

# The FAIR Principles – an important challenge to variational linguistics

Christina Mutter

<https://www.verba-alpina.gwi.uni-muenchen.de/>

ICLaVE | 11 2022

11-14 April 2022 - online





## 1. Project Overview of VerbaAlpina

research aims, area under investigation, conceptual domains, data

## 2. The FAIR principles

## 3. FAIR@VerbaAlpina

## 4. How the FAIR principles are challenging for variational linguistics



- *VerbaAlpina. Der alpine Kulturraum im Spiegel seiner Mehrsprachigkeit* (VerbaAlpina. The Alpine cultural region reflected through its multilingualism)
- Funded by the German Research Foundation (DFG)
- 1<sup>st</sup> term: 10/2014-10/2017, 2<sup>nd</sup> term: 11/2017-10/2020, 3<sup>rd</sup> term: 11/2020 – 10/2023 (perspective until 2026)
- Investigation of the multilingual Alpine region
- Combination of (geo-)linguistics and Digital Humanities (DH)



## Research Aims

- Selective and analytical investigation of the linguistically and dialectally highly fragmented alpine space in its historico-cultural and historical-linguistic unity
- Overcoming of the traditional limitation of geolinguistic investigation to nation-states
- recognition of connections regarding the etymology of the individual dialectal words
- Setting up a portal by using modern media technology: documentation, data collection, collaborative development
- cooperation with other projects is fundamental for VerbaAlpina



## Area under investigation: The Alpine region

- Area of investigation is limited to the territorial borders defined by the Alpine convention
- surface area of 190,600 km<sup>2</sup>, encompasses parts of six different countries (D, A, CH, I, F, SLO) and two entire countries (FL, MC)
- ethnographic and topographic homogeneity and strong linguistic heterogeneity → 3 language families





## Three conceptual domains

project years

1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---

calendar year

2014	2015	2016	2017	2018	2019	2020	2021	2022	2023
------	------	------	------	------	------	------	------	------	------

quarter

i, ii, iii, iv	i, ii, iii, iv	i, ii, iii, iv	i, ii, iii, iv	i, ii, iii, iv	i, ii, iii, iv	i, ii, iii, iv	i, ii, iii, iv	i, ii, iii, iv	i, ii, iii, iv
----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------

project phase

I	II	III
---	----	-----

focus

<p><b>traditional life</b></p> <ul style="list-style-type: none"> <li>• alpine pasture farming</li> <li>• milk processing</li> </ul>	<p><b>nature</b></p> <ul style="list-style-type: none"> <li>• landscape formations</li> <li>• weather</li> <li>• fauna</li> <li>• flora</li> </ul>	<p><b>modern life</b></p> <ul style="list-style-type: none"> <li>• ecology</li> <li>• tourism</li> </ul>
--	--	--



