predicates aggregated from all CPSS and y is the total number of unique predicates used to structure the knowledge about contributions in ORKG. The recommendation task attempted in this work can then be defined as, given a new paper P as its title T and abstract A , to semantify or describe its contribution C with an automatically discovered CPG from the ORKG KB of CPSS such that the predicates in $CPG \in P$.

4 Task Dataset

For our novel task as defined above in section 3, a novel dataset needed to be created. Our raw data source was the ORKG RDF data dump <https://orkg.org/orkg/api/rdf/dump> dated 2021-11-10. Our objective with creating the dataset

is to capture instances of the constructs of CPGs and CPSS with their respective scholarly articles’ title T and abstract A . Instantiated CPGs in a dataset can serve two purposes: 1) a supervision signal for machine learning if building a supervised recommendation system; and 2) as gold-standard data to evaluate system automatic recommendations. While obtaining CPSS and their corresponding articles may be a relatively straightforward process – for CPSS, we query the ORKG RDF data dump of contributions and for articles’ T and A , we query external services like Crossref – the process of obtaining CPGs is not.

Recall in section 2, our definition of ORKG *Comparisons*. Briefly, they are a downstream application enabling computing surveys over collections of machine-actionable, structured contributions with roughly similar CPSS. Given this and our need to obtain CPGs heuristically, we ask ourselves: could the CPSS aggregated in *Comparisons* be considered a CPG? The answer is “yes,” but with a caveat. We cannot consider the aggregated CPSS as CPGs from just about any *Comparison*. We want to generate CPGs that are strong candidates for templates. For this, we deem that the candidate CPSS need to demonstrate a strong repetition pattern of structuring several contributions as a determiner that they would apply to new contributions that have not yet been structured as well. Note here the connection between our heuristic and templates is that templates are defined with the intention of standardizing the process of structuring similar contributions on the same research problem across papers where they occur. Thus we concretely implement our heuristic as follows to generate CPGs for our dataset. The aggregated CPSS in *Comparisons* containing at least 10 contributions were considered as CPGs in our task dataset. To offer the reader better insights to the kind of *Comparisons* that were finally considered, we show in Table 1 some comparisons from whose aggregated CPSS, CPGs could be obtained. These *Comparisons* include between 10 to at most 55 structured contributions. Since structured contributions in ORKG can span Science at large, the comparisons shown in Table 1 have a diverse coverage of research fields: 1 is from the Information Science, 2 belongs to Semantic Web, 3 belongs to Bioinformatics, 4 is from Urban Studies and Planning, 5 is from Software Engineering, and 6 belongs to Natural Language Processing.

Finally, our task dataset contains a set of structured contributions as CPSS with their paper title T and abstract A . Further, these structured contributions only pertain to those from which CPGs could be obtained heuristically from *Comparisons* with the contribution included. The CPGs are also mappings to the original paper they roughly structure.

Dataset Statistics. We now offer the reader some concrete statistical insights into our task dataset. Table 2 shows the total unique papers, contributions, predicates and their research fields’ coverage. The minimum, maximum, and average numbers are aggregated by the selected *Comparisons* from which CPGs could be obtained. Thus our task dataset includes 3941 papers and 1681 unique predicates. The selected *Comparisons* have included 23.25 papers on average, with a minimum of 2 and maximum of 202. Further, contributions were structured by as few as 2 predicates and as many as 112 predicates at an average

Table 1. Example ORKG Comparisons, which are aggregations of collections of structured contributions, and correspondingly some contribution predicates in their respective contribution predicates’ groups (CPG).

	Comparison title	Predicates in <i>contribution predicate groups</i>
1.	Design and implementation of epidemiological surveillance systems https://orkg.org/comparison/R146851	epidemiological surveillance approach, epidemiological surveillance architecture, epidemiological surveillance software, epidemiological surveillance users, statistical analysis techniques
2.	Ontology learning from Folksonomies https://orkg.org/comparison/R144121	learning method, and binary valued properties as: terms learning, concepts learning, individual learning, axioms learning
3.	Review of the existing research applying Deep Learning related to mental health conditions https://orkg.org/comparison/R139050	study cohort, used models, data, outcomes, outcome assessment method, performance
4.	Enterprise architecture applications for managing digital transformation of smart cities https://orkg.org/comparison/R146458	has methodology, issues addressed, study purpose, technology deployed
5.	Overview of Approaches that Classify User Feedback as Feature Request https://orkg.org/comparison/R112387	tf-idf precision, tf-idf recall, tf-idf f1, bag-of-words precision, bag-of-words recall, bag-of-words f1
6.	NLP Datasets for Scientific Concept and Relation Extraction https://orkg.org/comparison/R150058	data domains, data coverage, dataset name, concept types, relation types, number of concepts, number of relations

Table 2. A tabular view of our task dataset statistics where the information in columns expressed w.r.t. the ORKG *Comparisons* that were considered for the dataset.

