# Knowledge extraction and footprint generation using the GeoSpace Body of Knowledge

Upeksha Indeewari Edirisooriya Kirihami Vidanelage[1], Mark van Vliet[2],

Sven Casteleyn[1], Carlos Granell[1], Stanislav Ronzhin[2], Rob Lemmens[2]

[1] Institute of New Imaging Technologies (INIT), Universitat Jaume I (UJI), Spain
[2] Faculty of Geo-Information Science and Earth Observation (ITC), University of Twente, The Netherlands

Correspondence: Rob Lemmens (r.l.g.lemmens@utwente.nl)

**Abstract.** Knowledge footprints are visualisations of personal and organisational expertise. They can be used for capturing, sharing and matching expertise in the context of promotion, research exposure, project collaboration, etc. We have developed a method for creating footprints in the geospatial domain, based on resources created by persons and organisations and using the GeoSpace Body of Knowledge as a shared, standardised vocabulary. We deployed an NLP-based keyword extraction method to annotate resources with GeoSpace BoK concepts and constructed a knowledge graph to connect these BoK concepts to personal or organisational profiles. Footprints are then created by querying the knowledge graph and visualizing the results. Initial tests have been carried out to validate the generated footprints.

**Submission Type.**

model, infrastructure

**BoK Concepts.**

[CF2] Cognitive and social foundations, [WB4] Resource Discovery, [CV4] Graphic representation techniques

**Keywords.** knowledge representation, body of knowledge, natural language processing, visualisation

## 1 Introduction

Knowledge representation forms the basis for knowledge sharing and reasoning in AI. It can be understood in different ways. According to Davis et al. (1993) it can range from a set of inferences to the medium of human expression. In this paper we focus on knowledge representation as ontology-based annotations of organisational and personal expertise through knowledge graphs. We visualise them in so-called *knowledge footprints,* which can be used for promotion, research exposure, project collaboration, etc. In this respect, we embark upon the ontological model, the method of annotation and the usage of knowledge representations for expertise mapping and matching between annotated resources. Our ontological model is the GeoSpace BoK, a Body of Knowledge, developed within subsequent projects on geoinformation technology, earth observation, satellite navigation and satellite communication. Related, though more generic, is the work done in structuring scholarly knowledge through the Open Research Knowledge Graph (ORKG)[1].

The GeoSpace Body of Knowledge (GeoSpace BoK) is a conceptual representation of the domain of geoinformation (GI) science, earth observation (EO) and satellite systems. Its development goes back to 2006 with the UCGIS GIS&T BoK (Wilson, 2014), focusing on GIS concepts, and was split off into a European version and further extended in the GI-N2K project (Vandenbroucke & Vancauwenberghe, 2016). Thereafter, during 2018-2022, the EO4GEO project embarked upon using the BoK

---

[1] https://orkg.org/

as a standard to support an evaluation of the European job market in GI-EO. Earth Observation concepts were added by domain experts and tools were developed to use the BoK concepts to annotate curricula, occupational profiles and job profiles. Further development on BoK content and tools that create and use the BoK are currently taking place in the SpaceSUITE project[2], in collaboration with the SPACE4GEO alliance[3] and the Association of Geographic Information Laboratories in Europe (AGILE)[4]. In this paper we describe the work on knowledge representation based on the GeoSpace BoK as a related ongoing research work. We first describe the BoK background, followed by creating and using BoK-based annotations for knowledge representation, ultimately visualised as knowledge footprints. As a use case we focus on the expertise of AGILE members by taking AGILE conference proceedings as a basis.

# 2 SpaceSUITE

Started in January 2024, SpaceSUITE is a 4-year Erasmus+ Blueprint project that aims to develop a program for upskilling and reskilling European professionals in the space downstream sector, encompassing the domains of earth observation, geoinformation technology and satellite systems, in order to bridge the gap between educational offer and professional demand. Its main objectives are to map the European educational and professional domains to uncover existing skill gaps, based on which key curricula and training actions are to be designed, developed and deployed to fill these gaps. Within this context, the maintenance and evolution of the GeoSpace BoK, which serves as a semantic backbone for an ecosystem of tools and resource annotation, is a key task.

