



ICTeSSH2021



SSHOC

social sciences & humanities open cloud

SSH Vocabulary Initiative What users want

June 28, 9:00 - 10:30 CEST

Agenda

1. Welcome - Marieke Willems (Trust-IT)
2. SSH Vocabulary Initiative - Daan Broeder (CLARIN)
3. How can researchers use vocabulary in SSH tools & use cases
 - a. **USAGE:** Using Vocabularies in the SSH community - Iulianna van der Lek (CLARIN)
 - b. **FINDABILITY:** Current ways of locating suitable vocabularies and using the SSH Open Marketplace - Matej Ďurčo (ACDC-CH & DARIAH)
 - c. **INTEROPERABILITY:** The Vocabulary Matching Tool - Holly Wright (Archaeology Data Service)
 - d. **INTEROPERABILITY:** Using MT to create multilingual vocabularies - Monica Monachini (CNR-ILC & CLARIN-IT)
4. The SSH Vocabulary Initiative - What users want. **Panel Discussion**
5. Wrap up



Monica Monachini
CLARIN



Holly Wright
Archaeology Data Service



Daan Broeder
CLARIN

#SSHOCaVocabulary
#SSHOCifyCLARIN



Taina Jääskeläinen
Finnish Social Science
Data Archive



Matej Durco
ACDH-CH



**Iuliana van
der Lek**
CLARIN



CLARIN



We are interested in
learning more about
your experience
with vocabularies.
Help us out by
filling in the
following polls!

SSH VOCABULARY INITIATIVE

Daan Broeder, CLARIN ERIC
ICTeSSH, June 2021





SSHOC

social sciences & humanities open cloud



Horizon 2020
European Union Funding
for Research & Innovation

Type of action & funding:
Research and Innovation action
(INFRAEOSC-04-2018)

Partners: 48

(23 beneficiaries + 25 LTPs)

SSH ESFRI Landmarks and Projects
& international SSH data infrastructures

Project budget:

€ 14,455,594.08

Duration: 40 months

(January 2019 – 30 April 2022)

Project website:

www.SSHOpenCloud.eu



Objectives:

- creating the social sciences and humanities (**SSH**) part of European Open Science Cloud (**EOSC**)
- maximising **re-use** through **Open Science** and **FAIR** principles (standards, common catalogue, access control, semantic techniques, training)
- interconnecting existing and new infrastructures (clustered cloud infrastructure)
- establishing appropriate **governance model** for SSH-EOSC



Diversity in describing phenomena

- A considerable part of research concerns describing and analyzing phenomena using descriptive schemas and concepts
- Typical for the SSH is a high variety of such schema and concepts, caused by
 - wide variety of data types, sub-community specifics, schools of thought,...
 - divergent purposes and available effort
- Suitable well crafted vocabularies are essential for
 - accurate descriptions and classifications, countering interpretative vagueness (reduce ambiguity)
 - efficient retrieval and search



Concordance Word List Thesaurus Find X Sketch-Diff Sketch-Eval Corpus Info

Save View options KWIC Sentence Sort Left Right Node References Shuffle Sample Filter Overlaps 1st hit in doc Frequency Node tags Node forms Doc IDs

Query colour 16,486 (147.0 per million)

Page 1 of 825 Go Next Last

JZL It would be tedious to list the types and colours of stone, ceramic etc. used at each site

JZL types of stone used for various shades of colour are predictable and limited in number.

JZL Birdcombe Avon Here, sandstone furnished a buff colour, pennant stone a blue, liar the white for

JZL most mosaics comprise three to six basic colours, a work of good quality will include many

JZL therefore, to note ten or twelve different colours of tesserae in one pavement. In some, such

JZL the Woodchester Orpheus mosaic. <-p> 3.2 The colour of Tesserae <-p> Sensitive use of shading

JZL 1976, 9). Elsewhere, intelligent use of colour is responsible for the blue shading which

JZL are notable. <-p> <-p> Whilst considering the colour of tesserae it is also pertinent to mention

JZL 0.5 cm. sq. and 1.5 cm. sq. <-p> <-p> Like colour, the size of the tesserae affects the perspective

JZL fairly dark tesserae (deep red is a favourite colour), so producing a stronger proximity effect

JZL panels (pl. 5b). At Leicester the rosettes - coloured (from the edges inwards) red, yellow and

JZL be cramped (although "loose"). There are colour contrasts however: the simple guilloche

JZL former. However, the more subtle use of colour in the latter also produces a less contrived

JZL angular appearance. An overall poverty of colour, and the use of slightly larger (but still

JZL mosaic A). Although including the same basic colours, as well as tesserae of a similar size,

JZL blending of many tones of five or six basic colours, is notable in both designs. It is a sensitivity

JZL shows a generally consistent interface of colour, one in every four tongues of the latter

JZL Oceanus panel (contrast the confusion of colour around the heads of the lion and stag)

JZL However, on balance, the use here of similar colours (red, yellow, grey, pale-blue, brown) and

JZL Street mosaic, the presence there of a richly coloured figured panel (enclosed by a chain-guilloche

Page 1 of 825 Go Next Last

social sciences & humanities open cloud

Lexical Computing

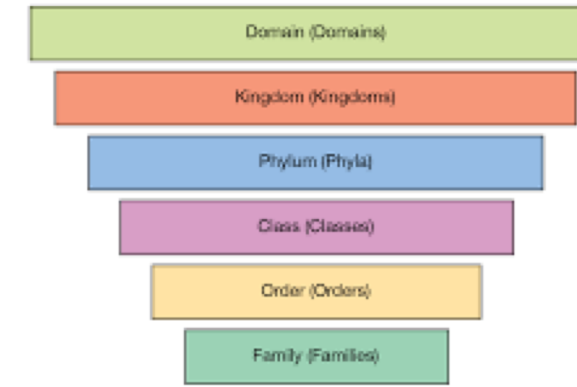


Controlled Vocabularies

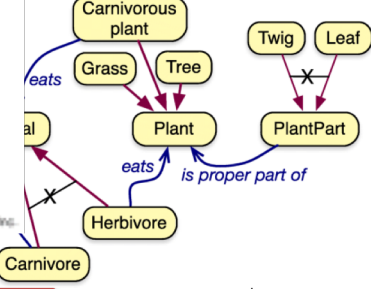
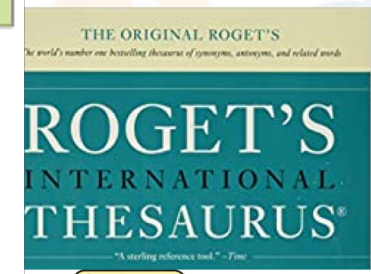
- Lists
- Glossaries
- *Thesauri*
- *Taxonomies*
- *Ontologies*

W3C std. RDF, SKOS, OWL:
widely accepted and supported formal way to represent vocabularies