-	Papers	Contributions	Predicates	Research Fields
Minimum per Comparison	2	10	2	1
Maximum per Comparison	202	250	112	5
Average per Comparison	23.25	35.47	12.86	1.19
Total	3941	5123	1681	44

rate of 12.86. Our corpus covers contributions from across 44 different research fields. Figure 1 depicts the trendline patterns of the distribution of contributions in ORKG *Comparisons* and the distribution of predicates to structure the contributions. We see they are respectively a long-tailed distribution, i.e. some comparisons are outliers in our dataset and include a large number of contributions, and on the other hand, some predicates are used most frequently to structure nearly most of the contributions. Our task dataset is publicly available at <https://doi.org/10.5281/zenodo.6513499>.

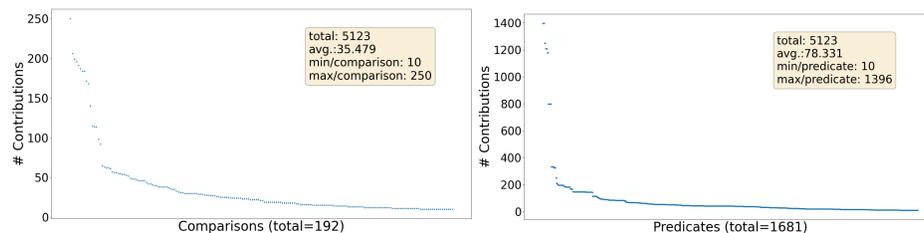


Fig. 1. Distribution of contributions in ORKG Comparisons in our task dataset (left) shown as a trendline. And the distribution of contributions over individual predicates in our task dataset shown as a trendline (right). This latter figure, in other words, shows the predicate repetition pattern for structuring contributions.

5 Background and Related Work

The Semantic Web [48] was first introduced as an extension to the World Wide Web as a comprehensive framework to create meaningful (semantic) description of information on the web as subject-predicate-object triples logical equivalent of unstructured text. Different Semantic Web languages (such as RDF, RDFS and OWL) were introduced and are now widely used standards to create structured data for intelligent computational agents [25].

Over the past decade, scholarly knowledge has started gaining traction for representation in machine-interpretable form leveraging Semantic Web standards. In other words, to semantically describe research knowledge gained among researchers from around the world. To this end, various initiatives are focused on constructing *Knowledge Graphs* (KGs) over different aspects of scholarly knowledge. One line of work concerns metadata-interlinking (e.g. authors, citations or keywords) [24,51,1,33]. Another line of work [4,9,36] supports the construction of KGs based on interlinking of research artifacts (e.g. source code, datasets or figures). In contrast, the ORKG [7] aims for contribution-centric interlinking of research resources.

Constructing a KG remains a challenge, in general, w.r.t. ensuring the graph quality and graph knowledge completeness. Four main groups of construction methods were classified by [41] as *i*) expert-based manual curation [32,38,13], *ii*) open community-based manual collaborative curation [14,49], *iii*) automated semi-structured approaches involving information extraction from structured tabulated data or using rule-based systems [6,14], and *iv*) automated unstructured approaches using machine learning to mine text [21,40,39,42,17]. In fact, the ORKG reflects all four of the information curation approaches. For instance, the creation of ORKG templates by domain experts as reusable graph patterns is an expert-based manual curation approach; the crowdsourcing of research contributions in the ORKG is a community-based collaborative approach [34,37]. As automated semi-structured methods, often rule-based approaches are also experimented with in the context of the ORKG, such as the system to acquire scientific entities from scholarly article titles [22]. Finally, the ORKG’s fully automated

text mining include the Leaderboard extraction system [30] and Computer Science Named Entity Recognition [23].

The predicates recommendation system described in this paper would then be a service offered for community-based curation of the ORKG to ensure repetition of the contribution graph patterns toward terminology convergence and comparability of the structured descriptions. Relatedly then, in terms of approaches, [45] introduced two simple algorithms. One based on generic similarity metrics [15,20], containment and resemblance, used to classify predicates similar to the query resource as a ranked list. And another approach based on computing co-occurrence matrix of predicates of resources. Despite the linear runtime performance of the algorithms, they rely on a structured query resource, while we are in need of a predicate suggestion service based on unstructured texts. To this end, [2] offered a clustering methodology using K-means [28] to semantify descriptions of Biological Assays. Our approach is modelled after this system. Finally, while at first glance, topic modeling [12] may seem similar to a clustering method, we observe that topic modeling would help us obtain distribution of a paper over topics, whereas our objective is instead to assign a paper to a single semantic group informed by its contribution, which we achieve via clustering where each paper is assigned to only a single cluster. Nevertheless, to ensure experimental completeness, we show results from a topic modeling baseline.

