# The Social Sciences and Humanties (SSH) FAIR Vocabulary Registry Demo

Liliana Melgar (<a href="https://orcid.org/0000-0003-2003-4200">https://orcid.org/0000-0003-2003-4200</a>)

Menzo Windhouwer (<a href="https://orcid.org/0000-0002-2204-4018">https://orcid.org/0000-0002-2204-4018</a>)

Kerim Meijer (<a href="https://orcid.org/0009-0000-3419-4307">https://orcid.org/0009-0000-3419-4307</a>)

Royal Netherlands Academy of Arts and Sciences (KNAW) - Humanities Cluster (HuC)

NKOS WORKSHOP at TPDL Tampere, September 23, 2025

#### Liliana Melgar

- Data specialist at the Digital Infrastructure department of the Royal Netherlands Academy of Arts and Sciences (KNAW) - Humanities Cluster (HuC)
- Focus on metadata standards, data modelling, data formats, data interoperability, controlled vocabularies
- Curator of the SSH FAIR vocabulary registry

#### Also on behalf of:

#### Menzo Windhouwer

 Lead Engineer Team Structured Data at the Digital Infrastructure department of the KNAW Humanities Cluster (HuC); engineer at research infrastructures: CLARIN, CLARIAH



#### Kerim Meijer

 Senior software engineer at the Team Structured Data at the Digital Infrastructure department of the KNAW Humanities Cluster (HuC)



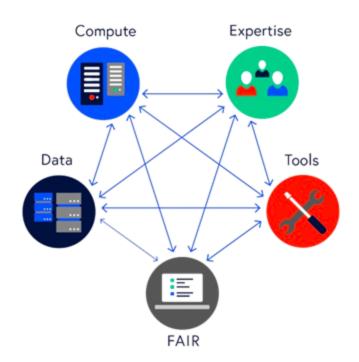
#### In the context of:

#### SSHOC-NL

Social Science and Humanities Open Cloud for the Netherlands (SSHOC-NL).

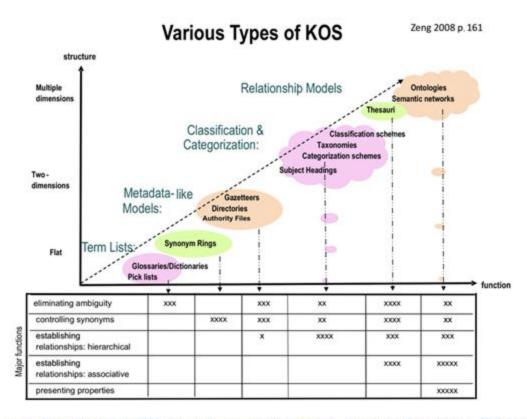
The project aims to enable groundbreaking interdisciplinary research on pressing societal questions.

SSHOC-NL is financed by the Dutch Research Council (NWO) Large-scale Research Infrastructure Grant.





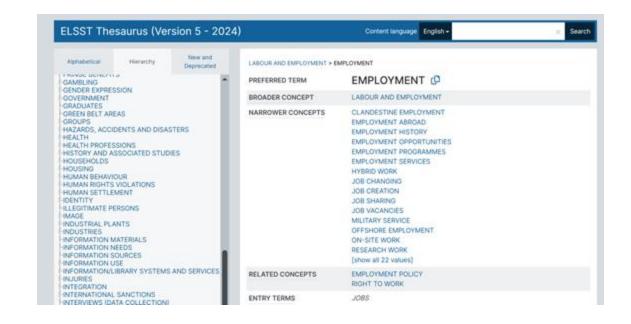
## Knowledge Organization Systems (KOS)



# Example: European Language Social Science Thesaurus (ELSST) by CESSDA

3,400 concepts covering the core social science disciplines

Translated into EU languages



# Banen en Ionen op basis van de Polisadministratie Central Dureuz voor Statolek, 2024, "Baren en lonen op basis van de Polisadministratie", Idma Usoc org 10,5793 406014410007:r7adf, OctoStit Portet, V1.

