

Science and Technology Ontology: A Taxonomy of Emerging Topics*

Mahender Kumar^{1,*}, Ruby Rani¹, Mirko Bottarelli¹, Gregory Epiphaniou¹ and Carsten Maple¹

¹Secure Cyber Systems Research Group, WMG, University of Warwick, Coventry, United Kingdom, CV4 7AL

Abstract

Ontologies play a critical role in Semantic Web technologies by providing a structured and standardized way to represent knowledge and enabling machines to understand the meaning of data. Several taxonomies and ontologies have been generated, but individuals target one domain, and only some of those have been found expensive in time and manual effort. Also, they need more coverage of unconventional topics representing a more holistic and comprehensive view of the knowledge landscape and interdisciplinary collaborations. Thus, there needs to be an ontology covering Science and Technology and facilitate multidisciplinary research by connecting topics from different fields and domains that may be related or have commonalities. To address these issues, we present an automatic Science and Technology Ontology (S&TO) that covers unconventional topics in different science and technology domains. The proposed S&TO can promote the discovery of new research areas and collaborations across disciplines. The ontology is constructed by applying BERTopic to a dataset of 393,991 scientific articles collected from Semantic Scholar from October 2021 to August 2022, covering four fields of science. Currently, S&TO includes 5,153 topics and 13,155 semantic relations. S&TO model can be updated by running BERTopic on more recent datasets.

Keywords

Science and Technology Ontology, Unconventional Topics, BERTopic, Scientific Knowledge Graph


1. Introduction


Ontologies are a valuable tool for representing and organising knowledge about a specific topic or set of topics, using a set of concepts, relationships, and rules within the domain [1, 2]. They have many applications, including data annotation and visualisation [3], forecasting new research areas [4], and scholarly data discovery [5]. Some topic ontologies created in different domains include ACM Computing Classification System¹, Physics and Astronomy Classification Scheme (PACS)², replaced in 2016 by the Physics Subject Headings (PhySH)³, Mathematics

TEXT2KG 2023: Second International Workshop on Knowledge Graph Generation from Text, May 28 - Jun 1, 2023, co-located with Extended Semantic Web Conference (ESWC), Hersonissos, Greece

✉ mahender.kumar@warwick.ac.uk (M. Kumar); ruby.rani@warwick.ac.uk (R. Rani); mirko.bottarelli@warwick.ac.uk (M. Bottarelli); gregory.epiphaniou@warwick.ac.uk (G. Epiphaniou); CM@warwick.ac.uk (C. Maple)

ORCID 0000-0001-7082-0475 (M. Kumar)

 © 2023 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

¹The ACM Computing Classification System: <http://www.acm.org/publications/class-2012>.

²Physics and Astronomy Classification Scheme: <https://publishing.aip.org/publishing/pacs>

³PhySH - Physics Subject Headings: <https://physh.aps.org/about>.

Subject Classification (MSC)⁴, the taxonomy used in the field of Mathematics, and Medical Subject Heading (MeSH)⁵. Creating these large-scale taxonomies is a complex and costly process that often requires the expertise of multiple domain experts, making it a time-consuming and resource-intensive endeavour. Consequently, these taxonomies are often difficult to update and maintain, quickly becoming outdated as new information and discoveries emerge. As a result, the practicality and usefulness of these taxonomies are significantly limited. One of the most notable advancements in ontology generation is the development of a large-scale automated ontology known as Computer Science Ontology (CSO) [6]. CSO ontology defines a significant breakthrough in the representation of research topics in the computer science domain, providing a structured and comprehensive framework for organising and integrating knowledge but limited to computer science concepts only.

Research Challenge: Understanding the dynamics associated with unconventional topics, which present a more comprehensive and holistic perspective of the knowledge landscape and interdisciplinary collaborations, poses a considerable challenge. Constructing an ontology for such unconventional topics necessitates recognising and collecting essential concepts and relationships from multiple domains. Furthermore, unconventional topics may necessitate multidisciplinary study, necessitating the integration of information from many fields. By overcoming these challenges, there is an opportunity for researchers and academicians to study new and developing areas of science and technology, as well as facilitate interdisciplinary collaboration across varied fields.

Contribution. This paper presents preliminary work to construct a S&TO ontology that automatically generates a taxonomy of unconventional S&T topics. S&TO ontology is built by applying BERTopic to a dataset of 393,991 scientific articles collected from Semantic Scholar from October 2021 to August 2022, covering four fields of science: computer science, physics, chemistry and Engineering. Currently, S&TO includes 5,153 topics and 13,155 semantic relations. Unlike existing ontology, S&TO ontology can provide many benefits for knowledge representation and discovery, facilitating interdisciplinary research and enabling dynamic updates.