6 Our Approach

6.1 Clustering of Contribution Predicates' Groups

Earlier in the task dataset section (section 4), we first described our heuristic reliance on ORKG *Comparisons* to obtain CPGs. The next question is: how can we develop a recommender of CPGs given a corpus of papers as their titles T and abstracts A structured for their contributions with CPSS? We propose an AI-based unsupervised clustering strategy of papers as the solution. With this approach, we aim to automatically obtain CPGs by aggregating all CPSS in a particular cluster of similar papers. Our hypothesis is that papers describing similar contributions are also similar to each other in terms of their unstructured text descriptions as T and A . Thus, the role played by the construct of *Comparisons* in generating CPGs are now replaced, in the context of an automated recommender, by a clustering algorithm.

6.2 Grouped Predicates Recommender System Workflow

Our automated recommender system workflow with clustering is as follows.

1. A user provides the paper's title and/or DOI they wish to add.
2. We fetch the paper's abstract from external service APIs as discussed in the dataset section (section 4).
3. The paper's title and abstract are concatenated and vectorized.

4. The vector representation will be fed to a pre-trained clustering model for most relevant cluster prediction. Note, each candidate cluster was constructed based on prior semantified papers already in the ORKG KB.
5. We fetch the predicates, i.e. CPSS, of all the structured paper contributions in the predicted cluster.
6. All CPSS are combined to a set to produce a CPG which is then recommended to the user for their query paper.

This workflow is illustrated in Figure 2 below.

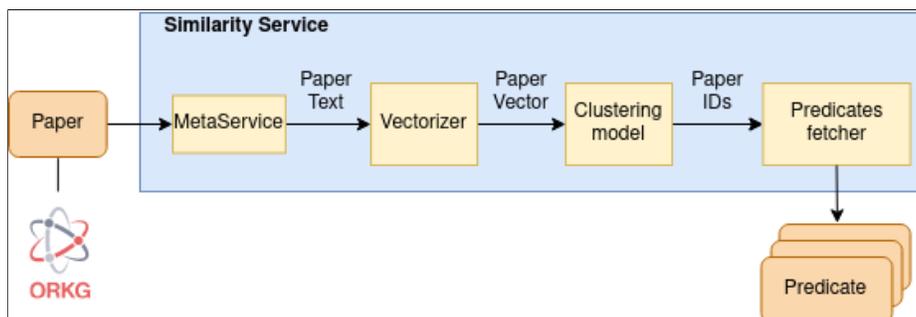


Fig. 2. Our grouped predicates recommender workflow. Arrows indicate the data flow.

For the vectorizer module, we experimented with two different vectorization functions. And, for the clustering model, we tried two different clustering algorithms. These experimental details are provided next.

6.3 Vectorization Functions

We rely on two vectorization methods. 1) TF-IDF - vectors are directly computed from our task corpus. 2) SciBERT [10] - vectors are computed from a pretrained model of embeddings over large-scale publications' data.

TF-IDF embeddings - We use the scikit-learn [16,46] library to convert our corpus of paper T and A into TF-IDF [47] vectors. TF-IDF vectors are n -dimensional real-valued vectors representing a given text with the term frequency-inverse document frequency (TF-IDF) value for each possible term in the corpus. 260,016 unique terms were found in our corpus.

SciBERT embeddings - We feed forward the pre-trained AllenNLP SciBERT [10] uncased model with our text corpus of paper T and A to output its final hidden state, which is then averaged via sentence transformers (<https://huggingface.co/sentence-transformers>) resulting in a vectorized text of dimension 768. We obtain the embeddings using a max sequence length of 512.

6.4 Clustering Algorithms

We rely on two complementary variants of clustering methods: one based on non-hierarchical clustering, specifically K-means; and another based on a hierarchical clustering, specifically agglomerative clustering.

K-MEANS. Following [2], we apply the centroid-based clustering algorithm K-means [28] to group similar scholarly contributions represented by their paper T and A . The scikit-learn implementation (<https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>) was leveraged and the models were trained on the *Google Colab Pro+* platform due to the complex time and space requirements of K-means.

AGGLOMERATIVE. The agglomerative bottom-up hierarchical-based clustering algorithm [52] with *ward* linkage was applied. This method, like the K-means objective function, minimizes the variance within a cluster. Again, the scikit-learn (<https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html>) implementation was used and the models were trained on *Google Colab*.

6.5 Experimental Setup

In this section, we describe our experimental setup to find the optimal vectorization method and clustering algorithm combination.

Dataset. First, we created training and test dataset splits of our task corpus. From each comparison, we split its papers in the 70:30 ratio for creating training and test datasets, respectively. The test dataset was reserved as a blind set with which the trained algorithm was queried for its predictions of clustered predicate groups. In total, our training set consisted of 3,696 contributions distributed over 192 comparisons, whereas our test set had 1,427 contributions distributed over 167 comparisons. The training and test sets contain mutually unique instances.

