

Markus Kunzmann, Florian Zacherl
Datengewinnung mittels Crowdsourcing
im Dienste der Sprachwissenschaft

Die virtuelle Forschungsumgebung
VerbaAlpina

Digital Humanities Austria 2017
04. – 06. Dezember 2017, Innsbruck



- **Projektbeschreibung**
- **Datenerhebung**
- **Eingliederung in den Datenbestand**
- **Visualisierung**
- **Nachhaltigkeit**

Projektbeschreibung



- *VerbaAlpina. Der alpine Kulturraum im Spiegel seiner Mehrsprachigkeit*
- Förderung durch die Deutsche Forschungsgemeinschaft
- seit 2014, bewilligt bis 2020 (Perspektive bis 2025)
- Untersuchung des sprachenreichen Alpenraums:
D, A, CH, I, F, FL, SLO, MC
- Umsetzung mittels moderner Konzepte und Methoden in der Verbindung mit den digitalen Möglichkeiten
(*Digital Humanities*)





Projektleitung

- Prof. Dr. Thomas Krefeld (Institut für Romanische Philologie)
- Dr. Stephan Lücke (IT-Gruppe Geisteswissenschaften)

MitarbeiterInnen

- Christina Mutter (wiss. Koordination)
- Markus Kunzmann (Germanistik)
- Aleksander Wiatr (Romanistik, Slowenisch)
- Florian Zacherl (Informatik)
- David Englmeier (Informatik)
- Monika Hausmann, Filip Hristov, Katharina Knapp, Marina Pantele, Daniela Warras (wissenschaftliche Hilfskräfte)



Projektjahre

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

Kalenderjahre

2014	2015	2016	2017	2018	2019	2020	2021	2022	2023
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

Quartale

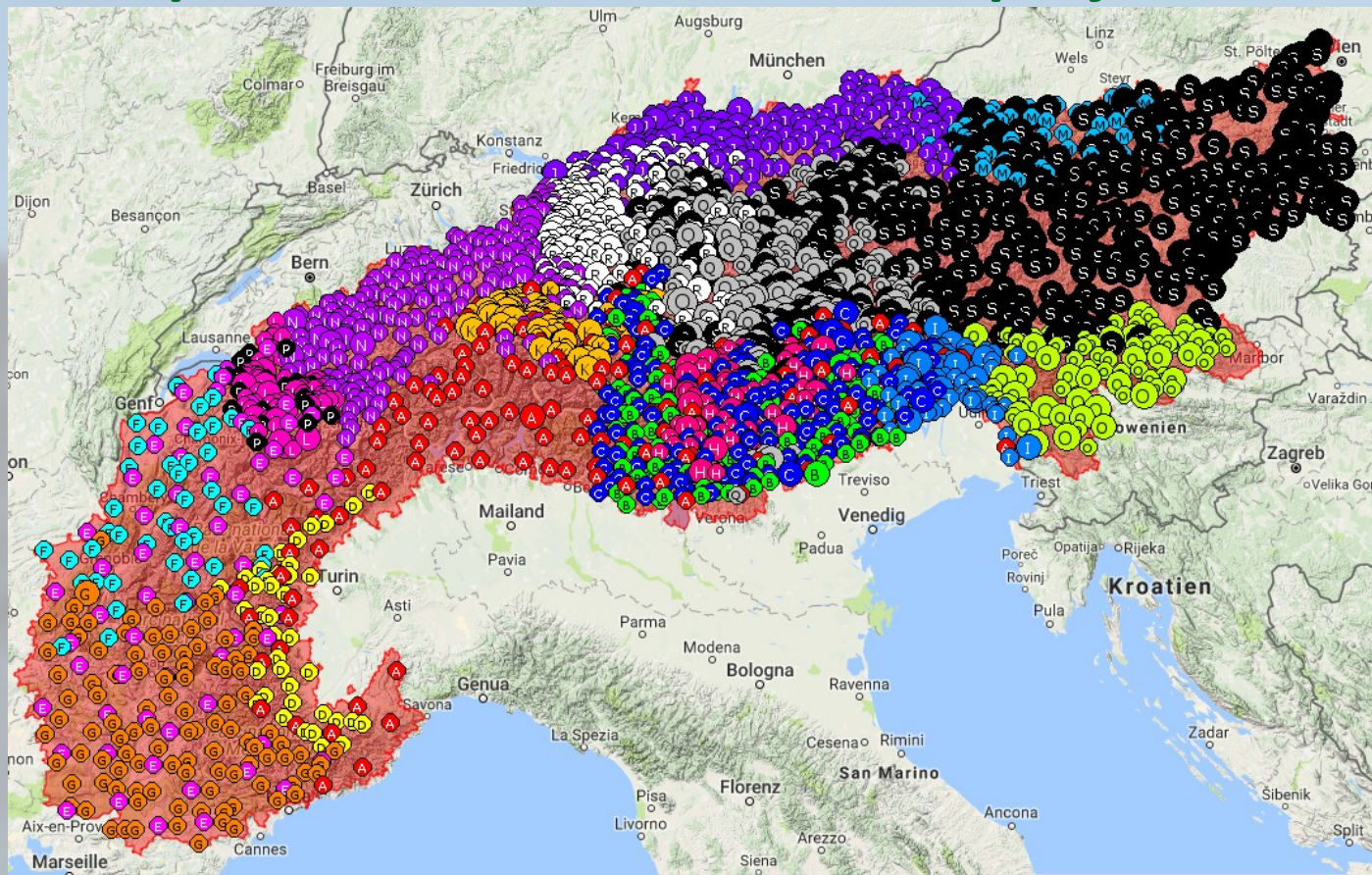
I			II			III		
---	--	--	----	--	--	-----	--	--

Projektphase

I	II	III
Kultur	Natur	moderne Lebenswelt
<ul style="list-style-type: none"> • Almwirtschaft • Volkstümliche Medizin • Traditionelle Küche 	<ul style="list-style-type: none"> • Landschaftsformen • Wetter • Fauna • Flora • traditionelle Küche 	<ul style="list-style-type: none"> • Ökologie • Tourismus

Schwerpunkt

Auflösung der Grenzen in Bezug auf sprachwissenschaftliche Teilprojekte



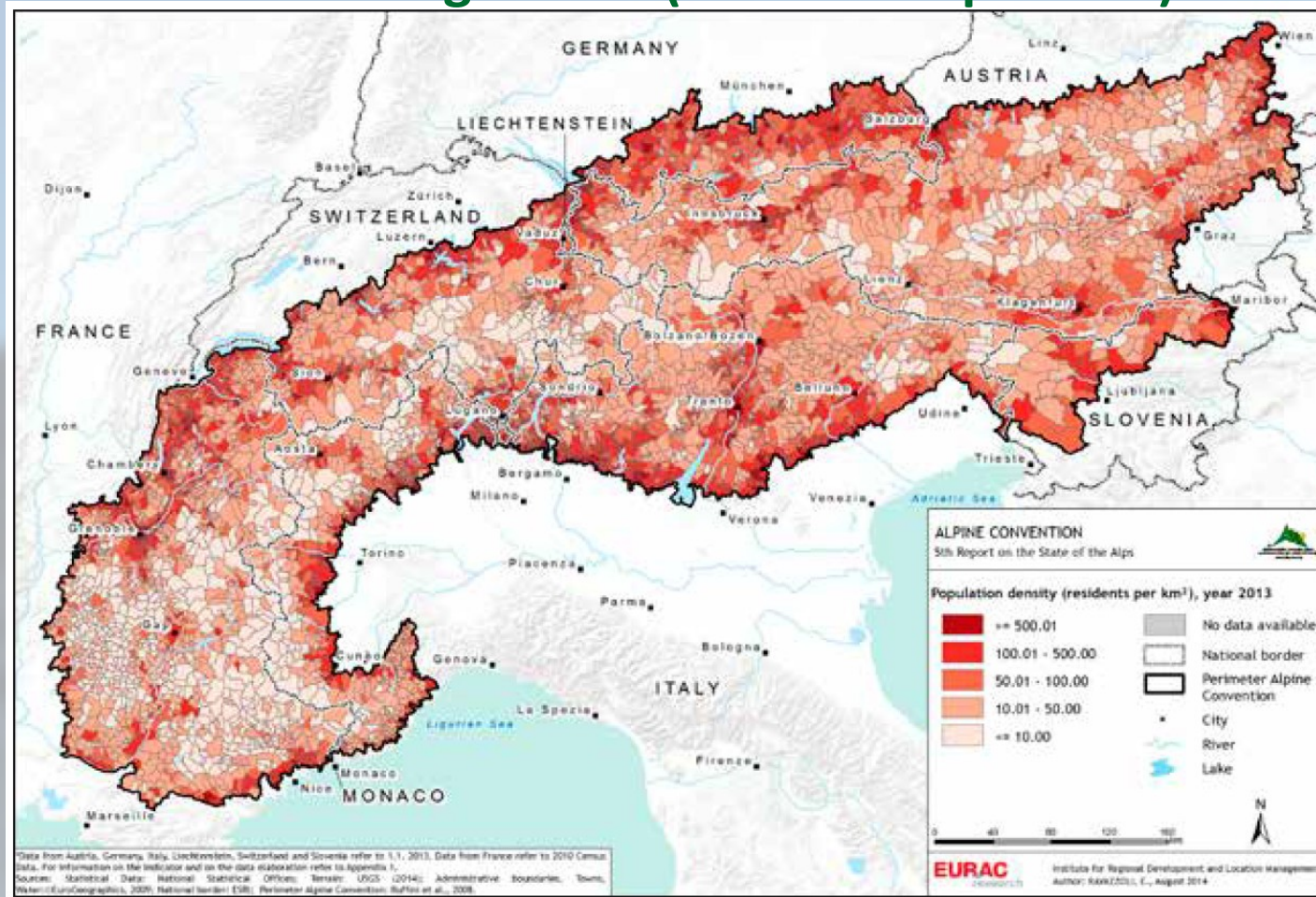


Auflösung der Grenzen hinsichtlich wissenschaftlicher Disziplinen

- Sprachwissenschaftler
- Historiker
- Ethnographen und Volkskundler
- Sprach- und Dialektinteressierte

