

Digital humanities on the Semantic Web: Sampo model and portal series

Eero Hyvönen

Department of Computer Science, Aalto University, Finland

Helsinki Centre for Digital Humanities (HELDIG), University of Helsinki, Finland

Semantic Computing Research Group (SeCo)

E-mail: eero.hyvonen@aalto.fi; URL: <https://seco.cs.aalto.fi/>

Editor: Christoph Schlieder, University of Bamberg, Germany

Solicited reviews: Christoph Schlieder, University of Bamberg, Germany; Kai Eckert, Stuttgart Media University, Germany

Abstract. Cultural heritage (CH) contents are typically strongly interlinked, but published in heterogeneous, distributed local data silos, making it difficult to utilize the data on a global level. Furthermore, the content is usually available only for humans to read, and not as data for Digital Humanities (DH) analyses and application development. This application report addresses these problems by presenting a collaborative publication model for CH Linked Data and six design principles for creating shared data services and semantic portals for DH research and applications. This *Sampo model* has evolved gradually in 2002–2021 through lessons learned when developing the *Sampo series* of linked data services and semantic portals in use, including MuseumFinland (2004), CultureSampo (2009), BookSampo (2011), WarSampo (2015), Norssit Alumni (2017), U.S. Congress Prosopographer (2018), NameSampo (2019), BiographySampo (2019), WarVictimSampo 1914–1922 (2019), MMM (2020), AcademySampo (2021), FindSampo (2021), WarMemoirSampo (2021), and LetterSampo (2022). These Semantic Web applications surveyed in this paper cover a wide range of application domains in CH and have attracted up to millions of users on the Semantic Web, suggesting feasibility of the proposed Sampo model. This work shows a shift of focus in research on CH semantic portals from data aggregation and exploration systems (1. generation systems) to systems supporting DH research (2. generation systems) with data analytic tools, and finally to automatic knowledge discovery and Artificial Intelligence (3. generation systems).

Keywords: Semantic Web, Digital Humanities, Linked Open Data, data services, portals

1. Breaking data silos of cultural heritage

Cultural Heritage content is published independently by different memory organizations, such as museums, libraries, archives, galleries, and media companies. The traditional web publishing model, where everybody can publish easily content for everybody to read, facilitates fast and flexible publication on the Web. However, using related local contents from separate data sources on a global level is difficult because of the incompatible *data silos*: the local databases and online systems of the publishers are associated in content, but heterogeneous in terms of incompatible data models, annotated using different thesauri and vocabularies, distributed geographically, based on different natural languages, and used with different kind of user interfaces. An even more fundamental problem is that the contents are typically published only for humans to read and not as data for computational analyses and application development. This means that the end users typically have to learn and use several different applications to cater their information needs about a topic. For the data publishers, lots of costly redundant work is needed in

creating the data silos, e.g., in developing the vocabularies, data services, and user interfaces. The availability of the data in a usable open form is a prerequisite of the work for the application developers.

To mitigate these problems, various massive international data aggregation systems have been created, such as Europeana¹ in Europe and the Digital Public Library of America² in the U.S. There are lots of similar systems around on a national and regional level (e.g., Deutsche Digitale Bibliothek³ in Germany and K-samsök service in Sweden) and within various thematic communities⁴ (e.g., ARIADNEplus⁵ in archaeology). Similar data aggregation systems have also been created within single organizations that may already have lots of siloed but related databases around, like in the case of the BBC in the U.K. [42]. There are lots of international and national standardization efforts for creating harmonized data models (e.g., Dublin Core,⁶ CIDOC CRM,⁷ and FRBRoo⁸ [72]), shared thesauri for annotating contents (e.g., AAT, TGN, and ULAN vocabularies of the Getty Research Institute⁹), as well as generic frameworks, such as the Semantic Web standards of the W3C.¹⁰

This paper concerns using Semantic Web (SW) technologies [14] and Linked Open Data (LOD) publishing [11,18] to address the data silo and data publishing problems above. A general model, called *Sampo Model*, is presented for the purpose. As empirical evidence of feasibility of applying the model in practise, the *Sampo series* of data services and semantic portals is presented.¹¹ They have had millions of users on the Semantic Web in total. The fundamental idea and goal of Linked Data is to create an interoperable interlinked Web of Data [11]. The novelty of the Sampo model lays in its attempt to address this goal using a set of re-usable design principles or guidelines for creating semantic portals, especially for Cultural Heritage applications and Digital Humanities research [6]. The Sampo model is a kind of consolidated approach for creating LOD services and semantic portals, something that the field of the Semantic Web is arguably still largely missing [13].

This paper is organized as follows. Section 2 presents the principles of the Sampo model. In Section 3, a survey of Sampo systems is presented as a proof-of-concept, illustrating use cases of the model and how it has evolved in 2002–2021. In conclusion, related works are discussed, contributions of the paper are summarized, and challenges and directions for further research are outlined. This paper extends substantially the earlier short paper [19] about the Sampo model at the DHN 2020 conference.

2. Sampo model principles

The Sampo Model is an informal collection of principles for LOD publishing and designing semantic portals listed in Table 1, supported by an ontology and data infrastructure and software tools for user interface design and

Table 1
Sampo model principles P1–P6

P1	Support collaborative data creation and publishing
P2	Use a shared open ontology infrastructure
P3	Make clear distinction between the LOD service and the user interface (UI)
P4	Provide multiple perspectives to the same data
P5	Standardize portal usage by a simple filter-analyze two-step cycle
P6	Support data analysis and knowledge discovery in addition to data exploration

¹<https://europeana.eu>

²<https://dp.la/>

³<https://www.deutsche-digitale-bibliothek.de/?lang=en>

⁴See <https://pro.europeana.eu/page/aggregators> for a list.

⁵<https://ariadne-infrastructure.eu/>

⁶<https://dublincore.org/>

⁷<https://cidoc-crm.org/>

⁸<https://www.cidoc-crm.org/frbroo/home-0>

⁹<https://www.getty.edu/research/tools/vocabularies/>

¹⁰<https://www.w3.org/standards/semanticweb/>

¹¹See <https://seco.cs.aalto.fi/applications/sampo/> for a complete list of “Sampo portals”, videos, and further information.

data publication. Principles P1–P3 can be seen as a foundation for developing data services; principles P4–P6 are related to creating semantic portals.¹² The model is called “Sampo” according to the Finnish epic Kalevala, where Sampo is a mythical machine giving riches and fortune to its holder, a kind of ancient metaphor of technology¹³ according to the most common interpretation of the concept. The principles P1–P6 of Table 1 are described and motivated in more detail in the following subsections, one after another.

P1. Support collaborative data creation and publishing The model is based on the idea of collaborative content creation. The data is aggregated from local data silos into a global service, based on a shared ontology and publishing infrastructure [18]. The local data are harmonized and enriched with each other by linking and reasoning, based on Semantic Web standards. In this model everybody can win, including the data publishers by enriched data and shared publishing infrastructure, and the end users by richer global content and services. However, collaborative publishing also complicates the publication process, as more agreements are needed within the community.

This model addresses the problems of semantic data interoperability and distributed content creation at the same time. A shared semantic ontology infrastructure that includes shared metadata schemas [87] and domain ontologies for populating the data models are used for harmonizing and interlinking data from separate silos. If the content providers provide the system with metadata about their contents using the shared infrastructure, the data is automatically linked and enriched with each other and forms a knowledge graph [7]. For example, if metadata about a painting created by Picasso comes from an art museum, it can be enriched (linked) with, e.g., biographies from Wikipedia and other sources, photos taken of Picasso, information about his wives, books in a library describing his works of art, related exhibitions open in museums, and so on. At the same time, the contents of any organization in the portal having Picasso related material get enriched by the metadata of the new artwork entered in the system.

Figure 1 depicts as an example how the collaborative Sampo publication model (P1) was used in the Mapping Manuscript Migrations (MMM) system [24,43]. MMM includes three key datasets about ca. 220 000 medieval and Renaissance manuscripts in total that originate from the U.S. (Schoenberg Institute (T1)), U.K. (Oxford University Libraries (T2)), and France (the Institut de recherche et d’histoire des textes (IRHT) (T3)). The data T1–T3 are transformed into the unified harmonizing data model used in the MMM Linked Data Service [43] that is depicted in the middle of the figure. The data service is used by the MMM portal (bottom) but can also be used in for research

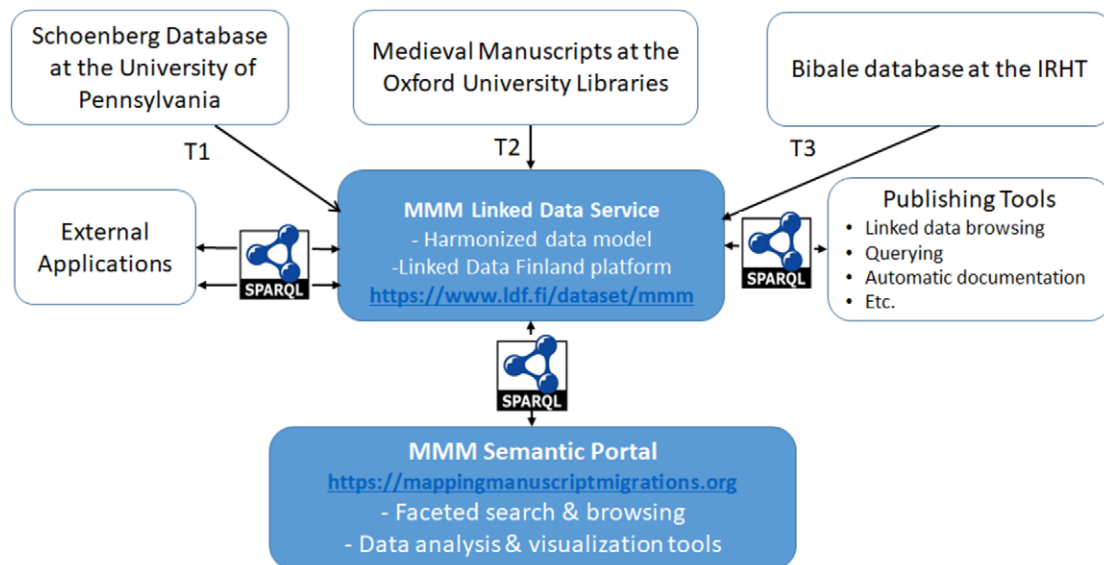


Fig. 1. Publishing and using heterogeneous distributed data in the MMM Sampo system.