Evaluation Metrics. Per the standard evaluation practice of information retrieval systems, we employed the macro- as well as the micro-average [5] of the precision (P), recall (R) and F-score ($F1$).

Selecting K . K was strategically chose in the range $|C| \leq k \leq |P|$ with a step size of 50, where $C = 200$ is the set of ORKG comparisons and $P = 2050$ is the set of training papers. 38 total models were obtained per vectorization method.

Predictions. Some considerations need to be taken w.r.t. evaluating our clustering models. We put emphasis on the absence of the prediction function in the agglomerative algorithm compared to its presence in K-means that can simply assign a new incoming data instance to one of the clusters based on the distance to the centroid. In hierarchical clustering on the other hand, assigning a new data instance can entirely change the clusters because it can trigger several

mergings based on the linkage measure. In order to avoid re-building the hierarchical clusters for each test instance, we build them only once on the entire dataset and evaluate by comparing the comparisons’ predicates of the training papers included in the cluster to which a test instance is assigned with the expected ones.

7 Results and Discussion

In this section, we discuss our experimental results for selecting the optimal vectorization and clustering model pair.

7.1 Quantitative Evaluations

Baselines. We implemented two baselines each driven by a research question (RQ). **Baseline 1 RQ:** what happens if the problem were reduced to a trivial solution where clusters of contributions are created simply based on the research field? For this baseline, the contribution CPSS in our training data were grouped to form a CPG per research field of the training data contributions. The 44 different research fields in our dataset thus resulted in 44 CPGs. Thus a new incoming paper from the test dataset would be assigned the CPG of its research field created from the training dataset contributions. Row 1 in Table 5 shows the results from this baseline. We find that while a perfect recall can be obtained, such an approach is not precise. This reveals an important characteristic of our dataset: i.e., *the structure of contributions within each research field can differ significantly across papers in the same field*. **Baseline 2 RQ:** what happens when 192 topics are generated from our dataset by topic modeling [12] analogous to the 192 comparisons? To implement this baseline, topic distributions were obtained for all papers in the training dataset and each paper was assigned to the best topic. Thus CPSS were obtained per topic from which CPGs were generated. A new incoming test paper was then classified to best topic and assigned its CPG. Row 2 in Table 5 shows the results from a topic modeling based approach. The results prove to not be promising in terms of both precision and recall. This is contrary to our initial assumption that topic groups could be a correlated semantic construct of comparisons. We find no correlation can be established.

Clustering Results. Tables 3 and 4 show the results from applying K-means and Agglomerative clustering, respectively, with both tables showing results of the two vectorization methods. The best results are highlighted as bold and underlined in the respective tables.

The evaluation results point out that each clustering method prefers a different vectorization strategy. The K-means clustering algorithm (see Table 3) show that SciBERT embeddings are the preferred vectorization method obtaining 0.726 micro $F1$ and 0.781 macro $F1$ ($k = 2050$). The Agglomerative clustering algorithm (see Table 4) show that TF-IDF embeddings is the preferred vectorization method obtaining 0.804 micro $F1$ and 0.834 macro $F1$ ($k = 1300$). Thus,

Table 3. Results of automatically generating contribution predicates groups using K-Means clustering. K was chosen in the range from 200 to 2050 in step sizes of 50. The table shows the most significant results obtained in terms of P , R , and $F1$.

-	Macro-Average						Micro-Average					
-	TF-IDF			SciBERT			TF-IDF			SciBERT		
Clusters K	P	R	$F1$	P	R	$F1$	P	R	$F1$	P	R	$F1$
200	0.344	0.957	0.506	0.380	0.921	0.538	0.057	0.953	0.108	0.242	0.913	0.383
350	0.453	0.917	0.607	0.466	0.924	0.620	0.272	0.906	0.419	0.320	0.919	0.475
1100	0.632	0.868	0.731	0.625	0.881	0.731	0.447	0.838	0.583	0.480	0.886	0.623
1650	0.650	0.821	0.726	0.681	0.855	0.758	0.535	0.799	0.641	0.588	0.832	0.689
1850	0.649	0.779	0.708	0.704	0.849	0.770	0.593	0.732	0.655	0.603	0.834	0.700
2050	0.609	0.748	0.672	0.728	0.844	0.781	0.486	0.696	0.572	0.659	0.808	0.726

Table 4. Results of automatically generating contributions predicates groups using Agglomerative clustering. K was chosen in the range from 200 to 2050 in step sizes of 50. The table shows the most significant results obtained in terms of P , R , and $F1$.