Datenerhebung

Bevölkerungsdichte (Einwohner pro km²)





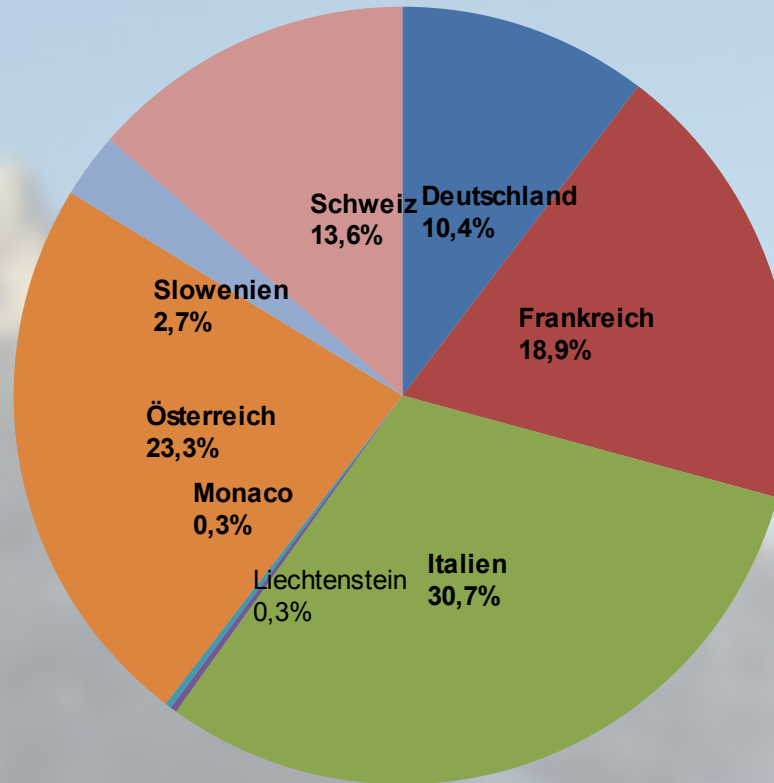
Bevölkerung, Fläche und Bevölkerungsdichte im Alpenraum

	Bewohner des alpinen Raumes	Fläche (km ²) des alpinen Raumes	Bevölkerungsdichte im alpinen Raum	Bevölkerungsdichte national
Deutschland	1.476.519	11.160	132,3	225,3
Frankreich	2.683.801	40.801	65,8	103,4
Italien	4.364.538	51.995	83,9	201,8
Liechtenstein	36.838	160	230,2	230,2
Monaco	36.950	2	18.475	18.475
Österreich	3.318.045	54.592	60,8	100,8
Slowenien	385.973	6.796	56,8	101,6
Schweiz	1.929.424	25.211	76,5	201,0
Alpen	14.232.088	190.717	74,6	-

ALPENKONVENTION (2015), S. 17



Bewohner des alpinen Raumes nach Staaten



EW Alpenraum:
14.232.088



Crowdsourcing Anforderungen

Welche Daten werden erhoben?

- Konzept
- Gemeinde
- Sprachbeleg (aktuell über Texteingabe)
- Implizit Datum

Zukunft?

- Genus bei Substantiven
- Sprachaufnahme
- phon. Ebene des Materials



Crowdsourcing Anforderungen

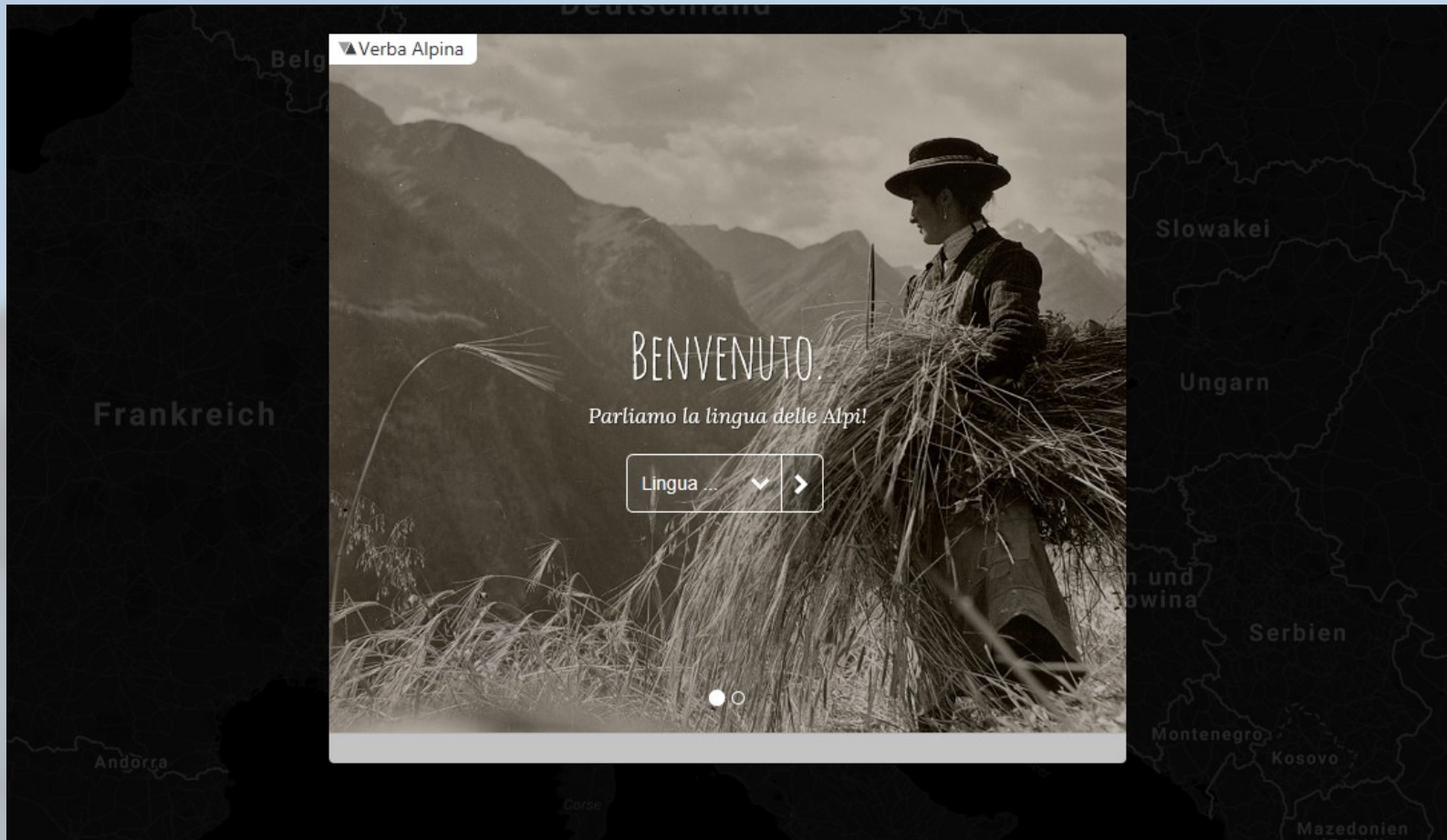
Zielgruppe

- Wissenschaftler und Laien
- Verzicht auf wissenschaftliche Terminologie

Bedienbarkeit

- Möglichst einfache Bedienbarkeit für mobile und Desktop-Geräte
- Simples responsives Design
- Registrierung nicht notwendig, Nutzer wird aber auf erweiterte Möglichkeiten hingewiesen