¹²The numbering of the principles P3 and P6 is switched in the table with respect to [24] to clarify this.

¹³<https://en.wikipedia.org/wiki/Sampo>

and in other external applications via the SPARQL endpoint (on the left) [4]. The global data is documented and can be studied using SPARQL and publishing tools (on the right), too. The aggregated global data can be used for solving research questions that cannot be answered by studying the local datasets separately.

P2. Use a shared open ontology infrastructure The Sampo model is based on a shared LOD ontology infrastructure with which the local datasets are made compatible. Re-using the same infrastructure, and developing it further step by step in each Sampo portal and application, saves a lot of effort for the developers of next Sampos and other applications. For example, the linked data-based geogazetteer of contemporary placenames in Finland, using data from the National Survey and introduced in NameSampo [40] for open use, contains some 800 000 geocoded places, and there are other ontologies for historical places, maps, and persons.

The infrastructure includes harmonising shared metadata models (schemas) for representing individuals as well as domain ontologies (thesauri, vocabularies) that are used in populating (instantiating) the metadata models. This can be done by using data transformations and by aligning ontologies, as described in detail in [43,46] for the WarSampo and MMM systems, respectively. The Sampo portals use in practise both Dublin Core-based models and the dumb-down principle¹⁴ for documents, and event-based models conforming and extending the CIDOC CRM ontology and FRBRoo. In addition to sharing same infrastructure components, different Sampos enrich each other's contents by mutual data linking, creating a gradually evolving network of Sampos, a kind of "SampoSampo" and data cloud. Also data from the international data infrastructure is used for this purpose, e.g., Wikidata¹⁵ and GeoNames.¹⁶ The WarSampo knowledge graph [46], for instance, is part of the international LOD Cloud.¹⁷

Many Sampo systems make use of the national FinnONTO ontology infrastructure [38]. Its development started in 2003 and the work is carried on today by the National Library of Finland as the Finto.fi ontology service,¹⁸ and under the research initiative Linked Open Data Infrastructure for Digital Humanities in Finland (LODI4DH)¹⁹ [20] that is part of the FIN-CLARIAH initiative²⁰ fostering work on the pan-European CLARIN²¹ and DARIAH²² infrastructures in a national context.

P3. Make clear distinction between the LOD service and the user interface (UI) The Sampo Model argues for the idea of separating the underlying Linked Data service *completely* from the user interface via a SPARQL API. The rationale for this is: Firstly, this simplifies the portal architecture. Secondly, the data service can be opened for data analysis research in Digital Humanities. For example, the YASGUI²³ [71] editor for SPARQL querying and visualizing the results can be used, or Python scripting in Google Colab²⁴ and using Jupyter²⁵ notebooks—see, e.g., the analyses of the BiographySampo data in [80].

On top of a SPARQL API it is possible to define more dedicated and simpler APIs depending on the application and user needs, such as *crlc* [61]. SPARQL querying can be computationally inefficient and has lead some developers to use other tools for developing search engines, such as Elasticsearch²⁶ and Solr.²⁷ However, SPARQL servers, such as Fuseki²⁸ used in later Sampo systems, include efficient tooling for text indexing, such as Lucene²⁹ and Solr.

¹⁴https://www.dublincore.org/resources/glossary/dumb-down_principle/

¹⁵<https://www.wikidata.org/>

¹⁶<https://www.geonames.org/>

¹⁷<https://lod-cloud.net/>

¹⁸<https://finto.fi/>

¹⁹<https://seco.cs.aalto.fi/projects/lodi4dh/>

²⁰<https://seco.cs.aalto.fi/projects/fin-clariah/>

²¹<https://www.clarin.eu/>

²²<https://www.dariah.eu/>

²³<https://yasgui.triply.cc>

²⁴<https://colab.research.google.com/notebooks/intro.ipynb>

²⁵<https://jupyter.org>

²⁶<https://www.elastic.co/elasticsearch/>

²⁷<https://solr.apache.org/>

²⁸<https://jena.apache.org/documentation/fuseki2/>

²⁹<https://lucene.apache.org/>

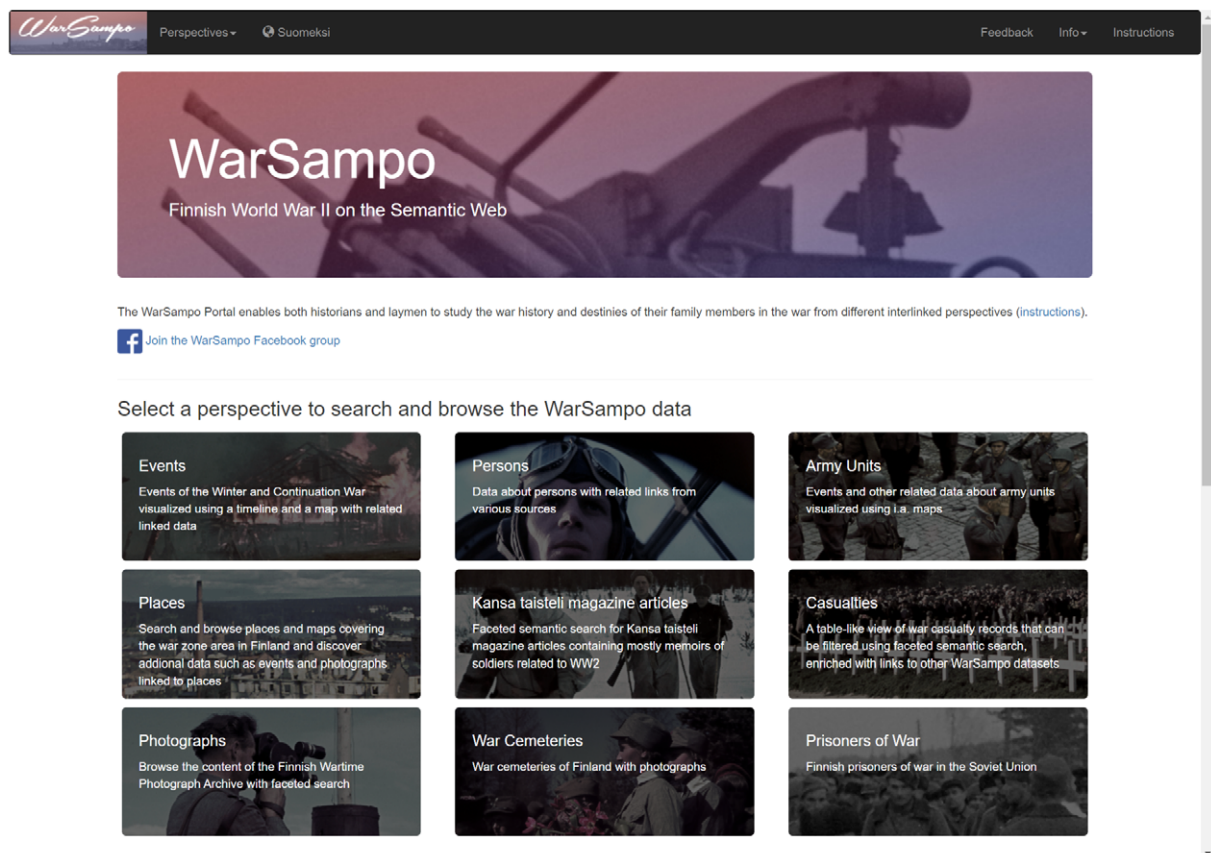


Fig. 2. Landing page of WarSampo with nine application perspectives.

P4. Provide multiple perspectives to the same data The Sampo model fosters the idea that on top of a LOD service different thematic *application perspectives* to the data can be created by re-using the data service. This means that the underlying data can be re-used without modifying it, which is typically costly [28] when dealing with Big Data.

The application perspectives are provided on the landing page of the Sampo portal, and they enrich each other by data linking. By selecting a perspective the corresponding application is opened. In addition, completely separate applications can be created on top of the data service by third parties, which is of help to memory organizations that typically are not strong in IT application development but are often willing to share the content openly through multiple channels.