-	Macro-Average						Micro-Average					
-	TF-IDF			SciBERT			TF-IDF			SciBERT		
Clusters K	P	R	$F1$	P	R	$F1$	P	R	$F1$	P	R	$F1$
200	0.521	0.970	0.678	0.309	0.031	0.056	0.160	0.979	0.275	0.198	0.032	0.055
250	0.550	0.967	0.701	0.354	0.031	0.057	0.189	0.977	0.317	0.265	0.032	0.057
350	0.592	0.955	0.731	0.390	0.030	0.056	0.239	0.964	0.383	0.312	0.032	0.058
1100	0.811	0.869	0.839	0.621	0.022	0.042	0.736	0.875	0.799	0.648	0.023	0.045
1300	0.823	0.845	0.834	0.751	0.021	0.041	0.760	0.853	0.804	0.761	0.023	0.044
1950	0.869	0.743	0.801	0.823	0.013	0.026	0.830	0.735	0.779	0.828	0.012	0.024
2000	0.874	0.733	0.797	0.823	0.013	0.026	0.835	0.727	0.777	0.828	0.012	0.024

Table 5. Overall results - Comparison between $Base_{RF}$ (Baseline Research Fields), $Base_{LDA}$ (Baseline Latent Dirichlet Allocation), K-Means and Agglomerative.

-	Macro-Average			Micro-Average		
Approach	P	R	$F1$	P	R	$F1$
$Base_{RF}$	0.186	1.0	0.250	0.028	1.0	0.055
$Base_{LDA}$	0.040	0.662	0.090	0.023	0.615	0.046
K-Means	0.728	0.844	0.781	0.659	0.808	0.726
Agglomerative	0.823	0.845	0.834	0.760	0.853	0.804

Agglomerative clustering surpasses K-means by nearly 10 points. While macro scores are evaluations at the Comparisons level, micro scores report evaluations at a more fine-grained predicates level. Based on this, our optimal model is at $k = 1300$ with the highest micro $F1$ using TF-IDF vectorization and Agglomerative clustering.

7.2 Qualitative Evaluations

We qualitatively analyze the ability of our system to regenerate ORKG *Comparisons* via clustering. In other words, are the contributions in ORKG *Comparisons*, even if contained in different clusters, distributed over pure or impure clusters? A pure cluster is one that contains contributions from only a single *Comparison*; an impure cluster is one that contains contributions from multiple *Comparisons*. We define this measure for regenerating the ORKG comparisons automatically as follows.

$$ReGen(comp) = \frac{|\{c \in C \mid c \text{ groups papers only from } comp\}|}{|C|} \quad (1)$$

We pick a representative example in our qualitative analysis. The ReGen value for “Smart cities and cultural heritage” ORKG *Comparison* which has 4 contributions originally (<https://www.orkg.org/orkg/comparison/R140131/>) is 50%. As a result of our method, this Comparison spanned 4 clusters (2 pure and 2 impure). However, observing the impure clusters closely, we noted that they included contributions from other *Comparisons* (e.g., “Smart city governance research categories analysis by references articles” or “Enterprise architecture applications for managing digital transformation of smart cities”, etc.) which were on the same research theme of “Smart Cities” and therefore had more or less very similar predicates. Thus, having not perfectly regenerable *Comparisons* by our method does not necessarily imply inaccurate predicted clusters. But this finding points to the fact that ORKG *Comparisons* are not all necessarily too semantically distinct from each other.

8 Conclusion and Future Work

Our experiments on the hierarchical Agglomerative algorithm have shown a quantitative result of 80.4% F1 and a qualitative result of similar recommendations of comparison predicates to those predefined in ORKG templates. Thus, the content-based recommender system based on clustered predicate units satisfies the templating concept of the ORKG. Overall, we offer among our methodology a semantification system for research contributions in the Semantic Web that does not limit the user autonomy, but instead directs the user to choose from an existing vocabulary, and hence prevent terminology divergence during later phases of graph construction.

As future work, the method will continue to be retrained for its clusters based on the ever-growing ORKG KB. Also, it is planned to implement a better association that is not heuristic-based to determine if indeed clustering related predicates produces templates.

Supplemental Material Statement: Dataset leveraged for constructing the clusters is available from <https://doi.org/10.5281/zenodo.6513499>. The code base for both training and evaluation is available from <https://doi.org/10.5281/zenodo.6514139>. Please check the publication descriptions for further details.