Crowdsourcing-Portal





Crowdsourcing Mailingaktion

Verlage

- Zeitschriften Almwirtschaft und Milchverarbeitung

Ausbildung

- Berufsschulen und sonstige Ausbildungsstätten

Gewerbe

- Käsereien
- Molkereien

Almen, Kultur

- Vereine
- Museen

Portale

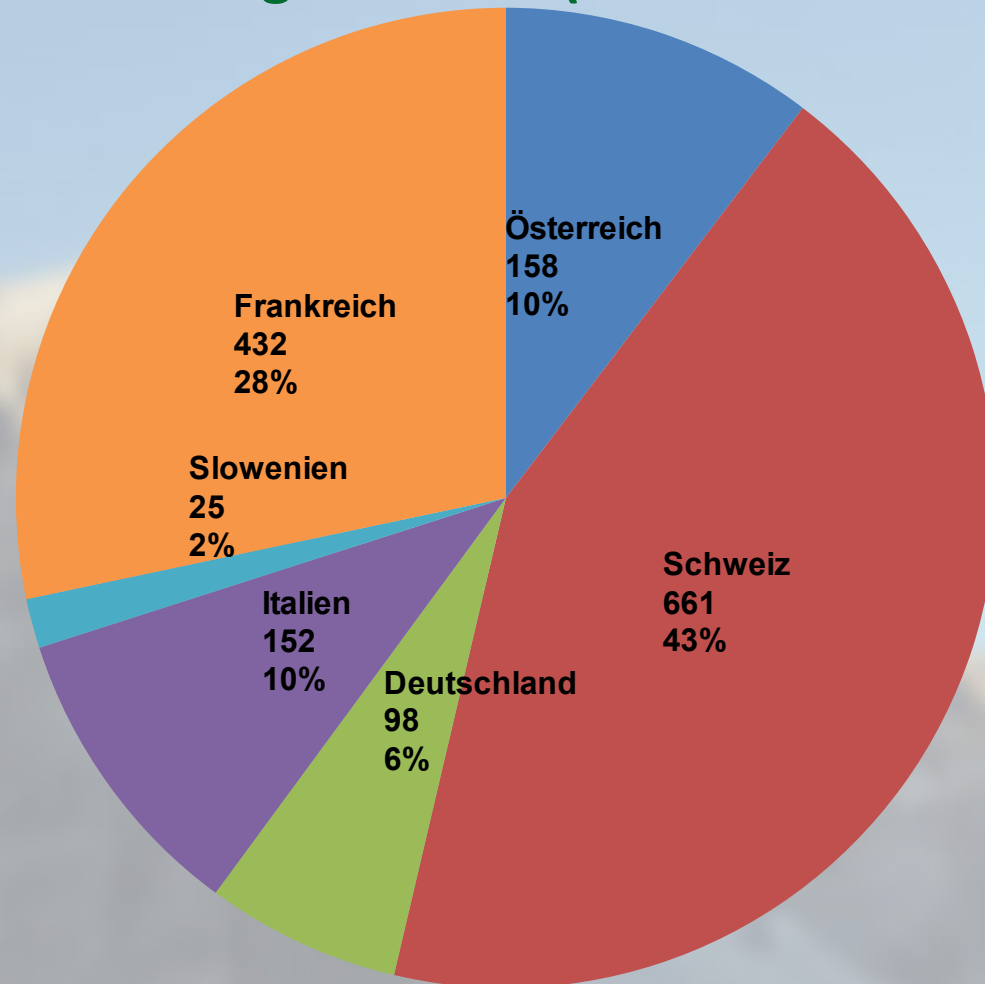
- Fachportale der Milchbranche (z.B. z'alp.ch)

Medien

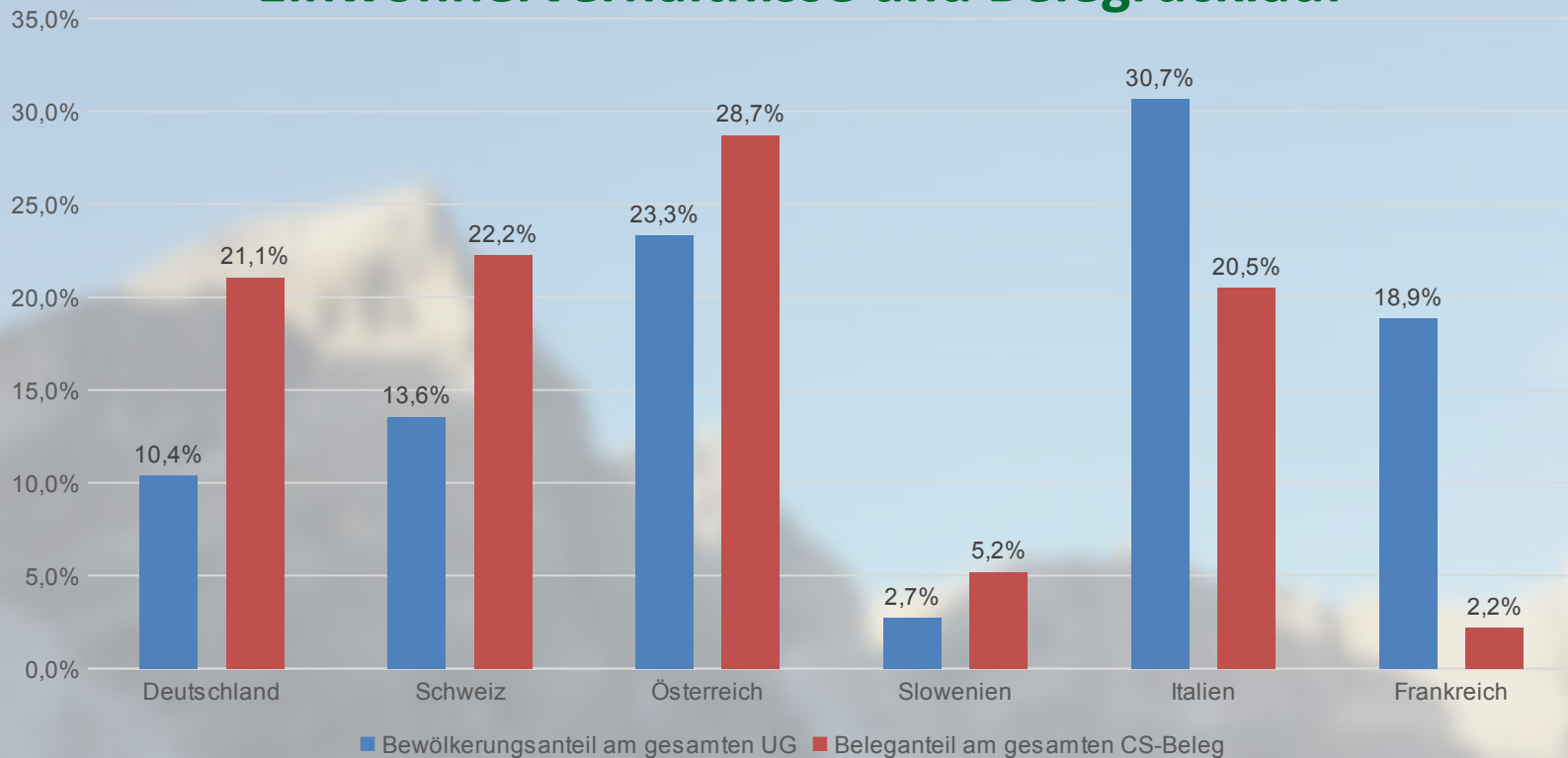
- Zeitungen
- Fernsehen
- Radio
- Verteiler

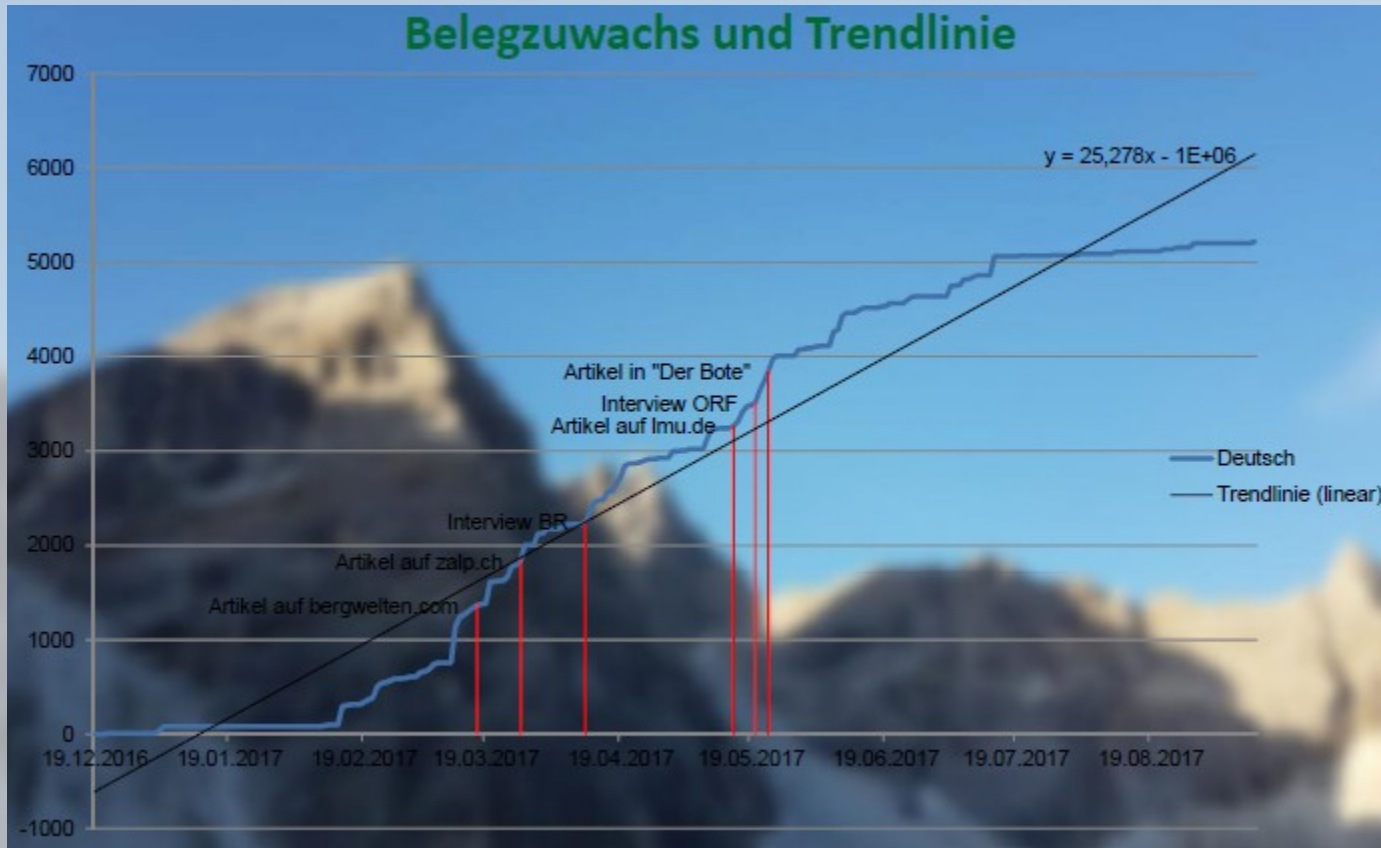
Private Kontakte

Crowdsourcing Kontakte (Februar – Mai 2017)