Organisation. The rest of the paper is organized as follows. Section 2 discusses the dataset and methods for constructing the proposed S&TO. Section 3 gives the proposed S&TO. The experimental results are discussed in section 4. Section 5 presents the applications and Usecases of S&TO, and the limitations of the current version are discussed in Section 6. Finally, the conclusion is given in section 7.

2. Data and Methods

2.1. Semantic Scholar

Semantic Scholar has many academic publications from various fields, including medical sciences, agriculture, geoscience, biomedical literature, and computer science. We used the RESTful Semantic Scholar Academic Graph (S2AG) API to retrieve a sample of these articles [7]. This

⁴2010 Mathematics Subject Classification: <https://mathscinet.ams.org/msc/msc2010.html>.

⁵MeSH - Medical Subject Headings: <https://www.nlm.nih.gov/mesh>.

API offers users on-demand knowledge of authors, papers, titles, citations, venues, and more. We obtained 393,991 Science and Technology articles from Semantic Scholar using the S2AG API. The API provides a dependable data source that allows users to link directly to the related page on [semanticscholar.org](https://www.semanticscholar.org), making it a convenient and accessible way to retrieve information about academic papers.

2.2. Methodologies



Figure 1: Data Flow

After downloading the dataset, we used the BERTopic method to obtain topics from the articles—some articles representing the multi-discipline need to be included as an outlier. To have these unconventional articles and reduce the outlier percentage, we adjusted parameters with BERTopic during topic clustering [8]. Table 3 lists the critical BERTopic parameters used in the taxonomy generation.

As shown in Figure 1, our suggested topic modelling workflow consists of five important steps: sentence embedding, dimension reduction, clustering, topic quality, and topic representation. The Sentence Embedding stage, in particular, involve turning textual input into numerical vectors that capture the underlying semantics of the text. Dimension Reduction is then applied to the vectors to lower their dimensionality and improve the effectiveness of the clustering procedure. Clustering is the process of combining similar vectors to generate coherent clusters of related material. The vectorisation ensures that the extracted topics are high quality, whilst topic Representation creates an interpretable summary of each topic. Overall, we provide a robust and effective approach to extracting meaningful unconventional topics from vast and heterogeneous datasets, utilising the power of BERTopic and careful parameter optimisation to assure optimal outcomes.

1. **Sentence Embedding:** We first transformed the input articles into numerical representations before analysing them. For this purpose, we utilised sentence transformers, the default embedding model used by BERTopic. This model can determine the semantic similarity of different documents. As default, BERTopic provides many pre-trained models among them we tried the following two: “all-MiniLM-L6-v2” and “paraphrase-MiniLM-L12-v2”. While various sentence embedding models are available, we opted for the “paraphrase-MiniLM-L12-v2” model in this work. This model effectively balances performance and speed, making it a good fit for our requirements. Thus, we can effectively translate textual data into numerical form and obtain relevant insights from large and diverse datasets using sentence transformers in conjunction with BERTopic.

Table 1
BERTopic Parameters

Parameters	Values
Embedding parameter	
embedding_model	sentence_model
UMAP parameters	
n_neighbors	2
n_components	5
low_memory	true
HDBSCAN parameters	
5min_cluster_size	10
metric	euclidean
cluster_selection_method	eom
prediction_data	true
min_samples	1
CountVectoriser parameters	
Vocabulary	Vocabulary
Min_df	10
c-TF-IDF parameters	
top_n_words	30
Verbose	True
min_topic_size	20
vectorizer_model	CountVectorizer
low_memory	True
calculate_probabilities	False
Diversity	0.4

Table 2
BERTopic clustering model selection for outlier reduction

Characteristics	HDBSCAN	K-means	HDBSCAN with Probability
Topic Quality	Good	Less	Good
Outlier	High	No	Low
Risk of missing unconventional topics	Less	Little	Least

2. **Dimension Reduction:** Clustering can be complex since embeddings are often high-dimensional. To solve this problem, the dimensionality of the embeddings is frequently reduced to a more practical level. We used the UMAP (Uniform Manifold Approximation and Projection) technique, representing local and global high-dimensional features in a

lower-dimensional domain [9]. “n_neighbors” and “n_components” are two important parameters in the UMAP method. These parameters have a considerable impact on the size of the generated clusters. Larger values for these factors, in particular, result in the formation of more important clusters. We got optimal clustering results and extracted important topics from the input data by carefully tweaking these parameters.