Learn about Data Citation Standards

Contact Owner Share

Dataset Metrice ©

0 Downloads ©

Description ()

In diff databedand sin kazantitatieve en kwalitatieve gopvenn opgenomen over banen en foran varwerknammer bij Michariandies bedrijven over een bepaald verstiggigaar of deel van meen verstiggaar. Die dafinities van "bean" die ten grondstag ligt aan dit bestand is een inkonsterverhooding (WVIII) in verband deut arbeid was ean werkgever heet een penacion. Een penacion kale per verliggever investigere seconsterverhooding integrijkenigt hebben 100 peza emprovent, waarbij de arbeid voor een baan de ekonomisterverhoodingen tegrijkenigt hebben 100 peza emprovent, waarbij de arbeid voor een baan de ekonomisterverhoodingen tegrijkenigt hebben 100 peza emprovent, waarbij de arbeid voor een baan de ekonomisterverhooding (WVIII) is in 2 teenhooding versitigser 2010. De patriangiseerd is polities ekonomisterverhooding (WVIII) is in 2 teenhooding versitigser 2010. De patriangiseerd is polities.

Read full Description [+]

Subject O

Social Sciences

Keyword 
Werkgelegenheid, Loon, Banen, Sociaaleconomische en numtelijke statistieken, Arbeid en tonen, Werkgelegenheid en tonen. Banen en toonsprimen

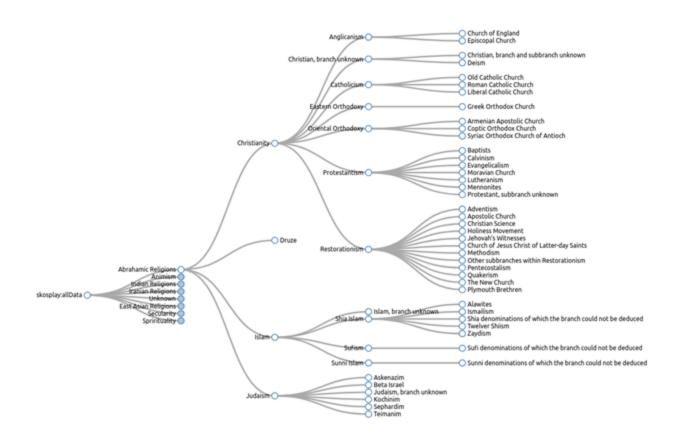
License/Data Use Agreement Custom Dataset Terms

City Dataset -

Vocabularies also used in data repositories (e.g., Dataverse)



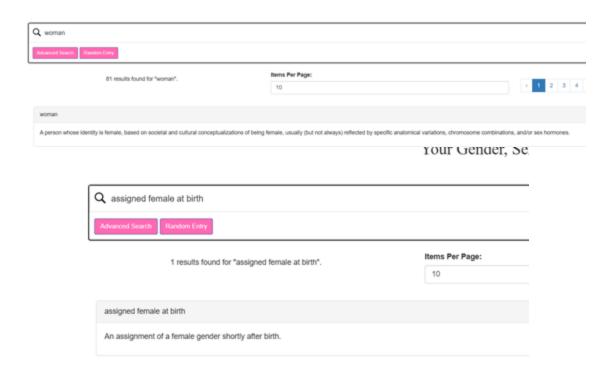
## Example (taxonomy of religions)



Taxonomy of religious denominations

(by Rick Mourits)

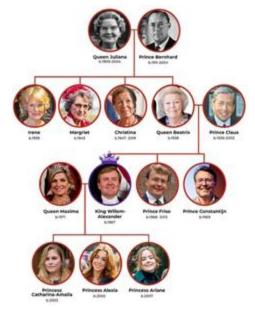
## Example: Gender, Sex, and Sexual Orientation (GSSO)