How animals are classified

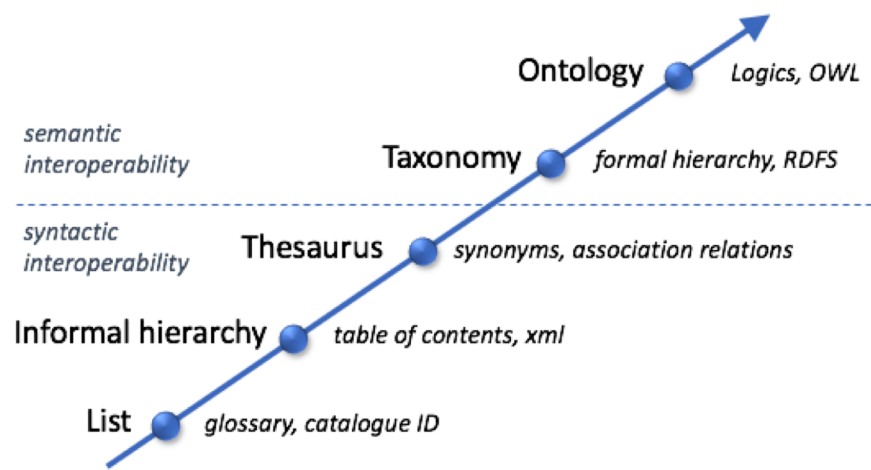


Symbol	Variable	Definition
Basic Series		
NR	S&P 500 Nominal Return	Quarterly return (in percentage) of S&P 500 index
ER	S&P 500 Excess Return	Quarterly return (in percentage) of S&P 500 index minus one quarter of 3-month T-bill rate (in percentage)
NCSPR	Change of credit spread	First order difference of the spread between Moody's BAA rate and 10 year Treasury Bond rate
ATSPPR	Change of term spread	First order difference of term spread between 30 year Treasury Bond rate-3 month Treasury Bill rate
ATBIL	Change of Treasury bill rate	First order difference of 3-month Treasury Bill rate
Change of Trading Activity and Liquidity		
ΔECM	Exchange Commission	First order log difference of quarterly Exchange Commission Revenue
ΔOCM	OTC Commission	First order log difference of quarterly OTC Commission Revenue
ΔMF	Mutual Fund Commission	First order log difference of quarterly Mutual Fund Sales Revenue
ΔNYSE	NYSE Volume	First order log difference of quarterly NYSE share volume
ΔMRG	Margin Trading	Quarterly Margin Interest divided by quarterly T-bill rate, then take the first order log difference
ΔECR	Effective Commission Rate	Quarterly Exchange Commission Revenue divided by quarterly NYSE share volume, then taken first order log difference.
Variation of Trading Activity and Liquidity		
VECM	Variation of Exchange Commission	Logarithm squared distance to mean of Exchange Commission Revenue
VOCM	Variation of OTC Commission	Logarithm squared distance to mean of OTC Commission Revenue
VMF	Variation of Mutual Fund Commission	Logarithm squared distance to mean of Mutual Fund Sales Revenue
NYSECV	Variation of NYSE Volume	Coefficient of variation of daily NYSE trading volume in quarter t
VECR	Variation of ECR	Logarithm squared distance to mean of Effective Exchange Commission Rate



An authority record

- O'Brien, Flann, 1911-1966
- Na Gopaleen, Myles, 1911-1966
- Knowall, George
- Na gCopaleen, Myles, 1911-1966
- His At Swim-Two-Birds ... 1939.
- His The best of Myles, 1983; CIP t.p. (Myles na Gopaleen (Flann O'Brien))
- His Myles away from Dublin, 1985; t.p. (Myles na Gopaleen (Flann O'Brien) selection written from the column written for ... under the name George Knowall)



Vocabularies in the SSHOC project

- **Coordination wrt. vocabularies: originally a limited effort**
 - Investigation of a common recommended platform for publishing and sharing vocabularies
 - Testing machine translation for vocabularies
 - Flexible integration of vocabularies in tools: e.g. SSHOC Dataverse
 - Identifying & creating proper vocabularies for SSH Marketplace and others
- **Identified more opportunities during the project**
 - Inventory and registration of relevant SSH vocabularies
 - Recommendations for further common approaches e.g. CV authoring tools
 - Opportunity & need to represent SSH interests with other stakeholders e.g. software & service providers

Vocabulary visibility and discovery

- Vocabularies not always FAIR yet; they should be properly registered and published, researchers & infrastructure providers should be able to find and reuse -> **FAIR semantic artefacts**
- SSH Vocabulary registry or a general one that supports sufficient discipline specificity e.g. Bartoc (3300 entries whereof 1200 SSH)
- Vocabulary search facility, that searches in vocabulary metadata **but also** the vocabulary terms themselves.
- *Note that providing optimal recommendations for researchers can be complicated e.g. also aspects of context and user profile play a role*

Vocabularies & Interoperability

- **Technical / Format interoperability.** SKOS and OWL are broadly accepted
 - but many projects use spreadsheets and tables and are locked in silos using highly specific software to manage and use these
 - Specific recommendations for vocabulary versioning are needed
- **Semantic interoperability.** Coming from different traditions different organizations and projects have developed different vocabularies to describe similar data. Normalization or conversion needed; the vocabularies involved can be huge and expertise expensive (see the VMT presentation).
- **Cultural & Human interoperability** aspects. Multi-lingual vocabularies, localization aspects.

SSH Vocabulary Initiative

We Identified the following topics for working towards common recommendations and implementations

- Vocabulary versioning (also in relation with recommendation for vocabulary authoring tools)
- SSH Vocabulary Registry; board a general registry or run our own
- Vocabulary recommendations; discovery via minimal metadata, term/token search, context via the SSH open marketplace?
- Recommended management platforms: SKOSMOS, VocBench, ..., and the APIs

Presentations

- **USAGE:**

- Use of Vocabularies in the SSH community - Iulianna van der Lek

- **FINDABILITY**

- Current ways of locating suitable vocabularies - Matej Ďurčo
- What can the SSH Open Marketplace add to that?

- **INTEROPERABILITY**

- Vocabulary Matching Tool - Holly Wright
- Use of MT for creating multilingual vocabularies - Monica Monachini



Using vocabularies in the SSH community

Iulianna van der Lek
CLARIN ERIC



SSH Vocabulary survey

- Timeline: Jan – June 2020
- Authors: CNRS/ HUMA-NUM
- Goal:
 - To learn about vocabulary usage and practices in SSH
 - To improve discoverability in SSHOC Open Marketplace



Clara Petitfils, Suzanne Dumouchel, Nicolas Larrousse, Laure Barbot, Klaus Illmayer, Matej Ďurčo and Tomasz Parkola. (2021). *SSHOC D7.6 Resources for Marketplace content description*. Zenodo.
doi:10.5281/zenodo.4558339

<https://zenodo.org/record/4558339#.YNi4nugzZPY>

SSH Vocabulary survey

52 out of 72 complete responses - >used vocabularies

- **SSH disciplines:** Linguistics, Archaeology & Prehistory; Sociology, Communication Sciences
- **Organizations:** Universities and research organisations (62,%), researchers, research libraries & archives, policy-making organisations
- **Countries:** France (55%), Spain, the Netherlands and Germany
- **Languages:** English, French, German and Spanish

SSH Vocabularies Survey

Vocabularies are used for:

- 58,49% “Concepts: Disciplines”
- 54,72% “Concepts: General concepts”
- 52,83% “Named entities : geographical entities”
- 45,28% “Named entities : persons”
- 45,28% “Concepts: Scholarly activities”
- 28,3% “Named entities: Institutions”
- 30,19% Other



SSH Vocabulary Survey

Type of used vocabularies:

Linguistics
Philosophy
Archaeology &
Prehistory
History & Sociology
Art & Art History

- Controlled vocabulary ([DDI](#), [CESSDA](#))
- Metadata schemas ([Dublin Core](#))
- Generic SSH thesaurus ([ELSST](#))
- Specialised thesaurus in SSH fields ([Pactols](#))
- Domain-specific dictionaries and glossaries
- Ontologies representing and structuring data ([SKOS](#))
- Gazetteers, i.e. authority files of particulars, persons, places ([Geonames](#), [Getty TGN](#))
- Research identifiers ([ORCID](#))