3. **Clustering:** BERTopic splits the input data into clusters of similar embeddings after the dimensionality reduction process. The clustering techniques’ accuracy directly impacts the quality of the generated topics. K-means [10], Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN) [11], and Agglomerative Clustering [12] are among the clustering techniques provided by BERTopic. The advantages and drawbacks of these clustering algorithms are summarised in Table 2, emphasising their capacity to generate high-quality topics, manage outlier percentages, and limit the danger of missing unconventional topics. HDBSCAN is a density-based clustering algorithm used to find clusters of varying densities in a dataset. It works by constructing a hierarchical tree of clusters based on the density of the data points. It starts by identifying the points with the highest density and forming a cluster around them. Then, it gradually adds lower-density points to the cluster until a natural cutoff is reached, indicating the end of the cluster. According to our findings, the HDBSCAN with the prediction_data parameter set to “True” was the best option. Our method efficiently balances the above elements, allowing us to obtain meaningful and valuable insights from large and complex datasets.
4. **Vectorisation:** The CountVectorizer technology turns text documents into vectors of phrase frequencies. However, it has significant drawbacks, such as failing to consider a specific topic’s relative relevance in an article. To fix this issue, we adopted C-TF-IDF (Class-Based TF-IDF), a variant of the classic TF-IDF (Term Frequency-Inverse Document Frequency) method that allocates weights to terms depending on their relevance to a specific class of documents. We could fine-tune the model’s performance in BERTopic by adjusting its parameters to optimise the clustering process using CountVectorizer with C-TF-IDF. This method made it possible to create higher-quality topics that are more fascinating and pertinent to the input data.
5. **Topic Representation:** BERTopic can adjust TF-IDF to work at the cluster level instead of the document level to obtain a concrete representation of topics from the bag-of-words matrix. This modified TF-IDF is called c-TF-IDF. For word x in class c , the c-TF-IDF value is:

$$w_{x,c} = |tf_{x,c}| \times \log(1 + A/f_x) \quad (1)$$

Where $tf_{x,c}$ denotes the frequency of word x in class c , f_x denotes word x across all classes, and A denotes the average number of words per class.

3. Science and Technology Ontology Generation

To create a topic network, also known as a knowledge graph, the metadata provided by the semantic scholar is utilised. The construction of topic ontologies involves the definitions of the following components:

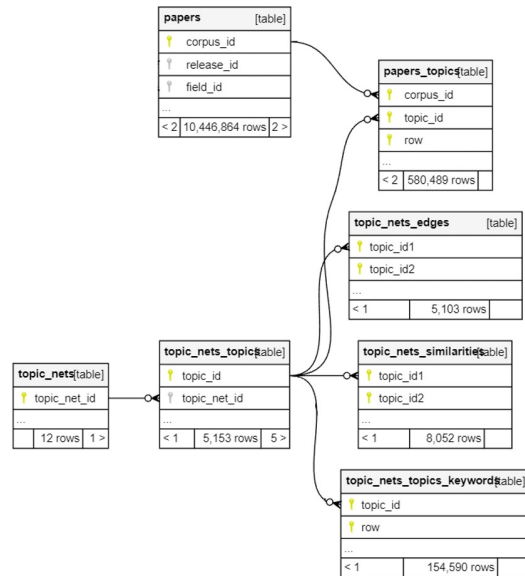


Figure 2: Schema of Topic Extraction and Topic Network Creation

- *Topics*: concepts of the topic ontology (e.g. Sports, Arts, Politics).
- *Predicates*: kinds of relationships that define the semantic link established between the ontology concepts. Many predicates can be defined in topic ontologies: hierarchical (e.g. superTopicOf) and non-hierarchical (e.g., part of, contribute to).
- *Relationships*: according to predicates and the set of elements they link, relationships are distinguished. They can be used to characterise the paths in the graphs and denoted as a triplet (T1, P, T2), where T1 and T2 denote the topics, and P denotes the predicate that links T1 and T2.

3.1. Topics

The KeyBERT tool is used on the associated publications to extract keywords representing the essential concepts and topics within each document to produce the vocabulary for BERTopic. These keywords are then sent into BERTopic, which generates a complete collection of topics that capture the overarching topic found in the dataset. The extracted topics are saved in a database's "topic_nets_topics" and "topic_nets_topics keywords" (see Figure 2). Each topic's weight denotes the number of papers for which it serves as the main association, showing its relevance within the dataset.