## 2.1 GeoSpace Body of Knowledge

As of February 2025, the GeoSpace BoK[5] Version 8.0 contains 1279 concepts in a hierarchical structure, including concepts on geoinformation science, earth observation, satellite navigation and satellite communication. Concept connections are based on Simple Knowledge Organisation System (SKOS) (World Wide Web Consortium, 2009) relationships: to constitute the hierarchy, *broader* and *narrower* relations between concepts are used. The *related* relationship is used to indicate a temporary relationship that needs to be specified further.

Furthermore, each concept contains a short description and links to external resources (which were used by the

experts to describe the concept). During the EO4GEO project, each concept has been also associated with skills. Skills can be used to further specify characteristics of educational offers, occupational and job profiles (see Section 2.2).

## 2.2 GeoSpace BoK tools

Having the GeoSpace BoK as a shared, common vocabulary, an ecosystem of tools was developed to address the needs of the educational sector on the one hand (i.e., skill offer), and the professional sector on the other hand (skill demand). As an educational tool, the Curriculum Design Tool allows the definition of educational offers at various levels of granularity (e.g., study program, course, lecture). Tools aimed at the professional sector include the Occupational Profile Tool, which allows the description of representative profiles in the field, and the Job Offer Tool to define concrete job offers. All these resources are internally described in terms of the GeoSpace BoK and can be exposed as annotated resources. To further support actors addressing the field's skill gaps, two auxiliary tools were also made available, one to annotate any pdf file using BoK concepts (i.e., utilizing RDFa descriptions), and the Matching Tool, to compare any two BoK-annotated resources (e.g., a job offer and a curriculum vitae). It is worthwhile to note that the BoK Annotation Tool only allows manual annotation through an intuitive graphical interface; no automatic generation is currently provided.

# 3 Semi-automatic BoK annotation of GI resources

Despite the fact that the BoK Annotation Tool is useful for manual annotation, it requires significant effort and domain expertise, which limits applicability and scalability. To address these challenges, Natural Language Processing (NLP) techniques were utilised to automate the annotation of GI and EO resources.

We started with a review of NLP techniques for concept annotation, focusing on key phrase extraction and semantic similarity measurements. Key phrase extraction identifies important terms which represent a document's core content (Sun et al., 2020), whereas, similarity measurements ensure alignment between extracted terms and BoK concepts (Gomaa & Fahmy, 2013). These techniques were used to develop an automated annotation approach, whereby various keyphrase extraction algorithms and similarity measures were applied and compared.
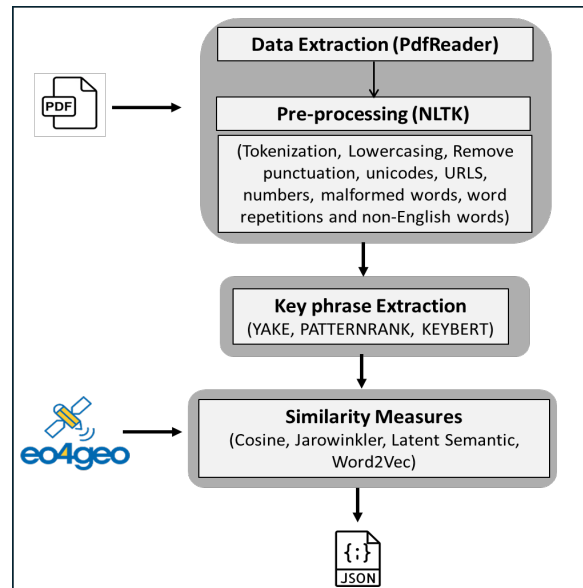
---

As shown in Fig. 1, the methodology consisted of three main phases: data extraction and pre-processing, keyphrase extraction, and similarity measure application. The data extraction consisted of extracting text from a pdf file using the PdfReader[6] Python tool. Subsequent data pre-processing tasks performed tokenization and data cleaning (i.e., lowercasing; remove punctuation, unicodes, URLS, numbers, malformed words, word repetitions and non-English words). On the resulting data, three popular keyphrase extraction algorithms, nl. YAKE[7], PatternRank[8] and KeyBert[9] were applied. Finally, the similarity between each keyphrase and each BoK concept was calculated, using three popular similarity measures, namely Cosine, Jaro-Winkler (Coken et al., 2003) and Word2Vec[10] similarity. Results were then ranked and the n (parameterisable) best matching BoK concepts were retained as annotations and stored in JSON format (see simplified example in Fig. 2). This approach resulted in 9 variants of the automated annotation process based on the 3 x 3 keyphrase extraction - similarity measure combinations.