Vocabulary inventory

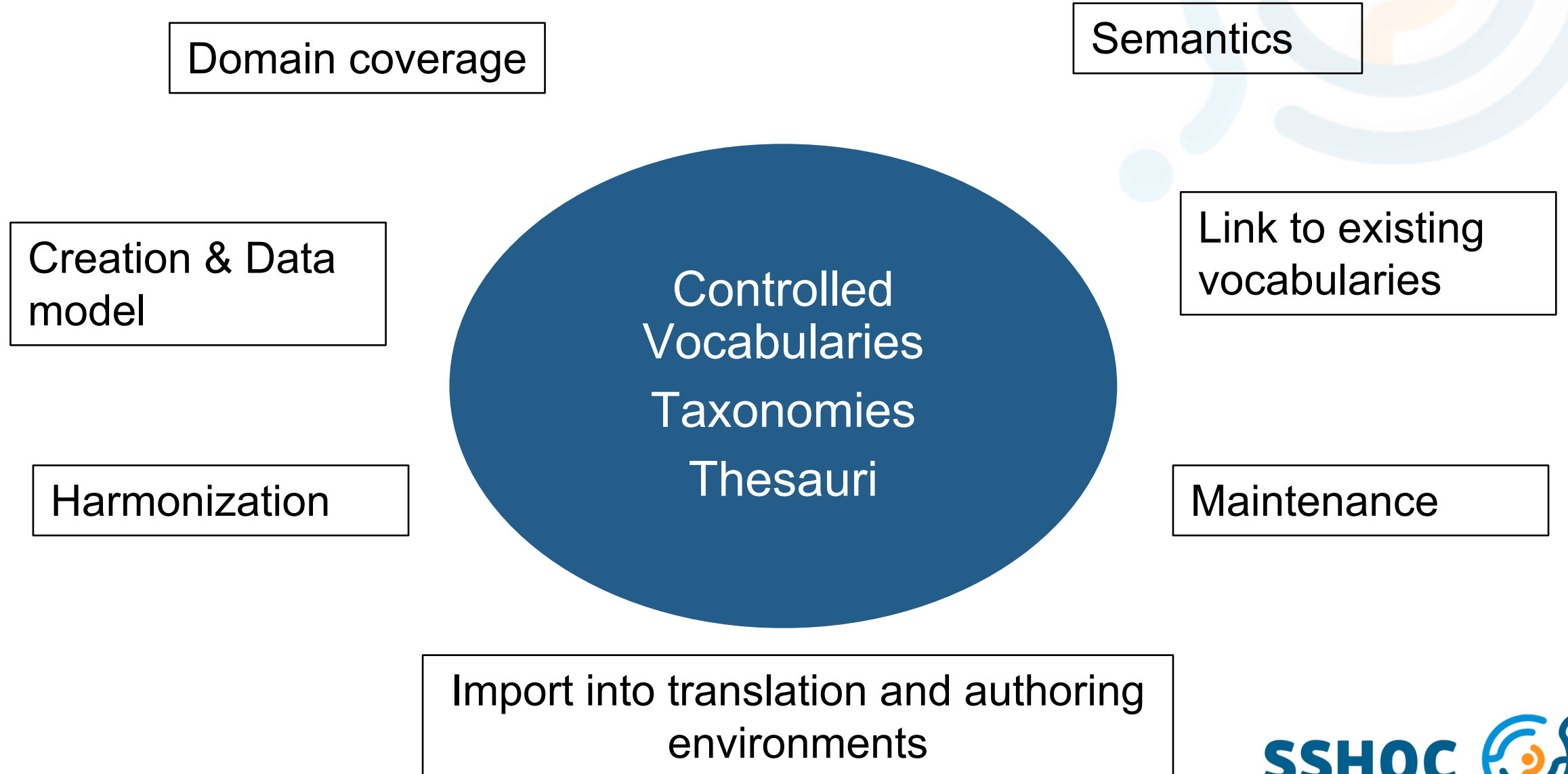
Inventory of vocabularies and thesauri for SSHOC

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	
1	Title	Homepage	Description	Creator	Publisher	URI	Domain	Subjects	Type of data out	Type (registry, repository, service)	Link to repository	Size	Format	User group	How it is used	Languages	License
2	AGROVOC - Food and Agriculture Organization of UN	http://www.fao.org/agrovoc/research	The largest Linked Open Data set about agriculture available for public use and facilitates	FAO	FAO	http://idm.fao.org/agrovoc	Food and Agriculture		Thesaurus	Repository	http://www.fao.org/agrovoc	38442 concepts	RDF, SKOS, XML	Researchers, librarians, information managers	AGROVOC is used by researchers, librarians and	Multilingual	CC-BY
4	Beckbone Thesaurus DARIAH	https://vocab.ac.fy.onew.ac.at/backbone_thesaurus/	It is a thesaurus for the humanities developed and maintained by the I Thesaurus Maintenance	EBT maintenance WG	Austrian Centre for Digital Humanities	http://idm.fao.org/agrovoc	Arts & Humanities	Computer science Information and general	Thesaurus	Repository		30 concepts, 10 collectors	RDF, SKOS, XML	Researchers in Humanities	It can be used to navigate the URI hierarchies and	French German	https://creativecommons.org/licenses/by/4.0/
5	BARTOC Register of Thesauri, Ontologies and Classifications	http://bartoc.skosmos.nl/bas/cv/	The Basel Register of Thesauri, Ontologies & Classifications (BARTOC) is a database of	Andreas Ledl as the Basel Register of Thesauri, Ontologies & Classifications	Common Library Network (CLN) in Göttingen		Computer Science, Social Sciences, Science		KOS	KOS Registry		5000 terminologies and 100 terminology registries		Researchers in Social Sciences and	Skosmos is a web based user source ontology	Any	BARTO available
6	BARTOC Software inventory for controlled vocabularies	http://bartoc.org/software/	A web page to collect basic info about software for controlled vocabularies		Common Library Network (CLN) in Göttingen				Controlled vocabulary	Open registry	https://github.com/globallibrarysoftwarefor			Researchers in Social Sciences and			
7	BIARTOC Vocabularies	https://bartoc.org/vocabularies	The Basel Register of Thesauri, Ontologies & Classifications (BIARTOC) is a database of	Different creators	Different publishers	Each entry has its own URI	SSH	D.J.C., EuroVoc, ILC	Vocabulary	Hierarchy		3289 vocabularies	KOS type	Researchers in Social Sciences and		Multilingual	Various
8	BBTalk	https://www.beckbone-thesaurus.eu/BBTalk/login	It is a submission and correction management tool that facilitates communication between the	EBT maintenance WG	DARIAH		Arts & Humanities		Thesaurus management platform	Online service		n.a.		Researchers in Humanities	Users can submit requests for possible changes	French German	
9	Bibliographic Ontology (BIBO)	http://biblographic-ontology.org/	The Bibliographic Ontology Specification provides main concepts and properties for	Bruce D'Arcus, Frédéric Glasson	Structured Dynamics LLC	http://purl.org/ontology/bibo/	Semantic Web		Ontology	Repository	https://github.com/structuredynamics/bibogloss		RDF, XML	Researchers in Social Sciences and	Can be used as a citation ontology, as a document	EN	The Bib. Copyrig. CC0 1.1
10	Bibliosima, the Observatory for Medieval and Renaissance Written Cultural Heritage	https://bibliosima.it/	A virtual library of libraries to discover the history of various texts and books that were written.		Bibliosima		Cultural Heritage		Observatory	Repository				Researchers in Social Sciences and			
11	BIO A vocabulary for legacy digital information	https://vocabulary.cosofa.org/discover	This document describes a vocabulary for describing biographical information about	Ian Davis, David Galbraith	CESSDA/ERIC	http://purl.org/vocab/bio/1/	Arts & Humanities	Biographical information (items are grouped by theme)	Ontology	Repository			OWL	Researchers in Social Sciences and		EN	No info
12	CESSDA Vocabulary service	https://vocabulary.cosofa.org/discover	CESSDA Vocabulary Service enables users to discover, browse, and download controlled		CESSDA/ERIC				Vocabulary	Service		28 vocabularies		Researchers in Social Sciences and			
13	COAR Controlled Vocabularies	https://www.coar-repositories.org/terms-and-conditions-for-coar-controlled-vocabularies	A collection of controlled vocabularies developed and maintained by the Federal Board	COAR	COAR	http://purl.org/coar/	Repositories	Resource type vocabulary, Access Rights vocabulary	Controlled vocabulary	Repository	https://www.coar-repositories.org/terms-and-conditions-for-coar-controlled-vocabularies		SKOS	Researchers in Social Sciences and	To enhance the interoperability across	Multilingual	
14	CRM Instruments	https://www.crm-centre.org/instruments	A terminology database of musical instruments		CRM Centre de Recherche en Chronologie	https://archives.crm-centre.org/instruments/	Cultural Heritage	Instruments	Thesaurus	Repository			Dublin core	Researchers in Social Sciences and		EN	
15	Cybergeo - EU Journal of Geography	http://journals.openedition.org/cybergeo/33	A catalogue of 552 journals that allow search by keywords.		Open Edition Journal - a Journals platform for the humanities and				Catalogue	Catalogue		552 journals		Researchers in Social Sciences and		Multilingual	
16	DRPedia	https://wiki.drpedia.org/	A project aiming to extract structured content from the information created in the Wikipedia						Open Knowledge Graph	Repository	https://github.com/drpedi/			Researchers in Social Sciences and	Alexis users to semantically query relationships and		GNU FDL
17	cdmrc	https://standaard.onroerend-erfgoed.nl/docs/terminologie	Metadata standards based on Dublin Core.		Dutch government				Controlled vocabulary	Repository			XML, OWMS 4.0	Researchers in Social Sciences and		NL	
18	DDalliance Controlled Vocabularies	https://dalliance.org/controlled-vocabularies	A set of controlled vocabularies that can be used with DD and other applications.		DDalliance		Social & Behavioral Sciences		Controlled vocabulary	Repository				The vocabularies are published in an XML	Researchers in Social Sciences and		
19	Devey Decimal Classification (DDC)	https://www.oclc.org/en/learn/decimalclassification/	The Dewey Decimal Classification—conceived by Melvil Dewey in 1873 and first	Melvil Dewey in 1873	OCLC				Classification system	Service				Libraries	The DDC is the most widely used classification	EN	Since 1913
20	Dictionnaire de sociologie clinique	https://www.cerim.fr/fr/dictionnaire-de-sociologie-clinique-17974525744.htm	The Dictionary of Clinical Sociology describes the central methods and questions of research	Agnès Yadeviciou-Rougale, Pascal Fugier, Vincent de Gaudemar	CARIN INCO		Social Sciences	Clinical sociology	Dictionary	?		245 terms	Deck	Researchers in Social Sciences and		FR	
21	Dublin Core categories DCM Metadata Terms	https://www.dublincore.org/specifications/dublin-core/dcm-terms/	An up-to-date specification of all metadata terms maintained by the Dublin Core Metadata			http://idm.fao.org/agrovoc			Controlled vocabulary	Repository		16 terms of the Dublin Core™ Metadata Element	RDF (creators of non-RDF metadata can	Researchers in Social Sciences and	These terms are intended to be used in combination	EN	http://creativecommons.org/licenses/by/4.0/
22	DYAS Humanities Thesaurus	https://humanities-thesaurus.academyofathens.gr/	The Thesaurus is the intellectual property of the Academy of Athens (AA), the National and	Academy of Athens (AA), the National and Kapodistrian University	BST Backbone Thesaurus	https://humanities-thesaurus.academyofathens.gr/	Arts & Humanities	Anthropology and ethnology	Thesaurus	Repository		1539 concepts	RDF, XML, X-Turtle	Humanities researchers		English Greek	https://creativecommons.org/licenses/by/4.0/
23	EIONET Data Dictionary - Vocabularies	http://eionet.emppan.eu/vocabularies	A vocabulary repository with search functionalities.		European Environment Agency		Environment	Air quality directives, Birds directives, Habitats directives.	Vocabulary	Repository				Researchers in Social Sciences and			
24	Encoded Archival Description - Iag library	https://www.loc.gov/read-lau-3a/gh/	Encoded Archival Description (EAD) is the international metadata transaction standard for		Developed by Society of American Archivists, published by Library of Congress	https://www.loc.gov/read-lau-3a/gh/			Schemas	Repository	https://github.com/SAA/SII-EAD3/tree/1.1			Researchers in Social Sciences and			Creative Commons
25	ET Thesaurus	https://www.et-thesaurus.com/et-thesaurus/	This series of ET additional terms			https://www.et-thesaurus.com/et-thesaurus/	Education		Thesaurus	Repository		7710 terms	RDF, XML, other	Education		Multilingual	No info