#### Example: PersonLink

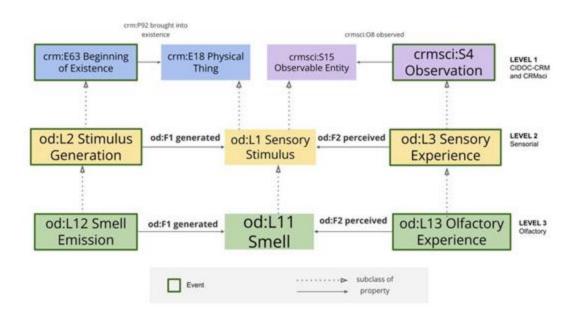
A Multilingual and Multicultural Ontology Representing Family Relationships





#### Example: Odeuropa

Example of an ontology produced in a research context



## Major challenge: lack of interoperability

Mourits et al., 2025 inventoried the vocabulary (schemas) used by bigger data centers in Asia, Europe, and North America.

#### They found that

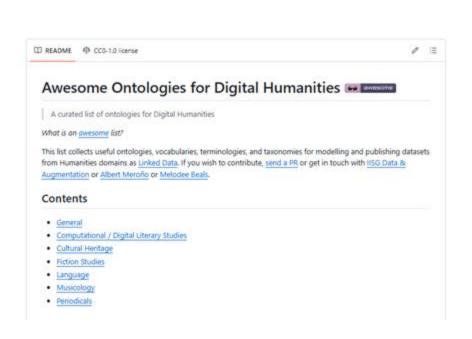
- "All schemas within historical demography deal with many similar concepts, but have very limited interoperability."

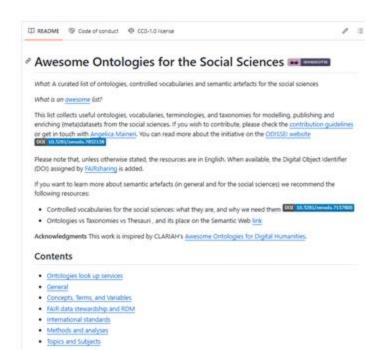
#### And they started harmonizing!

Mourits, R. J., Riswick, T., & Stapel, R. (2025). Common Language for Accessibility, Interoperability, and Reusability in Historical Demography. In B. Steffen (Ed.), Bridging the Gap Between Al and Reality (pp. 10–29). Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-73741-1\_2

CLAIR-HD (https://iisg.amsterdam/en/blog/clair-hd)

#### Reusing is needed, but: how to find existing vocabularies in SSH?





## FAIR principles

In 2016, the 'FAIR Guiding Principles for scientific data management and stewardship' were published in Scientific Data, a Nature journal.

- The authors intended to provide guidelines to improve the Findability, Accessibility,
   Interoperability, and Reuse of digital assets.
- The principles emphasise *machine-actionability* (i.e., the capacity of computational systems to find, access, interoperate, and reuse data with none or minimal human intervention)

Still what is FAIR and what that means from a technical perspective is still heavily debated and under construction (in the EU). Nothing is set yet, but there is progress in:

- community specific (implementation) strategies (FIPs, e.g. <u>CLARIN & ODISSEI</u>)
- OSTrails.eu

#### Interoperable & vocabularies

I2: (Meta)data use vocabularies that follow the FAIR principles

When we are describing data or metadata, we often use **vocabularies** that provide the terms or concepts that are adequate to represent their content. However, if we use vocabularies in our data or metadata, we should make sure that they are also FAIR in their own right so that others, humans or machines, can find, access, interoperate and reuse them. The controlled vocabulary used to describe datasets needs to be documented and resolvable using globally unique and persistent identifiers. This documentation needs to be easily findable and accessible by anyone who uses the dataset.