An experiment was set up to measure the accuracy of the set of variants for automating BoK annotations. As a data source, 8 selected research papers from the AGILE: GIScience Series were used as input for the data extraction and pre-processing, and the nine keyphrase extraction - similarity measure combinations were applied. Fig. 2 shows a sample output for one such research paper. The generated BoK annotations were then compared with (manual) annotations by the respective authors of the selected papers, hereby considering parent/child BoK concept annotations as equivalent to avoid granularity of annotations to negatively impact the results, and the amount of manually provided annotations was used as a parameter (n) to determine the number of automated annotations.

Based on the F1-score metric, which balances precision and recall to evaluate each variant 's accuracy, the combination YAKE – Jaro-Winkler achieved the highest performance (F1-score: 0,2828). Fig. 3 shows all evaluation results. Considering this is a multiclass classification problem with a large number of classes (i.e., 952 BoK concepts), and given the manual annotations were varying greatly due to different interpretations of how to annotate (i.e., granularity, number of annotations), this result is promising.



**Figure 1**. Pipeline of the proposed NLP-based tool.



**Figure 2**. The best matching BoK concepts, provided by the YAKE – Jaro-Winkler tool, represented in JSON format.



**Figure 3**. Precision (P), recall (R) and F1 (F) scores for the 9 NLP-based annotation automation variants, using keyphrase-concept matching allowing child-parent flexibility.

## 4 Knowledge graph data retrieval and visualisations

Next, we applied the best automated BoK annotation method (YAKE - Jaro Winckler; as described in Section 3) for the purpose of knowledge representation, i.e., personal and organisational expertise, based on papers in

---

[6] https://pypi.org/project/pdfreader/
[7] https://pypi.org/project/yake/
[8] https://pypi.org/project/keyphrase-vectorizers/

[9] https://maartengr.github.io/KeyBERT/
[10] https://www.tensorflow.org/text/tutorials/word2vec

**Figure 4.** A visual representation of the ontological model. OBOK (green) and BOKA (blue). The prefixes indicate the source ontology of reused concepts; SKOS is used for constructing hierarchy.

the proceedings of the last three AGILE conferences, and associated to the respective authors/organisations. As such, authors/organisations are annotated with expertise, in terms of BoK concepts. hereafter referred to as *expertise annotations*.

**4.1 Towards the GeoSpace knowledge graph**

To properly leverage the extracted expertise annotations, expressed in terms of GeoSpace BoK content, semantic web technologies were utilised. For this purpose, an upper ontology, called OBOK (namespace: obok), was created and specified as a standardized RDF data model, to generally define the semantics of bodies of knowledge. The ontology was created with the ontology editor Protégé[11]. This upper ontology was subsequently used to express the content in the GeoSpace BoK (version 7.0)[5]. The already semantically rich content and hierarchical structure present in the GeoSpace BoK allowed for an easy transformation to RDF triples, re-using existing vocabularies (i.e., SKOS[12], FOAF[13], BIBO[14], DC[15]) and aligning to the defined semantics in the OBOK upper ontology. This was done by mapping the JSON keys and values, provided by the BoK API to corresponding subjects, predicates, and objects based on the semantic

structure in the OBOK upper ontology. The transformation to RDF triples was done with the RDFLib python library[16]. The RDF triples are stored in GraphDB[17].

The next step was to extend the OBOK with the semantics needed to store and link expertise annotations to concepts in the GeoSpace BoK (namespace: boka), re-using the ORG[18] vocabulary. To effectively link expertise annotations to the RDF triples from the GeoSpace BoK, the existing persistent unique identifiers in the BoK (Lemmens et al., 2022) were leveraged as URI's. Fig. 4 shows a visual representation of the various constructs present in both ontologies and shows that expertise annotations are directly linked to BoK concepts.