91 entries

A few potential orphan vocabularies

Vocabularies in SSH: Challenges



Locating suitable vocabularies

How the SSH Open Marketplace can help?

Matej Ďurčo, ACDH-CH

ICTeSSH Conference
SSHOC Vocabulary Initiative – What Users Want
28 June 2021



Vocabularies use & reuse

Vocabularies serve:

- Formalize/conceptualize a specific dimension/aspect of an application domain
- Make semantics explicit
 - Through verbose descriptions/definitions
 - Through relations between concepts
- Semantic interoperability across project and dataset boundaries
 - However only if they are being reused



Where to look for useful/ already used (SSH) vocabularies?

- [BARTOC](#) - Basic Register of Thesauri, Ontologies & Classifications
 - > 3300 Vocabularies
 - No (partial) concept-based search/index
- [LOV](#) - Linked Open Vocabularies
 - Broad meaning of “vocabularies”!
- [EU vocabularies](#)
 - Good authority
 - multilingual!
- [D3.1 Report on SSHOC \(meta\)data interoperability problems](#)
 - Most used vocabularies in SSH: CESSDA Topic Classification, CLARIN Concept Registry, ISO 639-1 language list or TADIRaH



Vocabularies

Search

Filter

Search

search

Full-text search across vocabulary description, ranked by relevance

1...500

501...1000

1001...1500

1501...2000

2001...2500

2501...3000

3001...3306



Name

Links

[BARTOC data formats \(bartoc-formats\)](#) Data formats for Knowledge Organization Systems as listed in BARTOC.org...

[API](#)

[BARTOC access modes \(bartoc-access\)](#) This vocabulary is used to specify the access modes of a knowledge organization system listed in BARTOC.org: free (the vocabulary is freely available), registered (access requires registration), or licensed (access requires a personalized license)....

[API](#)

[Dewey Decimal Classification \(DDC\)](#) "The Dewey Decimal Classification (DDC) system, devised by library pioneer Melvil Dewey in the 1870s and owned by OCLC since 1988, provides a dynamic structure for the organization of library collections...

[API](#)

[ISO Language Codes \(639-1 and 693-2\) and IETF Language Types \(ISO639\)](#) Comprehensive language code information, consisting of ISO 639-1, ISO 639-2 and IETF language types...

[API](#)

[KOS Types Vocabulary](#) "Related to this DC Application Profile is a KOS Type vocabulary..."

[API](#)

[BARTOC ...](#)

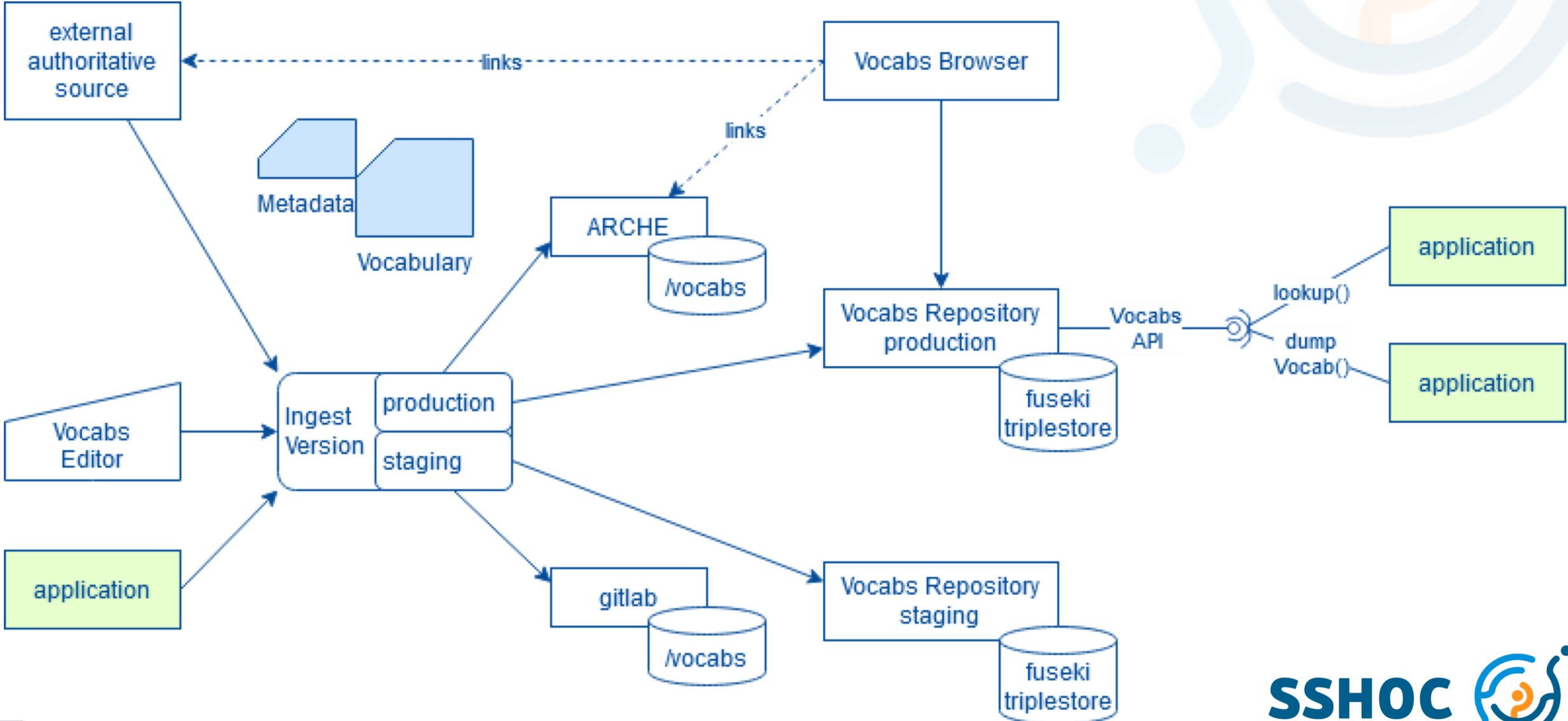
[API](#)

