

# A Measure for Linguistic Coherence in Spatial Language Variation

Alfred Lameli, Andreas Schönberg

Research Center Deutscher Sprachatlas, Germany

lameli@uni-marburg.de, andreas.schoenberg@uni-marburg.de

## Abstract

Based on historical dialect data we introduce a local measure of linguistic coherence in spatial language variation aiming at the identification of regions which are particularly sensitive to language variation and change. Besides, we use a measure of global coherence for the automated detection of linguistic items (e.g., sounds or morphemes) with higher or lesser language variation. The paper describes both the data and the method and provides analyses examples.

## 1 Introduction

Dialectometric work typically focuses on the co-occurrence of the distribution of variants in different sites (see Goebel 1984). From these co-occurrences, reasonably coherent regions of linguistic similarity can be identified. These regions then provide, for example, clues to the aggregated structuring of higher-level linguistic areas (e.g., within a nation). Alternatively, they show to what extent individual sites of a given corpus are integrated into the region under discussion in terms of their similarity or distance to other sites (e.g., Heeringa 2003). Such analyses, which at the same time constitute the classical field of dialectometry, thus benefit from the aggregation of all linguistic phenomena of a given corpus.

However, if the interest is not in the overall structuring of a region, but in the distribution

patterns of individual variants, non-aggregating procedures must be applied. For a single phenomenon, spots of variation may be identified in most cases by visual inspection (see Ormeling 2010 for a critical account). However, in order to capture this variation quantitatively, more recent studies have considered a number of solutions, for example based on resampling techniques (e.g., Wieling & Nerbonne 2015), Kernel Density Estimation (e.g., Rumpf et al. 2009) or the concept of entropy (e.g., Prokić et al. 2009).

This paper presents a diagnostic measure for the detection of coherence or heterogeneity in spatial language variation aimed at identifying those regions that are particularly prone to variation or particularly sensitive to language change. We perform an approach based on nearest neighbor comparison and exemplify the used measure.<sup>1</sup>

In the remainder, we provide information on the data and introduce both a local and a global measure of linguistic coherence and diversity. In what follows we present example analyses based on historical dialect data from southwestern Germany and discuss the introduced procedure.

## 2 Data

The study makes use of a data set collected by the German linguist Friedrich Maurer during the year 1941 in the Upper German dialect region within the boundaries of the national territory at the time. The survey was based on a questionnaire with 113

---

<sup>1</sup> The study builds on R programming (R Core Team 2021), using the packages `spatstat` (Baddeley & Turner 2005) and `Rvision` (Garnier et al. 2021) mainly. In order to perform our coherence measure more efficiently it has been implemented

into a R-package (LinguGeo). The current version of the LinguGeo package can be found at: <https://github.com/SchoenbergA/LinguGeo>

individual words (most of them nouns, but also adjectives and verbs) and 10 sentences together with biographic information of the participants. In contrast to both the earlier survey by Wenker (Wenker 2013) and the contemporaneous investigation by Mitzka (cf. Wrede et al. 1926–1956), Maurer focused more strongly on social and biographic information. Thus, in addition to the age of the participants, for example, their gender as well as the origin of their parents or their preferred market towns are documented.

We focus on the Alemannic part of the Maurer data which is mainly related to the southwestern part of nowadays Germany (the Baden region) and the Alsace in France (see Strobel 2021 for further information). In total, the data document 2344 locations, providing a quasi-total coverage of the region under discussion (Figure 1). The handwritten questionnaires of this area have been typewritten and therefore digitalized by student assistants. The data is stored in \*.csv files and will be publicly accessible in the future in the data repository of the Research Center Deutscher Sprachatlas.