For combining and storing expertise annotations, GeoSpace BoK content and the semantics defined in the ontologies for bodies of knowledge (OBOK) and applications (BOKA), the GeoSpace knowledge graph was created. This RDF graph dataset became the source for two new applications; a tool for generating knowledge footprints and a tool for knowledge footprint matching.

---

## 4.2 GeoSpace knowledge footprints

Once personal and organisational expertise annotations can be linked to the hierarchical structure in the GeoSpace BoK, it becomes possible to map expertise onto the geospatial domain using visualisation techniques and create a so-called knowledge footprint (a visual representation of the breadth of knowledge accumulated by a person or organisation based on information in the GeoSpace knowledge graph).

Knowledge footprints were created by leveraging SPARQL[19] to extract data from the knowledge graph and various JavaScript scripts to further parse queried data into visualisations using the D3.js library[20]. An example SPARQL query is provided in Appendix A. Knowledge footprints come in various forms:

- A paper-specific knowledge footprint shows all the concepts that are matched with a specific paper.
- An individual knowledge footprint represents an aggregation of all the knowledge displayed in papers a specific individual authored.
- An organisational knowledge footprint aggregates the knowledge footprints of all individuals within a specific organisation.

Fig. 5 provides an example of an organisational knowledge footprint. For its design it was chosen to create a D3 radial cluster tree including all the GeoSpace BoK concepts (yellow nodes) and their relations (blue lines) and combine that with a doughnut chart. The latter was inspired by Elsevier's Topic Wheel (Elsevier, 2021). This combined visualisation serves as a basemap of knowledge of our geospatial domain. The segments of the outer doughnut chart visualise the knowledge areas in the GeoSpace BoK. These segments aim to tell the viewer in which knowledge area an entity has knowledge without having to look at matched node labels.

On top of this basemap the entity's expertise is visualised. Highlighting matched concepts through red nodes, and highlighting so-called "knowledge paths" with green lines, make the hierarchical structure and thereby all the parent concepts of matched concepts visible through traversing this path from matched concept to the root concept. In this example, it shows that this department has expertise at different levels of granularity in Geocomputation, Cartography and Visualisation, and Image Processing and Analysis. There is not much expertise in Physical Principles.

The second application involves footprint matching. As shown in Fig. 6 the output of footprint matching is quite similar, but now incorporates the footprints of two entities. When both entities have knowledge of the same concept, their knowledge paths are drawn parallel to each other, adopting a visualisation style often used to visualise metro lines. The example shows some differences between the expertise of the two departments, especially in Cartography and Visualisation and in Web-based GI.

Note that for both examples, the footprints only represent the published papers in the proceedings of the last three AGILE conferences.

## 5 Evaluation of results

Knowledge footprints were introduced at the AGILE 2024 conference, to present the content of specific sessions (including a session keynote and the profile of the speaker) and the expertise profile of new AGILE member organisations. During a dedicated workshop, participants were introduced to the creation and use of footprints. Workshop participants perceived that the footprints provided an adequate representation of the expertise profiles, given the input available to the tools that produced the footprints (as described in Sections 3 and 4). This confirms the encouraging results of the automated annotation approach evaluation. However, the following aspects limit the input, and consequently, the accuracy and completeness of the footprints. First, the BoK hierarchy is not interpreted uniformly by all users; expertise about a concept is either seen as the combination of subconcepts or as the abstract level of expertise of subconcepts. Second, as papers have co-authors in most cases, the expertise of one person is not per se on every aspect of the paper, although some affinity with the paper topics is assumed. Third, in this work only the papers published in the last three Agile conference were considered. Authors publish in multiple outlets, and therefore a complete personal research profile should take into account *all* the author's publications.