For example, Fig. 2 depicts the landing page of WarSampo [22] with the following nine interlinked application perspectives for accessing the underlying LOD service data:

1. Major events (1050) of the Second World War (WW2) in Finland visualized on a timeline and maps with related linked data
2. People (100 000) with biographical data and links to related perspectives
3. Army Units (15 900) including events, war diaries, and people related to the units
4. Places perspective for searching the war zone events using contemporary and historical maps
5. *Kansa taisteli* magazine articles (3360) (1957–1986) containing memoirs of the soldiers after the war
6. Casualties data (95 000 death records) of all soldiers killed in action during the war
7. Authentic photographs (160 000) from the war zone by the Defence Forces of Finland, interlinked, e.g., with people, army units, and places
8. War Cemeteries (630) of the casualties of WW2 with 3000 photographs
9. Finnish Prisoners of War (4500) in the Soviet Union in 1939–1945

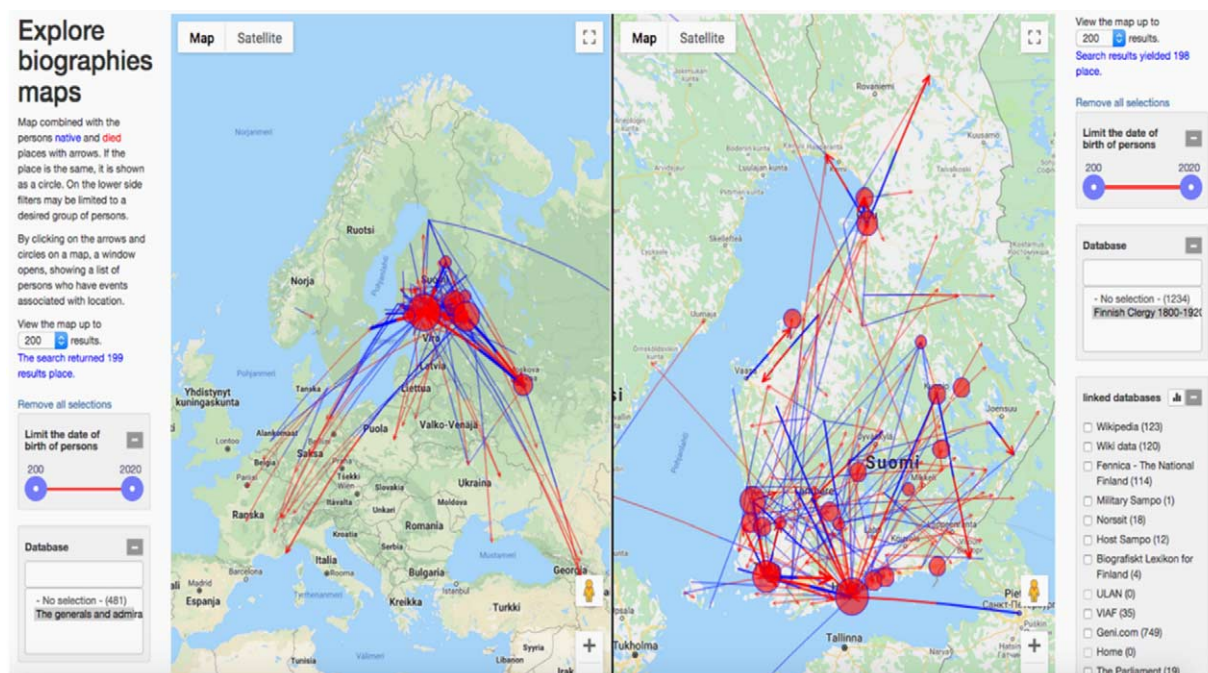


Fig. 3. Comparing the life charts of two target groups in BiographySampo, admirals and generals (left) and clergy (right) of the historical Grand Duchy of Finland (1809–1917).

P5. Standardize portal usage by a simple filter-analyze two-step cycle In later Sampos, the application perspectives can be used by a two-step cycle for research: First the focus of interest, the target group, is filtered out easily using faceted semantic search [32,81,83]. Second, the target group is visualized or analyzed by using ready-to-use DH tools of the application perspectives. The general idea here is to try to “standardize” the UI logic so that the portals are easier to use for the end users [39].

For example, Fig. 3 depicts a situation in BiographySampo where the user compares the life charts of two prosopographical groups in 1809–1917 when Finland was an autonomous Grand Duchy within the Russian Empire: 1) Finnish generals and admirals in the Russian armed forces (on the left). 2) Members of the Finnish clergy (1800–1920) (on the right). With a few selections from the facets the user can filter out the two target groups and see that, for some reason, quite a few officers moved to Southern Europe when they retired, like retirees today, while the Lutheran ministers usually stayed in Finland until their death.

P6. Support data analysis and knowledge discovery in addition to data exploration Three generations of semantic portals for Digital Humanities can be identified according to the vision [21] underlying the work on Sampos. Ten years ago the research focus in semantic portal development was on data harmonization, aggregation, search, and browsing (1. generation systems). However, the Sampo model aims not only at data publishing with search and data exploration [59]. The rise of DH research has started to shift the focus to providing the user with integrated tools for solving research problems in interactive ways (2. generation systems). The next step ahead is based on Artificial Intelligence: future portals not only provide tools for the human to solve problems but are used for finding research problems in the first place, for addressing them, and even for solving them automatically under the constraints set by the human researcher and explaining the results (3. generation systems). A step towards this is the relational search application perspective in BiographySampo where the machine tries to find “interesting” semantic connections in linked data and also explain them in natural language [30]. Principles P4–P6 are related to creating 2. and 3. generation systems.

FAIR Linked Data The widely used modern FAIR principles³⁰ for creating Findable (F1–F4), Accessible (A1–A2), Interoperable (I1–I3), and Re-usable (R1) data are:

- *Findability*: F1. (Meta)data are assigned a globally unique and persistent identifier; F2. Data are described with rich metadata (defined by R1 below) F3. Metadata clearly and explicitly include the identifier of the data they describe; F4. (Meta)data are registered or indexed in a searchable resource
- *Accessible* A1. (Meta)data are retrievable by their identifier using a standardised communications protocol; A2. Metadata are accessible, even when the data are no longer available
- *Interoperable* I1. (Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation; I2. (Meta)data use vocabularies that follow FAIR principles; I3. (Meta)data include qualified references to other (meta)data
- *Reusable* R1. (Meta)data are richly described with a plurality of accurate and relevant attributes

Obviously, the Sampo model principles are compatible with the FAIR principles above. This follows from the fact that the model was developed using the standards, linked data principles,³¹ and best practices³² of the W3C.

The Sampo model originates from applications in the CH domain but is generic and is not bound to its origin. The model has been applied in other domains, too. An example of this is the HealthFinland system [37,78] for health promotion information, that was deployed by the National Institute for Health and Welfare in Finland. HealthFinland got at the ISWC 2008 conference the international Semantic Web Challenge Award.

In [27] the idea of developing the Sampo model into domain-specific Sampo *frameworks* is presented using three epistolary datasets as case studies. The LetterSampo framework presented includes data models for representing metadata about historical letters as well as some domain-specific configurations for the Sampo-UI tool for user interfaces. In this way the LetterSampo framework could be applied easily to create portals for three different datasets related to the Republic of Letters [15] from the University of Oxford (Early Modern Letters Online³³), the Dutch Huygens Institute (CKCC corpus underlying ePistolarium³⁴), and the Berlin-Brandenburg Academy of Sciences (correspSearch (meta)dataset³⁵). The Sampo framework is applied also the in FindSampo system [31,69] in the domain of archaeological finds using collections of the National Museum of Finland, the PAS database of the British Museum, and collections of the Fitzwilliam Museum in Cambridge [64].

Linked data publications on the SW are typically evaluated with the W3C “5-star model”,³⁶ using a quality scale analogous to evaluating hotels. In the Linked Data Finland service³⁷ hosting most Sampos, the model is extended to a “7-star model” [36]. The 6th star is given to a data publication if it includes not only 5-star data but also the schemas of the data with documentation. This makes re-use of data easier. The 7th star is given to a data publication, if the publication includes some kind of evaluation that the data actually conforms to the provided schemas using, e.g., the SHACL Shapes Constraint Language³⁸ or ShEx Shape Expressions³⁹ [48]. The idea here is to encourage publishers to publish high quality data as data quality of LD is a severe issue on the SW.

3. Sampo series of semantic portals and LOD services

The Sampo model has evolved gradually over time in 2002–2021 via lessons learned in developing the Sampo series of semantic portals and related LOD services in various projects. This section overviews shortly a selection of these systems, listed in Table 2, in order to provide a proof-of-concept of the model and to give some examples

³⁰The FAIR principles are listed here and their numbering are based on <https://www.go-fair.org/fair-principles/>.

³¹<https://www.w3.org/DesignIssues/LinkedData.html>

³²<https://www.w3.org/TR/dwbp/>

³³<http://emlo.bodleian.ox.ac.uk>

³⁴<http://ckcc.huygens.knaw.nl/epistolarium/>

³⁵<https://correspsearch.net/>

³⁶<https://www.w3.org/community/webize/2014/01/17/what-is-5-star-linked-data/>

³⁷<https://ldf.fi>

³⁸<https://www.w3.org/TR/shacl/>

³⁹<https://shex.io/>

Table 2

A selection of sampo portals and LOD services for digital humanities; user counts by Google analytics 2021 October

Portal	Year	Domain	# Users	# Triples	Primary data owners
MuseumFinland	2004	Artefact collections	39 000	211 000	National Museums, City Museums of Espoo and Lahti, Finland
CultureSampo	2008	Finnish culture	107 000	11M	Memory organizations and the Web, ca 30 data sources
BookSampo	2011–	Fiction literature	2M/year	4.36M ^a	Public libraries in Finland (Kirjastot.fi)
WarSampo	2015–2019	World War II	857 000	14M	National Archives, Defense Forces, and others, Finland
Norssit Alumni	2017	Person registry	unknown	469 000	Norssi High School alumni organization Vanhat Norssit
U.S. Congress Prosopographer	2018	Parliamentary data	unknown	830 000	U. S. Congress Legislator data ^b
NameSampo	2019	placenames	37 000	26.0M ^c	Institute for the Languages of Finland (Kotus), National Land Survey of Finland, and the J. Paul Getty Trust TGN Thesaurus
BiographySampo	2019	Biographies	50 000	5.56M	Finnish Literature Society
WarVictimSampo 1914–1922	2019	Military history	29 000	9.96M	National Archives of Finland
Mapping Manuscript Migrations (MMM)	2020	Pre-modern manuscripts	9100	22.5M	Schoenberg Inst. for Manuscript Studies (U.S.), Oxford University Libraries (U.K.), and Inst. for Research and History of Texts (France)
AcademySampo	2021	Finnish Academics	8200	6.55M	University of Helsinki and National Archives, Finland
FindSampo	2021	Archaeology, finds	3100	1.0M	Finnish Heritage Agency, Finland
WarMemoirSampo	2021	World War II	3100	323 000	Veteran organization <i>Tammenlehvän Perinneytö</i> , National Archives, Finland
LetterSampo	2022	Historical Letters	TBA	3.75M	Huygens Institute, Berlin-Brandenburg Academy of Sciences, et al.