Crowdsourcing Mailingaktion (Februar – Mai 2017) Einwohnerverhältnisse und Belegrücklauf

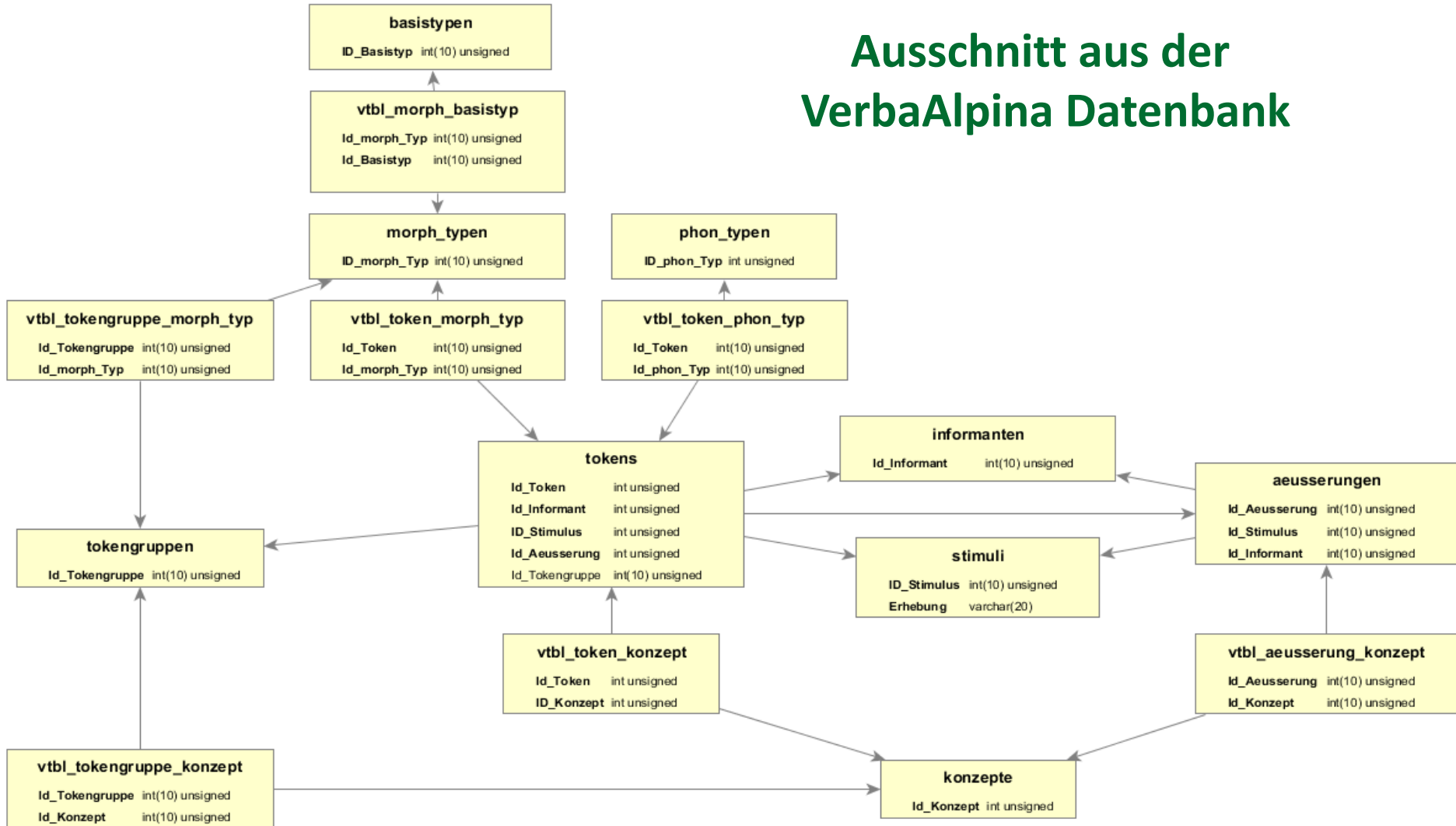




Eingliederung in den Datenbestand



Ausschnitt aus der VerbaAlpina Datenbank



Die VerbaAlpina Datenbank

Normalisierte Datenbank-Struktur

Vorteil

- Größtmögliche Vermeidung von inkonsistenten Daten





Die VerbaAlpina Datenbank

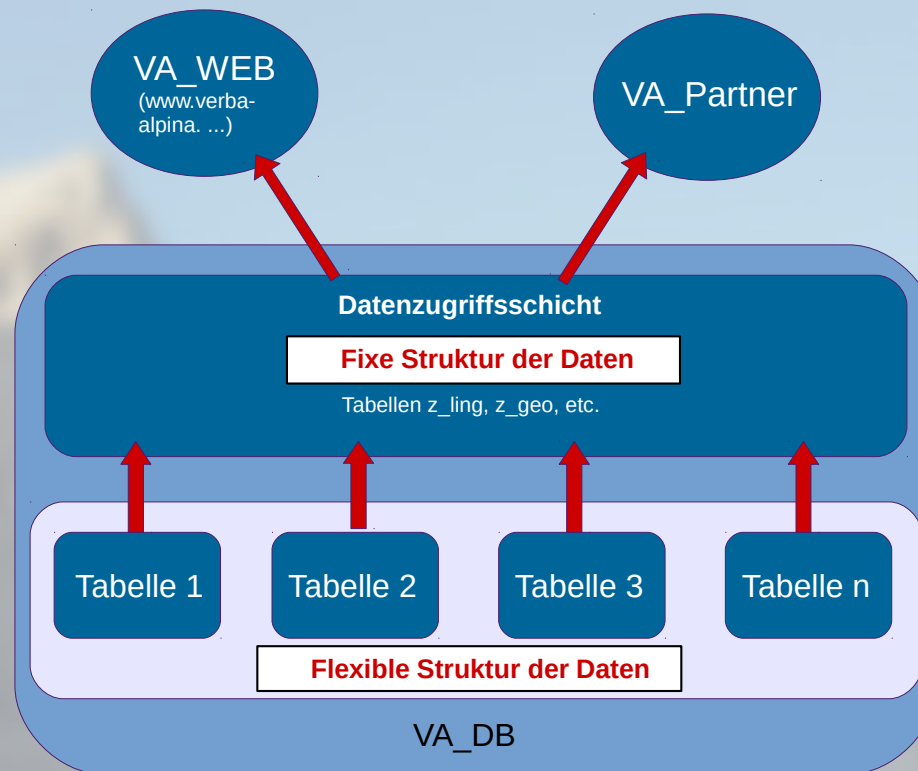
Normalisierte Datenbank-Struktur

Nachteile

- Komplexe SQL-Anfragen
- Langsame Datenbankabfragen