Linked Open Vocabularies

The screenshot shows the LOV website interface. At the top, there are navigation links: VOCABS, TERMS, AGENTS, and SPARQL/DUMP. Below this is a teal header with the text "Linked Open Vocabularies (LOV)". Underneath the header is a navigation bar with icons for "+ Suggest", "Documentation", "g+ Follow", a search box, and a "Feedback" icon. The main content area features a large bubble chart titled "605 Vocabularies in LOV". The chart displays various vocabulary acronyms as bubbles of different sizes and colors, with the largest bubbles being "dcterms", "foaf", "vann", "skos", "dce", and "cc". Below the chart is a "Category Tags" section with a grid of tags including: Methods, Metadata, Catalogs, Support, Geography, API, Society, Quality, RDF, Services, People, Industry, Vocabularies, Environment, IoT, General & Upper, Time, Multimedia, Events, Biology, Geometry, FRBR, Government, W3C Rec, SPAR, PLM, Academy, eBusiness, Tag, and Security.



Vocabulary management workflow (at ACDH-CH)



“Suitable”?

- Covering the “right” dimension
- Comprehensive
- Stable availability
- Stable reference to concepts
- Well-established
- Maintained (as opposed to orphaned)
- Never perfect match



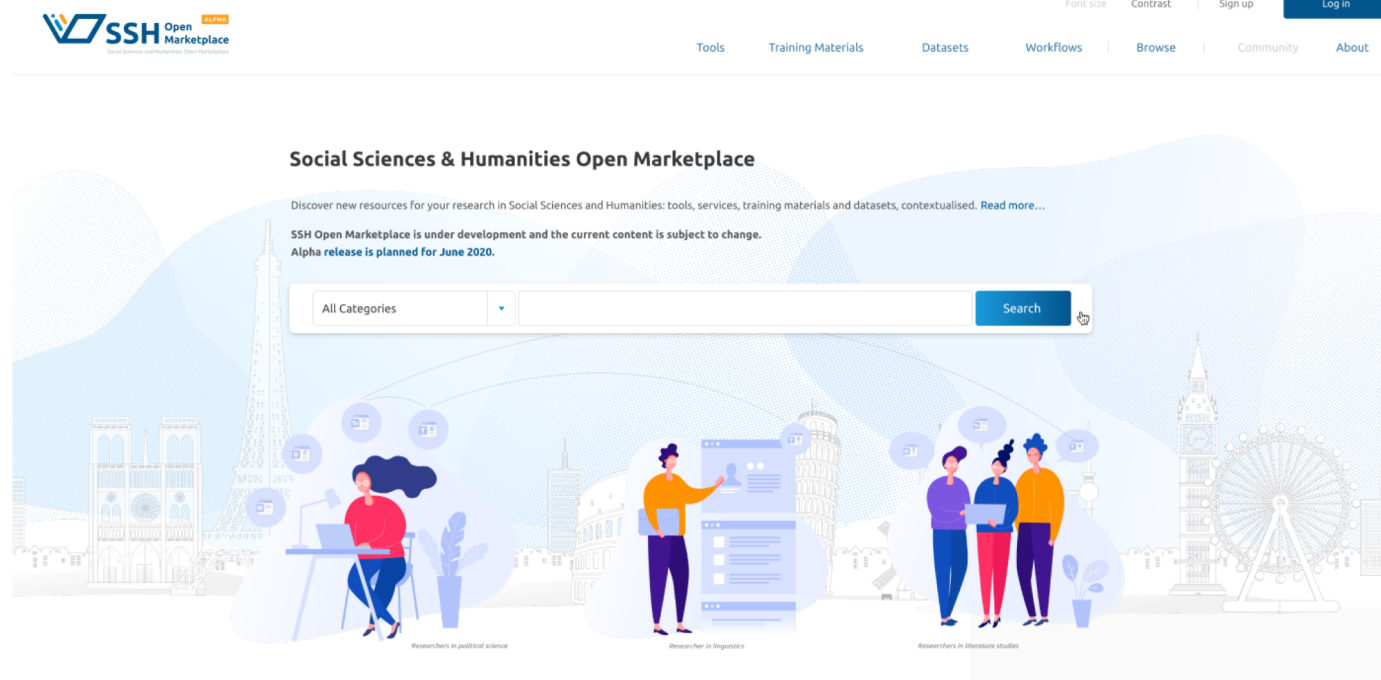
SSH Open Marketplace as one of the SSHOC services

- Discovery portal for SSH resources

- Tools & services
- Training materials
- Workflows
- Datasets
- Publications

- 3 guiding principles

- Contextualisation
- Curation
- Community



Beta version: marketplace.sshopencloud.eu

Vocabulary entries in the SSH Open Marketplace

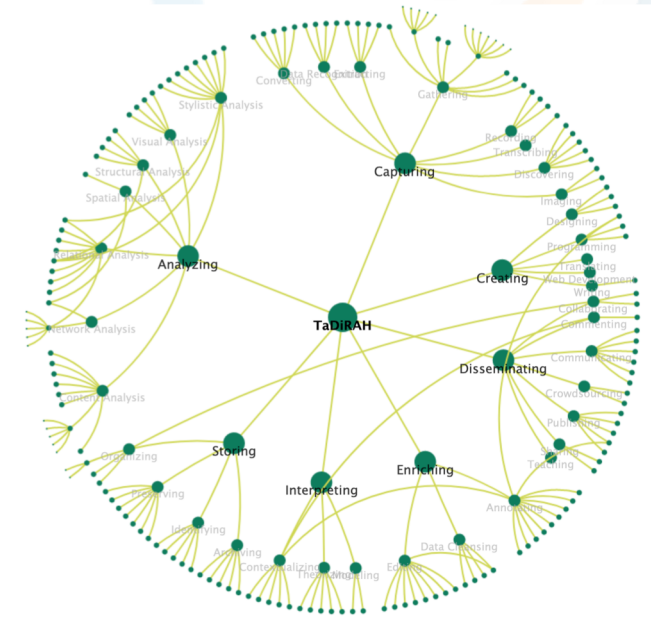
- Vocabularies as first-class citizens (represented as tools)

Contextualized (by related items):

- Training materials to support vocabulary use Controlled Vocabularies and SKOS
- Linked publications
- Examples of projects reusing the vocabularies
-

Referencing and contextualizing the most used SSH vocabularies - TaDIRAH example

TaDiRAH - Taxonomy of Digital Research Activities in the Humanities



[Github repo](#)

[Editors' Interview](#)

[Twitter](#)

[Linked publication](#)

Related services using TADIRaH

....

