

Language Documentation via Crowdsourcing

Birgit Alber

in cooperation with Joachim Kokkelmans and Emily Siviero

birgit.alber@unibz.it

Free University of Bozen/Bolzano

Workshop on Language Diversity, Data Elicitation and Multilingualism

Leibniz-ZAS, Berlin, 26 May 2023

Overview

- New ways of studying linguistic diversity: crowdsourcing
- The project *VinKo-AlpiLink*
 - Features of VinKo-AlpiLink (lights and shadows)
 - Quality of the data
 - Quantity of the data
 - Citizen science projects (*Vinkiamo*)

Crowdsourcing

3

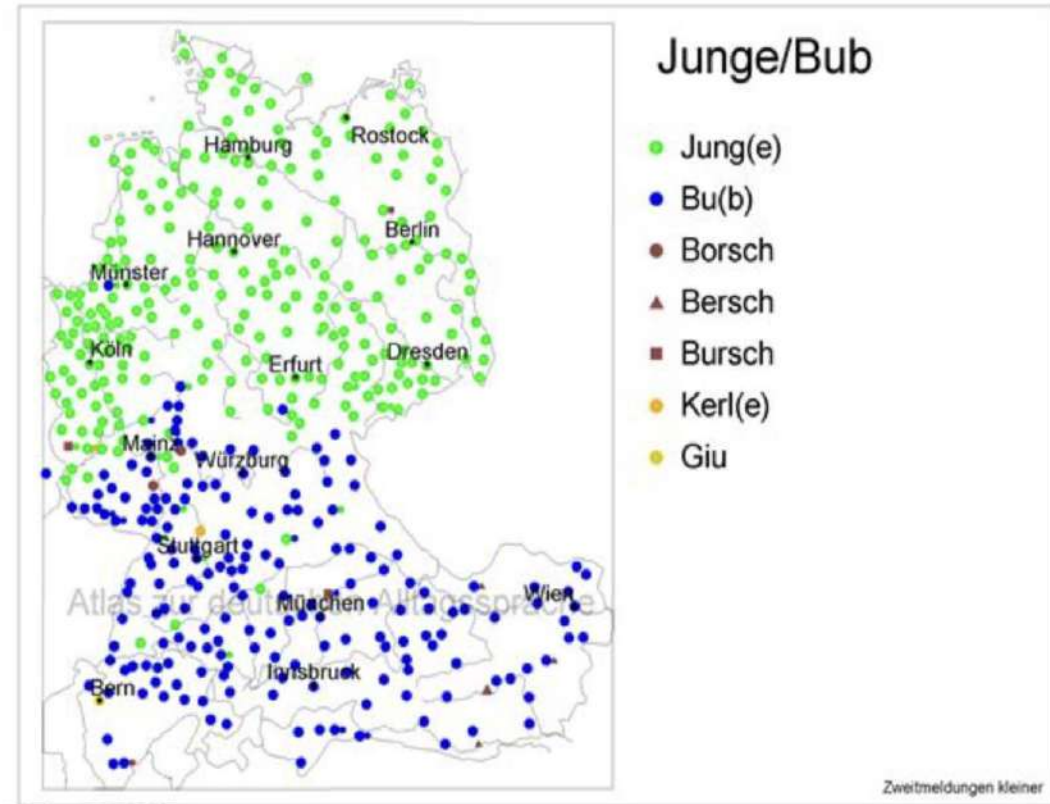
A new way of studying language diversity: crowdsourcing

AdA: Atlas zur deutschen
Alltagssprache

<http://www.atlas-alltagssprache.de/>

- *Fragerunden* (question rounds)
- Instructions: 'Bitte geben Sie bei den folgenden Fragen jeweils an, welchen Ausdruck man in Ihrer Stadt normalerweise hören würde – egal, ob es mehr Mundart oder Hochdeutsch ist. Antworten Sie bitte, ohne lange nachzudenken!'

Junge/Bub



Junge/Bub (Frage 1)

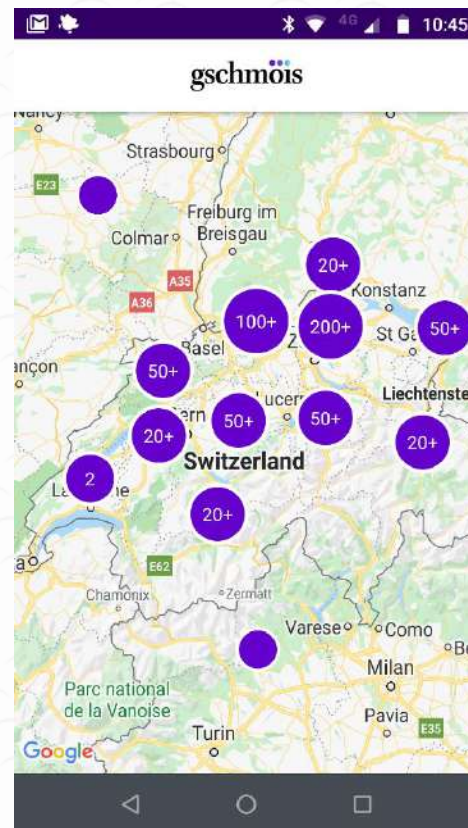
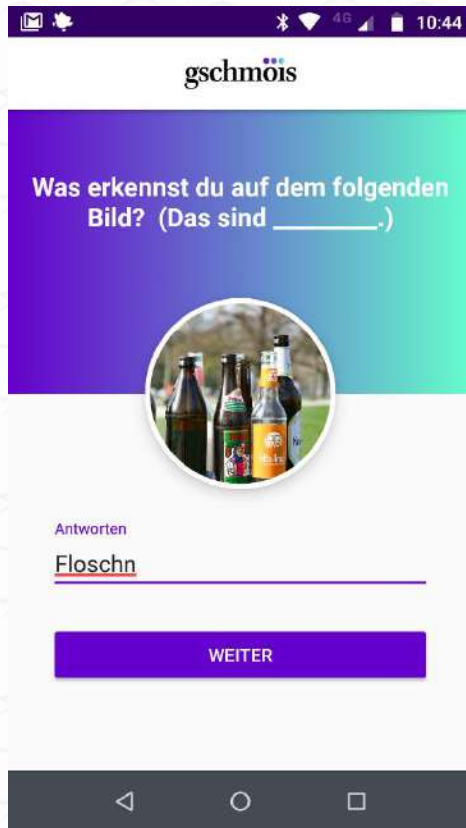
Zweitmeldungen Kleiner

Crowdsourcing

4

A new way of studying language diversity: crowdsourcing

Gschmöis: Swiss Dialects <https://www.gschmois.uzh.ch/de.html>



Typically

- focus on the lexicon
- written responses

Crowdsourcing

5

Advantages

- Informants: many
- Data: a lot, geographically fine-grained
- Restitution: high potential of involvement of informants (*citizen science*); immediate restitution in form of maps or similar