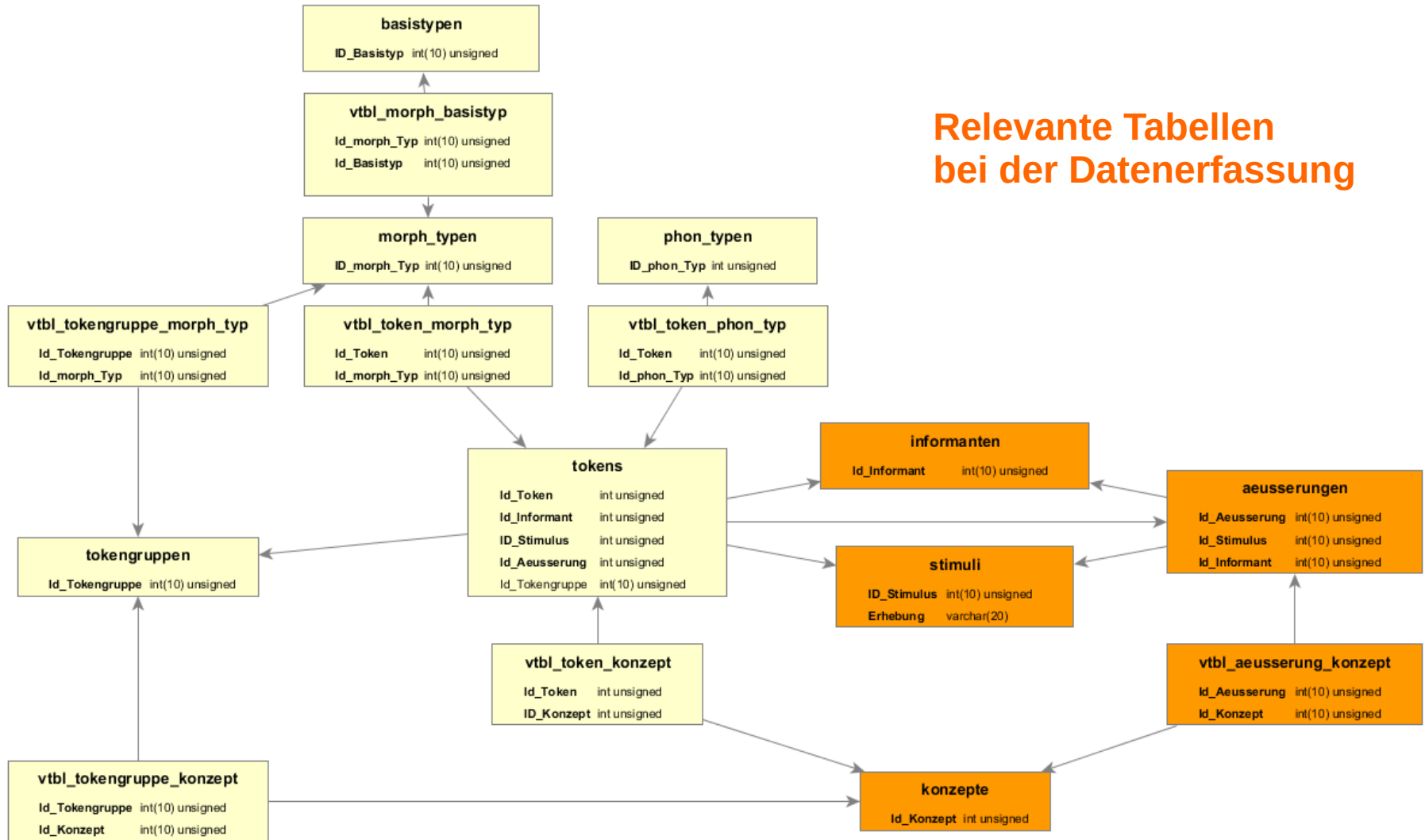
=> Datenzugriffsschicht

Datenzugriffsschicht



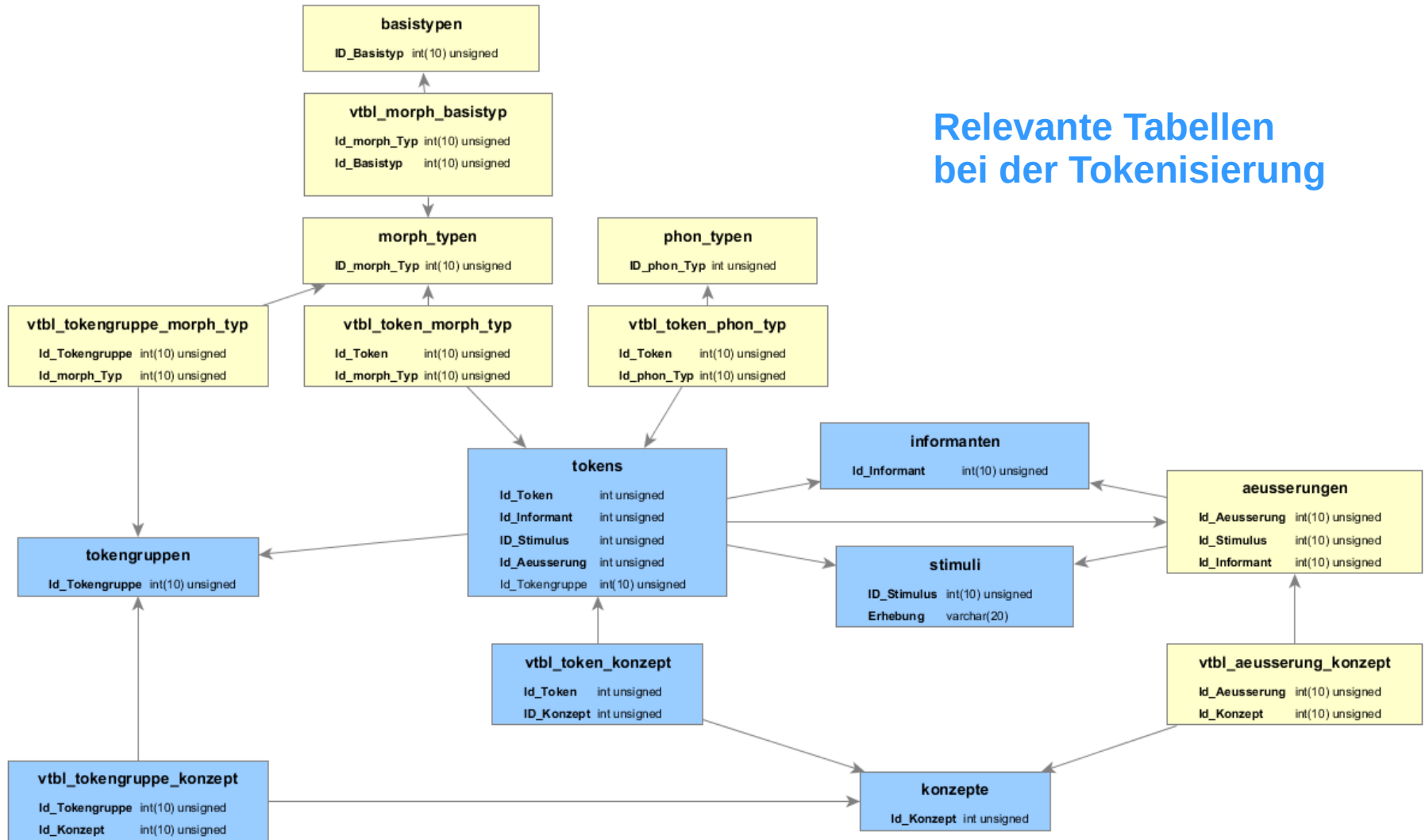


Relevante Tabellen bei der Datenerfassung





Relevante Tabellen bei der Tokenisierung





Tokenisierung

Äusserungen

Informant	Stimulus	Äusserung	Konzept
anonymousCrowder_58	CROWD_1_1	Milch ohseichn	MILCH SEIHEN



Tokens

Informant	Stimulus	Token	Konzept	Morph. Typ
anonymousCrowder_58	CROWD_1_1	Milch	MILCH	Milch
anonymousCrowder_58	CROWD_1_1	ohseichn	ABSEIHEN	Abseihen