There are two techniques to establish the relationship between papers and topics:

1. *Probability*: First, each paper's BERTopic/HDBSCAN probabilities are saved as entries in the "papers topics" table. These probabilities indicate how closely each document relates to each extracted topic.
2. *Embedding similarity*: Second, using the "main topic id" field, the major topic associated with each paper, as identified by BERTopic, is directly linked to the paper using a SQL

trigger. This allows for efficient querying and analysis of the topics and papers related to them within the corpus.



Figure 3: An instance of topic hierarchy

3.2. Relationship

The next step is to create topic networks using the relationship among topics. Currently, the S&TO ontology is built on 393,991 scientific papers collected from Semantic Scholar from October 2021 to August 2022. It covers four science fields: computer science, physics, chemistry and Engineering. S&TO ontology follows the data model SKOS⁶ and includes the following semantic relationships:

- “relatedIdentical“, It is a sub-metric of skos related, denotes that two topics can be viewed as identical for assessing research topics. The similarity between topics is calculated as cosine-similarity in the SQL stored procedure create_topic_nets. The relationship between topics is established if the similarity threshold is above 0.9.
- “superTopicOf“: It is a sub-metric of skos:narrower, which means that a topic is a super-area of another topic in the graph. For example, "streaming_rsi_retrieval_streaming_regression" is the super-topic of Topics with topic_ids 78 and 101, as shown in Figure 3.
- “CommonArticles“: It extracts common articles that appear in the two topics. The link between two topics is evaluated as the sum of the probability distribution by common articles assigned to the topics.
- “nSimilarTopics“: It returns the top x number of similar topics for an input keyword. For instance, the top 5 similar topics related to the keyword "motor" are shown in Table 3.

⁶SKOS Simple Knowledge Organization System - <http://www.w3.org/2004/02/skos>.

Table 3

Top 5 similar topics to keyword "motor"

Topic Id	Name
33826	33826_motorized_spindle_aerostatic_restrictors
50435	50435_multimotor_geartrain_rocker_truck
62709	62709_liner_motorbike_motored_honing
49677	49677_motorized_spindle_nanocatalysts_dragging
46619	46619_powertrains_powertrain_earthmoving_2025

4. Experimental Results and Discussion

In the literature [13], ontology has been evaluated by four methods: gold-standard based [14], corpus-based [15], application-based [16], and structure-based methods [17]. The gold-standard-based method compares the developed ontology with the referenced ontology developed earlier. The corpus-based method compares the significantly developed ontology with the contents of a text corpus that covers a given domain. The application-based approach considers applications and evaluations according to their performance across use cases. The structure-based approach quantifies structure-based properties such as size and ontology complexity.

Selecting the best evaluation approach and defining the rationale behind evaluating a developed ontology is necessary. In the proposed study, the science and technology ontology is in the early stages of development and will grow in future work. Thus, the application-based approach should not be a good evaluation approach because the proposed ontology is not proper for application purposes as it is currently in development. The proposed ontology is developed on a Semantic Scholar data set subset. Thus, the best reference ontology should be a semantic scholar. However, using Semantic Scholar as the gold standard ontology is impractical due to its unavailability.

Structure-based evaluation is performed on several measures, including knowledge coverage and popularity measures (i.e., number of properties and classes) and structural measures (i.e., maximum depth, minimum depth, and depth variance). These measures are adopted based on the belief that densely populated ontologies with high depth and breadth variance are more likely to result in meaningful semantic content. Structural metrics are related to the semantic accuracy of adaptively modelled knowledge in the ontology [19].

In the context of the proposed S&TO ontology, we quantified some structural measures by considering their taxonomic structure. S&TO ontology gives 5,153 topics and 13,155 semantic relations (of which 8052 topics are based on cosine similarity and 5103 topics are based on probability distribution). Figure 4 shows an example of a knowledge graph covering topics and semantic relationships. We used neo4j as the graph database to host the final ontology [18]. Links in green indicate semantic relationships assessed using cosine similarity, while links in yellow indicate relationships based on the probability distribution of papers assigned to topics. S&TO covers the maximum amount of articles for topic clustering and gives only 15.47% of articles as an outlier, enabling the extraction of topics belonging to unconventional articles.