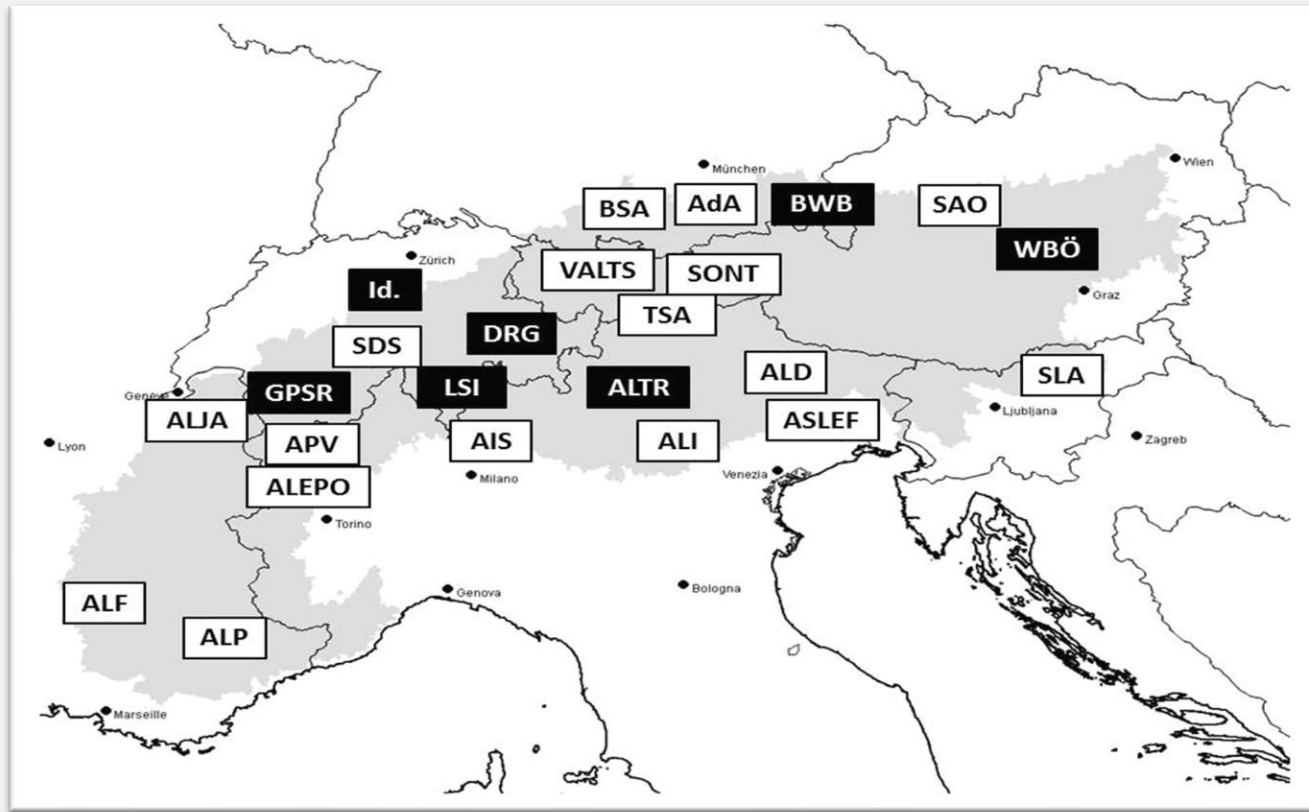
## Data

### Multiple different sources

- printed atlases/dictionaries (georeferenced)
- digital material from project partners
- crowdsourcing



## Atlases and Dictionaries in the Alpine region







## Crowdsourcing-Tool

[www.lmu.de/verbaalpina](http://www.lmu.de/verbaalpina)

Wählen Sie eine Gemeinde aus.

Wie sagt man zu *Begriff* in *Gemeinde*? Ihre Antwort



Research data have to be FAIR:

**F**\_indable

**A**\_ccessible

**I**\_nteroperable

**R**\_eusable

→ principles postulated by Wilkinson et al. 2016 as the guiding principles for scientific data management



**F**\_indable → via library catalogues and data aggregators

**A**\_ccessible → via open access licences

**I**\_nteroperable → via compatibility of databases and their interconnection

**R**\_eusable → results from F, A, I



## What VerbaAlpina does to make its data FAIR

### F\_indable

- Cooperation with the University Library of Munich University (VerbaAlpina data is available on [UB Discover](#), version 19/1 + 19/2) and the two RDM projects "[e-humanities-interdisziplinär](#)" (until 2021) and "[GeRDI](#)" (Generic Research Data Infrastructure) (until 2019)

### A\_ccessible

- Creative Commons licence (compatible with open access and open source) for all data managed by VerbaAlpina  
(up to version 18/1: CC BY SA 3.0, from 18/2: CC BY SA 4.0)



## I\_nteroperable

- through a fine granulation of the data stock via
  - structured data processing (transcription, tokenization, typification)
  - assignment of norm data (Q-ID, L-ID, GND, GeoNames etc.)
  - enrichment with metadata in DataCite and CIDOC CRM format
  - assignment of persistent identifiers (e.g. DOIs, Digital Object Identifiers)
- access to primary data and metadata (via [interactive map](#), [Lexicon Alpinum](#), [API](#))

## R\_eusable

- results from F, A, I



- requirements of F, A, R aim to be both human readable and machine readable
  
- apply to human-machine-human communication and to machine-machine communication
  
- I → only applies to machine-machine communication  
→ BUT: is crucial for the progress of research

## I\_nteroperable

- through a fine granulation of the data stock
  - structured data processing (transcription, tokenization, typification)
  - assignment of norm data (Q-ID, L-ID, GND, GeoNames etc.)
  - enrichment with metadata in DataCite and CIDOC CRM format
  - assignment of persistent identifiers (e.g. DOIs, Digital Object Identifiers)
- access to primary data and metadata (via interactive map, Lexicon Alpinum, API)

## R\_eusable

- results from F, A, I

## Assignment of norm data

### ▪ Norm data created by VerbaAlpina

For the 3 core entities

- morpho-lexical types → L
- concepts → C
- municipalities → A

e.g.: L1435, „babeurre (m.) (roa.)“  
C612, „ALMHÜTTE“ (*chalet*)  
A60171, „Sils in Engadin/Segl“

### ▪ Persistent identifiers of external institutions

(knowledge data bases/norm data bases/reference dictionaries)