But: Accompanied by a certain loss of control by the linguist/fieldworker →

Crowdsourcing

6

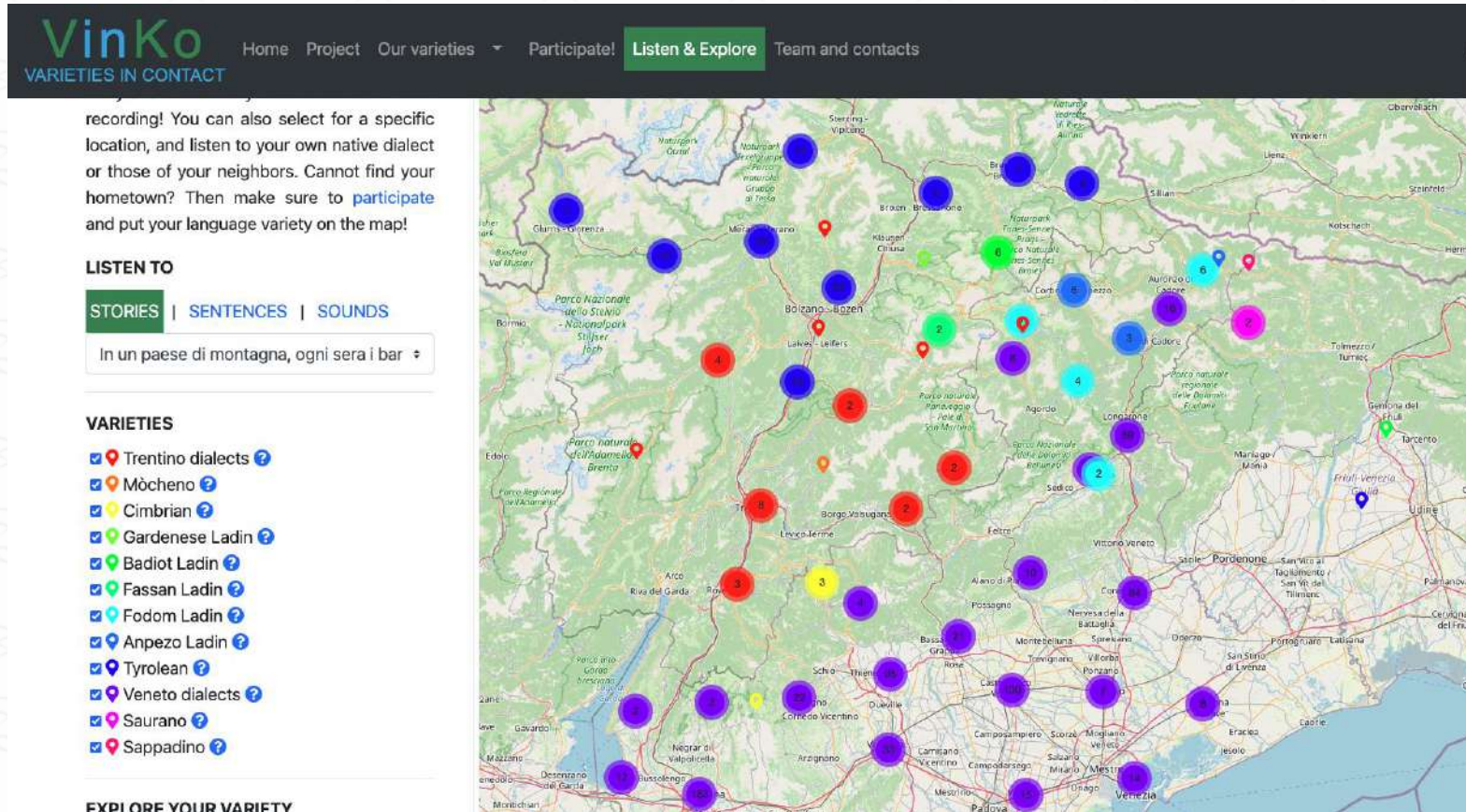
Critical points

- Which informants?
 - reliable speakers of the variety we want to elicit?
- Which data?
 - dialect, regional languages, regional standard?
 - form of the elicited data
 - written: can be used only for certain phenomena
 - audio data: technically challenging, may be noisy, has to be processed
 - impossibility of checking back with the informants
 - phenomena which can be studied: lexicon, phonology, morphology
syntax?

VinKo - AlpiLink

Main Features (<https://vinko.it/>) (Cordin et al. 2018, Rabanus et al. 2022, Kruijt et al. b. in press)

- **Crosslinguistic data elicitation:** Germanic and Romance varieties of Northern Italy



VinKo VARIETIES IN CONTACT

Home Project Our varieties Participate! **Listen & Explore** Team and contacts

recording! You can also select for a specific location, and listen to your own native dialect or those of your neighbors. Cannot find your hometown? Then make sure to [participate](#) and put your language variety on the map!

LISTEN TO

STORIES | SENTENCES | SOUNDS

In un paese di montagna, ogni sera i bar ▾

VARIETIES

- Trentino dialects
- Mòcheno
- Cimbrian
- Gardesane Ladin
- Badiot Ladin
- Fassan Ladin
- Fodom Ladin
- Anpezo Ladin
- Tyrolean
- Veneto dialects
- Saurano
- Sappadino

EXPLORE YOUR VARIETY

VinKo - AlpiLink



8

- **Crosslinguistic data elicitation:** Germanic and Romance varieties of Northern Italy

| Germanic | Romance |
|-------------------------------|-------------------------------------|
| Tyrolean dialects | Trentino dialects |
| Mòcheno (lg. island) | Veneto dialects |
| Cimbrian (lg. island) | Ladin (Badia, Fassa, Fodom, Anpezo) |
| Saurano (lg. island) | <i>Piedmontese dialects</i> |
| <i>Sappadino (lg. island)</i> | <i>Franco-Provençal</i> |
| <i>Walser (lg. island)</i> | <i>Occitan</i> |

Project partners:

Universities of

Verona: Stefan Rabanus (coordinator), Anne Kruijt

Bozen/Bolzano: Birgit Alber, Joachim Kokkelmans

Trento: Ermengildo Bidese

Torino: Livio Gaeta

Valle d'Aosta: Gianmario Raimondi

VinKo - AlpiLink



9

- **Crosslinguistic data elicitation:** Germanic and Romance varieties of Northern Italy

Same phenomena in different varieties:

- detection of contact-induced change
- possibility of typological comparisons

Example: Laryngeal features of stops

/p/ vs. /b/

| | | | |
|-----------|--------------|--------------|-----------------------|
| Tyrolean: | <i>Paarl</i> | <i>Baam</i> | 'type of bread, tree' |
| Veneto: | <i>pan</i> | <i>borsa</i> | 'bread, bag' |

VinKo - AlpiLink



10

Example: Laryngeal features of stops

- Vietti, Alber & Vogt (2018): acoustic analysis, Tyrolean dialect of Meran
 - 10 informants from Meran
 - laboratory setting (soundproof booth)
 - no aspiration in /p, t, k/
[differently from more Northern German dialects]
 - categorical neutralization of labials /p ~ b/, word-initially
 - intraspeaker variability in neutralizing /k ~ g/ and /t ~ d/, word-initially

- (1) a. [p]aarl
b. [p]aam

'Paarlbrod, type of bread'
'Baum, tree'

VinKo - AlpiLink



11

Example: Laryngeal features of stops

Qs

- can these results be replicated with a wider range of informants?
- is the contrast in terms of VOICE (not spread glottis) similar to that of the neighboring Romance languages?

These are questions we hope to answer with the help of crowdsourced data of the VinKo-Alplilink project

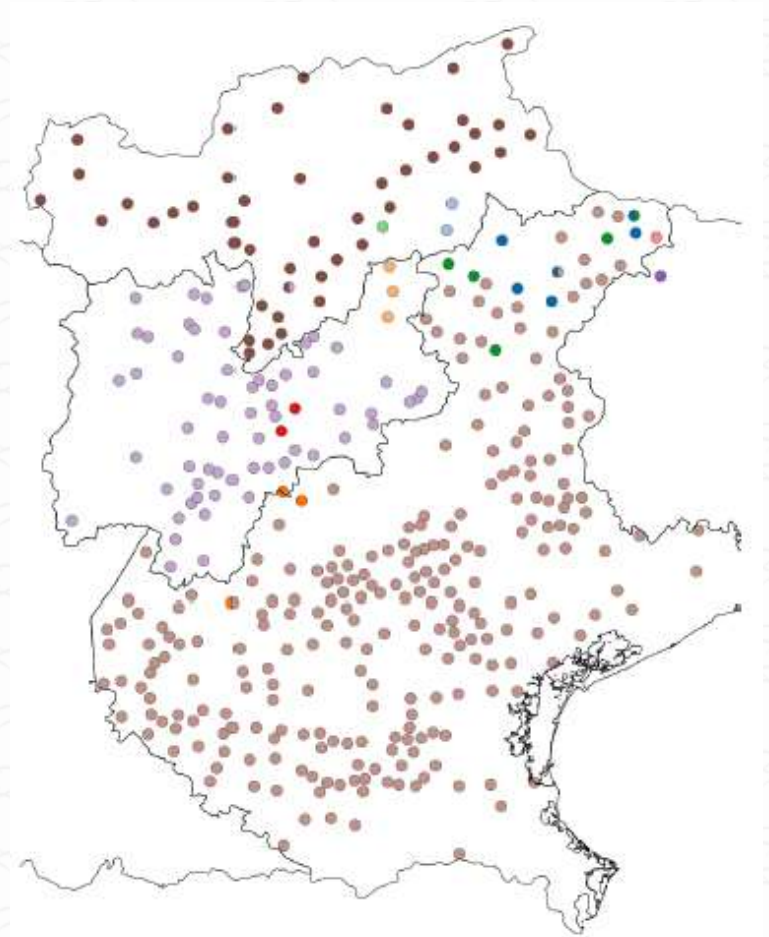
VinKo - AlpiLink



12

- **Data over five-year period (2018-2022) (Kruijt 2023, Kruijt et al., in press b)**
- 1392 questionnaires
- each questionnaire: 152 to 181 audio recordings.
- more than 125.000 audio files
- 377 locations
- corpus accessible at
Clarín repository, EURAC