^aOriginal KG size in 2011; the size in much larger today including also non-fiction works.

^b<https://github.com/unitedstates/congress-legislators>

^cThis count includes only data of Kotus; the total number of triples of all sources is 241M.

and historical background of the work. For each system, the year of publication, application domain, number of end users, size of the underlying triplestore, and primary data owners are listed. In below, each system is described shortly with a reference to its research homepage, to the portal online, and to research articles for more detailed information. These references provide links to related works and additional publications, and to the data services online.

MuseumFinland – Finnish Museums on the Semantic Web⁴⁰ (online since 2004) [29] was the first Sampo. It introduced principle P1 of Table 1 by aggregating and publishing heterogeneous, distributed artefact collection data from Finnish museums. No linked data infrastructure was available then, which sparked the idea of creating one in the next project CultureSampo. MuseumFinland application got the Semantic Web Challenge Award at the ISWC 2004 conference.

CultureSampo – Finnish Culture on the Semantic Web 2.0⁴¹ (online since 2009) [28,58] introduced the idea creating and using a data infrastructure (P2) [38] and thematic application perspectives [28] to the data (P4). It demonstrated how CH content of tens of different kinds, both tangible and intangible CH content, can enrich each other. CultureSampo includes, e.g., a semantic model of the Kalevala epic narrative [34,65], based on the national Finnish ontology infrastructure. The name “Sampo” originates from this connection to the epic and has been re-used as a “brand” name in most of the offspring systems after that.

⁴⁰Project: <https://seco.cs.aalto.fi/applications/museumfinland/index.fi.php>; portal: <http://museosuomi.fi>

⁴¹Project: <https://seco.cs.aalto.fi/applications/kulttuurisampo/>; portal: <http://www.kulttuurisampo.fi/>.

BookSampo – Finnish Fiction Literature on the Semantic Web⁴² (online since 2011) [55,56] publishes meta-data about virtually all Finnish fiction literature as a knowledge graph on top of which a portal was created. BookSampo data was originally part of CultureSampo but is today maintained independently by the Public Libraries of Finland. BookSampo has grown into one of their main web services and is used by ca. 2 million users in a year.

WarSampo – Finnish World War II on the Semantic Web⁴³ (online since 2015 with several new application perspectives published in 2016–2019) [22] is a popular Finnish service that has had nearly a million of distinct users by April 2022 according to Google Analytics. It introduced the principle of using a clearly separated SPARQL endpoint (P3) into the Sampo model. Also faceted search methods were developed in some of the application perspectives of the portal as well as some visualizations showing way towards the principles P5–6. The portal's core dataset provides information about the casualties and significant soldiers of the Second World War in Finland. There are also many other datasets in the knowledge graph, such as authentic photographs from the fronts, war diaries, historical maps, and memoir articles of soldiers constituting a LOD cloud of its own and an infrastructure for the Finnish WW2 data [46]. WarSampo application won in 2017 the LODLAM Open Data Prize in Venice.

Interest in WarSampo lead to a new Sampo in the same application domain of war history: **WarVictimSampo (1914–1922)**⁴⁴ (online since 2019) [68] publishes data about the deaths and battles of the Finnish Civil War 1918 and the Kindred Wars. Also this portal has become fairly popular, as many citizens are still looking for information about their lost relatives in the Civil War. In 2021, the WarSampo infrastructure was re-used to create **WarMemoirSampo**⁴⁵ that publishes video interviews of WW2 veterans. Based on timestamped semantic video annotations, it is possible to search for particular points inside the videos that are automatically interlinked to additional contextualizing information in WarSampo and Wikipedia [23]. For example, when a person in WarSampo is mentioned, his/her homepage in WarSampo can be provided during video viewing. In the same vein, places that were annexed to the Soviet Union after the war can be shown on historical Finnish maps on top of contemporary Google Maps with links to related data, such as events and photographs.

Both WarSampo and WarVictimSampo have a feedback channel by which the data can be commented, and indeed hundreds of comments and correction suggestions have been collected for the data owner, the National Archives of Finland, to consider. Based on this activity, a new citizen science project for collecting and maintaining Sampo data is underway.⁴⁶

A key idea in WarSampo is to reassemble the life stories of the soldiers based on data linking from different data sources. This biographical and prosopographical idea was a source of inspiration for several later biographical Sampos discussed below.

BiographySampo – Biographies on the Semantic Web⁴⁷ (online since 2018) [26] is yet another popular service with tens of thousands of users. It harnessed the principle P5 of using the filter-analyse two-step cycle as well as the idea of integrating data-analysis and knowledge discovery in the Sampo model (P6), with a focus on supporting biographical and prosopographical research [80]. The system is based on mining out a large knowledge graph from ca. 13 100 Finnish national biographies of the Finnish Literature Society, authored by some 940 scholars. The data is interlinked and enriched internally and by 16 external data sources and by reasoning, e.g., by inferring family relations [50] and connections of interest between people and places [30].

The idea of publishing textual biographies as structured LOD for data exploration and analysis was also developed in the Sampos **Norssit Alumni**⁴⁸ [25] (online since 2017) based on the student registry of a prominent Finnish high school (1867–1992) and **U.S. Congress Prosopographer**⁴⁹ [62] (online since 2018) using data of the United States Congress members from the 1st through 115th Congresses (1789–2018). **AcademySampo**⁵⁰ (online since 2021)

⁴²Research project: <https://seco.cs.aalto.fi/applications/kirjasampo/>; portal of Finnish public libraries: <https://www.kirjasampo.fi/>.

⁴³Project: <https://seco.cs.aalto.fi/projects/sotasampo/>; portal: <https://www.sotasampo.fi/>.

⁴⁴Project: <https://seco.cs.aalto.fi/projects/sotasurmat-1914-1922/>; portal: <https://sotasurmat.narc.fi/>.

⁴⁵Project: <https://seco.cs.aalto.fi/projects/war-memoirs/>; portal: <https://sotamuistot.arkisto.fi/>.

⁴⁶<https://seco.cs.aalto.fi/projects/sotasampo/citizens/>

⁴⁷Project: <https://seco.cs.aalto.fi/projects/biografiasampo/>; portal: <https://biografiasampo.fi/>.

⁴⁸Portal: <http://www.norssit.fi/semweb/>.

⁴⁹Portal: <https://semanticcomputing.github.io/congress-legislators/>.

⁵⁰Project homepage: url <https://seco.cs.aalto.fi/projects/yo-matriikkelit/>; portal: <https://akatemiasampo.fi/>.

[50,51] is yet another biographical in-use system based on 28 000 short biographies of all known Finnish academic people educated in Finland in 1640–1899. The system includes a rich set of data-analytic tools for DH research [53].

NameSampo – A Linked Open Data Infrastructure and Workbench for Toponomastic Research⁵¹ (online since 2019) [40] publishes data about placenames and places in Finland with old maps. It soon attracted tens of thousands of users on the Web. NameSampo core data originates from the Name Archive of the Institute of Languages of Finland, a database of over 2 million placenames collected in Finland over several decades. NameSampo also published the contemporary placename register (ca. 800 000 places) of the National Survey of Finland as Linked Open Data. Furthermore, the Thesaurus of Geographical Names (TGN)⁵² of the Getty Research Institute via its SPARQL endpoint is re-used, as well as various map services, including a collection of historical maps of Finland published as part of WarSampo.

The NameSampo project developed, based on the SPARQL Faceter tool [44] used in many earlier Sampos, the first version of the Sampo-UI framework [39] that has been used after this in all Sampos, supporting implementation of principles P3–P6 from an UI point of view. Sampo-UI has also been re-used in Norway by the Norwegian Language Collections for creating a national service similar to NameSampo: Norske stedsnavn.⁵³ The Sampo-UI framework, available openly in Github,⁵⁴ has also been re-used in commercial settings.

Mapping Manuscript Migrations (MMM)⁵⁵ (online since 2020) [24,43] is a Sampo, in spite of its name, based on metadata about some 220 000 pre-modern manuscripts from the Schoenberg Database of Manuscripts⁵⁶ in the U.S., Medieval Manuscripts in the Oxford University Libraries⁵⁷ in the U.K., and the Bibale⁵⁸ database in France (cf. Fig. 1). MMM is a result of the Trans-Atlantic Digging into Data research programme.⁵⁹ Both the portal and the underlying SPARQL endpoint is used by scholars for studying the manuscripts [4].

FindSampo⁶⁰ [31,69] (online since 2021) is a system and data service for supporting archaeology especially from a citizen science and metal detectorists' perspectives.