## Persistent identifiers of external institutions

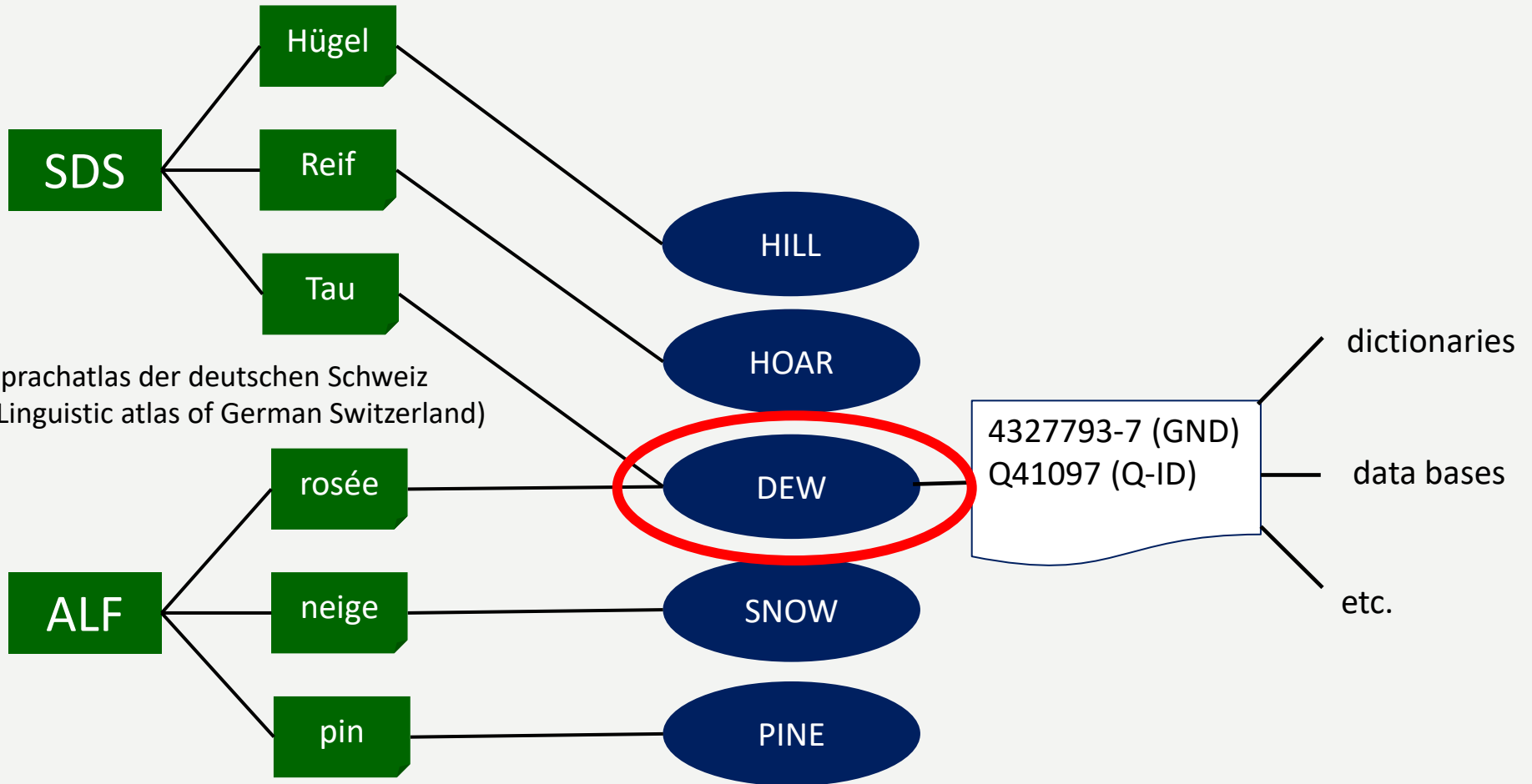
- **Q-IDs** of Wikidata (for concepts), partly also **L-IDs** of Wikidata (for morpho-lexical types)
- partly **GNDs** of the German National Library (for concepts) (*Gemeinsame Normdatei*, „Integrated Authority File“)
- **GeoNames** of [www.geonames.org](http://www.geonames.org) (for municipalities)
- **ISO Codes 639-3** (for languages)
- Identifiers of **reference dictionaries** (for morpho-lexical types + base types)
- **DOIs** (Digital Object Identifiers, assigned to every single data)

Limerick  
2962943

Q54050  
(Q-ID)

L73260  
(L-ID)

4135744-9 (GND)





- [Main page](#)
- [Community portal](#)
- [Project chat](#)
- [Create a new Item](#)
- [Create a new Lexeme](#)
- [Recent changes](#)
- [Random Item](#)
- [Query Service](#)
- [Nearby](#)
- [Help](#)
- [Donate](#)

- Tools
- [What links here](#)
  - [Related changes](#)
  - [Special pages](#)
  - [Permanent link](#)
  - [Page information](#)
  - [Concept URI](#)

English Not logged in Talk Contributions Create account

Item [Discussion](#)

Read [View history](#)



Wiki Loves Earth 2020 photo competition: take photos in nature and support Wikipedia.

**milk** (Q8495)

white liquid produced by the mammary glands of mammals


edit

[In more languages](#)

[Configure](#)

Language	Label	Description	Also known as
English	milk	white liquid produced by the mammary glands of mammals	
German	Milch	weißliche, undurchsichtige, als Milchfett-in-Wasser-Emulsion vorliegende, von Säugern produzierte Flüssigkeit	
French	lait	liquide biologique comestible produit par les mammifères femelles	
Bavarian	Muich	No description defined	





Kontakt | A-Z | Träger / Förderer | Datenschutz | Impressum | Hilfe | Mein Konto | English

**KATALOG DER DEUTSCHEN NATIONALBIBLIOTHEK**

Gesamter Bestand | Musikarchiv | Exilsammlungen | Buchmuseum

→ Suchformular zurücksetzen

sw all "Butter"   Expertensuche ?

eingeschränkt auf  
 - Materialarten: Elektronische Datenträger  
 - Normdaten: Sachbegriffe

**Ergebnis der Suche nach: sw all "Butter"**

[← Zurück zur Trefferliste](#)

Treffer 1 von 11

<small>GND</small>	
<b>Link zu diesem Datensatz</b>	<a href="http://d-nb.info/gnd/4009236-7">http://d-nb.info/gnd/4009236-7</a>
<b>Sachbegriff</b>	Butter
<b>Quelle</b>	M
<b>Oberbegriffe</b>	Milchprodukt Streichfett
<b>DDC-Notation</b>	637.2 641.372
<b>Systematik</b>	32.7 Milchwirtschaft ; 31.11 Lebensmitteltechnologie
<b>Typ</b>	Allgemeinbegriff (saz)
<b>Andere Normdaten</b>	LCSH: Butter RAMEAU: Beurre LCSH: Cooking (Butter) RAMEAU: Cuisine (beurre)