## 6 Data and Software Availability

The GeoSpace BoK and its tools are publicly available at the SpaceSUITE project website[21]. The BoK tools are licensed under GNU GPLv3. The tools and data of Sections 3 and 4 are available at the respective project's github repositories[22, 23].
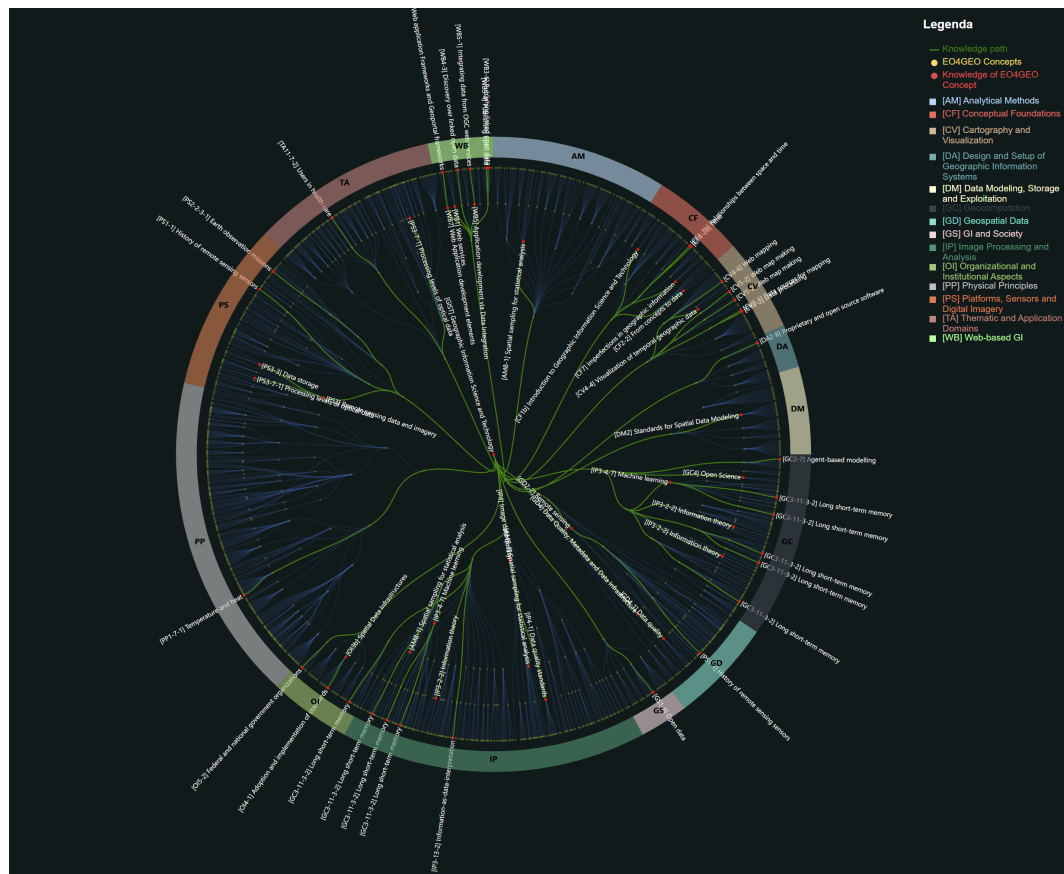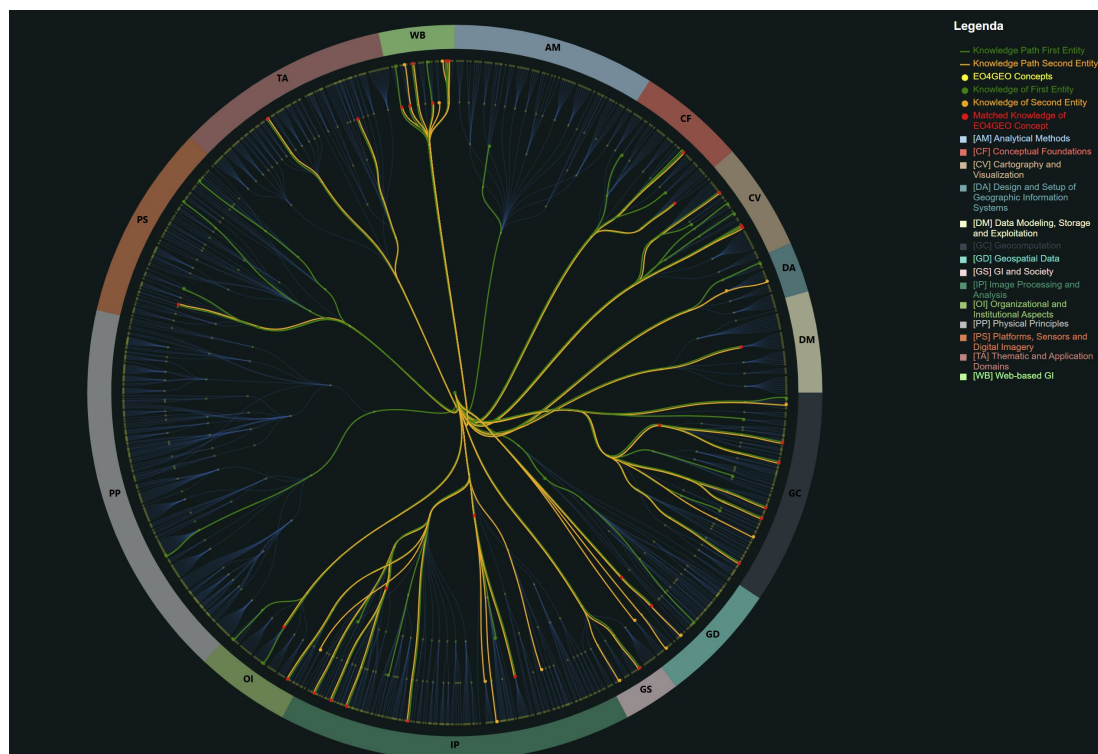
---