<http://hdl.handle.net/20.500.12124/46>



VinKo - AlpiLink

13

- **Data over five-year period (2018-2022) (Kruijt 2023)**

data not necessarily well balanced: Tyrolean

| | <i>Average age</i> | <i>Gender M/F</i> | | <i>Proficient</i> | <i>Frequent use</i> | <i>Use in family</i> | <i>Use with friends</i> |
|------------------|--------------------|-------------------|------|-------------------|---------------------|----------------------|-------------------------|
| <i>Anpezan</i> | 47,5 | 42% | 58% | 92% | 58% | 67% | 73% |
| <i>Badiot</i> | 34,7 | 14% | 86% | 100% | 86% | 100% | 100% |
| <i>Cimbrian</i> | 56,3 | 71% | 29% | 86% | 57% | 43% | 83% |
| <i>Fassan</i> | 23,3 | 25% | 75% | 100% | 100% | 100% | 100% |
| <i>Fodom</i> | 64,3 | 43% | 57% | 100% | 90% | 95% | 90% |
| <i>Gardenese</i> | 39,5 | 0% | 100% | 100% | 100% | 100% | 100% |
| <i>Mòcheno</i> | 39,1 | 43% | 57% | 100% | 100% | 100% | 100% |
| <i>Sappadino</i> | 50,0 | 0% | 100% | 100% | 100% | 100% | 0% |
| <i>Saurano</i> | 67,5 | 100% | 0% | 50% | 50% | 50% | 50% |
| <i>Trentino</i> | 39,3 | 56% | 44% | 89% | 84% | 86% | 64% |
| <i>Tyrolean</i> | 26,0 | 10% | 90% | 99% | 95% | 97% | 99% |
| <i>Venetan</i> | 48,2 | 45% | 55% | 82% | 72% | 79% | 76% |
| <i>Overall</i> | 45,3 | 43% | 57% | 85% | 76% | 82% | 78% |

Table 5: Speaker information per language variety (17 January 2023)

VinKo - AlpiLink

14

- **Phenomena investigated: (synchronic) structure**

Examples:

| | |
|------------|----------------------------------|
| Phonology | obstruent system |
| | processes involving sibilants |
| | /r/ |
| Morphology | pronoun and determiner paradigms |
| | <i>truncation patterns</i> |
| Syntax | pro-drop |
| | complementizers/subordination |
| | particle verbs |

VinKo - AlpiLink

15

- **Form of elicitation**

Phonology: word-list reading; separate questionnaires for each variety (with glosses)

Advantages: low impact of standard varieties

Critical points: ad hoc orthographies; sometimes items are not recognized

B71 ⏴ | ✕ ✓ fx Wies

| | A | B | C | D | E | F |
|----|----|----------|------------|-------|----|------|
| 33 | 32 | Paarl | Paarlbrot | obstr | wi | /p/ |
| 34 | 33 | Baam | Baum | obstr | wi | |
| 35 | 34 | Toog | Tag | obstr | wi | /t/ |
| 36 | 35 | Dâch | Dach | obstr | wi | /d/ |
| 37 | 36 | Kârtn | Karten | obstr | wi | /kx/ |
| 38 | 37 | Ggâggele | Ei | obstr | wi | /k/ |
| 39 | 38 | Ggugger | Fernglas | obstr | wi | /k/ |
| 40 | 39 | Gârtn | Garten | obstr | wi | /g/ |
| 41 | 40 | oper | schneefrei | obstr | wm | /p/ |
| 42 | 41 | lânnet | dumm | obstr | wm | /n/ |

⏴ ▶
Ampezzano
Badiot
Cimbrian (Lusern)
Fassan
Fodom
Gardenese
Mòcheno
Trentino
Tyrolean
Veneto

VinKo - AlpiLink

16

- **Form of elicitation**

Morphology: personal pronoun paradigms and articles (Kruijt 2022)

- translation task aided by pictures embedded in a story
- guided free production
- text in standard language; prompting through audio examples in a Germanic/Romance dialect

Die Hexe will die gefangenen Kinder kochen

La strega vuole cucinare i bambini presi



VinKo - AlpiLink



17

- **Form of elicitation**

Syntax: translation task from standard language (German or Italian)

Subordination:

- a. Ich weiß nicht, welcher Bus zuerst abfährt
- b. Non so quale autobus parta per primo
'I don't know which bus leaves first'

Critical points

- influence from the standard languages?
- which standard language to choose as a trigger for varieties such as Mòcheno or the Ladin varieties?
- perceived as the most complex and boring task by informants (Kruijt 2023)

VinKo - AlpiLink



18

General critical points (Kruijt 2023)

- length of the questionnaire (70% of feedback to VinKo experience by informants)
- complexity of the task (reading a word list is easier than translating)
- phonological questionnaire: proposed words don't match words in the local dialect
- tasks too repetitive: stories might be interesting than the traditional tasks

VinKo – Quality of the Data

19

Is crowdsourced data reliable?

- Kruijt (2022): Subject clitics in Veneto dialects described in the literature are present in VinKo (but not in the standard trigger sentences)

Input:

Tornano di corsa in paese e chiedono aiuto al cacciatore

'[They] return immediately to the village and [they] ask the hunter for help'

Veneto dialect from Bonavigo (Kruijt 2022: 48):

I= torna de corsa in paese e i= dimanda aiuto a=t
 3PL.M= return in run in village and 3PL.M= ask help to-DEF.M.SG
cacciatore.
 hunter(M)

More examples (morphosyntax) for the reliability of crowdsourced data: Rabanus (in press); Kruijt et al. (in press a.)

VinKo – Quality of the Data

20

Is crowdsourced data reliable?

- Alber & Kokkelmans (2022): **rhotics** /r/ in Tyrolean dialects
- Comparison of three sources
 - Tirolischer Sprachatlas: 1965-1970
 - Scheutz (2016): data elicited in 2010
 - VinKo: data elicited from 2017-22
- Results
 - similar structures are found
 - where structures diverge, this can be explained by the diachronic change that rhotics are undergoing

VinKo – Quality of the Data

Is crowdsourced data reliable? (Alber & Kokkelmans 2022)

- Comparability of data: Scheutz (2016) – VinKo (2017-22)
 - age ('young speakers' in Scheutz – VinKo mean of 24 years)
 - social background (non-agriculture-related job – university students)
 - location (points within reasonable travel time = 50 mins by bike)
 - sex (female: 63% in Scheutz – 91% in VinKo)