The **LetterSampo framework**,⁶¹ discussed in Section 2, has been applied to create three Sampo instances based on different epistolary datasets [27]. The CKCC corpus of the Dutch Huygens Institute, and the correspSearch corpus aggregated by the Berlin-Brandenburg Academy of Sciences have been published as SPARQL endpoints on the Linked Data Finland platform and as a semantic portal⁶² based on their joint knowledge graph.

In addition, new Sampos are already in prototype phase and planned to be published in 2022: **LawSampo**⁶³ [35] publishes Finnish legislation and case law based on data from the Ministry of Justice in Finland. **ParliamentSampo**⁶⁴ [33] publishes LOD extracted from the materials of the Parliament of Finland (1907–2021),⁶⁵ including over 960 000 Parliamentary debate speeches [75] and prosopographical data about the politicians' networks [52] in 1907–2022.

⁵¹ <https://seco.cs.aalto.fi/projects/nimisampo/en/>

⁵² <http://www.getty.edu/research/tools/vocabularies/tgn/>

⁵³ <https://toponymi.spraksamlingane.no/nb/app>

⁵⁴ <https://github.com/SemanticComputing/sampo-ui>

⁵⁵ <https://seco.cs.aalto.fi/projects/mmm/>

⁵⁶ See <https://sdbm.library.upenn.edu>.

⁵⁷ See <https://medieval.bodleian.ox.ac.uk>.

⁵⁸ The Bibale web service from the Institute for Research and History of Texts (IRHT) in Paris is described in <http://bibale.irht.cnrs.fr>.

⁵⁹ <https://diggingintodata.org/>

⁶⁰ <https://seco.cs.aalto.fi/projects/suall/>

⁶¹ <https://seco.cs.aalto.fi/projects/rrl/>

⁶² Portal: <https://lettersampo.demo.seco.cs.aalto.fi/>.

⁶³ <https://seco.cs.aalto.fi/projects/lawlod/>

⁶⁴ <https://seco.cs.aalto.fi/projects/sem parl/en/>

⁶⁵ <https://seco.cs.aalto.fi/projects/sem parl/en/>

4. Discussion

Design principles, models, and methods for software development are extensively studied and used in the field of Software Engineering [76]. The idea of trying to formulate general design principles for publishing and using linked data has turned out to be useful from a practical point of view. For example, the four Linked Data Principles and the 5-star model coined by Tim Berners-Lee have been quite influential, and ontology design patterns⁶⁶ are (re-)used as guidelines for data modelling. In the same vein, the FAIR principles for publishing data are widely used today. Also the Sampo model can be seen as a kind of effort for formulating a set of principles for publishing and using linked data in semantic portals. Our experiences in developing the Sampo series of data services and portals provide an empirical evaluation and evidence about the usability of the model in practical applications. The application domains of the model (cf. Table 2) are versatile including tangible and intangible cultural heritage collections, biography and prosopography, toponomastic research, manuscript studies, archaeology, legislation, and parliamentary studies. In many cases, language barriers have been crossed based on the language-agnostic ontology technology [77] of the Semantic Web.

Related work The Principles of Table 1 behind the Sampo model have been explored and developed before in different contexts:

1. The principle of collaborative content creation by data linking (P1) is a fundamental idea behind the Linked Open Data Cloud movement,⁶⁷ social media, and has been developed also in various other settings, e.g., in ResearchSpace⁶⁸ and WissKI.⁶⁹
2. The importance of developing shared open data models, thesauri, and ontologies for interoperability (P2) is a driving force behind the work of virtually all related standardization efforts. In our work, the ambitious goal has been to develop not only individual standards and datasets but an infrastructure in a national level effort [38] in terms of open ontology services [82,86] and LOD services [36].
3. The principle P3 of separating data related services from UI design is in line with modern software architectures, such as the Model-View-Controller (MVC) structure.⁷⁰ The Sampo model supports the idea of “separation of concerns” where each software layer can focus solely on its own role, and uses the Web idea of using the simple HTTP protocol for creating services based on distributed services.
4. The principle P4 of providing multiple analyses and visualizations for a set of filtered search results has been used in different contexts and also in other portals, such as the ePistolarium⁷¹ [70] for epistolary data. The idea of using multiple perspectives has also been studied as an approach in decision making [54].
5. Regarding principle P5, faceted search [10,32,81], also known as “view-based search” [67] and “dynamic taxonomies” [73], is a well-known paradigm for explorative search and browsing [59] in computer science and information retrieval, based on S. R. Ranganathan’s original ideas of faceted classification in Library Science in the 1930’s. The two-step usage model in the Sampo model is also used as a general research method in prosopographical research [85].
6. The principle of supporting data analysis and knowledge discovery (P6) based on Big Data is fundamental in, e.g., distant reading [63], Humanities Computing [60], and Digital Humanities [6] in general. However, what is still largely missing in the DH methodology and tools in semantic portals is the next conceptual level of automatic knowledge discovery from data [66]. The Sampo model aims to integrate such tools into a consolidated approach for creating portals and LOD services.

In addition to the Sampo portals, Linked Data and ontologies have been used as a basis for publishing collections in many museums [1,74], libraries [9], and archives [8,84]. Linked data has been used in building knowledge graphs

⁶⁶http://ontologydesignpatterns.org/wiki/Main_Page

⁶⁷<https://lod-cloud.net>

⁶⁸<https://www.researchspace.org>

⁶⁹<https://wiss-ki.eu/>

⁷⁰<https://en.wikipedia.org/wiki/Model-view-controller>

⁷¹<http://ckcc.huygens.knaw.nl>

and infrastructures, such as the Europeana linked open data [41] and ARIADNEplus⁷² for archaeology in Europe, Linked Art⁷³ in the U.S. [79], and in local efforts in Italy [5], the U.K. [49], and Spain [12] to list a few examples. Cultural heritage and DH have become a major application domain for Linked Data and the Semantic Web [3,88].

Contributions and Challenges The novelty of the Sampo model lies in the consolidated combination of the principles P1–P6 and in operationalizing them using an infrastructure and tooling for developing applications in Digital Humanities in a cost-efficient way. The approach aims at developing a gradually growing sustainable national LOD infrastructure: the work started with the Semantic Web Kick-off in Finland seminar [16] a few months after the seminal Semantic Web paper [2] was published in Scientific American and the W3C launched its Semantic Web Activity in 2001. The work presented demonstrates a shift of focus in research on CH semantic portals in three generations towards knowledge discovery and Artificial Intelligence [21]. The future work on the Sampo model aims at AI based DH tools that are able not only to present the data to the human researcher in useful ways but also to find and solve DH research problems with explanations. AI techniques are also useful when creating and enriching the knowledge graph underlying a semantic portal.

The model has also its limitations and challenges. For example, it assumes that the data is created by a separate pipeline. As suggested in [46], the transformation should be automatic and re-doable without a human in the loop, but optimally the RDF should be produced already when cataloging the data, not by correcting and aligning the data afterwards [17]. As Albert Einstein put it: “Intellectuals solve problems but geniuses prevent them”.

Neither does the Sampo model include principles for maintaining the knowledge graphs. This is a great challenge when using LD in general since the effects of a change may propagate all over the interlinked knowledge graph depending on the change. Especially changes in ontologies may have dramatic effects. Such knowledge management issues are discussed in [45] in relation to the WarSampo system, but more research is needed in this field.

A challenge of the semantic portal concept is related to the quality of the data produced typically using more or less automatic means, leading to problems of incomplete, skewed, and erroneous data [57]. This as well as conceptual difficulties in modeling complex real world ontologies, such as historical geogazetteers, become sometimes embarrassingly visible when using and exposing the knowledge structures to end users. In traditional systems the same problems are there, but are hidden in the unstructured presentations of the data. In general, more data literacy [47] is usually needed from the end user when using semantic portals and their data analytic tools. In spite of these challenges the linked data approach is according to our experiences useful in finding out interesting phenomena in Big Data using distant reading [63], but for interpreting the results traditional close reading is needed as before.

Acknowledgements

Tens of people⁷⁴ at the Semantic Computing Research Group (SeCo) have contributed to the Sampo model and systems, funded by over 50 organizations in Finland and beyond. Thanks to Marcia Zeng for inspirational discussions related to the notion of the Sampo model. Kai Eckert and Christoph Schlieder reviewed an earlier version of this paper and made insightful comments to it. The work of writing this paper is partly supported by the Semantic Parliament (ParliamentSampo) project of the Academy of Finland, the EU project InTaVia: In/Tangible European Heritage,⁷⁵ and the EU COST action Nexus Linguarum⁷⁶ on linguistic data science. CSC – IT Center for Science has provided computational resources for our projects.

References

- [1] L. Aroyo, N. Stash, Y. Wang, P. Gorgels and L. Rutledge, *CHIP Demonstrator: Semantics-Driven Recommendations and Museum Tour Generation*, in: *The Semantic Web*, Springer, 2007, pp. 879–886. doi:10.1007/978-3-540-76298-0_64.