Vocabs About Editor SPARQL API Help Interface language: English -

TaDIRAH: Taxonomy of Digital Research Activities in the Humanities

Content language: English - x Search

Alphabetical Hierarchy

A B C D E F G H I K L M N O P
Q R S T U V W

Abstract Thinking
Academic Publishing
Adding
Aggregating
Analyzing
Annotating
Archiving
Associating
Audio Annotation
Audio Conferencing
Audio Recording
Authorship Attribution

Vocabulary information

TITLE	TaDIRAH
DESCRIPTION	Taxonomy of Digital Research Activities in the Humanities
CREATOR	Luise Borek, Canan Hastik, Vera Khramova, Jonathan Geiger
CONTRIBUTOR	Elisabeth Burr, Francesca Tomasi, Tiziana Mancinelli, Monica Berti, Klaus Thoden, Luise Borek, Christof Schöch, Claudia Müller-Birn, Melanie Siemund, Saskia Lindner, Tamara Butigan, Vanja Savic, Toma Tasovac, Aurélien Berra, Thibault Cléricie, Martin Grandjean, Vincent Razanajao, Gimena del Rio Riande
LANGUAGE	English
VERSION	2.0.0
CREATED	2020
DATE ISSUED	2020-09-28
LICENSE	https://creativecommons.org/publicdomain/zero/1.0



What it looks like in the SSH Open Marketplace



TaDiRAH - Taxonomy of Digital Research Activities in the Humanities

The taxonomy of digital research activities in the humanities has been developed for use by community-driven sites and projects that aim to structure information relevant to digital humanities and make it more easily discoverable. The taxonomy is expected to be particularly useful to endeavors aiming to collect information on digital humanities tools, methods, projects, or readings.



[Go to Tool or service](#) ↗

Details

ACCESS

License: [Creative Commons Zero v1.0 Universal](#)

CATEGORISATION

Keyword: vocabulary

CONTEXT

See also: <https://openmethods.dariah.eu/2021/02/10/openmethods-spotlights-2-interview-with-luise-borek-and-canan-hastik-about-tadirah/>, <https://tadirah.info/>, https://twitter.com/tadirah_dh, https://epub.uni-regensburg.de/44951/1/isi_borek_et_al.pdf

TECHNICAL

Version: 2

EDITOR

Canan Hastik

[Website](#)

Jonathan Geiger

jonathan.geiger@adwmainz.de

Luise Borek

Vera Khramova

vera.khramova@stud.h-da.de

GitHub: <https://github.com/dhtaxonomy/TaDiRAH>

DOI: [10.5281/zenodo.32492](https://doi.org/10.5281/zenodo.32492)



open cloud



Vocabulary mapping tool for archaeology in **ARIADNE***plus*

Holly Wright, Archaeology Data Service

Ceri Binding & Douglas Tudhope

[University of South Wales](#), Trefforest

[tgn:7029392](#) World

[tgn:1000003](#) Europe

[tgn:7008591](#) United Kingdom

[tgn:7002443](#) Wales

[tgn:7018963](#) Rhondda Cynon Taf

[tgn:7441565](#) Trefforest



- **ARIADNE**

- 24 partners, 13 countries, 9 languages, 27 subject vocabularies
- 1.9 million data records aggregated/integrated
- Subject vocabularies coordinated via mapping to Getty AAT – total 6416 mappings produced

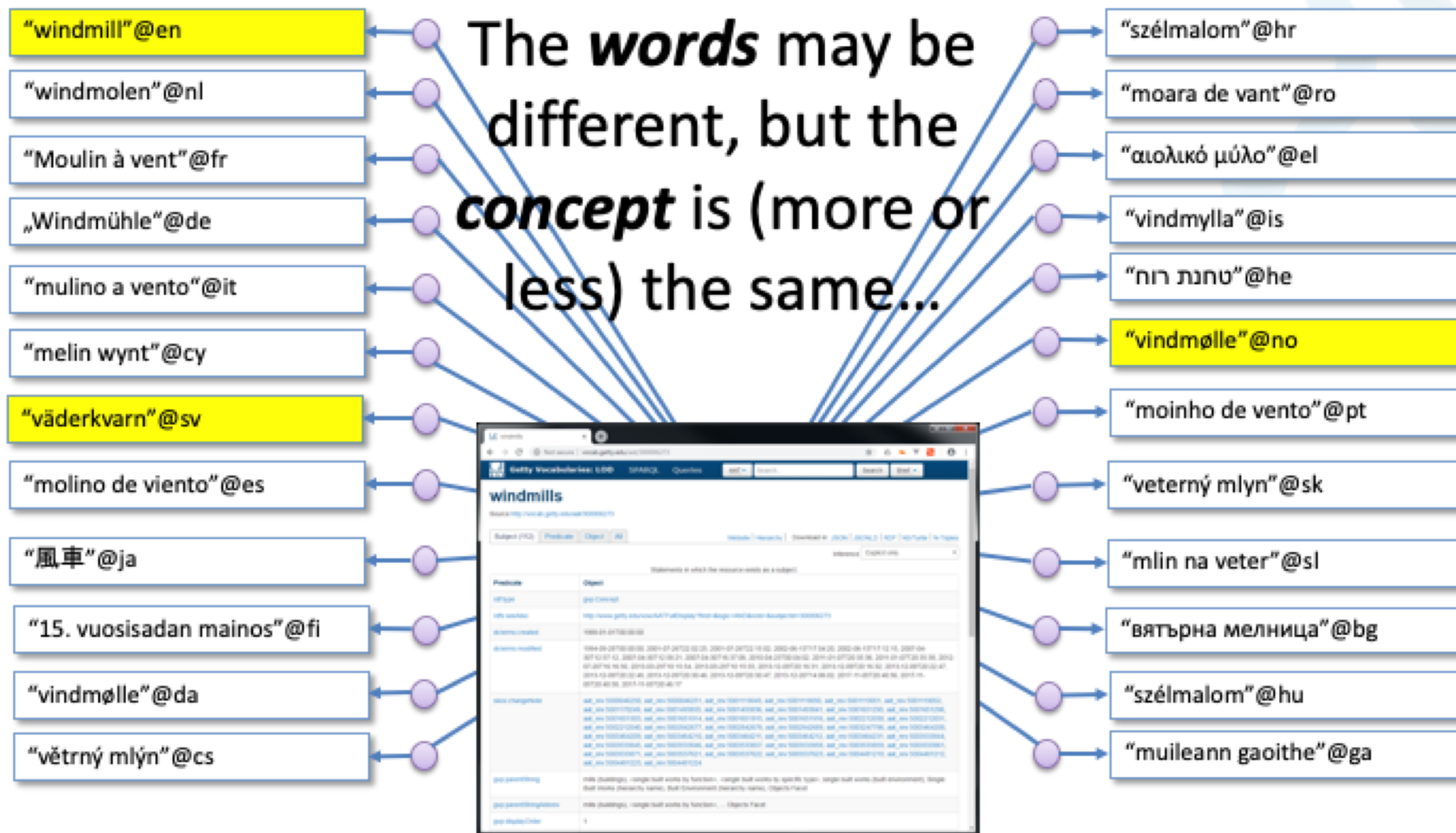
- **ARIADNE*plus***

- 41 partners, 29 countries, 22 languages, ?? subject vocabularies
- Data aggregation/integration work currently in progress
- Reusing, revising and supplementing previous mappings
- Adding vocabulary mappings from new data partners
- Adding Wikidata mappings (multilingual entry vocabulary)

Why do we need vocabulary matching in ARIADNE?

- Source datasets not necessarily produced with aggregation, consolidation, cross-search and reuse in mind
- I say “*potato*”, you say “*pomme de terre*”, she says “*maris piper*” – multiple barriers to cross-searching subject metadata: language, punctuation, spelling, homonyms, synonyms, level of specificity
- Text-based search is limited by any/all of these
- Need to establish mutually agreed meaning...