VinKo – Quality of the Data

22

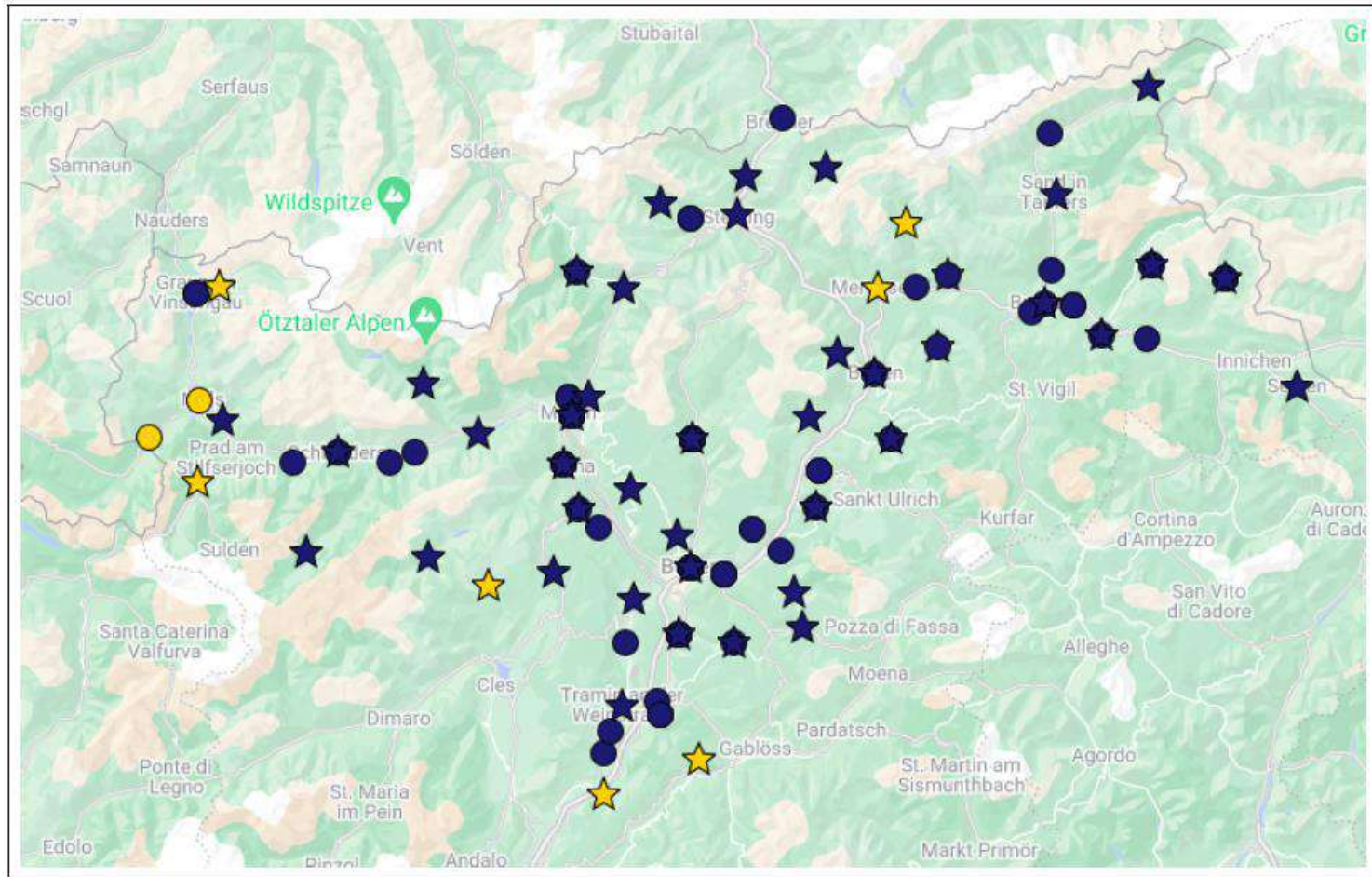
Is crowdsourced data reliable? (Alber & Kokkelmans 2022)

Rhotics in **Onset Position**: alveolar or uvular?

(2) *lea.[r]e.[r]in* vs. *lea.[R]e.[R]in* 'Lehrerin, teacher'

=> 84% match between Scheutz (2016) and VinKo (17-22)

VinKo – Quality of the Data



yellow = alveolar; blue = uvular; circles = *VinKo*; stars = *Insre Sproch*

VinKo – Quality of the Data

Is crowdsourced data reliable? (Alber & Kokkelmans 2022)

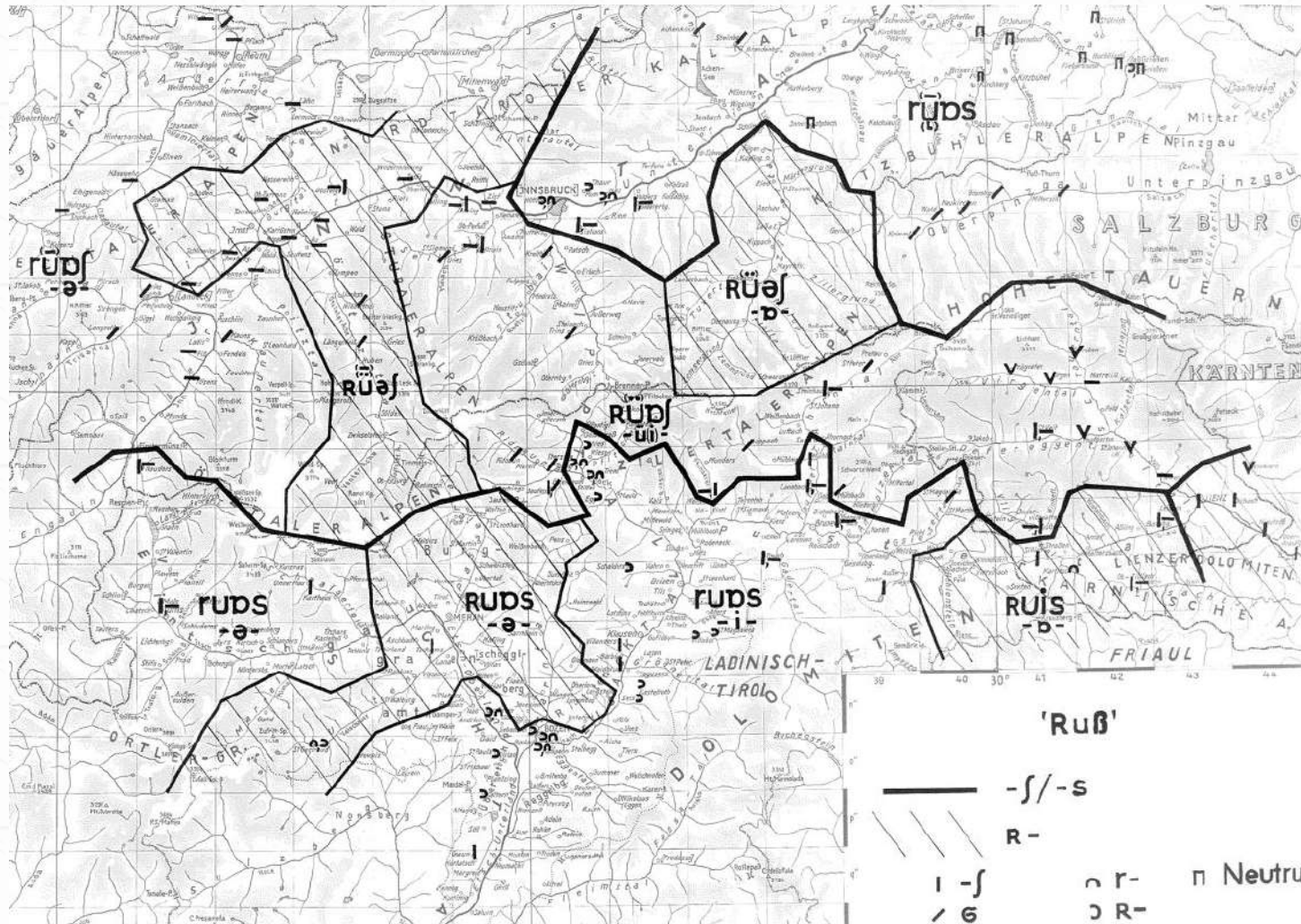
Interpretation

- 84% match for rhotics in **Onset position**: the diachronic shift from alveolar [r] to uvular [R] is next to complete

Compare with older sources: Tirolischer Sprachatlas/Kranzmayer (1956)

- alveolar [r] most widespread variant
- uvular [R] attested in urban contexts

VinKo – Quality of the Data



TSA

VinKo – Quality of the Data



26

Is crowdsourced data reliable? (Alber & Kokkelmans 2022)

Rhotics in Coda Position: consonantal or vocalized?

(3) Stressed syllables

a. after short vowel: *du[r]/[R]* vs. *du[ɐ]* 'dürr, meagre'

b. after long vowels: *joo[r]/[R]* vs. *joo[ɐ]* 'Jahr, year'

=> 68% and 66% match between Scheutz (2016) and VinKo (17-22)

VinKo – Quality of the Data

27

Is crowdsourced data reliable? (Alber & Kokkelmans 2022)

Rhotics in **Coda Position**: consonantal or vocalized?