## FAIR assessments, FAIR implementation profiles (FIPs)

#### What is your FIP?

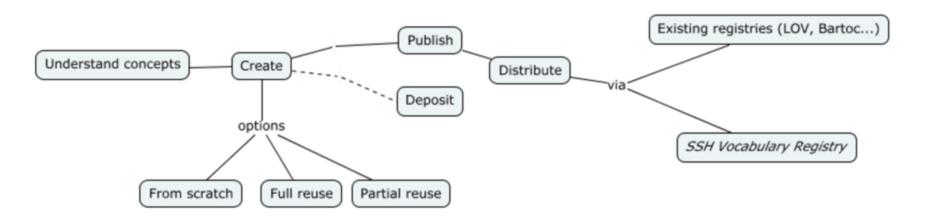
Check out the FIP mini-questionnaire which will lead you through the creation of your own FAIR Implementation Profile: https://bit.ly/yourFIP or download the questionnaire in PDF.

FAIR principle	Question	FAIR enabling resource types	Your answers
F1	What globally unique, persistent, resolvable identifiers do you use for metadata records?	Identifier type	e.g. PURL, DOI
F1	What globally unique, persistent, resolvable identifiers do you use for datasets?	Identifier type	
F2	Which metadata schemas do you use for findability?	Metadata schema	
F3	What is the technology that links the persistent identifiers of your data to the metadata description?	Metadata-Data linking mechanism	
F4	In which search engines are your metadata records indexed?	Search engines	
F4	In which search engines are your datasets indexed?	Search engines	
A1.1	Which standardized communication protocol do you use for metadata records?	Communication protocol	
A1.1	Which standardized communication protocol do you use for datasets?	Communication protocol	
A1.2	Which authentication & authorisation technique do you use for metadata records?	Authentication & authorisation technique	
A1.2	Which authentication & authorisation technique do you use for datasets?	Authentication & authorisation technique	
44	Military materials to come the plan of a constant	Administration of the Control of the	
11	Which knowledge representation languages (allowing machine interoperation) do you use for metadata records?	Knowledge representation language	
11	Which knowledge representation languages (allowing machine interoperation) do you use for datasets?	Knowledge representation language	
12	Which structured vocabularies do you use to annotate your metadata records?	Structured vocabularies	
12	Which structured vocabularies do you use to encode your datasets?	Structured vocabularies	
13	Which models, schema(s) do you use for your metadata records?	Metadata schema	
13	Which models, schema(s) do you use for your datasets?	Deta schema	
R1.1	Which usage license do you use for your metadata records?	Data usage toense	
R1.1	Which usage license do you use for your datasets?	Data usage license	
R1.2	Which metadata schemas do you use for describing the provenance of your metadata records?	Provenance model	
		Provenance model	

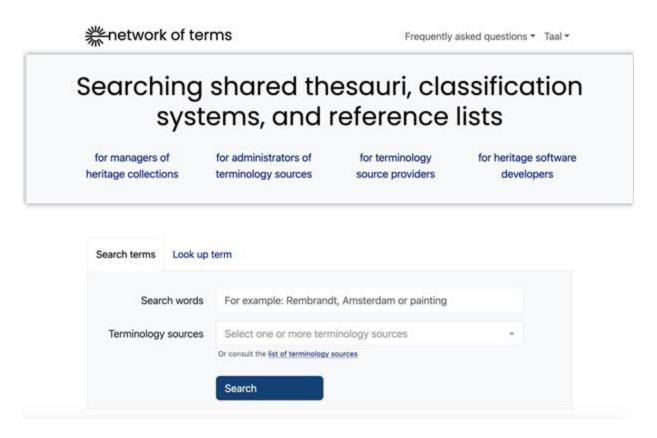
## FAIR (limitations)

- are semantic artifacts only limited to the Linked Data domain?
  - we (will) have lot of (legacy) data that is valuable and should be available in a FAIR manner
  - also these have semantic artefacts, e.g. XSD
  - FAIR tends to have a RDF/LD focus, but needs to go beyond that

## Challenge for end-user: vocabulary workflows



## Solution for curators (GLAM sector)



## Existing vocabulary (semantic artifacts) registries

- Linked Open Vocabularies
- Bartoc.org
- Ontoportal Alliance
- FAIRsharing.eu
- Limitations
  - Not all do common Apis
  - Not all do caching
  - Registries are general or per scientific domain, but not for SSH
- We need(ed) one registry for SSH