[19] https://www.w3.org/TR/sparql12-query/
[20] https://d3js.org/
[21] https://www.spacesuite-project.eu/body-of-knowledge/

[22] https://github.com/UpekshaIndeewari/geotec_thesis_EO4GEO/tree/main
[23] https://github.com/MPvliet/Thesis-GIMA-2023-2024.

**Figure 5**. An organisational knowledge footprint of the University of Twente. Interactive version available at: https://mpvliet.github.io/?footprintType=Organisational&footprintEntity=University of Twente



**Figure 6.** Footprint matching University of Twente (green) with Technische Universität Dresden (yellow). Interactive version available at: https://mpvliet.github.io/footprintMatching.html?footprintType=Organisational&footprintFirstEntity=University of Twente&footprintSecondEntity=Technische Universität Dresden

# 7 Conclusion and recommendations

The extraction of concepts from literature resources works fairly well. An NLP-based tool for semi-automatic BoK annotation of GI resources was successfully developed, and tested using three keyphrase extraction and three similarity measurement techniques. Among the methods evaluated, the YAKE – Jaro-Winkler combination performed the best. To achieve better performance, future research should focus on exploring supervised learning, refining keyphrase extraction and fine tuning the expert-based evaluation method. Encouraging comparative author verification and investigating abstract-based keyphrase extraction could further enhance accuracy and efficiency. These advancements will strengthen automated annotation tools and their applicability in the EO*GI research community.

The implementation of knowledge graphs facilitates querying GeoSpace BoK concepts and the NLP expertise annotations. A knowledge footprint, generated from a body of literature, can currently at best represent a snapshot of the aggregate knowledge of an individual or an association. To complete specific personal or organisational profiles, there is a need to include additional sources and/or add annotations manually. The current footprint visualisations are detailed and therefore require a large visual space. While detailed, footprints lack any form of displaying the level of expertise someone's profile holds. Leveraging visual indicators, e.g., size or coloring of nodes, could provide a visual clue in the future. Smaller visualisations would require an abstraction of the concepts and their relationships. To overcome the limitations due to the ambiguity in BoK hierarchy interpretations, better instructions are needed for annotators and end-users.

**Declaration of Generative AI in writing**

The authors declare that they have not used Generative AI tools in the preparation of this manuscript. All intellectual and creative work, including the analysis and interpretation of data, is original and has been conducted by the authors without AI assistance.