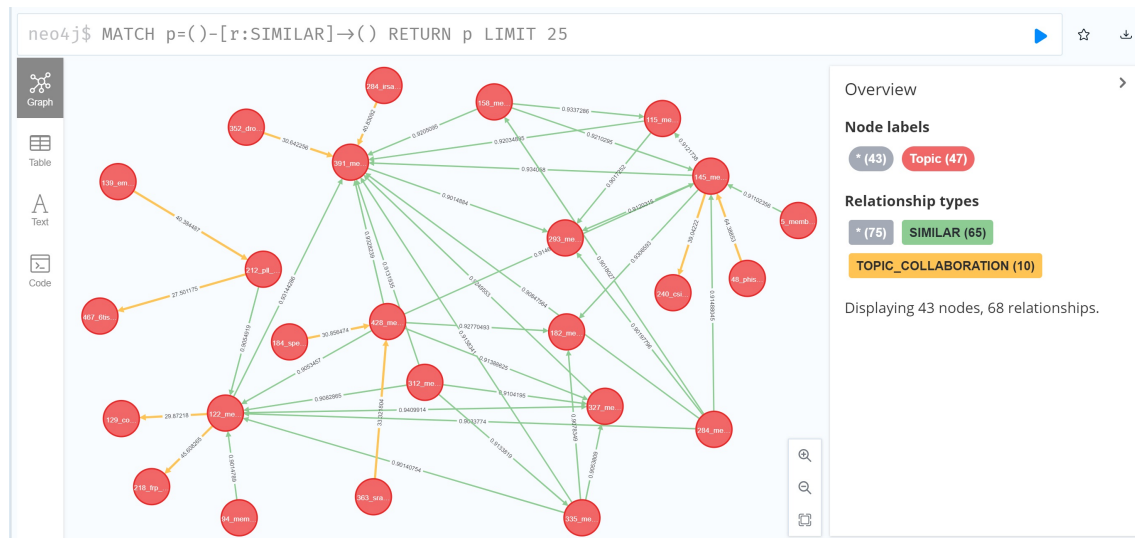


Figure 4: An instance of knowledge graph in Neo4j Browser [18]

Table 4
Structure of Table 'topic_nets'

Field	Description
topic_net_id	PK ⁷ , network unique identifier
created_on	Indicates when the network was created
status	Indicates the status of the network (NEW, DONE, etc)
year_month	Indicate the month from which this network was created

4.1. Topics Details

This section discusses the structure of topic networks, topics in topic networks and keywords related to topics in the network.

The "topic_nets" table (Table 4) gives information about the development and status of topic modelling networks. Each network is assigned a unique identifier known as a "topic_net_id," its primary key. The "created_on" parameter specifies the date and time the network was established. The "status" field offers information on the network's current status, which might take various values. This field assists in tracking the process's progress and ensuring that all networks are correctly generated and assessed. In addition, the "year_month" parameter provides the month the network was founded. This feature is beneficial for tracking the temporal evolution of themes within the corpus since it allows researchers to understand how topics and their associations change over time.

The "topic_nets_topics" table (Table 5) provides essential information about the topics associated with the Topic networks. Each topic has a unique identifier known as a "topic_id", the primary key. Each topic links to the corresponding network the "topic_nets" using a unique identifier known as a "topic_net_id", its foreign key. A descriptive label is assigned to the topic

Table 5
Structure of Table 'topic_nets_topics'

Field	Description
topic_id	PK, unique identifier of the topic
topic_net_id	FK ⁸ to the topic_nets table number Integer Topic number
label	Topic label
topic_weight	Number of papers associated with this topic
embedding	Topic embedding used for cosine similarity
similar_topics	Array of topic ids related to similar topics

Table 6
Structure of Table 'topic_nets_topics_keywords'

Field	Description
topic_id	PK, unique identifier of the topic
number	Topic number
row	Auto-incremented number, ordered by increasing score
keyword	The keyword
score	The score associated with the keyword for this topic

based on the most common terms in the associated papers. The topic has "topic_weight" which indicates the number of papers related to it, indicating its importance and relevance within the corpus. It also stores the "embedding", which will be used for cosine similarity calculations, and "similar_topics", an array of topic ids related to searched topics.

The "topic_nets_topics_keyword" table (Table 6) provides essential information about the keywords associated with the topic, represented by a unique identifier known as a "topic_id", which is the primary key. It stores the fields such as: "number" representing the topic number, "row" is an auto-incremented number, and "keyword" representing the name of the keywords. In addition, "score" represents a weight associated with the keyword for that topic.