Map local terms to a central concept

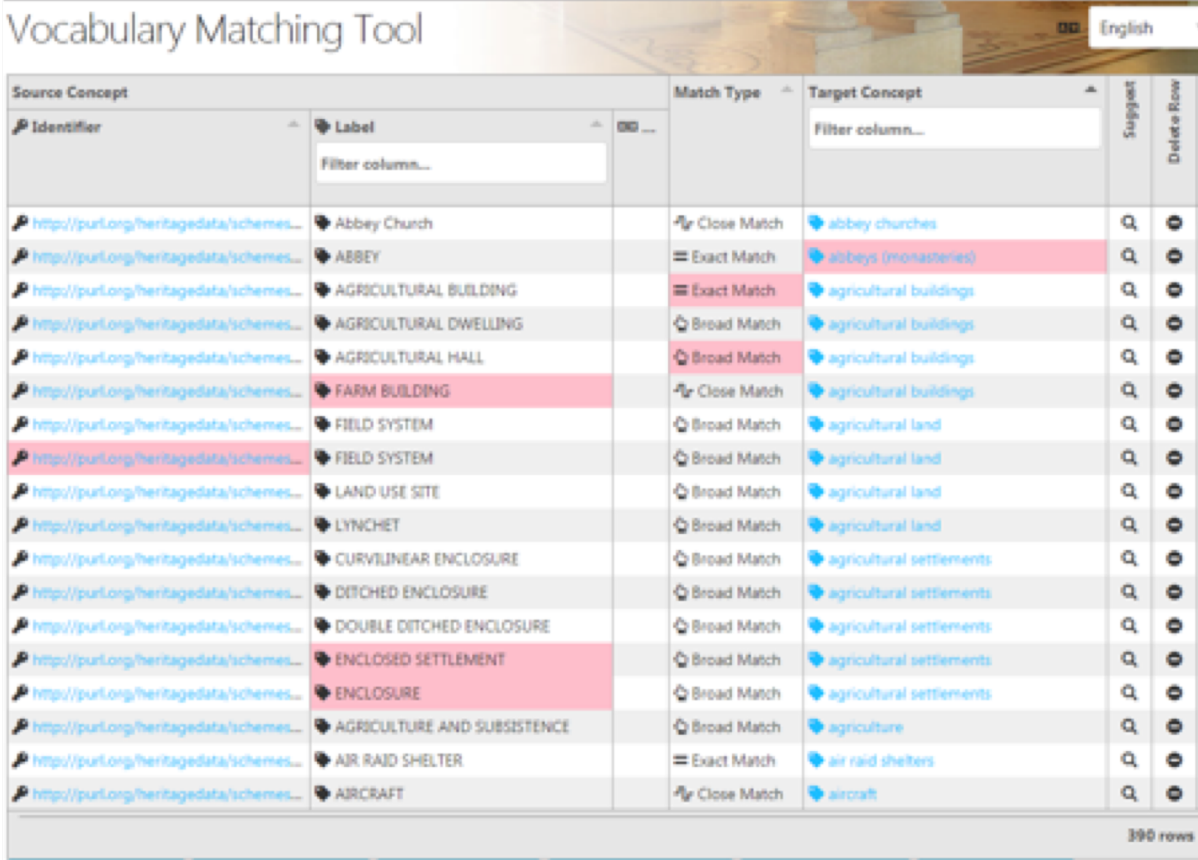


Now we can include any/all of these variants in a single query

Vocabulary Matching Tool

<https://vmt.ariadne.d4science.org/vmt/>

- For matching local subject terms / concepts to Getty AAT concepts
- Search & browse Getty AAT
- No auto match: examine scope and context of source / target concepts



Vocabulary Matching Tool

Source Concept	Match Type	Target Concept	Suggest	Delete Row
Identifier		Filter column...		
Label				
Filter column...				
http://purf.org/heritagedata/schemes... Abbey Church	Close Match	abbey churches	Q	⊘
http://purf.org/heritagedata/schemes... ABBEY	Exact Match	abbeys (monasteries)	Q	⊘
http://purf.org/heritagedata/schemes... AGRICULTURAL BUILDING	Exact Match	agricultural buildings	Q	⊘
http://purf.org/heritagedata/schemes... AGRICULTURAL DWELLING	Broad Match	agricultural buildings	Q	⊘
http://purf.org/heritagedata/schemes... AGRICULTURAL HALL	Broad Match	agricultural buildings	Q	⊘
http://purf.org/heritagedata/schemes... FARM BUILDING	Close Match	agricultural buildings	Q	⊘
http://purf.org/heritagedata/schemes... FIELD SYSTEM	Broad Match	agricultural land	Q	⊘
http://purf.org/heritagedata/schemes... FIELD SYSTEM	Broad Match	agricultural land	Q	⊘
http://purf.org/heritagedata/schemes... LAND USE SITE	Broad Match	agricultural land	Q	⊘
http://purf.org/heritagedata/schemes... LYNCHET	Broad Match	agricultural land	Q	⊘
http://purf.org/heritagedata/schemes... CURVILINEAR ENCLOSURE	Broad Match	agricultural settlements	Q	⊘
http://purf.org/heritagedata/schemes... DITCHED ENCLOSURE	Broad Match	agricultural settlements	Q	⊘
http://purf.org/heritagedata/schemes... DOUBLE DITCHED ENCLOSURE	Broad Match	agricultural settlements	Q	⊘
http://purf.org/heritagedata/schemes... ENCLOSED SETTLEMENT	Broad Match	agricultural settlements	Q	⊘
http://purf.org/heritagedata/schemes... ENCLOSURE	Broad Match	agricultural settlements	Q	⊘
http://purf.org/heritagedata/schemes... AGRICULTURE AND SUBSISTENCE	Broad Match	agriculture	Q	⊘
http://purf.org/heritagedata/schemes... AIR RAID SHELTER	Exact Match	air raid shelters	Q	⊘
http://purf.org/heritagedata/schemes... AIRCRAFT	Close Match	aircraft	Q	⊘

390 rows

IMPORT JSON EXPORT JSON EXPORT CSV + ADD NEW ROW CLEAR ROWS SHOW HELP

ARIADNE plus

Created by University of South Wales

ARIADNEplus is a Horizon 2020 project funded by the European Commission (Grant Agreement No 823184)

This application retrieves some information originating from Getty Art & Architecture Thesaurus (AAT)® which is made available under the ODC Attribution License. See <https://vocab.getty.edu/> for further details.

Selected references and links

References

- Binding, C, Tudhope, D & Vlachidis, A 2018, 'A study of semantic integration across archaeological data and reports in different languages' Journal of Information Science, vol 45, no. 3, pp. 364-386. [doi:10.1177/0165551518789874](https://doi.org/10.1177/0165551518789874)
- Binding, C & Tudhope, D 2016, 'Improving interoperability using vocabulary linked data' International Journal on Digital Libraries, vol 17, no. 1, pp. 5-21. [doi:10.1007/s00799-015-0166-y](https://doi.org/10.1007/s00799-015-0166-y)

Links

- ARIADNEplus project: <http://www.ariadne-infrastructure.eu/>
- ARIADNE portal: <https://ariadne-infrastructure.eu/portal/>
- Vocabulary Matching Tool (VMT): <https://vmt.ariadne.d4science.org/vmt/>
- USW Hypermedia Research Group: <https://hypermedia.research.southwales.ac.uk/>

Contact

- ceri.binding@southwales.ac.uk ORCID: 0000-0002-6376-9613
- douglas.tudhope@southwales.ac.uk

Use of MT for creating multilingual vocabularies

Monica Monachini, CNR-ILC & CLARIN IT

ICTeSSH, June 2021



Topic: Why Do Multilingual Vocabularies matter

- Are essential for a proper description of resources and phenomena
- Help people find resources and determine their value
- Are critical in digital environments, where humans rely on computer processing for reliable and timely results
- Some use-cases in SSH:

- users searching data in all languages



- multi-country social surveys



Challenges: Machine translations for multilingual vocabularies

- perform automatic translation of
 - **occupation ontologies** <https://www.surveycodings.org/occupation-measurement>
 - **metadata and definitions (CLARIN CCR)** <https://concepts.clarin.eu/ccr/browser/>
- using and testing different systems
 - “MT tool-suite” at UFAL <https://lindat.mff.cuni.cz/services/translation/>
 - “Google Translate” <https://translate.google.com/>
 - “DeepL” <https://www.deepl.com/en/translator>
 - “Reverso” www.reverso.net
- Languages: English French, Italian, German, Dutch, Greek, Czech, Slovene, Russian (evaluation running)