References

1. Ammar, W., Groeneveld, D., Bhagavatula, C., Beltagy, I., Crawford, M., Downey, D., Dunkelberger, J., Elgohary, A., Feldman, S., Ha, V., Kinney, R., Kohlmeier, S., Lo, K., Murray, T., Ooi, H.H., Peters, M., Power, J., Skjonsberg, S., Wang, L., Wilhelm, C., Yuan, Z., van Zuylen, M., Etzioni, O.: Construction of the literature graph in semantic scholar. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers). pp. 84–91. Association for Computational Linguistics, New Orleans - Louisiana (Jun 2018), <https://aclanthology.org/N18-3011>
2. Anteghini, M., D’Souza, J., dos Santos, V.A.P.M., Auer, S.: Easy semantification of bioassays (2021), <https://arxiv.org/abs/2111.15182>
3. Aryani, A., Poblet, M., Unsworth, K., Wang, J., Evans, B., Devaraju, A., Hausstein, B., Klas, C.P., Zapilko, B., Kaplun, S.: A Research Graph dataset for connecting research data repositories using RD-Switchboard. *Scientific data* **5**, 180099 (2018)
4. Aryani, A., Wang, J.: Research Graph: Building a Distributed Graph of Scholarly Works using Research Data Switchboard (3 2017), https://bridges.monash.edu/articles/preprint/Research_Graph_Building_a_Distributed_Graph_of_Scholarly_Works_using_Research_Data_Switchboard/4742413
5. Asch, V.V.: Macro-and micro-averaged evaluation measures [[basic draft]] (2013)
6. Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z.: Dbpedia: A nucleus for a web of open data. In: Aberer, K., Choi, K.S., Noy, N., Allemang, D., Lee, K.I., Nixon, L., Golbeck, J., Mika, P., Maynard, D., Mizoguchi, R., Schreiber, G., Cudré-Mauroux, P. (eds.) *The Semantic Web*. pp. 722–735. Springer Berlin Heidelberg, Berlin, Heidelberg (2007)
7. Auer, S., Oelen, A., Haris, M., Stocker, M., D’Souza, J., Farfar, K.E., Vogt, L., Prinz, M., Wiens, V., Jaradeh, M.Y.: Improving access to scientific literature with knowledge graphs. *Bibliothek Forschung und Praxis* **44**(3), 516–529 (2020)
8. Baas, J., Schotten, M., Plume, A., Côté, G., Karimi, R.: Scopus as a curated, high-quality bibliometric data source for academic research in quantitative science studies. *Quantitative Science Studies* **1**(1), 377–386 (2020)
9. Bechhofer, S., Buchan, I., De Roure, D., Missier, P., Ainsworth, J., Bhagat, J., Couch, P., Cruickshank, D., Delderfield, M., Dunlop, I., Gamble, M., Michaelides, D., Owen, S., Newman, D., Sufi, S., Goble, C.: Why linked data is not enough for scientists. *Future Generation Computer Systems* **29**(2), 599–611 (2013), <https://www.sciencedirect.com/science/article/pii/S0167739X11001439>, special section: Recent advances in e-Science
10. Beltagy, I., Lo, K., Cohan, A.: Scibert: A pretrained language model for scientific text. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 3606–3611 (2019)
11. Birkle, C., Pendlebury, D.A., Schnell, J., Adams, J.: Web of science as a data source for research on scientific and scholarly activity. *Quantitative Science Studies* **1**(1), 363–376 (2020)
12. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. *Journal of machine Learning research* **3**(Jan), 993–1022 (2003)
13. Bodenreider, O.: The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic Acids Res* **32**(Database issue), D267–70 (Jan 2004)