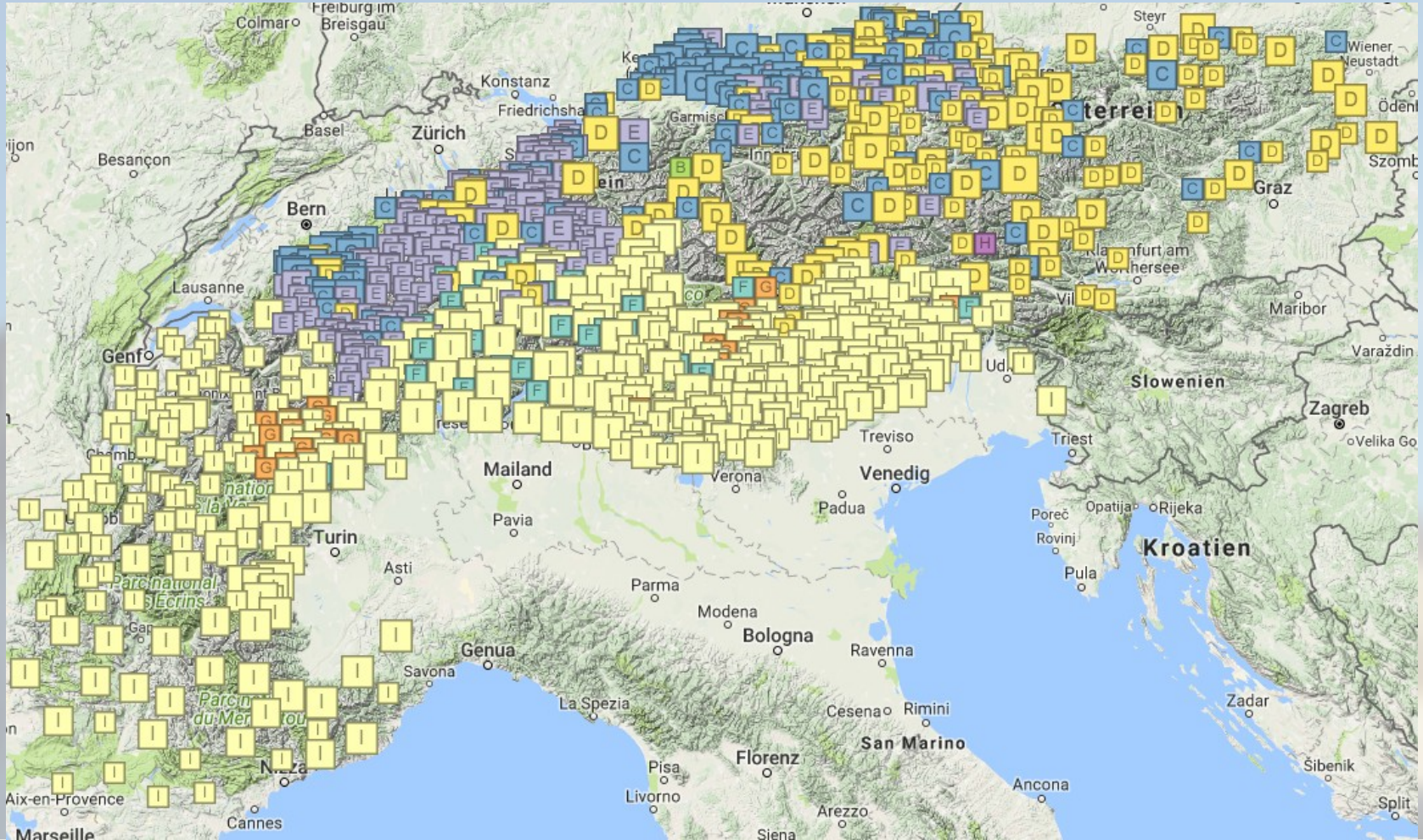
Tokengruppen

Informant	Stimulus	Tokengruppe	Konzept	Morph. Typ
anonymousCrowder_58	CROWD_1_1	Milch ohseichn	MILCH SEIHEN	Milch abseihen

Visualisierung



Qualitative Karte



Onomasiologisch

vs.

Semasiologisch

<input type="radio"/>	Konzept QUARK	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ Schotte(n) (ger.) (788 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ Topfen (ger.) (339 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ Ziger (ger. m.) (92 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ Spröss (ger.) (65 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ Ziger (ger.) (58 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ Press (ger.) (50 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ magnóca (rom. f.) (7 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ Topfenkäse (ger.) (13 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ Topfen (ger. m.) (7 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Morpho-lexikalischer Typ magnuc (rom. m.) (5 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

<input type="radio"/>	Morpho-lexikalischer Typ Schotte(n) (ger.)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Konzept QUARK (788 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Konzept MOLKE (60 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Konzept EIWEISSTEILCHEN, DIE NACH DER ZWEITEN SCHEIDUNG BEI ERHITZEN DER MOLKE AUFSTEIGEN (6 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Konzept FLÜSSIGKEIT NACH ENTNAHME DER KÄSEMASSE, ERSTE SCHEIDUNG (2 Belege)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="radio"/>	Konzept FLÜSSIGKEIT NACH ENTNAHME DES QUARKS (1 Beleg)	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

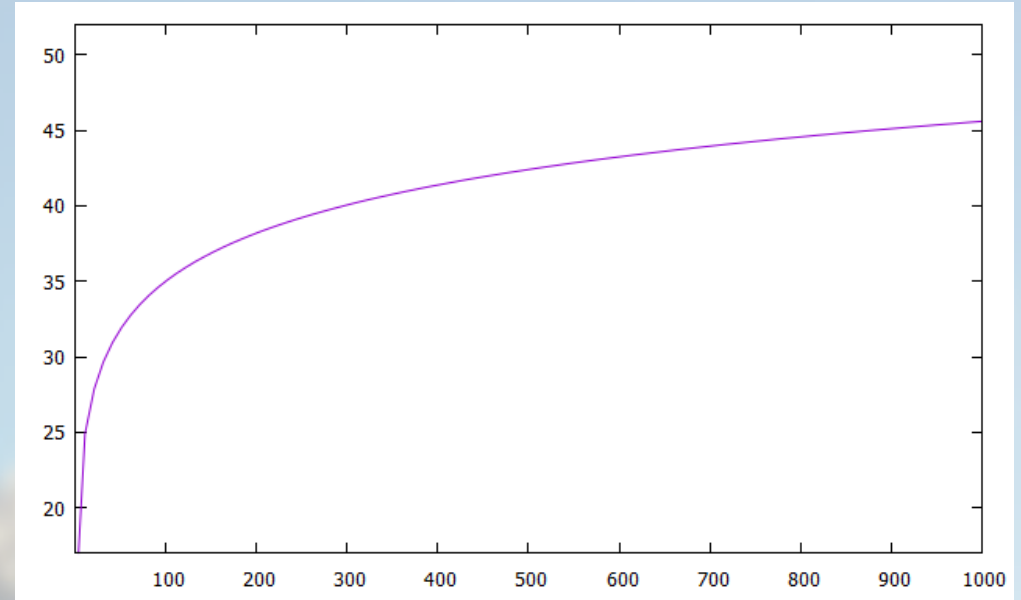


Symbole

Prinzipiell: Ein Symbol pro Beleg

Referenzgröße: Gemeinde
=> gleichartige Symbole innerhalb einer Gemeinde werden zusammengefasst

Logarithmisches Wachstum bei identischen Symbolen



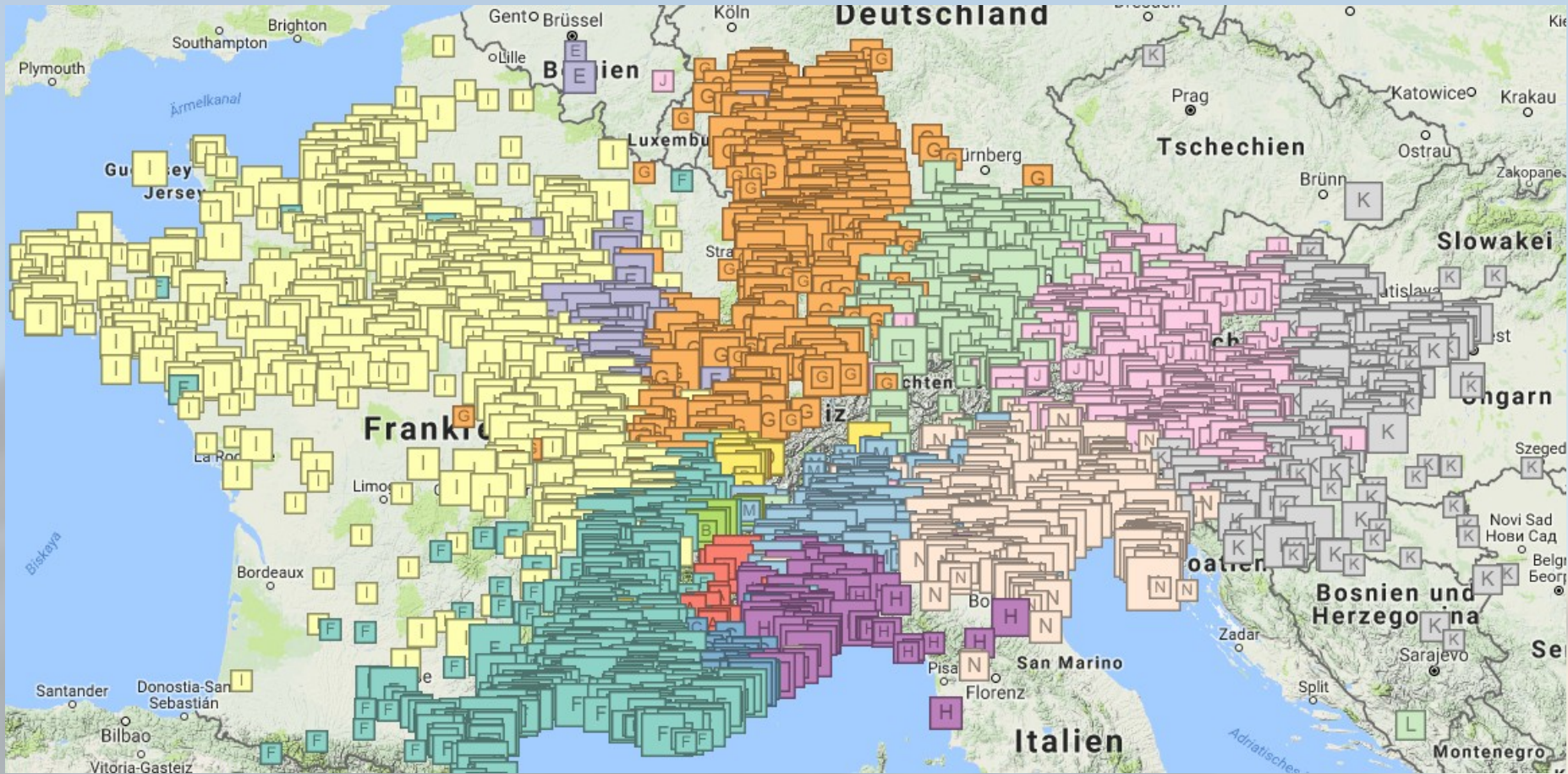
Anzahl Symbole	1	2	3	4	5	6	7	8	9	10	20	30	40	50	100	200	500	1000	
Größe																			

Kombinierte Symbole





Große Anzahl an Belegen möglich



Karte mit 92890 Belegen



Belegfenster

Innsbruck

topfn
(Einzelbeleg)

Phonetischer Typ	(nicht typisiert)	VA
Morpho-lexikalischer Typ	Topfen (ger.) D D K	VA

Quelle	Konzept
TSA III_99_1, D35_3 (Hötting)	QUARK

Topfen
(Einzelbeleg)