↓ Katalog

- Einfache Suche
- Erweiterte Suche
- Browsen (DDC)
- Suchverlauf
- Meine Auswahl
- Hilfe
- Datenschop
- Mein Konto
- Ablieferung von Netzpublikationen
- Informationsvermittlung

Login →

→ Über die Deutsche Nationalbibliothek

## Detail view of one specific data point

**IPA**  
Darstellung: IPA VA

**ISO Code**

**Source + link**  
Goebel, Hans (Wiesbaden): Atlant linguistich di ladin dolomitch y di dialec vejins I, vol. 1-7 (sprechend: http://ald.sbg.ac.at/ald/ald-i/index.php), 1998, vol. 1-7, Reichert  
[Link](#)

**GeoNames**

**Dictionaries**  
Treccani: latte  
CNRTL: lait

**Etymological dictionary**  
Georges: lac 2, 525

**Wikidata**

**l'ate** (Einzelbeleg)

Erbezzo

Phonetischer Typ (nicht typisiert) VA

Morpho-lexikalischer Typ lait / latte (roa m.) T C VA

Basistyp läcte(m) (lat.) G VA

Quelle	Konzept
ALD- 348#1 177 (Erbezzo)	MILCH (Wikidata)

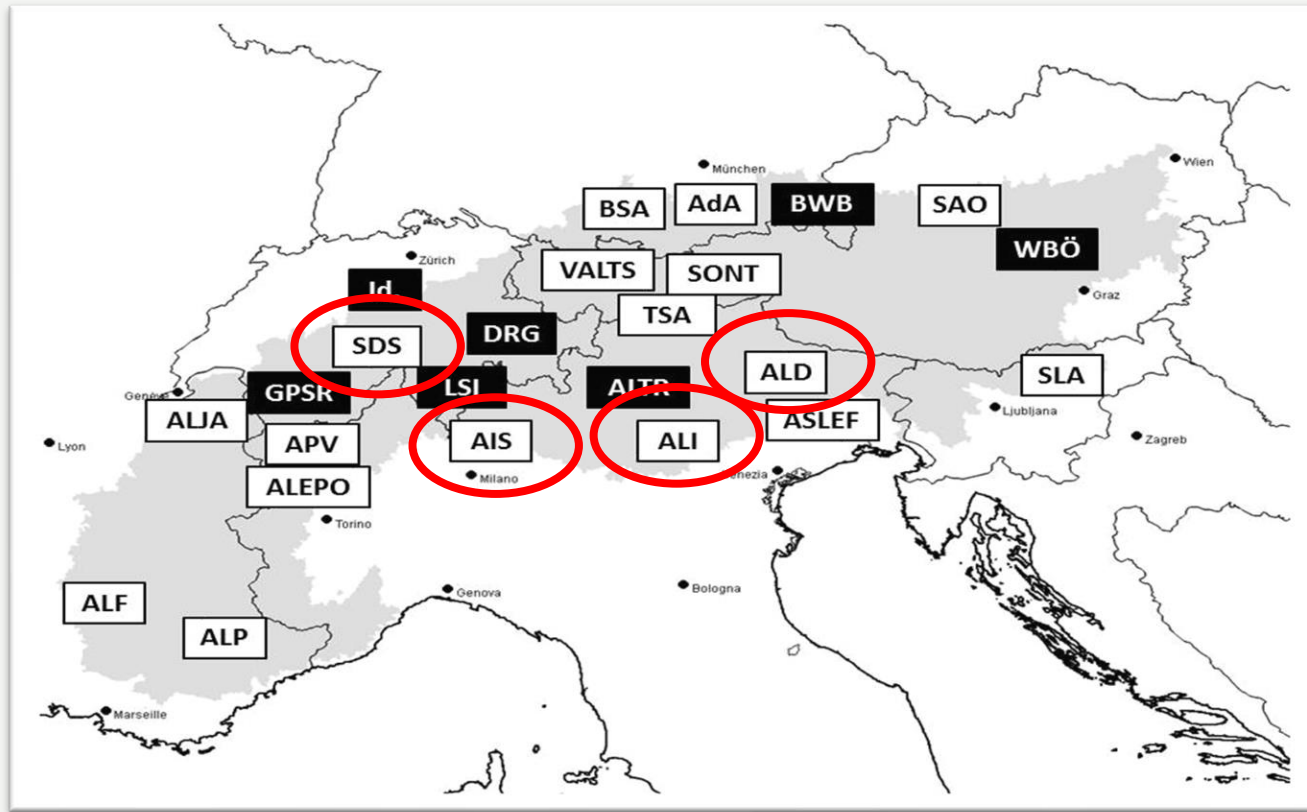


VerbaAlpina tries to create FAIR output

**BUT:** many data sources (the input) are far away from being FAIR

**Main reason:** lack of interoperability

## FAIRness of atlases in the Alpine region



## Comparison of the FAIRness of 4 geolinguistic atlases with VerbaAlpina

	Findable		Accessible		Inter- operable	Reusable	
	hum. read.	mach. read.	hum. read.	mach. read.		hum. read.	mach. read.
ALI	–	–	–	–	–	–	–
SDS	+	+	+	–	–	+	–
AIS	+	+	+	–	–	+	–
ALD	+	+	+	–	–	+	+
VA	+	+	+	+	+	+	+