## Trend exploration and analysis

#### Aggregate portal services

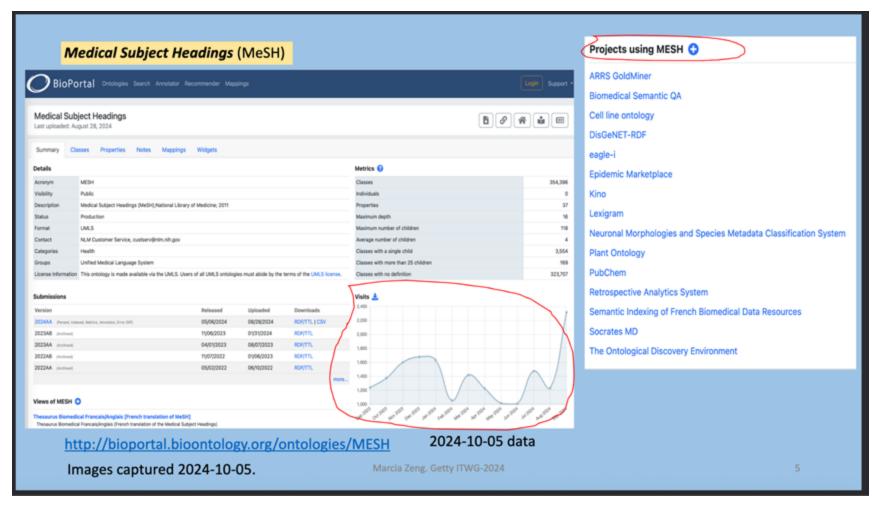
 provide support for searching, sharing, mapping, visualizing, and analyzing a large repository of ontologies, vocabularies, terminologies, and annotations.

#### **Individual KOS-based services**

## The KOS lifecycle encompassed by such services includes

- · creation,
- · transformation,
- mediation,
- · migration,
- exchange,
- integration, and
- · aggregation.

2

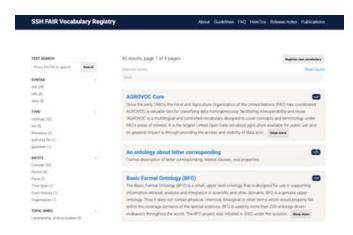


# Distributing vocabularies via the The SSH Vocabulary Registry

## FAIR vocabulary registry & recommender

Production version launched at DH Benelux (June, 2025) <a href="https://registry.vocabs.clariah.nl/">https://registry.vocabs.clariah.nl/</a>

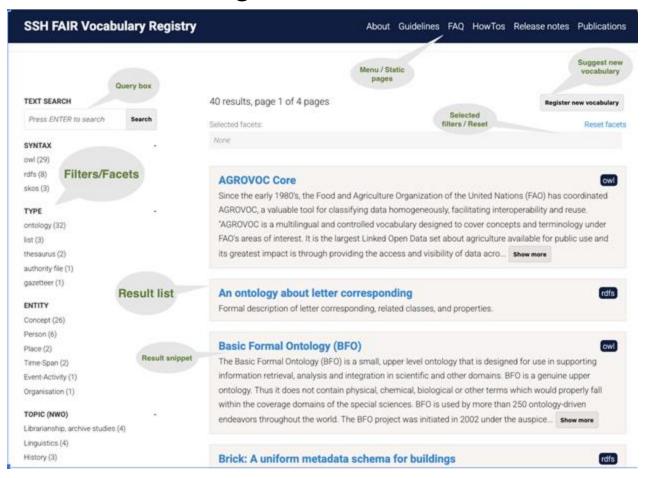
- Development started during CLARIAH-Core project (2022)
- Currently 41 vocabularies only
  - But ca. 200 are in review
  - Initially, it included most vocabularies from two curated lists:
    - <u>Awesome Ontologies for Digital Humanities</u> (ca. 40 vocabs)
    - YALC: a registry of Linked Open Datasets (ca. 126 vocabs)
- Some are highly FAIR-compliant, some are reference only