Phonetischer Typ	(nicht typisiert)	VA
Morpho-lexikalischer Typ	Topfen (ger.) D D K	VA

Quelle	Konzept
CROWD 1_1, anonymousCrowder_200 (Innsbruck)	QUARK

Bardonecchia

mnt' aja*
(Einzelbeleg)

Phonetischer Typ	(nicht typisiert)	VA
Morpho-lexikalischer Typ	montagne / montagna (rom. f.) C T L	VA
Basistyp	*montania	VA

Quelle	Konzept
AIS 1192_1, 140 (Rochemolles)	BERG

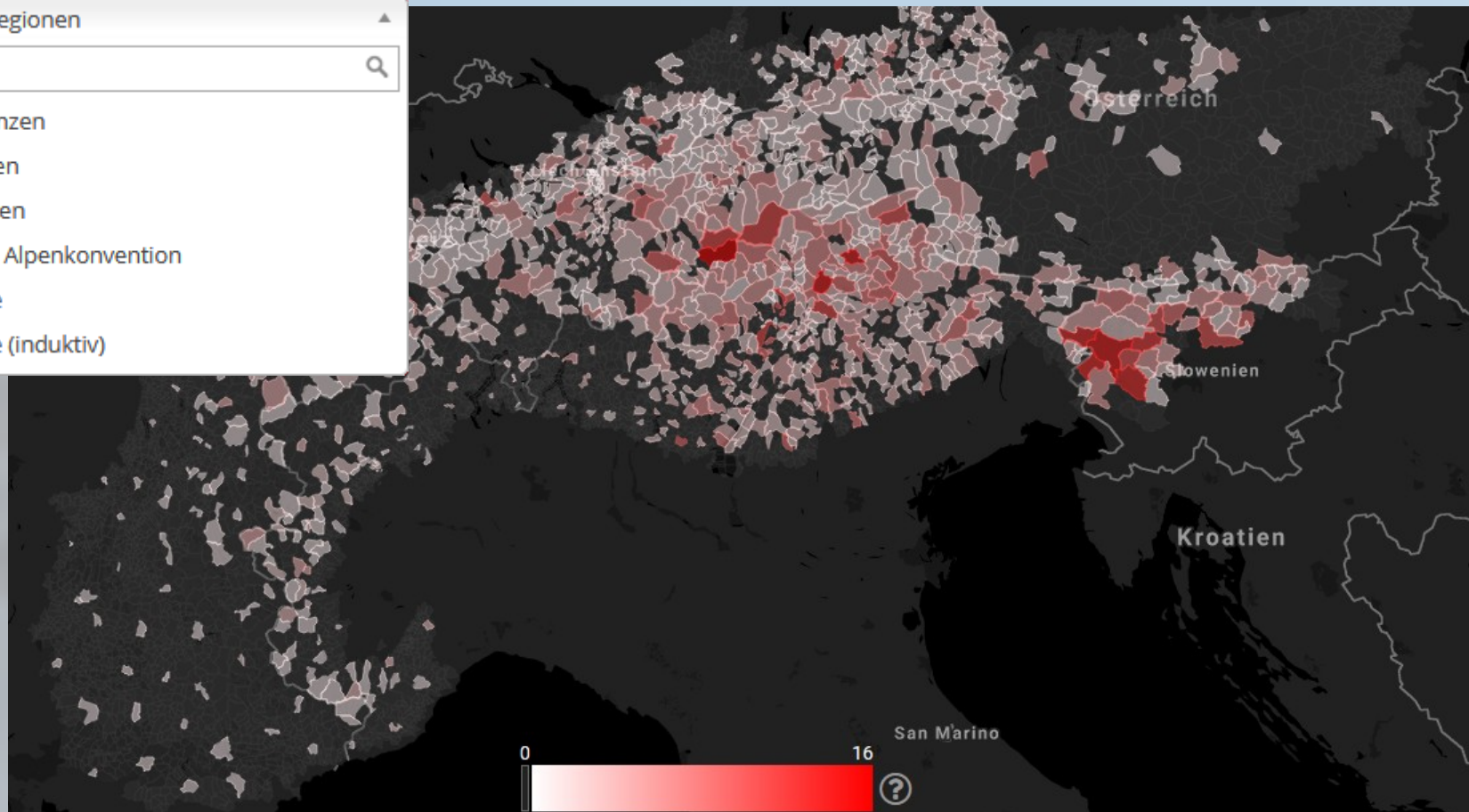
* belegt als Teil von me:z' un d mnt' aja

Quantitative Darstellung

Abbildung von Punktdaten auf Flächen

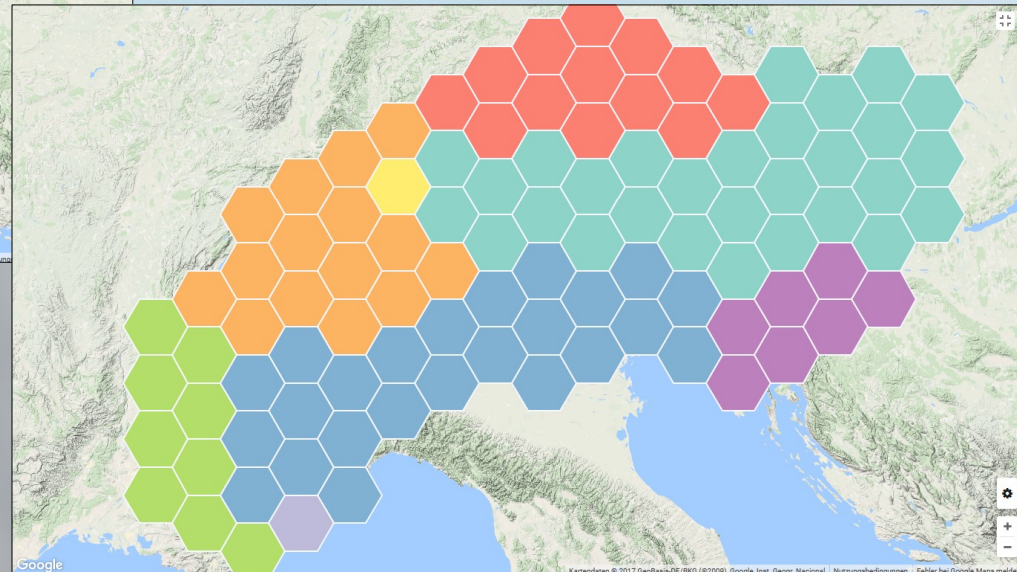
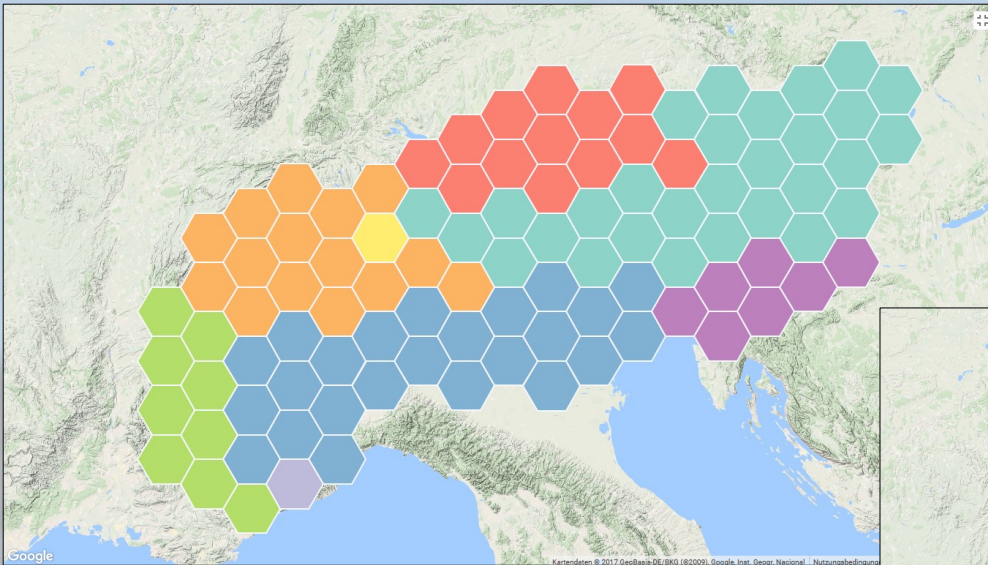
Flächen und Regionen

Gemeindegrenzen
Nationalstaaten
NUTS-3 Grenzen
Perimeter der Alpenkonvention
Sprachgebiete
Sprachgebiete (induktiv)