# References

Cohen, W. W., Ravikumar, P., & Fienberg, S. E. (2003). A Comparison of String Distance Metrics for Name-Matching Tasks. In *IIWeb* (Vol. 3, pp. 73-78).

Davis, R., Shrobe, H., & Szolovits, P. (1993). What Is a Knowledge Representation? AI Magazine, 14(1), Article 1. https://doi.org/10.1609/aimag.v14i1.1029

Elsevier. (2021). SciVal Topics Elsevier. www.elsevier.com. Retrieved 7 February 2025, from https://www.elsevier.com/products/scival/overview/topics

Gomaa, H. W. & Fahmy, A. A. (2013). A Survey of Text Similarity Approaches. International Journal of Computer Applications, 68(13), 13–18. https://doi.org/10.5120/11638-7118

Lemmens, R., Albrecht, F., Lang, S., Casteleyn, S., Stelmaszczuk-Górska, M., Olijslagers, M., Belgiu, M., Granell, C., Augustijn, E.-W., Pathe, C., Missoni-Steinbacher, E.-M., & Monfort Muriach, A. (2022). Updating and using the EO4GEO Body of Knowledge for (AI) concept annotation. AGILE: GIScience Series, 3, 1–6. https://doi.org/10.5194/agile-giss-3-44-2022

Sun, C., Hu, L., Li, S., Li, T., Li, H., & Chi, L. (2020). A review of unsupervised keyphrase extraction methods using within-collection resources. Symmetry, 12(11), 1–20. https://doi.org/10.3390/sym12111864

Vandenbroucke, D. & Vancauwenberghe, G. (2016). Towards a New Body of Knowledge for Geographic Information Science and Technology. Micro, Macro & Mezzo Geoinformation, 2016 (6), 7-19. ISSN: 1857-9000, EISSN: 1857-9019, 2016.

Wilson, J. P. (2014). Geographic Information Science & Technology Body of Knowledge 2.0 Project (Final report), 2014. Retrieved from https://ucgis.memberclicks.net/assets/docs/ucgis_bok2_wilson_report_dec 20 14.pdf

World Wide Web Consortium (2009). SKOS Simple Knowledge Organization System Reference. (W3C Recommendation) Miles, A., & Bechhofer, S. (Editors) Retrieved 1 February 2025, from https://www.w3.org/TR/skos-reference/#semantic-relations

# Appendix A

SPARQL query to create individual knowledge footprints.

```sparql
1   PREFIX eo4geo: <https://bok.eo4geo.eu/>
2   PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
3   PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
4   PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
5   PREFIX boka: <http://example.org/BOKA/>
6   PREFIX foaf: <http://xmlns.com/foaf/0.1/>
7
8   SELECT ?conceptName ?childName ?conceptID ?childID ?nodeColour ?showLabel ?labelSize ?nodeValue
9   FROM eo4geo:applications
10  FROM eo4geo:concepts
11  WHERE {
12    {
13      SELECT ?concept ?conceptName ?childName ?conceptID ?childID (IF(?knownByFirstEntity, 1 , 0 ) AS ?nodeValue) WHERE {
14        ?concept rdf:type skos:Concept;
15          rdfs:label ?conceptName;
16          skos:notation ?conceptID.
17        OPTIONAL {
18          ?concept skos:narrower ?child.
19          ?child rdfs:label ?childName;
20            skos:notation ?childID.
21        }
22        BIND(EXISTS {
23          ?expertURI rdf:type boka:Expert;
24            foaf:name ?expertName;
25            boka:hasKnowledgeOf ?concept.
26          FILTER(CONTAINS(LCASE(STR(?expertName)), LCASE("Sven Casteleyn")))
27        } AS ?knownByFirstEntity)
28      }
29    }
30    BIND(IF(?nodeValue = 1 , "#FF0000", "#FFFF00") AS ?nodeColour)
31    BIND(IF(?nodeValue = 1 , 16 , 0 ) AS ?labelSize)
32    BIND(IF(?nodeValue = 1 , 1 , 0 ) AS ?showLabel)
33  }
34  ORDER BY (?conceptName)
```