Challenges: Machine translations for multilingual vocabularies



	UFAL		DEEP-L (https://www.deepl.com/en/translator)		Google translate (https://translate.google.com)		REVERSO		
Initiation	source term translation	target definition translation (UFAL)	source term translation	target definition translation	source term translation	target definition translation	source term translation	target definition translation	
research experiment in (sou)	magicien de	N Expérience de recherche dans laquelle	magicien d'oz	N Une expérience de r	magicien de l'oz	N Une expérience de	magicien de l'oz	N	L
visual representation (sou)	systèmes d'é	Y La représentation visuelle de la langu	systèmes d'écriture	Y La représentation vis	systèmes d'écriture	Y La représentation vi	systèmes d'écriture	Y	L
specification of a persistent (sou)	identifiant pers	M Spécification d'un identificateur pers	identifiant persistant	M Spécification d'un id	identifiant persistant	M Spécification d'un ic	identifiant persistant	M	S
access to the communication (sou)	public	Y L'accès à l'événement de communication	public	Y L'accès à l'événement	publique	Y L'accès à l'événement	public	Y	L
address of an organization (sou)	adresse	Y L'adresse d'une organisation qui a pa	adresse	Y L'adresse d'une orga	adresse	Y L'adresse d'une org	adresse	Y	T
number of years that (sou)	âge	Y Le nombre d'années de vie de quelqu	âge	Y Le nombre d'années	âge	Y Le nombre d'année	âge	Y	L
investigator asks speaker (sou)	élimé	N L'enquêteur demande au ou aux locu	a suscité	M L'enquêteur deman	suscité	M L'enquêteur deman	obtenue	M	L
email address of a person (sou)	courriel	M L'adresse électronique d'une person	email	Y L'adresse électroniq	email	Y L'adresse e-mail d'u	courriel	M	L
indicates the structure of (sou)	structure évé	Y Indique la structure de l'événement	structure des évé	M Indique la structure	structure de l'événem	M Indique la structure	structure d'événem	M	li
identification of the location (sou)	lieu d'exécut	Y Identification de l'endroit où l'outil o	lieu d'exécution	Y Identification de l'en	lieu d'exécution	Y Identification de l'e	emplacement d'e	M	li
transmission of the content (sou)	cadre expéri	Y Une transmission du contenu se déro	milieu expériment	M Une transmission du	cadre expérimental	Y Une transmission d	cadre expériment	Y	T
transmission of the message (sou)	face à face	Y La transmission du message assure u	face à face	Y La transmission du n	face à face	Y La transmission du	face à face	Y	L
contrary to what is true (sou)	faux	Y Contrairement à ce qui est vrai, error	faux	Y Contraire à ce qui es	faux	Y Contrairement à ce	faux	Y	C
access to the communication (sou)	famille	Y L'accès à l'événement de communication	famille	Y L'accès à l'événement	famille	Y L'accès à l'événement	famille	Y	L
Fax number of a person (sou)	numéro de tél	M Le numéro de télécopieur d'une pers	numéro de fax	Y Le numéro de fax d'l	numéro de fax	Y Le numéro de fax d'Y	numéro de téléc	M	M
communication event with (sou)	monologue	Y Événement de communication avec u	monologue	Y Événement de comn	monologue	Y Événement de com	monologue	Y	É
study of the structure (sou)	morphologie	Y L'étude de la structure et de l'électo	morphologie	Y L'étude de la structu	morphologie	Y L'étude de la struct	morphologie	Y	L
speaker prepares in (sou)	prévu	N L'intervenant prépare en détail la str	prévu	N L'orateur prépare à l	prévue	N L'orateur prépare e	prévue	N	L
indicates in how far the (sou)	type de plan	Y Indique dans quelle mesure le consu	type de planification	Y Indique dans quelle	type de planification	Y Indique dans quelle	type de planificat	Y	li
content is composed (sou)	noé	Y Le contenu est composé en vers ou e	noé	Y Le contenu est comp	noé	Y Le contenu est com	noé	Y	li

Benefits for SSH Researchers

- **Multilingual occupation ontologies**
 - During the interviews, the respondent can self-select a job title from a list of occupations, and find the appropriate version (also the male vs. female form)
 - a multilingual occupational database where all titles are coded according to the International Standard Occupational Classification (ISCO)
- **Multilingual controlled vocabulary**
 - The user can perform a query in native language and retrieve data in all languages

Final considerations

- At CNR, we are creating and experimenting with MT for creating multi-lingual vocabularies
- Platform where to host vocabularies: in general, the hosting of vocabularies (that do not have a safe home yet) is needed
- There is not yet a default solution for hosting and publishing such vocabularies in general, but we are looking into it.
- Vocabularies are also important data that should be FAIR



Monica Monachini
CLARIN



Holly Wright
Archaeology Data Service



Daan Broeder
CLARIN

#SSHOCaVocabulary
#SSHOCifyCLARIN



Taina Jääskeläinen
Finnish Social Science
Data Archive



Matej Durco
ACDH-CH



**Iuliana van
der Lek**
CLARIN



CLARIN



We are interested in
learning more about
your experience
with vocabularies.
Help us out by
filling in the
following polls!

ICTeSSH2021 SSH Vocabulary Initiative - What users want

17 - 29 Jun 2021

Poll results

Table of contents

- What is your role in your organisation?
- Are you familiar with vocabularies?
- How and where did you find the vocabulary/ies you decided to use?
- How do you determine the quality of a vocabulary?
- How do you handle vocabulary that fails to cover your descriptive needs?
- Mapping (conversion) between entries of different vocabularies, how do you manage that?
- Do you use vocabularies in your own language?
- What kind of support do you expect from research infrastructures/ support services for vocabularies?

What is your role in your organisation?

0 2 2



Are you familiar with vocabularies?

0 2 4

Not familiar at all



Beginner



Familiar



Very familiar



Expert



I manage vocabularies



How and where did you find the vocabulary/ies you decided to use?

0 2 5

Provided as part of the tool or workflow



Recommendations from peers



Recommendations from research infrastructures



Dedicated vocabulary registry



Other



How do you determine the quality of a vocabulary?

(1/2)

0 2 3

Intensive use



Provider



Recommendations by peers



Proscription



Standardisation



How do you determine the quality of a vocabulary?

(2/2)

0 2 3

Tradition

 4 %

Other

 4 %

How do you handle vocabulary that fails to cover your descriptive needs?

0 2 2

I add additional descriptions in a description field.



I propose new terms to the vocabulary manager.



I choose the closest alternative.



I don't use that vocabulary.



Other



Mapping (conversion) between entries of different vocabularies, how do you manage that?

0 1 3

- not applicable
- In a shared spreadsheet.
- Sincerely, I don't convert. Reason is, I am not exposed yet to such tools in the part of the world I'm in.
- Cross walks
- I don't.
- I don't
- Ontology mapping tools (e.g., LOOM, AML, YAM++) and then the use of ontology repositories to share mappings.
- Start comparison on excel spreadsheet
- I don't
- With a lot of frustration ;)
- No good solution yet found, need to try it out with different tools and select the best for my needs
- Automated or manual?
- I don't.

Do you use vocabularies in your own language?

0 1 7

Yes



No



No, I use English versions



Would like to, if they would be available



What kind of support do you expect from research infrastructures/ support services for vocabularies?

- Guidance on which vocabulary to use for a particular purpose. Quality, authority, sustainability. Backed by the community. Fit for purpose. Domain relevance. Machine access.
- Ease of access, guidance on which vocabulary to use and support in terms of mapping
- Hosting a mapping tool
- to put a quality stamp to a certain vocabulary
- Support for as many vocabularies as feasible, crosswalks.
- Registries and vocabs
- services like mapping
- I expect them to hide all interoperability complexity
- advice how to create vocabs
- Ease of access.
- Suitable repos
- Technical, advisory support



SSHOC

social sciences & humanities open cloud

SSH Vocabulary Initiative What users want

June 28, 9:00 - 10:30 CEST

Join us



sshopencloud.eu/register

sshvocabularyinitiative@sshopencloud.eu

Thank you for your attention!

Join our community



sshopencloud.eu



[@SSHOpenCloud](https://twitter.com/SSHOpenCloud)



info@sshopencloud.eu



[/in/company/sshoc](https://www.linkedin.com/company/sshoc)