(3) Unstressed syllables

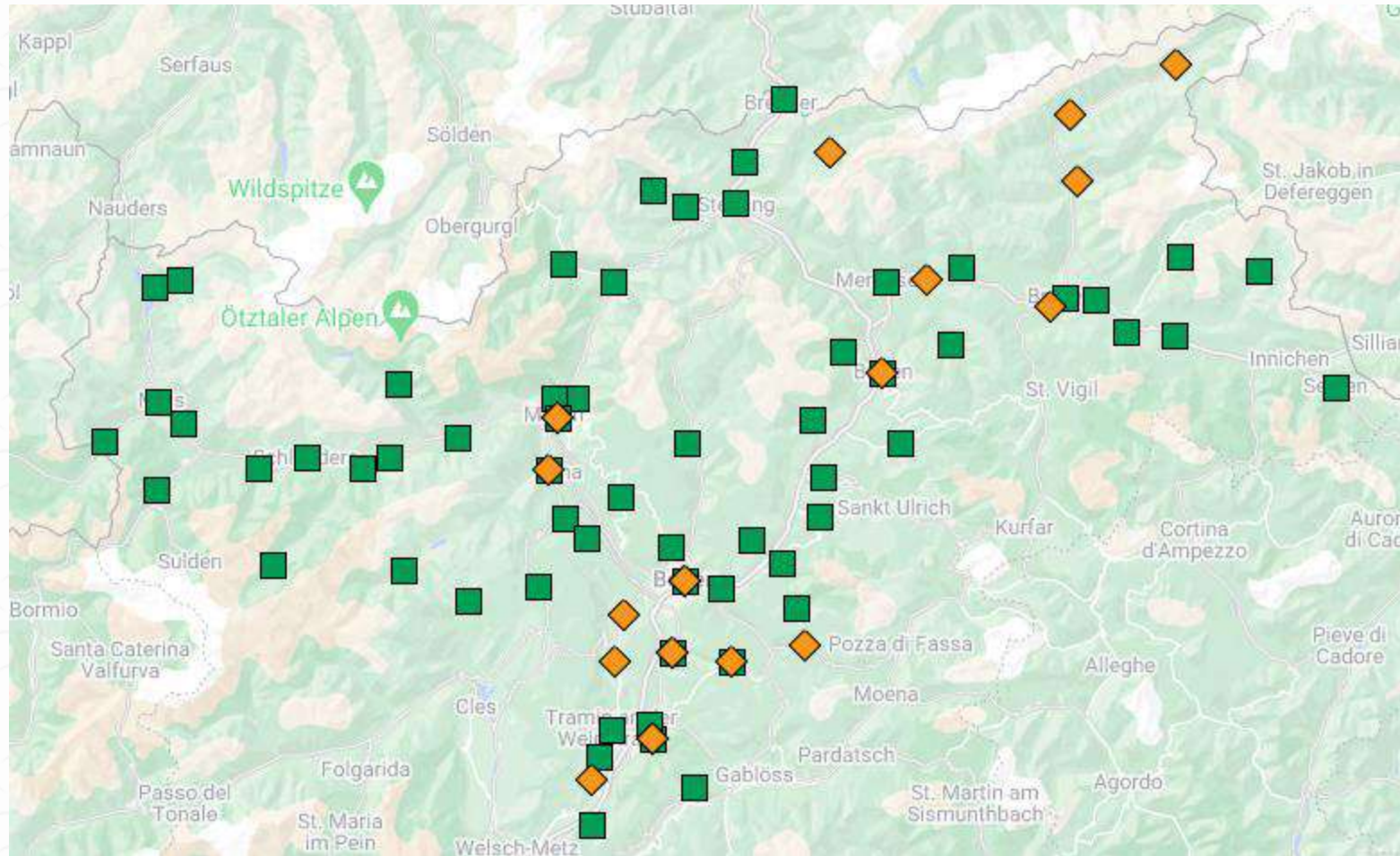
wosse[r]/[R] vs. *woss*[ɐ] 'Wasser, water'

=> 75% match between Scheutz (2016) and VinKo (17-22)

VinKo – Quality of the Data



28



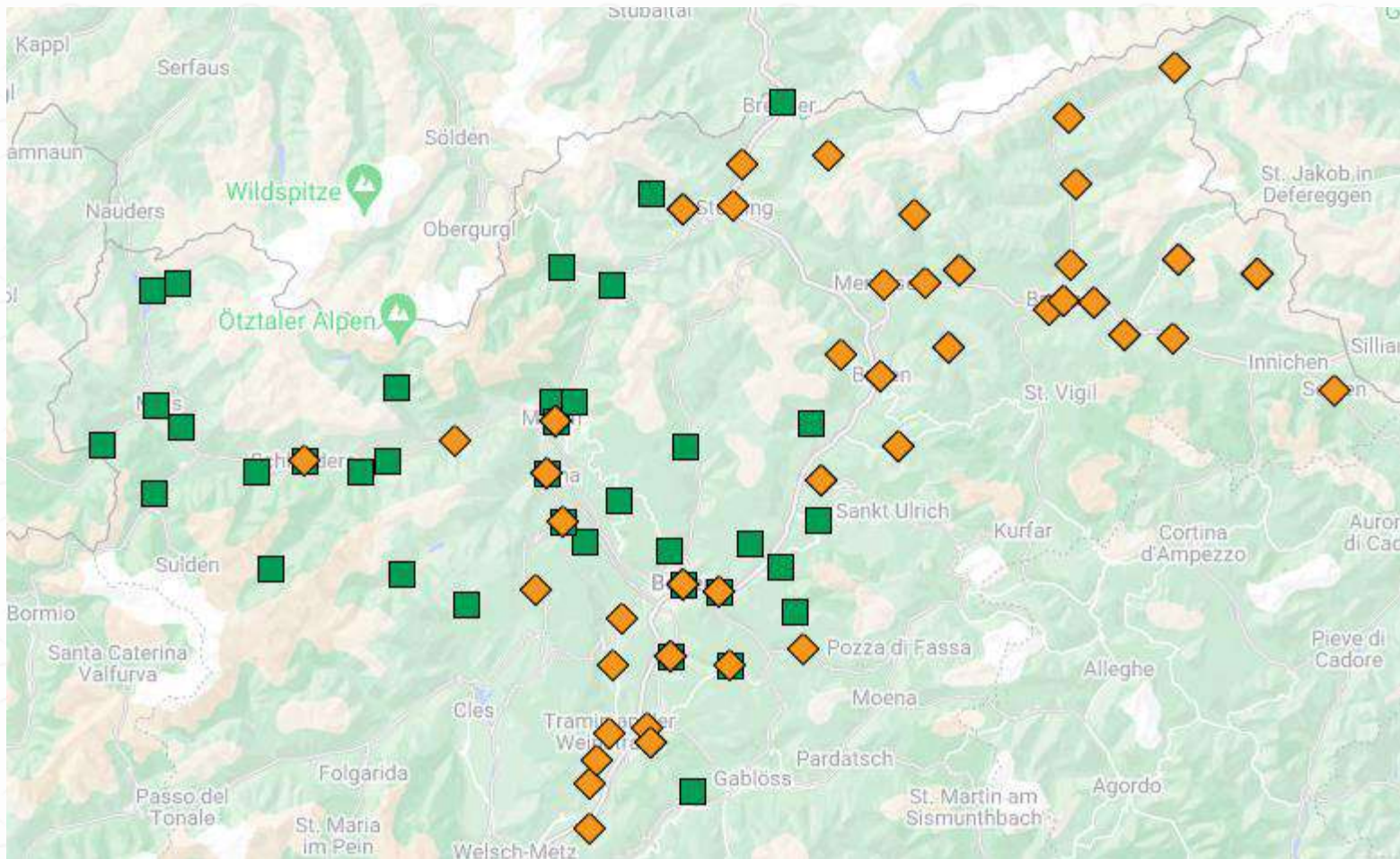
jaar

green:
consonantal

orange:
vocalized

Combined
data Scheutz
(2016) +
VinKo

VinKo – Quality of the Data



wosser

green:
consonantal

orange:
vocalized

combined
data Scheutz
(2016) +
VinKo

VinKo – Quality of the Data

30

Is crowdsourced data reliable? (Alber & Kokkelmans 2022)

Interpretation

- 66-68-75% match for rhotics in **Coda position**: the vocalisation process is still ongoing, varying according to contexts

durr < joor < wosser

Compare with older sources: Tirolischer Sprachatlas

- vocalisation only attested for the *wosser*-Kontext, and only in the East (Pustertal)

VinKo – Quality of the Data

32

Is crowdsourced data reliable?