**ALI = Atlante linguistico italiano**

**SDS = Sprachatlas der Deutschen Schweiz**

**AIS = Sprach- und Sachatlas Italiens und der Südschweiz**

**ALD = Atlant linguistisch dl ladin dolomitich y di dialec vejins**

**VA = VerbaAlpina**

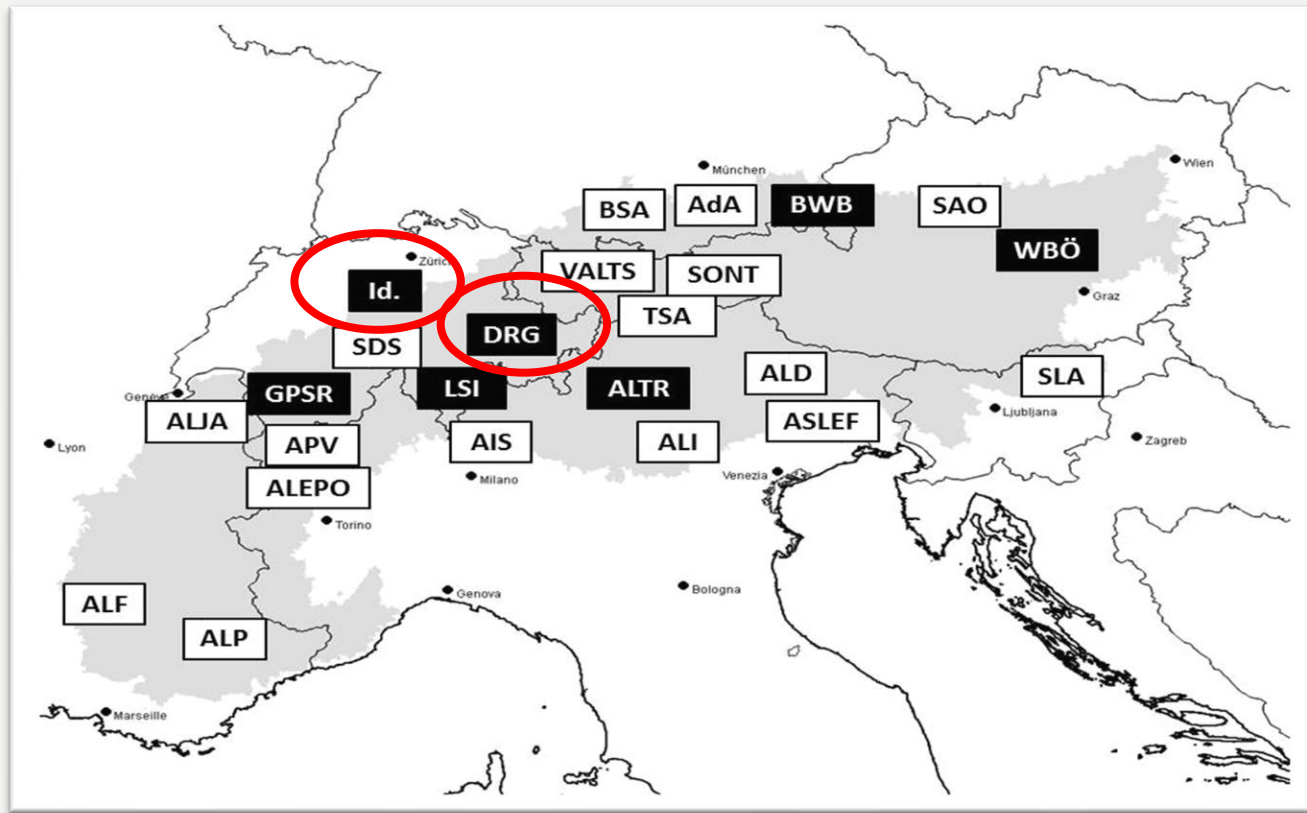




## Lack of FAIRness because...

- linguistic atlases are usually only accessible as printed works
- only a few offer at least the most basic level of digitisation i. e. digital photos (scans), e.g. AIS ([NavigAIS](#)), SDS for original material
- no older atlas has yet been prepared in the form of a structured corpus that also allows export of the data
- only the ALD is based on a non- interoperable digital format that proved to be machine-readable after certain adaptations

## Interoperability of dictionaries in the Alpine region





## ■ Schweizerisches Idiotikon

### Many Cons

- referencing is possible, but not precise enough
- structure of URLs is cryptic/semantically intransparent

Example: *Teie(n)* (Alpine chalet)

<https://digital.idiotikon.ch/idtkn/faksimile.php?band=12&spalte=31&hl=229458>

- referencing is based on the pages of the printed edition, not on the specific entries (naming of volume and column in the URL)
- all entries on one page are addressed via the same URL (no individual URLs)
- URLs refer to a scan of the printed volume and retain its page/column logic as reference system (see URL)
- entries cannot be copied/pasted (image and no HTML code)



**Bd XII 31** [gedruckt 1952]: *enerdei* ▾ · *sël<sup>b</sup>dei* ▾ · *Tei N.* ▾ · *dei II* ▾ · *Deia N.* ▾ · *deianen* ▾ · *deidei I* ▾ · *deidei II* ▾ · *deider<sup>ch</sup>* *deidur<sup>ch</sup>en* ▾ · *Teien* *Teielen*, *Ti(e)je<sup>n</sup>* ▾ · *deiënne<sup>n</sup>* ▾ · *deihëre<sup>n</sup>* ▾ · *deihin* ▾ · *deihin(d)en* ▾ · *deioben* ▾ · *deiun(d)en* ▾