⁷²<https://ariadne-infrastructure.eu/>

⁷³<https://linked.art/>

⁷⁴<https://seco.cs.aalto.fi/people/>

⁷⁵<https://intavia.eu/>

⁷⁶<https://nexuslinguarum.eu/the-action>

- [2] T. Berners-Lee, J. Hendler and O. Lassila, *The Semantic Web*, *Scientific American* **284**(5) (2001), 34–43.
- [3] A. Bikakis, E. Hyvönen, S. Jean, B. Markhoff and A. Mosca (eds), 2021, Special Issue on Semantic Web for Cultural Heritage, *Semantic Web – Interoperability, Usability, Applicability* **12**(2). doi:[10.3233/SW-210425](https://doi.org/10.3233/SW-210425).
- [4] T. Burrows, L. Cleaver, D. Emery, E. Hyvönen, M. Koho, L. Ransom, E. Thomson and H. Wijsman, Medieval manuscripts and their migrations: Using SPARQL to investigate the research potential of an aggregated Knowledge Graph, *Digital Medievalist* (2022), forthcoming, <https://seco.cs.aalto.fi/publications/2021/burrows-et-al-digital-medievalist-2021.pdf>.
- [5] V.A. Carriero, A. Gangemi, M.L. Mancinelli, L. Marinucci, A.G. Nuzzolese, V. Presutti and C. Veninata, ArCo: The Italian cultural heritage knowledge graph, in: *The Semantic Web – ISWC 2019*, Springer, 2019, pp. 36–52. doi:[10.1007/978-3-030-30796-7_3](https://doi.org/10.1007/978-3-030-30796-7_3).
- [6] E. Gardiner and R.G. Musto, *The Digital Humanities: A Primer for Students and Scholars*, Cambridge University Press, New York, NY, USA, 2015. doi:[10.1017/CBO9781139003865](https://doi.org/10.1017/CBO9781139003865).
- [7] C. Gutierrez and J.F. Sequeda, Knowledge graphs, *Communications of the ACM* **64**(3) (2021), 96–104. doi:[10.1145/3418294](https://doi.org/10.1145/3418294).
- [8] M. Hallo, S. Luján-Mora, A. Maté and J. Trujillo, Current state of linked data in digital libraries, *Journal of Information Science* **42**(2) (2016), 117–127. doi:[10.1177/0165551515559479](https://doi.org/10.1177/0165551515559479).
- [9] B. Haslhofer, A. Isaac and R. Simon, Knowledge graphs in the libraries and digital humanities domain, 2018, arXiv preprint [arXiv:1803.03198](https://arxiv.org/abs/1803.03198).
- [10] M. Hearst, Design recommendations for hierarchical faceted search interfaces, in: *ACM SIGIR Workshop on Faceted Search*, Seattle, WA, 2006, pp. 1–5.
- [11] T. Heath and C. Bizer, *Linked Data: Evolving the Web into a Global Data Space*, 1st edn, Morgan & Claypool, Palo Alto, California, 2011. doi:[10.2200/S00334ED1V01Y201102WBE001](https://doi.org/10.2200/S00334ED1V01Y201102WBE001).
- [12] F. Hernández, L. Rodrigo, J. Contreras and F. Carbone, Building a cultural heritage ontology for Cantabria, in: *Annual Conference of CIDOC*, 2008, pp. 1–14, https://cidoc.mini.icom.museum/wp-content/uploads/sites/6/2018/12/64_papers.pdf.
- [13] P. Hitzler, A review of the semantic web field, *Commun. ACM* **64**(2) (2021), 76–83. doi:[10.1145/3397512](https://doi.org/10.1145/3397512).
- [14] P. Hitzler, M. Krötzsch and S. Rudolph, *Foundations of Semantic Web Technologies*, Springer, 2010.
- [15] H. Hotson and T. Wallnig (eds), *Reassembling the Republic of Letters in the Digital Age*, Göttingen University Press, 2019. doi:[10.17875/gup2019-1146](https://doi.org/10.17875/gup2019-1146).
- [16] E. Hyvönen (ed.), *Semantic Web Kick-Off in Finland – Vision, Technologies, Research, and Applications*, in *HIIT Publications 2002-01*, 2002, <http://www.seco.hut.fi/publications/2002/hyvonen-semantic-web-kick-off-2002.pdf>.
- [17] E. Hyvönen, Preventing interoperability problems instead of solving them, semantic web, *Interoperability, Usability, Applicability* **1**(1–2) (2010), 33–37. doi:[10.3233/SW-2010-0014](https://doi.org/10.3233/SW-2010-0014).
- [18] E. Hyvönen, *Publishing and Using Cultural Heritage Linked Data on the Semantic Web*, Morgan & Claypool, Palo Alto, California, 2012. doi:[10.2200/S00452ED1V01Y201210WBE003](https://doi.org/10.2200/S00452ED1V01Y201210WBE003).
- [19] E. Hyvönen, “sampo” model and semantic portals for digital humanities on the semantic web, in: *DHN 2020 Digital Humanities in the Nordic Countries. Proceedings of the Digital Humanities in the Nordic Countries 5th Conference*, CEUR Workshop Proceedings, Vol. 2612, 2020, pp. 373–378, <http://ceur-ws.org/Vol-2612/poster1.pdf>.
- [20] E. Hyvönen, Linked open data infrastructure for digital humanities in Finland, in: *DHN 2020 Digital Humanities in the Nordic Countries*, Proceedings of the Digital Humanities in the Nordic Countries 5th Conference, CEUR Workshop Proceedings, Vol. 2612, 2020, pp. 254–259, <http://ceur-ws.org/Vol-2612/short10.pdf>.
- [21] E. Hyvönen, Using the semantic web in digital humanities: Shift from data publishing to data-analysis and serendipitous knowledge discovery, semantic web, *Interoperability, Usability, Applicability* **11**(1) (2020), 187–193.
- [22] E. Hyvönen, E. Heino, P. Leskinen, E. Ikkala, M. Koho, M. Tamper, J. Tuominen and E. Mäkelä, WarSampo data service and semantic portal for publishing linked open data about the second world war history, in: *The Semantic Web – Latest Advances and New Domains (ESWC 2016)*, H. Sack, E. Blomqvist, M. d’Aquin, C. Ghidini, S.P. Ponzetto and C. Lange, eds, Springer, 2016, pp. 758–773. doi:[10.1007/978-3-319-34129-3_46](https://doi.org/10.1007/978-3-319-34129-3_46).
- [23] E. Hyvönen, E. Ikkala, M. Koho, R. Leal, H. Rantala and M. Tamper, How to search and contextualize scenes inside videos for enriched watching experience: Case stories of the second world war veterans, in: *Proceedings of the 19th Extended Semantic Web Conference (ESWC 2022)*, *Poster and Demo Papers*, 2022, forthcoming, <https://seco.cs.aalto.fi/publications/2022/hyvonen-et-al-wms-2022.pdf>.
- [24] E. Hyvönen, E. Ikkala, M. Koho, J. Tuominen, T. Burrows, L. Ransom and H. Wijsman, Mapping manuscript migrations on the semantic web: A semantic portal and linked open data service for premodern manuscript research, in: *The Semantic Web – ISWC 2021*, A. Hotho, E. Blomqvist, S. Dietze, A. Fokoue, Y. Ding, P. Barnaghi, A. Haller, M. Dragoni and H. Alani, eds, Lecture Notes in Computer Science, Springer, 2021, pp. 615–630. ISBN 978-3-030-88360-7. doi:[10.1007/978-3-030-88361-4_36](https://doi.org/10.1007/978-3-030-88361-4_36).
- [25] E. Hyvönen, P. Leskinen, E. Heino, J. Tuominen and L. Sirola, Reassembling and enriching the life stories in printed biographical registers: Norssi high school alumni on the semantic web, in: *Proceedings, Language, Technology and Knowledge (LDK 2017)*, Springer, 2017, pp. 113–119. doi:[10.1007/978-3-319-59888-8_9](https://doi.org/10.1007/978-3-319-59888-8_9).
- [26] E. Hyvönen, P. Leskinen, M. Tamper, H. Rantala, E. Ikkala, J. Tuominen and K. Keravuori, BiographySampo – publishing and enriching biographies on the semantic web for digital humanities research, in: *The Semantic Web. 16th International Conference, ESWC 2019*, Springer, 2019, pp. 574–589. doi:[10.1007/978-3-030-21348-0_37](https://doi.org/10.1007/978-3-030-21348-0_37).
- [27] E. Hyvönen, P. Leskinen and J. Tuominen, LetterSampo – Historical Letters on the Semantic Web: A Framework and Its Application to Publishing and Using Epistolary Data of the Republic of Letters, 2022, Submitted for peer review, <https://seco.cs.aalto.fi/publications/2020/hyvonen-et-al-lettersampo-2020.pdf>.