14. Bollacker, K., Evans, C., Paritosh, P., Sturge, T., Taylor, J.: Freebase: A collaboratively created graph database for structuring human knowledge. In: Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data. p. 1247–1250. SIGMOD '08, Association for Computing Machinery, New York, NY, USA (2008). <https://doi.org/10.1145/1376616.1376746>, <https://doi.org/10.1145/1376616.1376746>
15. Broder, A.Z., Glassman, S.C., Manasse, M.S., Zweig, G.: Syntactic clustering of the web. *Computer Networks and ISDN Systems* **29**(8), 1157–1166 (1997). [https://doi.org/https://doi.org/10.1016/S0169-7552\(97\)00031-7](https://doi.org/https://doi.org/10.1016/S0169-7552(97)00031-7), <https://www.sciencedirect.com/science/article/pii/S0169755297000317>, papers from the Sixth International World Wide Web Conference
16. Buitinck, L., Louppe, G., Blondel, M., Pedregosa, F., Mueller, A., Grisel, O., Niculae, V., Prettenhofer, P., Gramfort, A., Grobler, J., Layton, R., VanderPlas, J., Joly, A., Holt, B., Varoquaux, G.: API design for machine learning software: experiences from the scikit-learn project. In: ECML PKDD Workshop: Languages for Data Mining and Machine Learning. pp. 108–122 (2013)
17. Carlson, A., Betteridge, J., Kisiel, B., Settles, B., Hruschka, E.R., Mitchell, T.M.: Toward an architecture for never-ending language learning. In: Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence. p. 1306–1313. AAAI'10, AAAI Press (2010)
18. Chiarelli, A., Johnson, R., Richens, E., Pinfield, S.: Accelerating scholarly communication: The transformative role of preprints (2019)
19. Dessi, D., Osborne, F., Reforgiato Recupero, D., Buscaldi, D., Motta, E., Sack, H.: AI-KG: an Automatically Generated Knowledge Graph of Artificial Intelligence. *Lecture Notes in Computer Science* (2020)
20. Dhyani, D., Ng, W.K., S Bhowmick, S.: A survey of web metrics. *ACM Comput. Surv.* **34**, 469–503 (12 2002). <https://doi.org/10.1145/592642.592645>
21. Dong, X.L., Gabrilovich, E., Heitz, G., Horn, W., Lao, N., Murphy, K., Strohmann, T., Sun, S., Zhang, W.: Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In: The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New York, NY, USA - August 24 - 27, 2014. pp. 601–610 (2014), <http://www.cs.cmu.edu/~nlao/publication/2014.kdd.pdf>, evgeniy Gabrilovich Wilko Horn Ni Lao Kevin Murphy Thomas Strohmann Shaohua Sun Wei Zhang Jeremy Heitz
22. D'Souza, J., Auer, S.: Pattern-based acquisition of scientific entities from scholarly article titles. In: Ke, H.R., Lee, C.S., Sugiyama, K. (eds.) *Towards Open and Trustworthy Digital Societies*. pp. 401–410. Springer International Publishing, Cham (2021)
23. D'Souza, J., Auer, S.: Computer science named entity recognition in the open research knowledge graph (2022). <https://doi.org/10.48550/ARXIV.2203.14579>, <https://arxiv.org/abs/2203.14579>
24. Färber, M.: The microsoft academic knowledge graph: A linked data source with 8 billion triples of scholarly data. In: Ghidini, C., Hartig, O., Maleshkova, M., Svátek, V., Cruz, I., Hogan, A., Song, J., Lefrançois, M., Gandon, F. (eds.) *The Semantic Web – ISWC 2019*. pp. 113–129. Springer International Publishing, Cham (2019)
25. Fernández, J.D., Martínez-Prieto, M.A.: RDF Serialization and Archival, pp. 1–11. Springer International Publishing, Cham (2018), https://doi.org/10.1007/978-3-319-63962-8_286-1
26. Fricke, S.: Semantic scholar. *Journal of the Medical Library Association: JMLA* **106**(1), 145 (2018)

27. Hou, Y., Jochim, C., Gleize, M., Bonin, F., Ganguly, D.: Identification of tasks, datasets, evaluation metrics, and numeric scores for scientific leaderboards construction. arXiv preprint arXiv:1906.09317 (2019)
28. Jin, X., Han, J.: K-Means Clustering, pp. 695–697. Springer US, Boston, MA (2017), https://doi.org/10.1007/978-1-4899-7687-1_431
29. Johnson, R., Watkinson, A., Mabe, M.: The stm report. An overview of scientific and scholarly publishing. 5th edition October (2018), https://www.stm-assoc.org/2018_10_04_STM_Report_2018.pdf
30. Kabongo, S., D’Souza, J., Auer, S.: Automated mining of leaderboards for empirical ai research. In: Ke, H.R., Lee, C.S., Sugiyama, K. (eds.) Towards Open and Trustworthy Digital Societies. pp. 453–470. Springer International Publishing, Cham (2021)
31. Kabongo, S., D’Souza, J., Auer, S.: Automated mining of leaderboards for empirical ai research. In: International Conference on Asian Digital Libraries. pp. 453–470. Springer (2021)
32. Lenat, D.B.: Cyc: A large-scale investment in knowledge infrastructure. Commun. ACM **38**(11), 33–38 (nov 1995). <https://doi.org/10.1145/219717.219745>, <https://doi.org/10.1145/219717.219745>
33. Lo, K., Wang, L.L., Neumann, M., Kinney, R., Weld, D.: S2ORC: The semantic scholar open research corpus. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 4969–4983. Association for Computational Linguistics, Online (Jul 2020), <https://aclanthology.org/2020.acl-main.447>
34. Lund, B., Hammond, T., Flack, M., Hannay, T., NeoReality, I.: Social bookmarking tools (ii). D-Lib magazine **11**(4), 1–1 (2005)
35. Manghi, P., Atzori, C., Bardi, A., Shirrwagen, J., Dimitropoulos, H., La Bruzzo, S., Foufoulas, I., Löhden, A., Bäcker, A., Mannocci, A., Horst, M., Baglioni, M., Czerniak, A., Kiatropoulou, K., Kokogiannaki, A., De Bonis, M., Artini, M., Ottonello, E., Lempesis, A., Nielsen, L.H., Ioannidis, A., Bigarella, C., Summan, F.: OpenAIRE Research Graph Dump (Dec 2019). <https://doi.org/10.5281/zenodo.3516918>, <https://doi.org/10.5281/zenodo.3516918>
36. Manghi, P., Bardi, A., Atzori, C., Baglioni, M., Manola, N., Schirrwagen, J., Principe, P.: The openaire research graph data model (Apr 2019). <https://doi.org/10.5281/zenodo.2643199>, <https://doi.org/10.5281/zenodo.2643199>
37. Marlow, C., Naaman, M., Boyd, D., Davis, M.: Ht06, tagging paper, taxonomy, flickr, academic article, to read. In: Proceedings of the seventeenth conference on Hypertext and hypermedia. pp. 31–40 (2006)
38. Miller, G.A.: Wordnet: A lexical database for english. Commun. ACM **38**(11), 39–41 (nov 1995). <https://doi.org/10.1145/219717.219748>, <https://doi.org/10.1145/219717.219748>
39. Nakashole, N., Theobald, M., Weikum, G.: Scalable knowledge harvesting with high precision and high recall. In: Proceedings of the Fourth ACM International Conference on Web Search and Data Mining. p. 227–236. WSDM ’11, Association for Computing Machinery, New York, NY, USA (2011). <https://doi.org/10.1145/1935826.1935869>, <https://doi.org/10.1145/1935826.1935869>
40. Nakashole, N., Weikum, G., Suchanek, F.: Patty: A taxonomy of relational patterns with semantic types. In: Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language