**Bd XII 32**: *Suppe(n)teili* ▾ · *Deuängger* ▾ · *Theodor N.*, *Tedi* ▾ · *Theodosius N.*, *Dedli N.* ▾ · *Theophil N.*, *Fili N.*, *Töfel N.* ▾ · *di* ▾ · *dī* ▾ · *Diabel N.* ▾ · *didi* ▾ · *Tita* ▾ · *titä* ▾ · *titata* ▾ · *titi* ▾ · *Drëckteie<sup>n</sup>* ▾ · *Tuppe<sup>n</sup>teili* ▾

verloren gehen. 's *ist dei-dei gange<sup>n</sup>* B<sup>a</sup> (St.<sup>2</sup>). *Aber derchtüsing Gotts Willen, soll denn alles deidei gän, was mer so gwues verdiene<sup>n</sup>!* MWALDEN 1884. Nicht verdoppelt. *Bist bald fertig mit dem Chrishacken? Es ruckt em, es ist wie däi ZTöf<sup>st</sup>al. Es dunkt mich, das Fäfl<sup>i</sup> sei bald lär. Ja, 's ist wie däi.* ebd.

**Tei, Dei:** Matthäus (vgl. Bd IV 551/2, sowie *Tewes*), auch Amadeus, Taddäus (vgl. *Tede*) W.

**Deia:** Andreas GRVal.

**Deieldē** s. *Däldē*.

**Teien**, in GRD. (Bed. 1) *Tieien*, *Teiele<sup>n</sup>* (Bed. 2) GRMu. — f.: 1. a) Alp-, Sennhütte GRKl. (veraltet! It Tsch.); s. die Anm. ‚Welche Liegends kauft hetten und daruff gebuwen, gezimmeret oder sonst gebesseret... es werend Hüser, Ställ, Tieyen, Wald oder was einer sonst daran gebesseret hette mit Dach oder nüwen Zünen oder Muren...‘ GRD. LB. — b) ‚Vorraum hinter der Tür mit Küche in den Alphütten‘, Wohnraum, an den Heustadel angebaut GRD. (B.). — 2. übertr., gebrechliche, beschränkte, schwerfällige Weibsperson GRChur, D. (B.), He. (Tsch.), L., Mu., Schs, Valz. (Tsch.); Synn. *Noggen* (Bd IV 710); *Tschänggelen*, *Tschappe(le)n*. *Das ist en*

**Ti-ta** n.: Kinderspr., Uhr ZZoll.; vgl. *Tigg-Tagg*.

**ti-ta-ta:** Nachahmung des Taktes des Webstuhles (?). ‚Ach, der Lohn wird immer kleiner, und die Seide immer reiner, 's *ist nüt mē mit dem Titata*‘, Vers aus der Zeit, da es mit dem häuslichen Seidenweben zurückging. AFV. 24, 256 (ZMaur).

**ti-tä:** Nachahmung des Geräusches der Dachtraufe. Rätsel: *Es gät öppis um's Hūs umer und cha<sup>n</sup>n nie inen, wänn mer's scho<sup>n</sup> wol<sup>t</sup> inen lär, und macht ti-tä* Z (Dän.).

**di-di:** Lockruf für Zicklein, Katzen BsL.; SCHR.; ZO. (Stutz); vgl. *de-de I. Hale<sup>n</sup> chumm! Hale<sup>n</sup> didi!* SCHR.; s. auch Bd II 1768 o. (ZRafz). *Das ischt doch en Ströf. Die Tüfels Chatz!* [die Würste angefressen hat] *Ich schlohne<sup>n</sup> si z' Töd uf dem Platz. Wo ischt si ech<sup>t</sup>?* *Chumm di di di!* STUTZ, Gem. Auch als Bezeichnung des Tieres BsL.

**ti-ti:** Beteuerungsformel, gewiß BHa.; vgl. *doi, dui*.

**Diabel** m.: 1. Name einer Kuh ‚mit spitzen Hörnern, die einer Gabel gleich aufgerichtet sind‘ WUlr. — 2. Name (einer ‚Farbe‘) im Kartenspiel WVt. — 3. leichter Reisewagen; s. Bd VII 360 M. (E. XVIII., SWbl.). — Frz. *diable*.



- [Dicziunari Rumantsch Grischun \(DRG\)](#)

## Pros

- Links do not refer to the book page as in the Idiotikon, but to the corresponding entry

Example: BARGUN (Hay barn in mountain pastures)

<http://online.drg.ch/main.aspx#35ac0b4c3795c91831f3534a40a0c5f0>

- Contents are available in HTML and can be copied/pasted and cited

## Cons

- Identifiers are too long and make intuitive and automatic generation of specific URLs impossible



[online.drg.ch/main.aspx#35ac0b4c3795c91831f3534a40a0c5f0](http://online.drg.ch/main.aspx#35ac0b4c3795c91831f3534a40a0c5f0)

FsK — Milch Home - Milch SQL Fiddle Ilion-Troja Deutsche Dialekte i... Italienisch kostenlo... italiano/LeLingue/L... NavigAIS - AIS Navi... Google Maps lat/lo... 710 TL



deutsch | rumantsch

Institut | Über das Wörterbuch | Angebot | Publikationen | Shop | Fototeca Online | DRG-Online

DRG-ONLINE ↓

Neue Suche

Suchtipps

Abkürzungen

Vorworte

Hinweise

A

B

C

D

E

«Zurück

**BARGUN**

(DRG 02 / 192) | Artikel als PDF laden

Markieren Sie einfach grosszügig eine Textstelle, um die Band- und Seitenzahl anzuzeigen

Suchen

**BARGUN** surselv., **MARGUN** oengad., m. 'Heustadel auf Bergwiesen, Alpstall, Alphütte auf der obersten Weidestufe' usw., siehe unten. Vgl. Karte B II, p. 179.