The "papers_topics" table (table 7) illustrates the relationship between academic papers and the topics they cover. The "corpus_id" denotes the unique identifier of the corpus the papers that were extracted from. The "topic_id" represents the unique identifier of the topic that the paper covers, making it easier to track papers within that topic. The "probability" represents the weight of the paper assigned to the topic. This weight indicates the degree to which the paper covers the topic. The probability value is typically normalized, which is scaled to a range between 0 and 1.

4.2. Relation Details

This section discusses the structure of topic relations, such as information on edges and similarities among topics.

The "topic_nets_topics_edges" stores (Table 8) the information of edges among the topics "topic_nets_topics" in the network "topic_nets". The edge is established between the two topics represented by their unique identifier (i.e., "topic_id1" and "topic_id2"). "edge_Weight" is the

Table 7
Structure of Table 'papers_topics'

Field	Description
corpus_id	Part of PK, and FK to the papers table
topic_id	Part of PK, and FK to the topic_nets_topics table
row	Auto-incremented number, ordered by increasing probability
probability	Probability of the paper to be assigned to this topic

Table 8
Structure of Table 'topic_nets_topics_edges'

Field	Description
topic_id1	Part of PK, and FK to the topic_nets_topics table
topic_id2	Part of PK, and FK to the topic_nets_topics table
edge_weight	Sum of probabilities of papers sharing the two topics
str_of_col	Field weight computed as harmonic mean and normalised weight

Table 9
Structure of Table 'topic_nets_topics_similarities'

Field	Description
topic_id1	Part of PK, and FK to the topic_nets_topics table
topic_id2	Part of PK, and FK to the topic_nets_topics table
similarity	Cosine similarity between the two topics' embeddings

sum of possibilities ("possibility" field of Table 7) of papers sharing between two topics. The strength of collaboration "str_of_col" represents the weight computed as harmonic mean and normalised based on topics weights, shown in Eq (2).

$$str_of_col = HarmonicMean\left(\frac{edge_weight(id1, id2)}{topic_weight(id1)}, \frac{edge_weight(id1, id2)}{topic_weight(id2)}\right) \quad (2)$$

The "topic_nets_topics_similarity" stores (Table 9) the information of edges among the topics based on the similarity.

5. Advantages and Use-cases

The proposed S&TO with unconventional topics could have the following advantages and Use cases.

5.1. Knowledge expansion

S&TO can broaden the scope of knowledge representation beyond existing ontologies by incorporating previously unconsidered topics. Offering a more holistic and comprehensive

view of the knowledge landscape can lead to new insights and discoveries. In medical research, unconventional topics like holistic therapies or mindfulness practice might be incorporated into the ontology to provide a more comprehensive view of the more extensive health and wellness landscape [20].

5.2. Interdisciplinary Collaboration

Second, by connecting topics from diversified fields that may have commonalities or be connected, an unconventional topics ontology encourage interdisciplinary collaboration. This can encourage the discovery of new research areas and cross-disciplinary cooperation, leading to novel solutions to complicated issues. For example, an ontology incorporating computer science and psychology issues could make it easier for academics in both domains to collaborate on human-computer interaction or affective computing [21].

5.3. Scalability and Adaptability

An unconventional topics ontology has the benefit of being easily updatable and adaptable to reflect the most contemporary developments and topics, resulting in a dynamic and flexible knowledge representation system. This capability is significant in fast-paced sectors like technology and healthcare, where new topics and concepts develop regularly. For example, an ontology that includes topics relating to emerging technologies such as artificial intelligence [22] or blockchain [23] can be easily updated to include new concepts and trends.

6. Limitations

The current version of the proposed S&TO ontology has the following limitations.

6.1. Limited dataset

The current version of S&TO ontology is built on the Semantic scholar dataset covering 393,991 S&T articles from October 2021 to August 2022. However, it could be built on more datasets.

6.2. Topic labelling

Since S&TO ontology utilises BERTopic, an unsupervised topic computation library, ontology suffers from the consequences of unlabeled topics. Due to a lack of labelled data, it may be challenging to determine the significance and relevance of unconventional topics and distinguish them from noise or irrelevant topics. This is incredibly challenging when working with massive, complicated datasets containing various topics.

6.3. Domain coverage

S&TO can capture various research topics and domains by covering these four domains: computer science, physics, chemistry and Engineering. However, many other disciplines and subfields within science and technology are not yet included in S&TO. For example, biology,

environmental science, and neuroscience are all essential areas of research that could be integrated into an ontology to create a more comprehensive and multidisciplinary framework for understanding scientific research. Expanding the coverage of S&TO to include additional domains would have several potential benefits.