YES!

Comparison with 'traditional' sources shows that

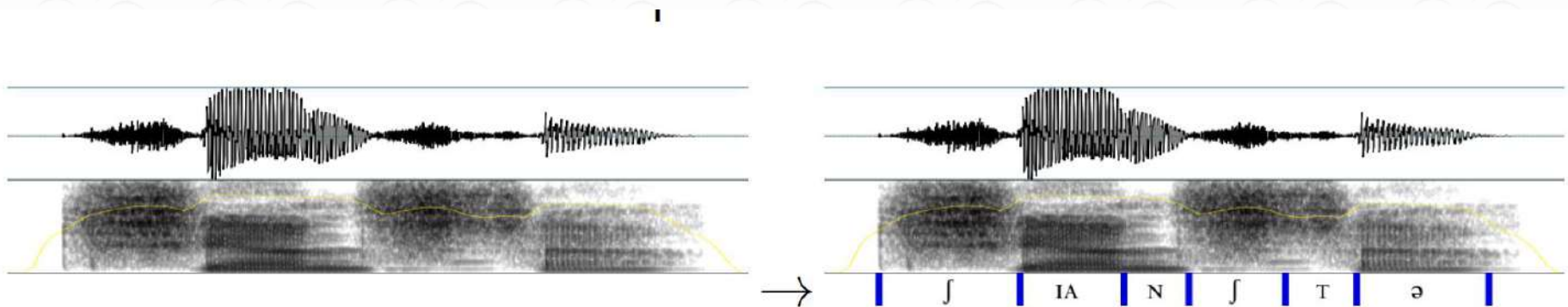
- similar structures can be found in the data
- divergences are in line with observed paths of language change

VinKo – Quantity of the Data

33

A lot of data – how to deal with it?

From audio files to (possibly automatically) annotated corpus (Kokkelmans 2023)



Phonology questionnaire

- designed to elicit certain words → we know what we can expect (more or less)
- goal: isolate single segments to study their properties (e.g. laryngeal properties of stops)
- prepares the ground for future work: machine learning, transcription of sentences

VinKo – Quantity of the Data



34

A lot of data – how to deal with it? (Kokkelmans 2023)

possibly in an automated way, employing the least manual work

1. Data preprocessing
2. Forced alignment of segments with the help of WebMAUS (LMU Munich)
3. Manual correction of transcriptions
4. Development of interface for corpus search

VinKo – Quantity of the Data

A lot of data – how to deal with it? (Kokkelmans 2023)

1. Data preprocessing

Association of audio file of word-list item to three types of transcription

- VinKoGraphy: item as presented to informants
- HistPhonGraphy: transcription useful to follow variation in language change through dialects
- SAMPA: assumed phonetic realization (in majority of dialects)

| | A | B | C | D | E |
|---|--------|------|-------------|----------------|--------|
| 1 | WordID | Lang | VinKoGraphy | HistPhonGraphy | SAMPA |
| 2 | W0001 | tre | servel | kErbeI | sErveI |
| 3 | W0002 | tre | salt | galt | zalt |
| 4 | W0003 | tre | sal | s`al | s`al |

VinKo – Quantity of the Data

36

A lot of data – how to deal with it? (Kokkelmans 2023)

2. Forced alignment via WebMAUS

WebMAUS General

Files

Please drag & drop the input signal + BPF file pairs here, e.g. 'file.wav' + 'file.par' (allowed formats are: aiff, au, avi, csv, flac, flv, mpg, mp3, mpeg, mp4, nis, nist, ogg, par, snd, sph, wav) or multiple signals all to be paired with the same annotation file `_TEMPLATE_FILE_` [par|csv].

Service options

| | | | |
|----------------|-----------------------------------------------|--------------------------------------|---|
| Language | <input type="button" value="Show inventory"/> | Language indep. (sampa) | ? |
| MAUS modus | | Forced alignment to input transcript | ? |
| Input Encoding | | X-SAMPA (ASCII) | ? |
| Output format | | Praat (TextGrid) | ? |

Expert Options (click to hide):

| | | | |
|-----------------|--|-------------|---|
| Output Encoding | | IPA (UTF-8) | ? |
|-----------------|--|-------------|---|

Schiel (1999, 2015); Kisler et al. (2017)

VinKo – Quantity of the Data

37

A lot of data – how to deal with it? (Kokkelmans 2023)

3. Manual correction

The screenshot displays the Praat interface with a 'Pause: Options for manual TextGrid correction' dialog box open. The background shows an audio spectrogram with a transcription grid. The dialog box contains the following options:

- SAVE:** The TextGrid as it is now is correctly aligned; this will save the corrected TextGrid instead of the old one, and go to the next uncorrected TextGrid.
- Align:** Rather than moving and adjusting all boundaries, you only need to adjust those of the MAU tier (bottom); the other boundaries will be adjusted accordingly. When you click SAVE, the boundaries are adjusted automatically.
- Noise:** Click this option if the audio is only irrelevant noise; the audio, TextGrid and .par files will be deleted.
- Lexic:** Click this option if the audio and the transcription do not match. E.g. the audio says 'ja' and the transcription 'nein'.
- SKIP:** Do not save the modifications and skip to the next audio.
- EXIT:** Close this window and stop the script without saving the modifications made to the last TextGrid.

The spectrogram shows a yellow highlighted region for the transcription 'A S A R'. The MAU tier (bottom) shows the transcription 'A: R (...) B A Z A: R (...)'. The visible part is 2.090000 seconds, and the total duration is 3.780000 seconds.

VinKo – Quantity of the Data

A lot of data – how to deal with it? (Kokkelmans 2023)

4. Interface for corpus search

here: <p> and word-initially, in VinKoGraphy

Filter *VinKoQuestionnaires.csv* with regex

Regexes applied to filter the data:

Note: The option 'Match VinKoGraphy and SAMPA at the same place' is activated, meaning that the regexes for these two columns must refer to the same segment in the string.

| WordID | Lang | VinKoGraphy | HistPhonGraphy | SAMPA | SAMPASAMPA |
|--------|------|-------------|----------------|------------|------------|
| .* | .* | ^(p b P B) | .* | ^(p b P B) | .* |

Generate list of files to analyse

Result of filtering:

| Include? | WordID | Lang | VinKoGraphy | HistPhonGraphy | SAMPA | SAMPASAMPA | NumWords | Nb. of files |
|-------------------------------------|--------|------|---------------------|----------------|---------------|-----------------------------|----------|--------------|
| <input checked="" type="checkbox"/> | W0005 | tre | p resa | prEs's'a | prEs's'a | prEs's'a prEs's'a | 2 | 0 |
| <input checked="" type="checkbox"/> | W0011 | tre | b us, busi | buk, buki | bUs', bu:z'i | bUs', bu:z'i bUs', bu:z'i | 4 | 0 |
| <input checked="" type="checkbox"/> | W0012 | tre | p esc, pesci | pisk, piski | pEs', pEs'i | pEs', pEs'i pEs', pEs'i | 4 | 0 |
| <input checked="" type="checkbox"/> | W0032 | tre | b islonch | bis'long | hiz'lonk | biz'lonk biz'lonk | 2 | 0 |
| <input checked="" type="checkbox"/> | W0033 | tre | p ensar | pens'a:r | pEnz'a:r | pEnz'a:r pEnz'a:r | 2 | 0 |
| <input checked="" type="checkbox"/> | W0042 | mhn | b èsp | wEs'p | bEs'p | bEs'p bEs'p | 2 | 3 |
| <input checked="" type="checkbox"/> | W0056 | mhn | p rotestiarn | protEsti:r@n | prOtEs'ti:r:n | prOtEs'ti:r:n prOtEs'ti:r:n | 2 | 2 |
| <input checked="" type="checkbox"/> | W0065 | mhn | b irt | wirt | birt | birt birt | 2 | 3 |
| <input checked="" type="checkbox"/> | W0088 | cim | p ost | pOs't | pOs't | pOs't pOs't | 2 | 4 |