- [28] E. Hyvönen, E. Mäkelä, T. Kauppinen, O. Alm, J. Kurki, T. Ruotsalo, K. Seppälä, J. Takala, K. Puputti, H. Kuittinen, K. Viljanen, J. Tuominen, T. Palonen, M. Frosterus, R. Sinkkilä, P. Paakkari, J. Laitio and K. Nyberg, CultureSampo – Finnish culture on the Semantic Web 2.0. Thematic perspectives for the end-user, in: *Museums and the Web 2009*, Archives & Museum Informatics, Toronto, 2009, <https://www.archimuse.com/mw2009/papers/hyvonen/hyvonen.html>.
- [29] E. Hyvönen, E. Mäkelä, M. Salminen, A. Valo, K. Viljanen, S. Saarela, M. Junnila and S. Kettula, MuseumFinland—Finnish museums on the Semantic Web, *Journal of Web Semantics* 3(2) (2005), 224–241. doi:10.1016/j.websem.2005.05.008.
- [30] E. Hyvönen and H. Rantala, *Knowledge-Based Relational Search in Cultural Heritage Linked Data, Digital Scholarship in the Humanities (DSH)*, Vol. 36, Oxford University Press, 2021, pp. 55–64. doi:10.1093/lhc/fqab042.
- [31] E. Hyvönen, H. Rantala, E. Ikkala, M. Koho, J. Tuominen, B. Anafi, S. Thomas, A. Wessman, E. Oksanen, V. Rohiola, J. Kuitunen and M. Ryyppö, *Citizen Science Archaeological Finds on the Semantic Web: The FindSampo Framework, Antiquity, a Review of World Archaeology* 95(382), 2021, p. E24. doi:10.15184/auq.2021.87.
- [32] E. Hyvönen, S. Saarela and K. Viljanen, Application of ontology-based techniques to view-based semantic search and browsing, in: *Proceedings of the First European Semantic Web Symposium*, Springer, 2004. doi:10.1007/978-3-540-25956-5_7.
- [33] E. Hyvönen, L. Sinikallio, P. Leskinen, M.L. Mela, J. Tuominen, K. Elo, S. Drobac, M. Koho, E. Ikkala, M. Tamper, R. Leal and J. Kesäniemi, *Finnish Parliament on the Semantic Web: Using ParliamentSampo Data Service and Semantic Portal for Studying Political Culture and Language*, in: *Digital Parliamentary Data in Action (DIPADA 2022), Workshop at the 6th Digital Humanities in Nordic and Baltic Countries Conference, CEUR Workshop Proceedings*, 2022, forth-coming, <https://seco.cs.aalto.fi/publications/2022/hyvonen-et-al-semparl-dhnb-2022.pdf>.
- [34] E. Hyvönen, J. Takala, O. Alm, T. Ruotsalo and E. Mäkelä, Semantic Kalevala – accessing cultural contents through semantically annotated stories, in: *Proceedings of the Cultural Heritage on the Semantic Web Workshop at the 6th International Semantic Web Conference (ISWC 2007)*, Busan, Korea, 2007, <https://seco.cs.aalto.fi/publications/2007/hyvonen-et-al-kalevala-2007.pdf>.
- [35] E. Hyvönen, M. Tamper, E. Ikkala, S. Sarsa, A. Oksanen, J. Tuominen and A. Hietanen, Publishing and using legislation and case law as linked open data on the semantic web, in: *The Semantic Web: ESWC 2020 Satellite Events*, Lecture Notes in Computer Science, Vol. 12124, Springer, 2020, pp. 110–114. doi:10.1007/978-3-030-62327-2_19.
- [36] E. Hyvönen, J. Tuominen, M. Alonen and E. Mäkelä, Linked data Finland: A 7-star model and platform for publishing and re-using linked datasets, in: *The Semantic Web: ESWC 2014 Satellite Events*, Springer, 2014, pp. 226–230. doi:10.1007/978-3-319-11955-7_24.
- [37] E. Hyvönen, K. Viljanen and O. Suominen, HealthFinland – Finnish health information on the Semantic Web, in: *The Semantic Web*, K. Aberer, K.-S. Choi, N. Noy, D. Allemang, K.-I. Lee, L. Nixon, J. Golbeck, P. Mika, D. Maynard, R. Mizoguchi, G. Schreiber and P. Cudré-Mauroux, eds, Springer, 2007. doi:10.1007/978-3-540-76298-0_560.
- [38] E. Hyvönen, K. Viljanen, J. Tuominen and K. Seppälä, *Building a National Semantic Web Ontology and Ontology Service Infrastructure – the FinnONTO Approach*, in: *The Semantic Web: Research and Applications, 5th European Semantic Web Conference, ESWC 2008*, Springer, 2008, pp. 95–109. doi:10.1007/978-3-540-68234-9_10.
- [39] E. Ikkala, E. Hyvönen, H. Rantala and M. Koho, Sampo-UI: A full stack JavaScript framework for developing semantic portal user interfaces, semantic web, *Interoperability, Usability, Applicability* 13(1) (2022), 69–84. doi:10.3233/SW-210428.
- [40] E. Ikkala, J. Tuominen, J. Raunamaa, T. Aalto, T. Ainiala, H. Uusitalo and E. Hyvönen, NameSampo: A linked open data infrastructure and workbench for toponomastic research, in: *Proceedings of the 2nd ACM SIGSPATIAL Workshop on Geospatial Humanities, GeoHumanities '18*, ACM, New York, NY, USA, 2018, pp. 2:1–2:9. ISBN 978-1-4503-6032-6. doi:10.1145/3282933.3282936.
- [41] A. Isaac and B. Haslhofer, Europeana linked open data – data.Europeana.eu, semantic web, *Interoperability, Usability, Applicability* 4(3) (2013), 291–297. doi:10.3233/SW-120092.
- [42] G. Kobilarov, T. Scott, Y. Raimond, S. Oliver, C. Sizemore, M. Smethurst, C. Bizer and R. Lee, *Media Meets Semantic Web – How the BBC Uses DBpedia and Linked Data to Make Connections*, in: *The Semantic Web: Research and Applications*, Springer, Berlin Heidelberg, Berlin, Heidelberg, 2009, pp. 723–737. ISBN 978-3-642-02121-3.
- [43] M. Koho, T. Burrows, E. Hyvönen, E. Ikkala, K. Page, L. Ransom, J. Tuominen, D. Emery, M. Fraas, B. Heller, D. Lewis, A. Morrison, G. Porte, E. Thomson, A. Velios and H. Wijsman, Harmonizing and publishing heterogeneous pre-modern manuscript metadata as linked open data, *Journal of the Association for Information Science and Technology (JASIST)* 73(2) (2022), 240–257. doi:10.1002/asi.24499.
- [44] M. Koho, E. Heino and E. Hyvönen, SPARQL faceter – client-side faceted search based on SPARQL, in: *Joint Proc. of the 4th International Workshop on Linked Media and the 3rd Developers Hackshop, CEUR Workshop Proceedings*, Vol. 1615, 2016, <http://ceur-ws.org/Vol-1615/semdevPaper5.pdf>.
- [45] M. Koho, E. Ikkala, E. Heino and E. Hyvönen, Maintaining a linked data cloud and data service for second world war history, in: *Digital Heritage. Progress in Cultural Heritage: Documentation, Preservation, and Protection. 7th International Conference, EuroMed 2018*, Springer, 2018. doi:10.1007/978-3-030-01762-0.
- [46] M. Koho, E. Ikkala, P. Leskinen, M. Tamper, J. Tuominen and E. Hyvönen, WarSampo knowledge graph: Finland in the second world war as linked open data, semantic web, *Interoperability, Usability, Applicability* 12(2) (2021), 265–278. doi:10.3233/SW-200392.
- [47] T. Koltay, Data literacy for researchers and data librarians, *Journal of Librarianship and Information Science* 49(1) (2015), 3–14. doi:10.1177/0961000615616450.
- [48] J.E. Labra Gayo, E. Prud'hommeaux, I. Boneva and D. Kontokostas, *Validating RDF Data, Synthesis Lectures on the Semantic Web: Theory and Technology*, Vol. 7, Morgan & Claypool Publishers LLC, 2017, pp. 1–328. doi:10.2200/s00786ed1v01y201707wbe016.
- [49] Y. Lei, V. Lopez, E. Motta and V. Uren, An infrastructure for semantic web portals, *Journal of Web Engineering* 6(4) (2007), 283–308, <https://journals.riverpublishers.com/index.php/JWE/article/view/4105>.