## Vocabularies in the SSH FAIR Vocabulary Registry

- Vocabulary metadata is stored in XML (using a CMDI profile)
  - <a href="https://www.clarin.eu/content/cmdi-component-metadata-infrastructure">https://www.clarin.eu/content/cmdi-component-metadata-infrastructure</a>
- Not all vocabularies are FAIR
- We have been working on:
  - Selecting most relevant ones
  - Adding more metadata
    - and creating/adapting the vocabulary metadata using existing vocabularies (e.g., VOID, DCAT, MOD)
  - Testing download / Improving workers

#### The search and browsing interface



## Detailed page per vocabulary



#### Vocabularies for vocabulary description

DCAT (and DC terms) -> Data Catalog Vocabulary

VOAF -> Vocabulary of a Friend

VANN -> A vocabulary for annotating vocabulary descriptions

REV -> RDF Review Vocabulary

PAV -> Provenance, Authoring and Versioning

Schema.org

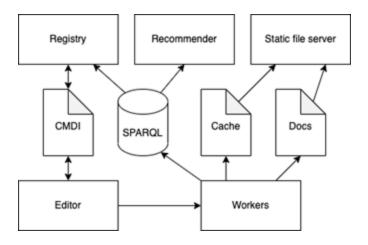
VoID -> Vocabulary of Interlinked Datasets

MOD -> Metadata for Ontology Description and publication (Ontoportal)

#### Basic approach

- the registry is metadata only, i.e., all vocabs "live" somewhere else ... has an owner/home outside of the registry ... unless it "died" than we have still a cache
- we add common APIs as far as possible, e.g., SKOSMOS for SKOS, SPARQL for any RDFbased vocab
- the registry is for more than just vocabularies, but for any semantic artifacts, i.e. including schema's models, ontologies, etc. (this is a main difference with Termennetwerk)
- as not all datasets will become LD, we also want to cater for semantic artifacts for other data representations like XML or CSV
- as there are unlimited ways a homepage for a vocab might look and lead to a download of the actual vocab, there is still a lot manual/semi-automatic labour going on, and, as many other similar registries do, we have to revisit the vocabulary from time to time to curate our entry ... we're working to make this move to automatic more and more ...

#### FAIR vocabulary registry & recommender (architecture overview)



#### Repositories:

- https://github.com/CLARIAH/vocab-workers
- <a href="https://github.com/CLARIAH/vocab-registry">https://github.com/CLARIAH/vocab-registry</a>
- https://github.com/CLARIAH/vocab-registry-editor
- <a href="https://github.com/CLARIAH/vocab-registry-docs">https://github.com/CLARIAH/vocab-registry-docs</a>
- (vocabulary recommender temporarily in development in private Gitlab)

#### Recent paper:

Meijer, K. and Windhouwer, M. (2024). The CLARIAH FAIR Vocabulary Registry. CLARIN Annual Conference Proceedings, page 154. <a href="https://www.clarin.eu/sites/default/files/CLARIN2024\_ConferenceProceedings\_final.pdf">https://www.clarin.eu/sites/default/files/CLARIN2024\_ConferenceProceedings\_final.pdf</a>

#### Guidelines and HowTos

SSH FAIR Vocabulary Registry

About **Guidelines** FAQ HowTos Release notes

Concepts Software Standards Creation Reuse Publication Example workflow 1

#### **SSH FAIR Vocabularies Guidelines**

Just the start of a collaborative effort!