I. *bargun*

C 1 *bargōn*, C 20 *bargēun*, C 21–23 *bargáun*, C 24, 3–4 *bargúy*, C 45 *barvúy*, C 50 *bragún*, C 60, 61 *bargúy*, C 62 *galbúy*, C 66 *balgún*, *galbúy*, C 68, 69 *bargúy*, C 70 *bargúñ*, C 81 *bargéun*, *barkéun*, C 93 *bargúy*, S 10, 13, 21, 25, 31–42, (64?) *bargún*. Zur Endung cf. Tab. 1 a. – Deutschbünden: Untervaz *pərgú*, Maienfeld *pərgú*, pl. *pərgú* (cf. auch MEINHERZ 116), Valzeina *pərgáu* neben *bōrgə*, Seewis, Fanas, Schiers, Furva Luzern *pərgáu*, Jenaz, Klosters *pərgún*, St. Antönien *pərgú* neben *bōrgə*, St. Gallen

 Source: DRG-Online <http://online.drg.ch/#35ac0b4c3795c91831f3534a40a0c5f0>



## Difficulties in the assignment of norm data

- ISO-Codes 639-3 are not detailed enough

Varieties according to Lia Rumantscha		Wikidata Q-ID	Glottolog	Constitutions CH, GR	ISO-693.3
"standard language"	Rumantsch Grischun	Q688873	Rumantsch Grischun	"Rhaeto-Romance" (Räto-romanisch)	roh
"written idioms"	Puter	Q688309	Upper Engadine		
	Vallader	Q690226	Lower Engadine		
	Surmiran	Q690216	Surmiran-Albula		
	Sursilvan	Q688348	Sursilvan		
	Sutsilvan	Q688272	Sursilvan-Oberland		
spoken dialect	Jauer	Q690181	Sutsilvan		



- **Concept assignment is sometimes unclear**  
Example: concept GHIACCIAIO (GLACIER)

**IL GHIACCIAIO**  
GLETSCHER – GLACIER  
ALF Suppl. 1. s. glacier  
A 91,14 = 62,16 = 0

**Legende:**

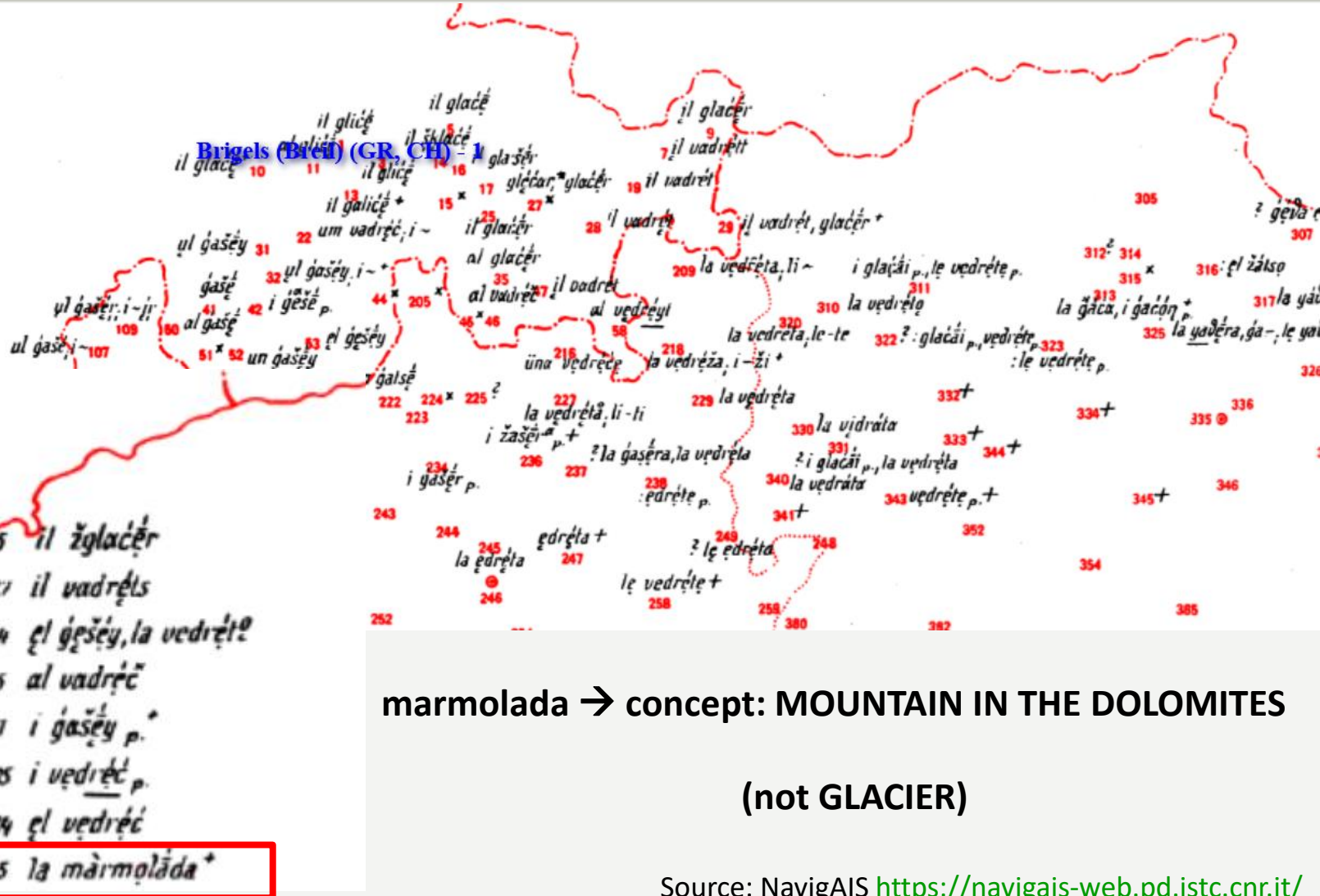
10 al vadrëc = aller Lawinenschnee über  
13 il vadrëc = ewiger Schnee, Firn (einem Bache.  
29 la murëna = Moräne.  
32 'Wenn der Gletscher donnert (kuvand u hüt ul gashëy), gibt es bald darauf Hagelwetter.'  
51 -vereiste Färhien an plattigen Stellen.  
120 = grass, im Frühling stets abschmelzende  
218 i lyëski = Gletscherspalten. fëisfläche  
313 z. B. i gacëj de la marmolëda.  
325 Bezieht sich nicht bloss auf die Marmolëda  
312 z. B. la yäda de l'antelëw [.. del Monte Antelao].