- [50] P. Leskinen and E. Hyvönen, Linked open data service about historical Finnish academic people in 1640–1899, in: *DHN 2020 Digital Humanities in the Nordic Countries. Proceedings of the Digital Humanities in the Nordic Countries 5th Conference*, CEUR Workshop Proceedings, Vol. 2612, 2020, pp. 284–292, <http://ceur-ws.org/Vol-2612/short14.pdf>.
- [51] P. Leskinen and E. Hyvönen, Reconciling and using historical person registers as linked open data in the AcademySampo knowledge graph, in: *The Semantic Web – ISWC 2021. 20th International Semantic Web Conference, ISWC 2021, Proceedings*, Springer, 2021, pp. 714–730. doi:10.1007/978-3-030-88361-4_42.
- [52] P. Leskinen, E. Hyvönen and J. Tuominen, Members of Parliament of Finland knowledge graph and its linked open data service, in: *Proceedings of SEMANTICS – in the Era of Knowledge Graphs*, 6–9, 2021, Studies on the Semantic Web, IOS Press, Amsterdam, 2021, pp. 255–269. doi:10.3233/SSW210049.
- [53] P. Leskinen, H. Rantala and E. Hyvönen, Analyzing the lives of Finnish academic people 1640–1899 in nordic and Baltic countries: AcademySampo data service and portal, in: *6th Digital Humanities in Nordic and Baltic Countries Conference, Proceedings, CEUR Workshop Proceedings*, 2022, forth-coming, <https://seco.cs.aalto.fi/publications/2022/leskinen-et-al-academysampo-dhnb-2022.pdf>.
- [54] H.A. Linstone, Multiple perspectives: Concept, applications, and user guidelines, *Systems practice* 2(3) (1989), 307–331. doi:10.1007/BF01059977.
- [55] E. Mäkelä, K. Hypén and E. Hyvönen, BookSampo—lessons learned in creating a semantic portal for fiction literature, in: *The Semantic Web – ISWC 2011 10th International Semantic Web Conference, Proceedings, Part II*, Springer, 2011, pp. 173–188. doi:10.1007/978-3-642-25093-4_12.
- [56] E. Mäkelä, K. Hypén and E. Hyvönen, Fiction literature as linked open data – the BookSampo dataset, *Semantic Web, Interoperability, Usability, Applicability* 4(3) (2013), 299–306. doi:10.3233/SW-120093.
- [57] E. Mäkelä, K. Lagus, L. Lahti, T. Säily, M. Tolonen, M. Hämäläinen, S. Kaislaniemi and T. Nevalainen, Wrangling with non-standard data, in: *Proceedings of the Digital Humanities in the Nordic Countries 5th Conference, CEUR Workshop Proceedings*, 2020, pp. 81–96, <http://ceur-ws.org/Vol-2612/paper6.pdf>.
- [58] E. Mäkelä, T. Ruotsalo and Hyvönen, How to deal with massively heterogeneous cultural heritage data—lessons learned in CultureSampo, *Semantic Web, Interoperability, Usability, Applicability* 3(1) (2012), 85–109. doi:10.3233/SW-2012-0049.
- [59] G. Marchionini, Exploratory search: From finding to understanding, *Communications of the ACM* 49(4) (2006), 41–46. doi:10.1145/1121949.1121979.
- [60] W. McCarty, *Humanities Computing*, Palgrave, London, 2005.
- [61] A. Meroño-Peñuela and R. Hoekstra, gric makes GitHub taste like linked data APIs, in: *The Semantic Web. Latest Advances and New Domains. The 13th International Conference, ESWC 2016*, Springer, 2016, pp. 342–353. doi:10.1007/978-3-319-47602-5_48.
- [62] G. Miyakita, P. Leskinen and E. Hyvönen, Using linked data for prosopographical research of historical persons: Case U.S. congress legislators, in: *Digital Heritage. Progress in Cultural Heritage: Documentation, Preservation, and Protection. 7th International Conference, EuroMed 2018*, Springer, 2018. doi:10.1007/978-3-030-01765-1.
- [63] F. Moretti, *Distant Reading*, Verso Books, 2013.
- [64] E. Oksanen, H. Rantala, J. Tuominen, M. Lewis, D. Wigg-Wolf, F. Ehrnsten and E. Hyvönen, Digital humanities solutions for pan-European numismatic and archaeological heritage based on linked open data, in: *Proceedings of the Digital Humanities in Nordic and Baltic Countries 2022, CEUR Workshop Proceedings*, 2022, forthcoming, <https://seco.cs.aalto.fi/publications/2022/oksanen-et-al-diginuma-dhnb-2022.pdf>.
- [65] T. Palonen, J. Hyvönen, J. Takala and E. Hyvönen, *Semanttinen Kalevala – Kulttuurisammon taontaa (Semantic Kalevala – Forging the CultureSampo)*, *eLore* 16(2), 2009. doi:10.30666/elore.78811.
- [66] M.J. Pazzani, Knowledge discovery from data?, *IEEE Intelligent Systems* 15(2) (2000), 10–13. doi:10.1109/5254.850821.
- [67] A.S. Pollitt, The key role of classification and indexing in view-based searching, Technical Report, Centre for Database Access Research, University of Huddersfield, 1998.
- [68] H. Rantala, E. Ikkala, I. Jokipii, M. Koho, J. Tuominen and E. Hyvönen, WarVictimSampo 1914–1922: A national war memorial on the Semantic Web for digital humanities research and applications, *ACM Journal on Computing and Cultural Heritage* 15(1) (2022), 1–18. doi:10.1145/3477606.
- [69] H. Rantala, E. Ikkala, V. Rohiola, M. Koho, J. Tuominen, E. Oksanen, A. Wessman and E. Hyvönen, FindSampo: A linked data based portal and data service for analyzing and disseminating archaeological object finds, in: *Proceedings of the 19th Extended Semantic Web Conference (ESWC 2022)*, Springer, 2022, forth-coming, <https://seco.cs.aalto.fi/publications/2022/rantala-et-al-findsampo-2022.pdf>.
- [70] W. Ravenek, C. van den Heuvel and G. Gerritsen, The ePistolarium: Origins and techniques, in: *CLARIN in the Low Countries*, A. van Hessen and J. Odijk, eds, Ubiquity Press, 2017, pp. 317–323. doi:10.5334/bbi.
- [71] L. Rietveld and R. Hoekstra, The YASGUI family of SPARQL clients, *Semantic Web, Interoperability, Usability, Applicability* 8(3) (2017), 373–383. doi:10.3233/SW-150197.
- [72] P. Riva, M. Doerr and M. Zumer, FRBRoo: Enabling a common view of information from memory institutions, in: *World Library and Information Congress: 74th IFLA General Conference and Council*, 2008, https://archive.ifla.org/IV/ifla74/papers/156-Riva_Doerr_Zumer-en.pdf.
- [73] G.M. Sacco, Dynamic taxonomies for intelligent information access, in: *Encyclopedia of Information Science and Technology*, M. Khosrow-Pour, ed., 3rd edn, 2015, pp. 3883–3892. doi:10.4018/978-1-4666-5888-2.ch382.
- [74] G. Schreiber, A. Amin, L. Aroyo, M. van Assem, V. de Boer, L. Hardman, M. Hildebrand, B. Omelayenko, J. van Osenbruggen, A. Tordai, J. Wielemaker and B. Wielinga, Semantic annotation and search of cultural-heritage collections: The MultimediaN E-culture demonstrator, *Journal of Web Semantics* 6(4) (2008), 243–249. doi:10.1016/j.websem.2008.08.001.

- [75] L. Sinikallio, S. Drobac, M. Tamper, R. Leal, M. Koho, J. Tuominen, M.L. Mela and E. Hyvönen, Plenary debates of the Parliament of Finland as linked open data and in parla-CLARIN markup, in: *Proceedings, Language, Data and Knowledge (LDK, 2021)*, Vol. 93, Schloss Dagstuhl – Leibniz-Zentrum für Informatik GmbH, 2021, pp. 1–17. doi:[10.4230/OASICS.LDK.2021.8](https://doi.org/10.4230/OASICS.LDK.2021.8).
- [76] I. Sommerville, *Software Engineering*, 10th edn, Pearson, 2016.
- [77] S. Staab and R. Studer (eds), *Handbook on Ontologies*, 2nd edn, Springer, 2009.
- [78] O. Suominen, E. Hyvönen, K. Viljanen and E. Hukka, HealthFinland – a national semantic publishing network and portal for health information, *Journal of Web Semantics* 7(4) (2009), 287–297. doi:[10.1016/j.websem.2009.09.003](https://doi.org/10.1016/j.websem.2009.09.003).
- [79] P. Szekely, C.A. Knoblock, F. Yang, E.E. Fink, S. Gupta, R. Allen and G. Goodlander, Publishing the data of the Smithsonian American art museum to the linked data cloud, *International Journal of Humanities and Arts Computing* 8(supplement) (2014), 152–166, <http://usc-isi-i2.github.io/papers/szekely14-ijhac.pdf>. doi:[10.3366/ijhac.2014.0104](https://doi.org/10.3366/ijhac.2014.0104).
- [80] M. Tamper, P. Leskinen, E. Hyvönen, R. Valjus and K. Keravuori, *Analyzing Biography Collection Historiographically as Linked Data: Case National Biography of Finland, Semantic Web – Interoperability, Usability, Applicability*, 2021, forth-coming, <https://seco.cs.aalto.fi/publications/2021/tamper-et-al-bs-2021.pdf>.
- [81] D. Tunkelang, *Faceted Search*, Morgan & Claypool, Palo Alto, California, 2009. doi:[10.2200/S00190ED1V01Y200904ICR005](https://doi.org/10.2200/S00190ED1V01Y200904ICR005).
- [82] J. Tuominen, M. Frosterus, K. Viljanen and E. Hyvönen, *ONKI SKOS Server for Publishing and Utilizing SKOS Vocabularies and Ontologies as Services*, in: *The Semantic Web: Research and Applications: 6th European Semantic Web Conference, ESWC 2009*, Springer, 2009, pp. 768–780. doi:[10.1007/978-3-642-02121-3_56](https://doi.org/10.1007/978-3-642-02121-3_56).
- [83] Y. Tzitzikas, N. Manolis and P. Papadakos, Faceted exploration of RDF/s datasets: A survey, *Journal of Intelligent Information Systems* 48(2) (2017), 329–364. doi:[10.1007/s10844-016-0413-8](https://doi.org/10.1007/s10844-016-0413-8).
- [84] S. Van Hooland and R. Verborgh, *Linked Data for Libraries, Archives and Museums: How to Clean, Link and Publish Your Metadata*, Facet Publishing, 2014. doi:[10.1080/00048623.2016.1162277](https://doi.org/10.1080/00048623.2016.1162277).
- [85] K. Verboven, M. Carlier and J. Dumolyn, A short manual to the art of prosopography, in: *Prosopography Approaches and Applications. A Handbook, Unit for Prosopographical Research*, Linacre, College, 2007, pp. 35–70. doi:[1854/8212](https://doi.org/10.1017/9780521882122.003).
- [86] K. Viljanen, J. Tuominen and E. Hyvönen, Ontology libraries for production use: The Finnish ontology library service ONKI, in: *The Semantic Web: Research and Applications: 6th European Semantic Web Conference, ESWC 2009*, Springer, 2009, pp. 781–795. doi:[10.1007/978-3-642-02121-3_57](https://doi.org/10.1007/978-3-642-02121-3_57).
- [87] M. Zeng and J. Qin, *Metadata*, 3rd edn, ALA Neal-Schuman, Chicago, 2022. ISBN 978-0-8389-4875-0.
- [88] M. Zeng, C. Sula, K. Gracy, E. Hyvönen and V.M.A. Lima, (eds), JASIST Special Issue on Digital Humanities (DH), *Journal of the Association for Information Science and Technology (JASIST)* (2021), 1–5. doi:[10.1002/asi.24584](https://doi.org/10.1002/asi.24584).