## Roadmap

#### Current work in progress:

- Inventorizing & adding existing vocabularies in SS and Humanities
- Making the vocabulary metadata schema also FAIR
- Defining workflows (accepting/rejecting suggested vocabs, curation)
- User evaluation
- Documentation pages / Guidelines

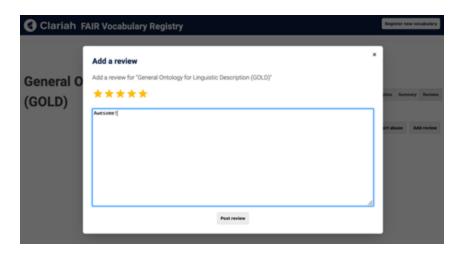
#### Future development (public roadmap in progress):

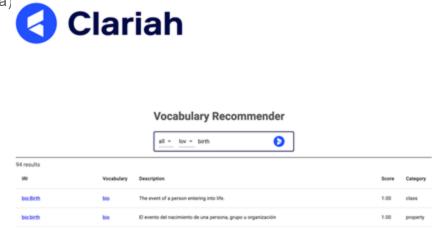
- Browser/website technical and usability improvements
- Continue developing the vocabulary recommender
- Implement APIs (OpenRefine, Lenticular Lens)
- Work on specificities for non-RDF vocabularies
- Integrate to INEO portal
- Integrate FAIR assessment scores for vocabularies

## FAIR vocabulary recommender (in the roadmap)

#### https://recommender.vocabs.dev.clariah.nl/

- scoring of a vocab to be based on:
  - Automatically checks on whether a vocabulary exists in other registries (e.g., Bartoc or LOV)
  - Registry gets user reviews (stars and comments)
  - Recommender Processes these inputs and give recommendations
  - Also recommendations based on user-provided csv (schema)





# Key points

## Key points

- Controlled vocabularies are an essential part (and/or product) of research work, but finding, reusing and sharing doesn't happen often.
- FAIR-ness of vocabularies in SSH is far from ideal.
- The SSH FAIR Vocabulary Registry aims to fill in these gaps: collecting, describing & facilitating reuse, and providing guidelines to researchers to encourage FAIR-ness of vocabularies from the creation stage.

## Use and contribute to SSH FAIR Vocabulary Registry

#### https://registry.vocabs.clariah.nl/

Suggesting new vocabularies to add



 Contributing content or questions to the Guidelines Participating of (user)
 evaluation



## extra

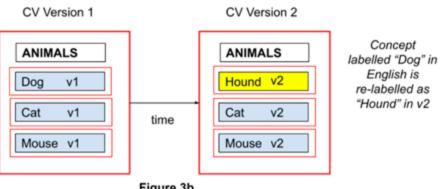


Figure 3b Versioning Propagated UPWARDS

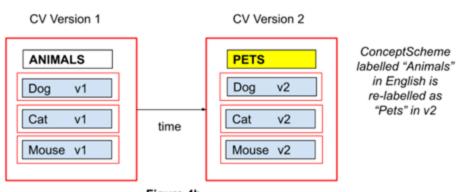
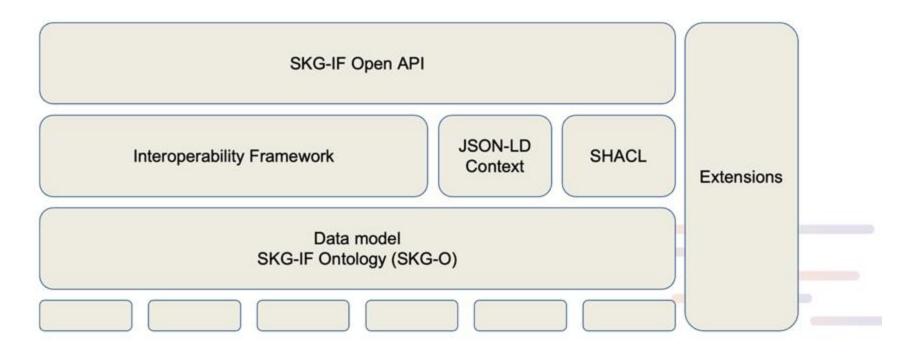