6.4. Topic quality

While S&TO depicts a significant effort towards organising and categorising S&T topics, there is still room for improvement regarding the quality of the topics included in the ontology.

7. Conclusion and Future Work

S&T Ontology, an automated ontology of science and technology that includes all scientific study topics, was introduced in this paper. We constructed an ontology encompassing four different science domains by utilising BERTopic on a collection of 393,991 scientific articles acquired from Semantic Scholar from October 2021 to August 2022. S&TO can be updated using BERTopic on recent datasets, offering a dynamic and flexible foundation for knowledge representation. S&TO ontology has the potential to broaden the scope of knowledge representation and stimulate interdisciplinary collaboration, making it a valuable resource for scientists and technologists.

The S&TO ontology constantly evolves and requires ongoing enhancements to meet the expanding knowledge landscape's demands. Currently, S&TO is growing, and we are employing topic labelling techniques to improve the organisation and comprehension of different topics by giving them meaningful "tags". This makes it easy for users to browse the ontology and derive valuable insights. Furthermore, we intend to improve the topic quality by investigating additional methodologies and algorithms for topic modelling and clustering. This will strengthen the ontology's accuracy and efficacy in describing the knowledge landscape. In addition, we intend to expand the ontology to a more extensive dataset, allowing for the inclusion of more unconventional categories and topics, which will improve and diversify the knowledge base.

References

- [1] H. Saif, Y. He, H. Alani, Semantic sentiment analysis of twitter, in: *The Semantic Web–ISWC 2012: 11th International Semantic Web Conference*, Boston, MA, USA, November 11-15, 2012, Proceedings, Part I 11, Springer, 2012, pp. 508–524.
- [2] F. Osborne, A. Salatino, A. Birukou, E. Motta, Automatic classification of springer nature proceedings with smart topic miner, in: *The Semantic Web–ISWC 2016: 15th International Semantic Web Conference*, Kobe, Japan, October 17–21, 2016, Proceedings, Part II 15, Springer, 2016, pp. 383–399.
- [3] M. Dudáš, S. Lohmann, V. Svátek, D. Pavlov, Ontology visualization methods and tools: a survey of the state of the art, *The Knowledge Engineering Review* 33 (2018) e10.
- [4] A. A. Salatino, F. Osborne, E. Motta, Augur: forecasting the emergence of new research topics, in: *Proceedings of the 18th ACM/IEEE on joint conference on digital libraries*, 2018, pp. 303–312.

- [5] S. Fathalla, S. Vahdati, S. Auer, C. Lange, Semsur: a core ontology for the semantic representation of research findings, *Procedia Computer Science* 137 (2018) 151–162.
- [6] A. A. Salatino, T. Thanapalasingam, A. Mannocci, F. Osborne, E. Motta, The computer science ontology: a large-scale taxonomy of research areas, in: *The Semantic Web–ISWC 2018: 17th International Semantic Web Conference*, Monterey, CA, USA, October 8–12, 2018, *Proceedings, Part II* 17, Springer, 2018, pp. 187–205.
- [7] A. D. Wade, The semantic scholar academic graph (s2ag), in: *Companion Proceedings of the Web Conference 2022*, 2022, pp. 739–739.
- [8] M. Grootendorst, Bertopic: Neural topic modeling with a class-based tf-idf procedure, *arXiv preprint arXiv:2203.05794* (2022).
- [9] L. McInnes, J. Healy, J. Melville, Umap: Uniform manifold approximation and projection for dimension reduction, *arXiv preprint arXiv:1802.03426* (2018).
- [10] J. A. Hartigan, M. A. Wong, Algorithm as 136: A k-means clustering algorithm, *Journal of the royal statistical society. series c (applied statistics)* 28 (1979) 100–108.
- [11] L. McInnes, J. Healy, S. Astels, hdbscan: Hierarchical density based clustering., *J. Open Source Softw.* 2 (2017) 205.
- [12] D. Müllner, Modern hierarchical, agglomerative clustering algorithms, *arXiv preprint arXiv:1109.2378* (2011).
- [13] M. Fernández, C. Overbeeke, M. Sabou, E. Motta, What makes a good ontology? a case-study in fine-grained knowledge reuse, in: *The Semantic Web: Fourth Asian Conference, ASWC 2009*, Shanghai, China, December 6-9, 2009. *Proceedings* 4, Springer, 2009, pp. 61–75.
- [14] A. Maedche, S. Staab, Measuring similarity between ontologies, in: *Knowledge Engineering and Knowledge Management: Ontologies and the Semantic Web: 13th International Conference, EKAW 2002 Sigüenza*, Spain, October 1–4, 2002 *Proceedings* 13, Springer, 2002, pp. 251–263.
- [15] C. Brewster, H. Alani, S. Dasmahapatra, Y. Wilks, Data driven ontology evaluation (2004).
- [16] M. Sabou, J. Gracia, S. Angeletou, M. d’Aquin, E. Motta, Evaluating the semantic web: A task-based approach, in: *The Semantic Web: 6th International Semantic Web Conference, 2nd Asian Semantic Web Conference, ISWC 2007+ ASWC 2007*, Busan, Korea, November 11-15, 2007. *Proceedings*, Springer, 2007, pp. 423–437.
- [17] P. Buitelaar, T. Eigner, T. D. OntoSelect, A dynamic ontology library with support for ontology selection in: *Proc. of the demo session at the international semantic web conference*, Hiroshima, Japan, Nov (2004).
- [18] D. Fernandes, J. Bernardino, Graph databases comparison: Allegrograph, arangodb, infinitegraph, neo4j, and orientdb., in: *Data*, 2018, pp. 373–380.
- [19] D. Sánchez, M. Batet, S. Martínez, J. Domingo-Ferrer, Semantic variance: an intuitive measure for ontology accuracy evaluation, *Engineering Applications of Artificial Intelligence* 39 (2015) 89–99.
- [20] J. Howard, Artificial intelligence: Implications for the future of work, *American journal of industrial medicine* 62 (2019) 917–926.
- [21] W. Xu, Toward human-centered ai: a perspective from human-computer interaction, *interactions* 26 (2019) 42–46.
- [22] C. Zhang, Y. Lu, Study on artificial intelligence: The state of the art and future prospects,