Abstrahierte Karte

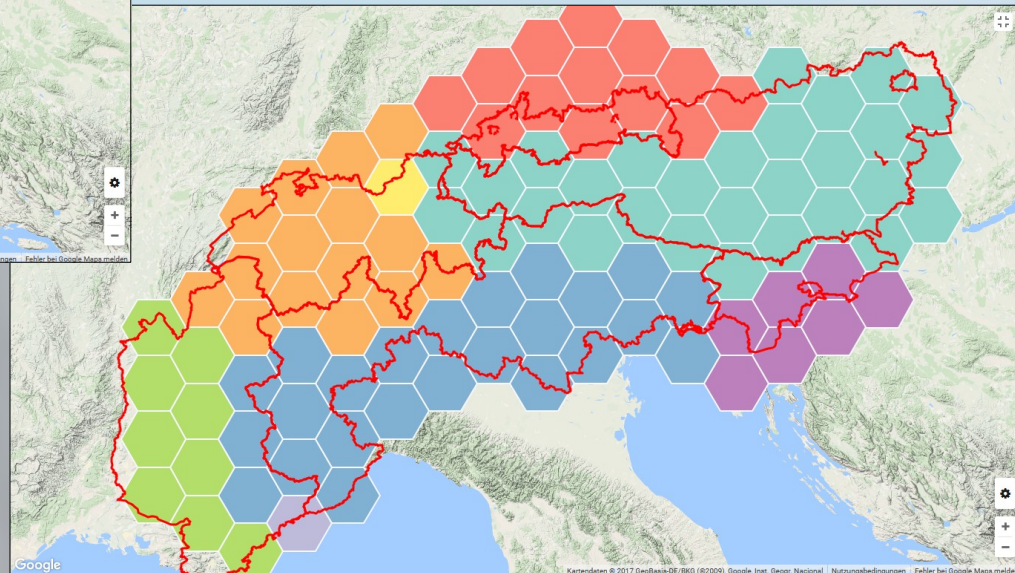
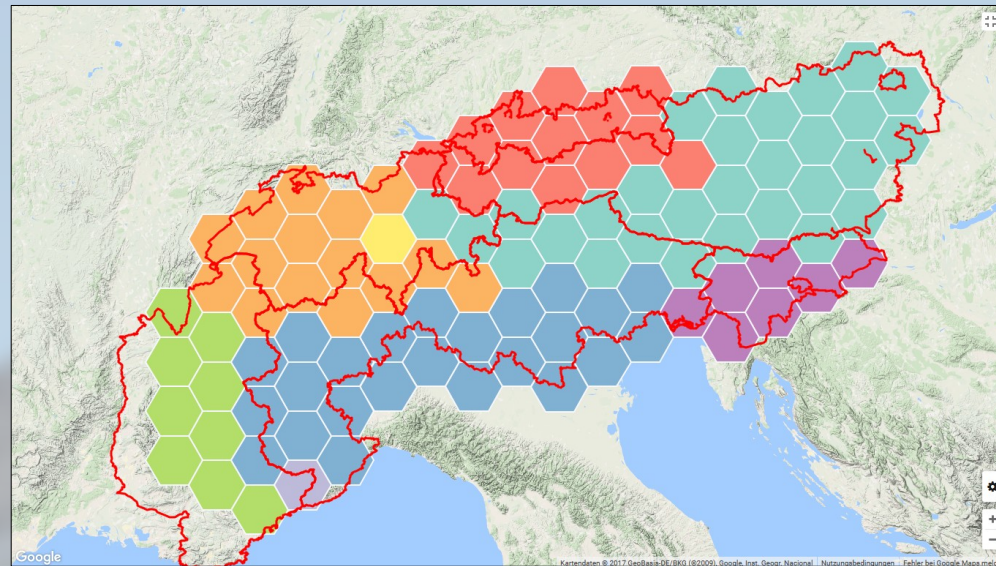
NUTS-3 Grenzen Schematisch



NUTS-3 Grenzen Geographisch

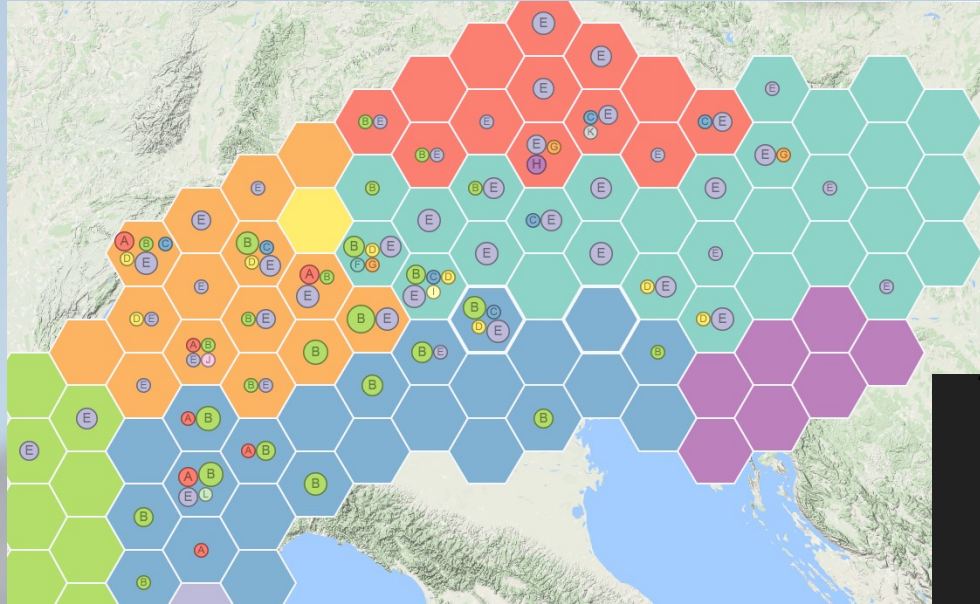
Abstrahierte Karte

NUTS-3 Grenzen Schematisch



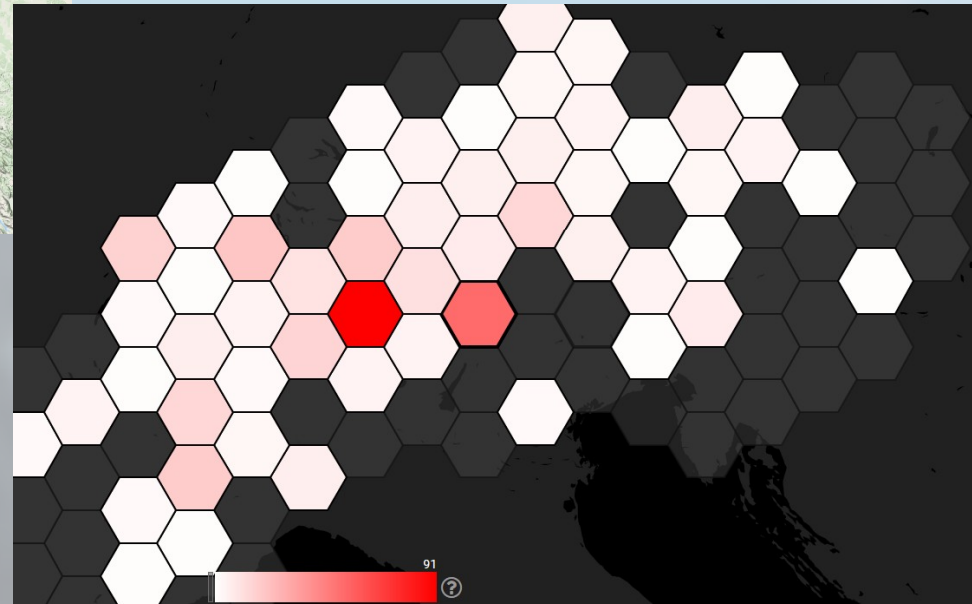
NUTS-3 Grenzen Geographisch

Abstrahierte Karte



Qualitativ

Quantitativ



91 ?

Nachhaltigkeit

Versionierung

Vollständiger Datenbestand wird regelmäßig versioniert (ca. halbjährlich)

Grund: Zitierbarkeit

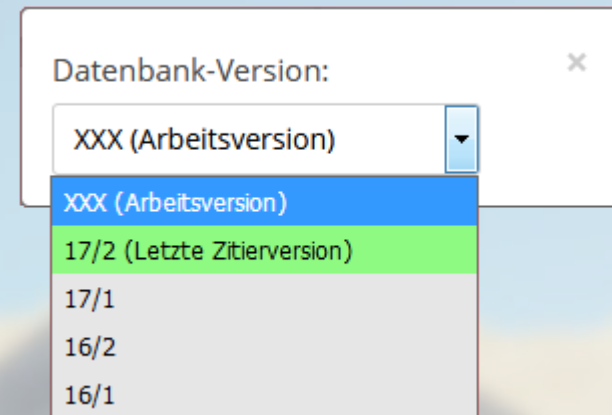
Technische Umsetzung: Kopie der vollständigen Datenbank

Vorteile

- Schneller Zugriff
- Geringe Fehleranfälligkeit

Nachteile

- Höherer Speicherbedarf
- Nachträgliche Strukturanpassungen





Versionierung

Außerdem Bereitstellung des vollständigen Programmcodes unter

<https://github.com/VerbaAlpina>

Plugin Interactive Map

Plugin VerbaAlpina

Plugin CS-Tool

Theme VerbaAlpina



Langfristige Nutzung

Probleme

Nicht rückwärtskompatible Software

Lösungsansatz: Zeitlich unbegrenzter Zugriff über einen Docker auf dem Server der Universitätsbibliothek: Eingefrorene Konfiguration ab Förderungsende (Prototyp Ende 2016)

Abhängigkeit von externen Diensten für die Kartierung

Lösungsansatz: Eigenes Hosten des Kartenmaterials über das Leibniz-Rechenzentrum (aktuell in Arbeit)



Archivierung der Datenbasis

Archivierung bei verschiedenen Institutionen

- IT-Gruppe Geisteswissenschaften
- Clarin-D
- Evtl. Generic Research Data Infrastructure (GeRDI)