- Learning. p. 1135–1145. EMNLP-CoNLL '12, Association for Computational Linguistics, USA (2012)
41. Nickel, M., Murphy, K., Tresp, V., Gabrilovich, E.: A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE* **104**(1), 11–33 (2016). <https://doi.org/10.1109/JPROC.2015.2483592>
 42. Niu, F., Zhang, C., Ré, C., Shavlik, J.: Elementary: Large-scale knowledge-base construction via machine learning and statistical inference. *Int. J. Semant. Web Inf. Syst.* **8**(3), 42–73 (jul 2012). <https://doi.org/10.4018/jswis.2012070103>, <https://doi.org/10.4018/jswis.2012070103>
 43. Oelen, A., D'Souza, J., Stocker, M., Vogt, L., Farfar, K.E., Haris, M., Fadel, K., Jaradeh, M.Y., Wiens, V.: Covid-19 reproductive number estimates (2020). <https://doi.org/10.48366/R44930>, <https://www.orkg.org/orkg/comparison/R44930>
 44. Oelen, A., Jaradeh, M.Y., Stocker, M., Auer, S.: Generate FAIR Literature Surveys with Scholarly Knowledge Graphs, p. 97–106. Association for Computing Machinery, New York, NY, USA (2020), <https://doi.org/10.1145/3383583.3398520>
 45. Oren, E., Gerke, S., Decker, S.: Simple algorithms for predicate suggestions using similarity and co-occurrence. In: Franconi, E., Kifer, M., May, W. (eds.) *The Semantic Web: Research and Applications*. pp. 160–174. Springer Berlin Heidelberg, Berlin, Heidelberg (2007)
 46. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011)
 47. Sammut, C., Webb, G.I. (eds.): *TF-IDF*, pp. 986–987. Springer US, Boston, MA (2010), https://doi.org/10.1007/978-0-387-30164-8_832
 48. Shadbolt, N., Berners-Lee, T., Hall, W.: The semantic web revisited. *IEEE Intelligent Systems* **21**(3), 96–101 (2006)
 49. Vrandečić, D., Krötzsch, M.: Wikidata: A free collaborative knowledgebase. *Commun. ACM* **57**(10), 78–85 (sep 2014). <https://doi.org/10.1145/2629489>, <https://doi.org/10.1145/2629489>
 50. Wang, K., Shen, Z., Huang, C., Wu, C.H., Dong, Y., Kanakia, A.: Microsoft academic graph: When experts are not enough. *Quantitative Science Studies* **1**(1), 396–413 (2020)
 51. Yaman, B., Pasin, M., Freudenberg, M.: Interlinking SciGraph and DBpedia Datasets Using Link Discovery and Named Entity Recognition Techniques. In: Eskevich, M., de Melo, G., Fäth, C., McCrae, J.P., Buitelaar, P., Chiarcos, C., Klimek, B., Dojchinovski, M. (eds.) *2nd Conference on Language, Data and Knowledge (LDK 2019)*. OpenAccess Series in Informatics (OASIS), vol. 70, pp. 15:1–15:8. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany (2019), <http://drops.dagstuhl.de/opus/volltexte/2019/10379>
 52. Zepeda-Mendoza, M.L., Resendis-Antonio, O.: *Hierarchical Agglomerative Clustering*, pp. 886–887. Springer New York, New York, NY (2013), https://doi.org/10.1007/978-1-4419-9863-7_1371