- [23] A. A. Monrat, O. Schelén, K. Andersson, A survey of blockchain from the perspectives of applications, challenges, and opportunities, *IEEE Access* 7 (2019) 117134–117151.

A. BERTopic Parameters

Here, we summarised the parameters we set throughout S&TO development. The following parameters have been set:

- “n_neighbors” refers to the number of neighbouring data points needed to estimate the manifold. Large sample point embeddings produce a more global perspective of the structure, while low values produce a narrower one. To get a good strike, we set n=2 as the result of the estimation.
- “n_components” refers to the number of components after the reduction in dominance. This value directly affects the clustering performance, so it is necessary to set an optimal value. By default, it is set to 5 to reduce the dimensionality as much as possible while maximizing the information in the generated embeddings.
- “low_memory”: It is set to TRUE because we use a huge dataset and need a lot of memory.
- “min_cluster_size”: The number of cluster generations highly relies on the cluster size. It is necessary to adjust the minimum size. After several experiments, a cluster size of 50 was found to be the optimal one. While high value gives few clusters of considerable size, and low value gives microclusters.
- “metric”: metric, like HDBSCAN, calculates the distances. Here, we went with Euclidean as, after reducing the dimensionality, we have low dimensional data, and not much optimisation is necessary. However, if you increase “n_components” in UMAP, it would be advised to investigate metrics that work with high dimensional data.
- “prediction_data”: Make sure you always set this value to True, as it is needed to predict new points later. You can set this to False if you do not wish to predict any unseen data points.
- “min_samples”: It is automatically set to “min_cluster_size” and controls the number of outliers generated. Setting this value significantly lower than “min_cluster_size” might help you reduce the amount of noise you will get. Do note that outliers are typical to be expected, and forcing the output to have no outliers may not properly represent the data.
- “top_n_words” refers to the number of words extracted per topic. In practice, we keep this value below 30, preferably between 10 and 20. The reasoning is that the more words representing a topic, the less relevant it may be. In this case, the top words are most representative of the topic and are the focus.
- “min_topic_size” specifies the minimum size of a topic. The lower the size value, the more topics are created. If the value is set too high, no topics may be created. We set this value too low, and we get many micro-clusters.
- “calculate_probabilities” give probabilities of all topics per document. This could slow down the extraction of topics for a large number of many documents.

- “low_memory“ set to true ensures that less memory is used in the calculations. This slows computation but allows UMAP to run on machines with little memory.
- “diversity“ reports a range of topic diversity from 0 to 1, where 0 indicates no diversity and 1 indicates a lot of diversity. Higher diverse topics mean less coherent topics in smaller cluster sizes. In our case, the diversity is assumed to be 0.4 or above.