VinKo – Quantity of the Data

39

A lot of data – how to deal with it? (Kokkelmans 2023)

Problems encountered on this path (corrected, in part, manually)

- Speakers don't follow instructions: 'read word twice'
- Speakers add comments: 'this word in my dialect really means y, not x ...'
- Speakers choose a lexical item different from the one proposed (*verruckt* instead of *norret* 'crazy')
- phonetic realization varies so much that WebMAUS has difficulties aligning:

Wirt → Wi[r]t, Wi[R]t, Wi[ɐ]t, Wi[ʃ]t

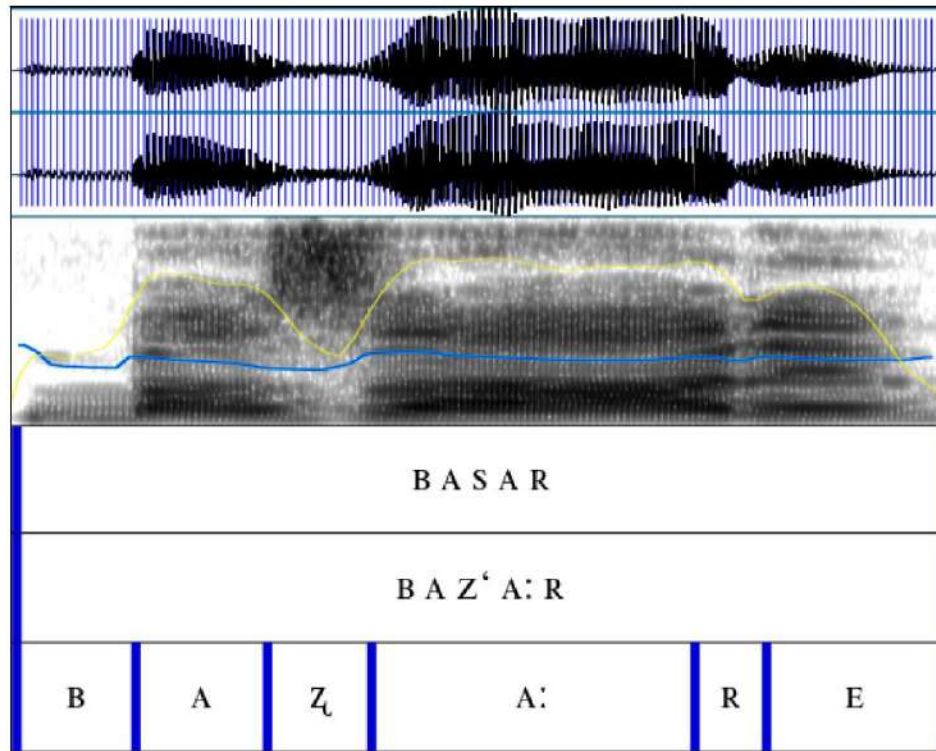
- some files consist only of noise or silence
- segments are added in some variety (differently from expected)

VinKo – Quantity of the Data

A lot of data – how to deal with it? (Kokkelmans 2023)

Problems encountered on this path (corrected, in part, manually)

- segments are added in some variety (differently from expected)



VinKo – Quantity of the Data

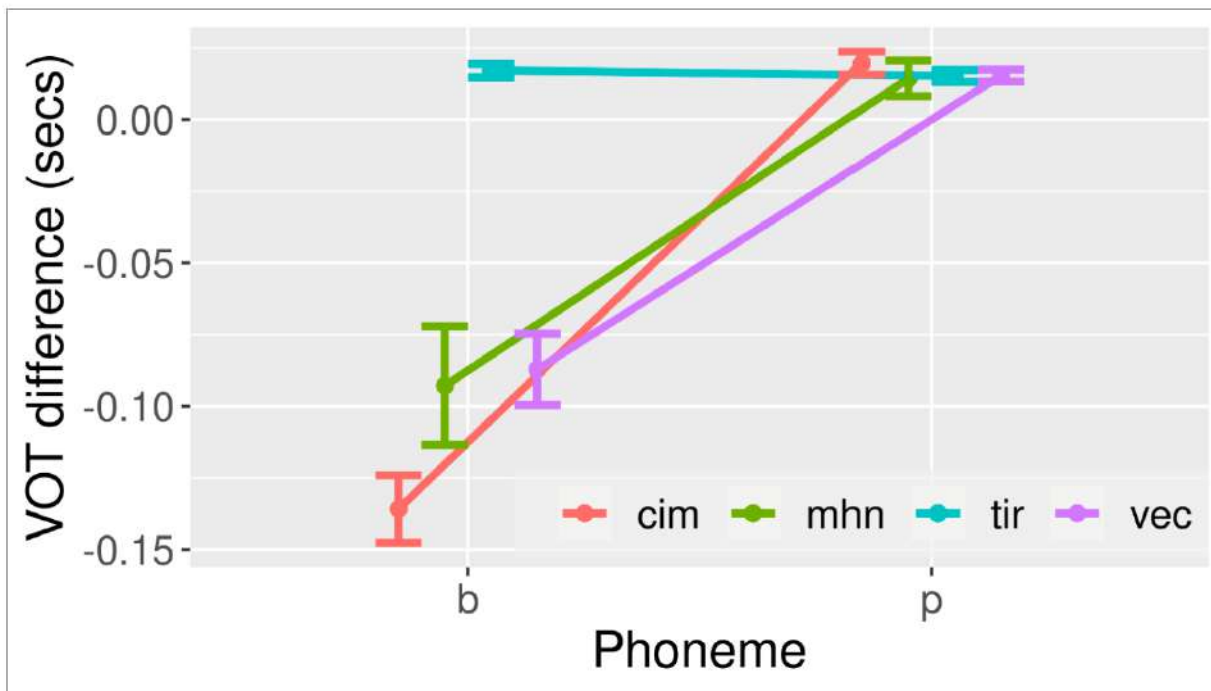
41

A lot of data – how to deal with it? (Kokkelmans 2023)

But eventually, analysis can start!

Here: historic /p/ vs. /b/ contrast, word-initially (118 items)

- Variable 1: *Voice Onset Timing*



First results

- neutralization of contrast in Tyrolean
- preservation of contrast in Cimbrian, Mochoeno and the Veneto dialects

VinKiamo Südtirol

42

Crowdsourcing and Citizen Science

Crowdsourcing platforms have a high potential of involving the communities of speakers in activities related to language documentation

VinKiamo

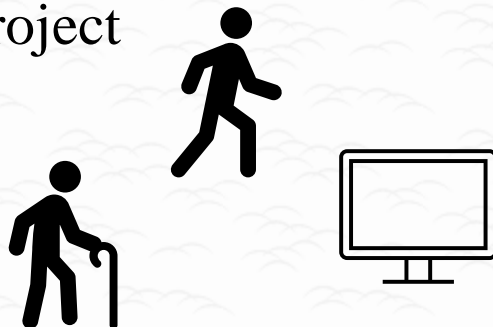
- an offspring of VinKo-Alpilink, first implemented in Veneto in 2021 (University of Verona: <https://sites.hss.univr.it/vinkiamo/>)
- 2023 start of VinKiamo Südtirol (University of Bozen/Bolzano: Birgit Alber, Joachim Kokkelmans, Emily Siviero)