Brigels (Bret) (GR, CH)

Source: NavigAIS <https://navigais-web.pd.istc.cnr.it/>



# How the FAIR principles are challenging for variational linguistics



Source: NavigAIS <https://navigais-web.pd.istc.cnr.it/>



- Alpine Convention. Contracting parties. <https://www.alpconv.org/en/home/> [accessed 26 June 2020].
- Force11 (eds.), “The Fair Data Principles”, <https://www.force11.org/group/fairgroup/fairprinciples> [accessed 26 June 2020]
- Kümmer, S. / Lücke, S. / Schulz, J. / Zacherl, F.: s.v. “Forschungsdatenmanagement”, in: VerbaAlpina-de 19/1 (Erstellt: 18/2, letzte Änderung: 18/2), Methodologie, [https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage\\_id%3D493%26db%3D191%26letter%3DF%23112](https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage_id%3D493%26db%3D191%26letter%3DF%23112)
- Thomas Krefeld & Stephan Lücke (2019): Kleinsprachen, Dialekte und die FAIR-Prinzipien (am Beispiel von VerbaAlpina), Version 1 (02.10.2019, 15:57). In: Korpus im Text, Serie A, 48389. url: <http://www.kit.gwi.uni-muenchen.de/?p=48389&v=1>
- Krefeld, Thomas / Lücke, Stephan (2020): 54 Monate VerbaAlpina – auf dem Weg zur FAIRness, in: Ladinia, vol. XLIII, 139-156
- Lücke, S.: s.v. “FAIR-Prinzipien”, in: VerbaAlpina-de 19/1 (Erstellt: 18/2, letzte Änderung: 18/2), Methodologie, [https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage\\_id%3D493%26db%3D191%26letter%3DF%23128](https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage_id%3D493%26db%3D191%26letter%3DF%23128)
- Lücke, S.: s.v. “Metadaten”, in: VerbaAlpina-de 19/1, Methodologie, [https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage\\_id%3D493%26db%3D191%26letter%3DM%23140](https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage_id%3D493%26db%3D191%26letter%3DM%23140)
- Lücke, S.: s.v. “Normdaten”, in: VerbaAlpina-de 19/1 (Erstellt: 18/2, letzte Änderung: 18/2), Methodologie, [https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage\\_id%3D493%26db%3D191%26letter%3DN%23114](https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage_id%3D493%26db%3D191%26letter%3DN%23114)
- Lücke, S. / Schulz, J.: s.v. “Digital Object Identifier (DOI)”, in: VerbaAlpina-de 19/1 (Erstellt: 16/1), Methodologie, [https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage\\_id%3D493%26db%3D191%26letter%3DD%2373](https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage_id%3D493%26db%3D191%26letter%3DD%2373)
- Mutter, C.: s.v. “Wikidata”, in: VerbaAlpina-de 19/1 (Erstellt: 18/1), Methodologie, [https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage\\_id%3D493%26db%3D191%26letter%3DW%23105](https://doi.org/10.5282/verbaalpina?urlappend=%3Fpage_id%3D493%26db%3D191%26letter%3DW%23105)
- M. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L. O. Bonino da Silva Santos, P. Bourne, J. Bouwman, A. J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C. Evelo, R. Finkers, and B. Mons. The fair guiding principles for scientific data management and stewardship. *Scientific Data*, 3, 03 2016. doi: 10.1038/sdata.2016.18

# Thank you for your attention!

<https://www.verba-alpina.gwi.uni-muenchen.de/>