Figure 4b Versioning Propagated DOWNWARDS

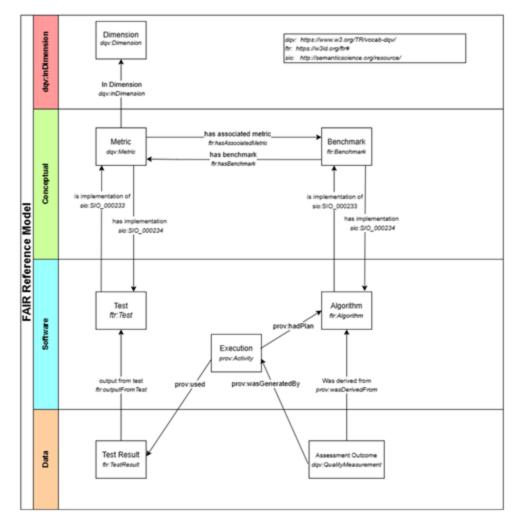
## **SKG-IF Overview**

#### https://zenodo.org/records/15309838







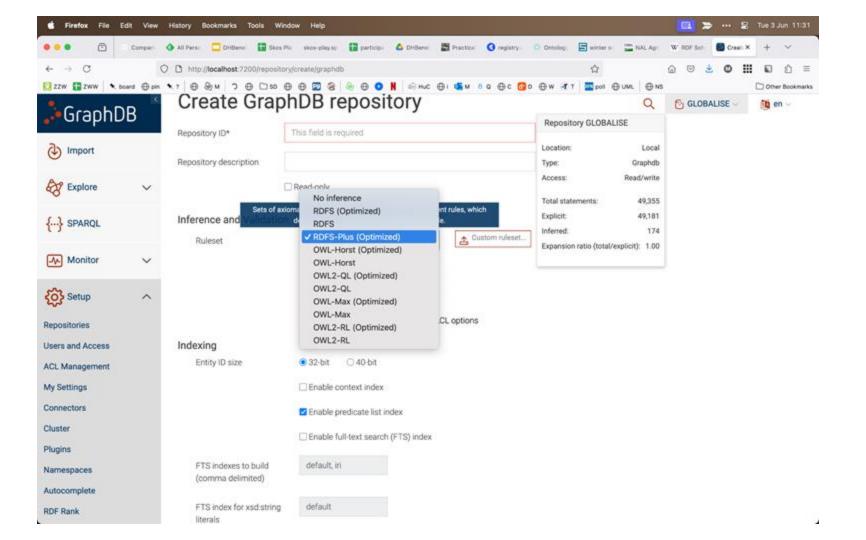


https://zenodo.org/records/14794901

#### reuse

- RDFS and beyond add reasoning capabilities
  - range <- prop -> domain
  - using the prop has impact on the subject and/or object
  - 0 ...
- reusing an ontology also takes over inference
- mixing & matching multiple ontologies might lead to incompatible inferences
- still try to reuse as much as possible ;-)

http://workingontologist.org/



#### Publish vocabularies: handson

- github pages
- skosplay
- skosmos

#### List of concepts to place in the taxonomy:

#### https://lod.nal.usda.gov/nalt/en/page/29712

squashes

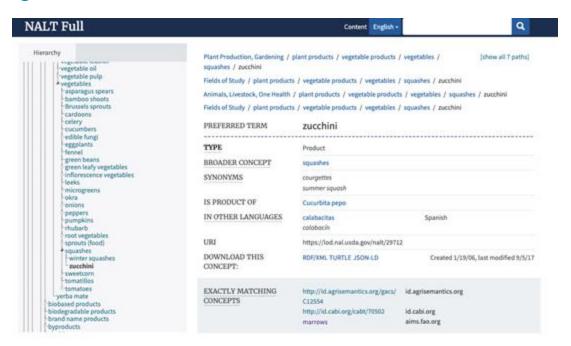
vegetables

zucchini

winter squashes

summer squash

courgette



## Coffee/tea break