VinKiamo Südtirol

43

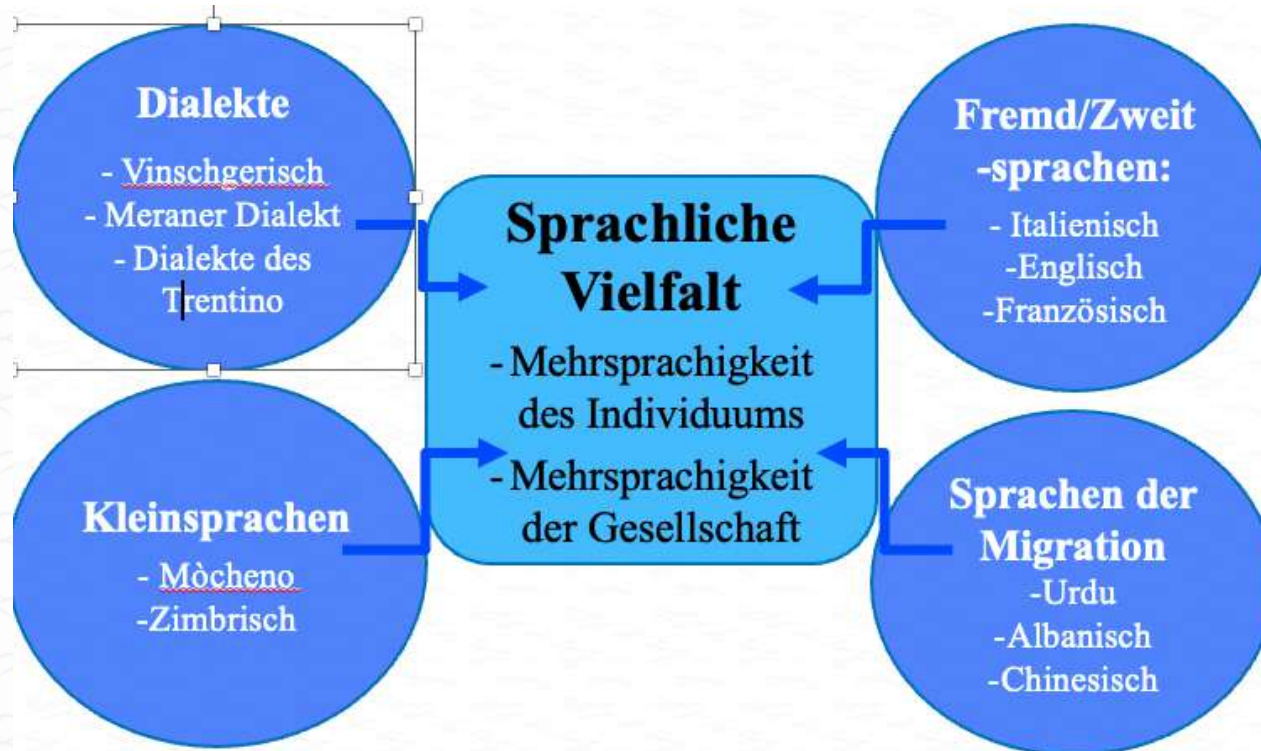
Structure

- target group: high-school students in South-Tyrol (about 17 years old)
- Presentation of basic concepts of linguistic diversity, multilingualism etc. [university staff]
- Introduction to the platform VinKo and to linguistic fieldwork [university staff]
- Fieldwork: students use VinKo to interview speakers of the older generation
- they write a report on their experience and leave a final overall feedback on the project



VinKiamo Südtirol

44



Was finden wir an Sprachenvielfalt schon alleine in diesem Klassenzimmer vor?

Welche Sprachen spricht ihr?

Wir sammeln die Antworten an der Tafel ...

VinKiamo Südtirol

45

Goals

- Enhancing awareness of linguistic diversity
- Enhancing competences of students: organization of fieldwork, teamwork, data elicitation and analysis
- documenting the linguistic knowledge of the older generation
- bridging the digital gap between the younger and the older generation
- initiating an intergenerational discourse on linguistic diversity

Language Documentation via Crowdsourcing

46

Summary

- Language documentation via crowdsourcing: potentials and challenges
- The project VinKo-AlpiLink, with a focus on
 - elicitation of linguistic structure (phonology, morphology, syntax)
 - collecting data across the Germanic-Romance language family
 - collecting spoken language
- Potentials and challenges emerging in VinKo-AlpiLink
 - in data elicitation
 - concerning the quality of the data (no problem, so far)
 - concerning the quantity of the data
 - possibilities of citizen science projects

Acknowledgment of support by:

PRIN 2020 2020SYSYBS_002



PNRR iNEST (Interconnected North-East
Innovation Ecosystem)

I43C22000250006

ECS 00000043

References

- Alber, Birgit & Joachim Kokkelmans. 2022. South Bavarian rhotics in crowdsourced linguistic data from Northeastern Italy: a diachronic and qualitative comparison. Talk presented at the conference 'Beyond Borders: German-speaking Minorities in Italy and around the World', Trento, 6.10.2022.
- Cordin, P., S. Rabanus, B. Alber, A. Mattei, J. Casalicchio, A. Tomaselli, E. Bidese & A. Padovan. 2018. VinKo, Versione 2 (20.12.2018, 09:20). In *Lo spazio comunicativo dell'Italia e delle varietà italiane. Korpus im Text*, 1–15. München: Ludwig-Maximilians-Universität
- Kisler, T., U. Reichel & F. Schiel. 2017. Multilingual processing of speech via web services. *Computer Speech & Language* 45. 326-347.
- Kokkelmans, Joachim. 2023. Halbautomatisierte Annotation von Dialektaufnahmen (VinKo & AlpiLink). Guest presentation at the Seminar 'Deutsche Dialektologie', MA in Applied Linguistics, Free University of Bozen/Bolzano.

- Kranzmayer, E. 1956. Historische Lautgeographie des gesamt-bairischen Dialektraumes. Wien: Böhlau.
- Kruijt, Anne. 2022. Crowdsourcing language contact: pronoun and article morphology in Trentino-South Tyrol and Veneto. PhD Dissertation: University of Verona.
- Kruijt, Anne. 2023. VinKo. Final Report. Technical Report, University of Verona.
- Kruijt, Anne, Patrizia Cordin, and Stefan Rabanus. (in press a.). 'On the validity of crowdsourced data'. In Elissa Pustka, Carmen Quijada Van den Berghe, and Verena Weiland (eds.): *Corpus Dialectology: from methods to theory* (French, Italian, Spanish). Amsterdam / Philadelphia: John Benjamins Publishing Company.
- Kruijt, Anne, Stefan Rabanus, and Marta Tagliani (in press b.). 'The VinKo-Corpus: Oral data from Romance and Germanic local varieties of Northern Italy'. In Marc Kupietz and Thomas Schmidt (eds.): *Neue Entwicklungen in der Korpuslandschaft der Germanistik: Beiträge zur IDS-Methodenmesse 2022. (= Korpuslinguistik und interdisziplinäre Perspektiven auf Sprache (CLIP) 11)*. Tübingen: Narr.

Rabanus, Stefan (in press). Nome di battesimo e articolo espletivo – crowdsourcing e cartografica linguistica nello studio della variazione linguistica in Trentino-Alto Adige e Veneto. In: Schöntag, Rober & Laura Linzmeier (Hrsg.): *Neue Ansätze und Perspektiven zur sprachlichen Raumkonzeption und Geolinguistik*

Rabanus, Stefan; Kruijt, Anne; Tagliani, Marta; Tomaselli, Alessandra; Padovan, Andrea; Alber, Birgit; Cordin, Patrizia; Zamparelli, Roberto; Vogt, Barbara Maria. 2022, VinKo (Varieties in Contact) Corpus v1.1, Eurac Research CLARIN Centre, <http://hdl.handle.net/20.500.12124/46>.

Scheutz, H. 2016. *Insre Sproch. Deutsche Dialekte in Südtirol. Mit dem ersten "sprechenden" Dialektatlas auf CD-ROM*. Bozen: Athesia.

Schiel, F. 1999. Automatic Phonetic Transcription of Non-Prompted Speech. In *Proceedings of the ICPHS*. San Francisco CA. 607-610.

Schiel, F. 2015. A Statistical Model for Predicting Pronunciation. In *Proceedings of the ICPHS*, Glasgow UK. Paper 195.

- TSA = Klein, K. K., L. E. Schmitt & E. Kühebacher. 1965–1971. Tirolischer Sprachatlas (3 Bände) Deutscher Sprachatlas. Regionale Sprachatlanten Nr. 3: Tirolischer Sprachatlas. Marburg: N. G. Elwert Verlag.
- Vietti, Alessandro, Birgit Alber & Barbara Vogt. 2018. Initial Laryngeal Neutralization in Tyrolean. *Phonology*. 35. 79-114