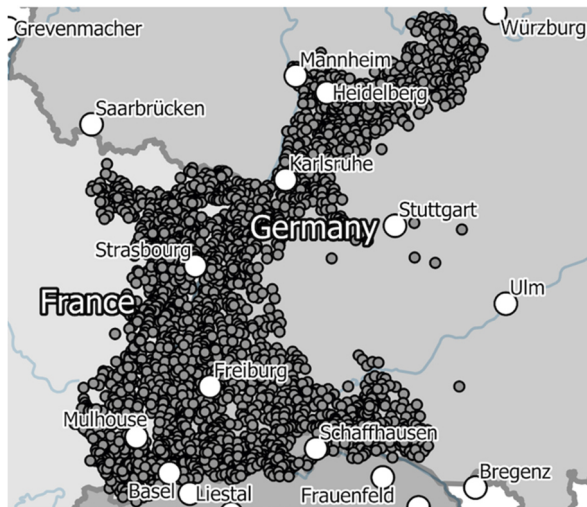


Figure 1: Study area.

### 3 Method

#### 3.1 Local Measure

In order to analyze the spatial variation of the area under discussion we compare the linguistic realizations of one site with the realizations of its geographic neighbors. Behind the selection of neighborhood relations is the assumption of the so-called “Fundamental Dialectological Postulate” (Nerbonne & Kleiweg 2007), which states that

closer objects are linguistically more similar than distant objects.

From a technical point of view, for every site  $r$  we compare the linguistic realization of an individual item  $i$  of the questionnaire (e.g., a word) with its geographic neighbor  $s$ .  $\text{Coh}_{rs|i}$  is then the number of identities between  $r$  and  $s$  with  $\text{Coh}_{rs|i} = 1$  in case of identity and  $\text{Coh}_{rs|i} = 0$  otherwise.

To obtain a better insight into how the individual sites fit into the language region, the number of compared sites should be  $S > 1$ . In the present paper, we consider up to 19 neighbors ( $0 \leq S \leq 19$ ), where 0 is used for the rendering of the original data.  $\text{Coh}_{rS}$  is then the average overlap between  $r$  and its set of neighbors  $S$  with  $0 \leq \text{Coh}_{rS} \leq 1$  and  $\text{Coh}_{rS} = 1$  indicating identity between  $r$  and  $S$  and  $\text{Coh}_{rS} = 0$  indicating no identity between  $r$  and  $S$ . In case a location has several variants for a linguistic variable (e.g., because of several participants or multiple responses), the number of matches between  $r$  and  $s$  is related to the number of local variants.

An example is provided by Figure 2. The centrally located site is opposed by a total of 5 nearest neighbors, which have a total of 2.5 matches with the central site, resulting in  $\text{Coh} = 2.5/5 = 0.5$ . The number of variants is irrelevant for this approach but is relevant for the global measure (cf. 3.2)

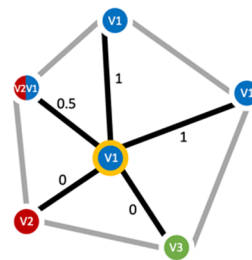


Figure 2: Model of distribution of variants.

Inverting the scale results in a measure of linguistic diversity instead of linguistic coherence which we refer to as  $\text{Div} = 1 - \text{Coh}$ . We use this  $\text{Div}$  measure in order to identify moments of particular dynamics on language maps.

Another point is worth mentioning. The nearest neighbor approach relies heavily on the definition of geographic coordinates and distances. In our approach, the geometric information of the spatial position for each survey site is thus originally stored in the WGS 84 format (longitude and latitude). Due to the ellipsoidal coordinate system, the distances are heavily distorted which directly

affects the selection of the nearest neighbors. To use the quasi-exact distances a cartesian coordinate system is required. Therefore, we projected our data to the UTM system related to the ETRS89 ellipsoid.

### 3.2 Global Measure

While the local measure indicates the integration of individual sites into its nearest spatial neighborhood, it says nothing about the coherence or heterogeneity of an overall map. Various options are available for this purpose. For example, the mean of all local Coh values could be taken as a global measure of coherence (CohG). However, as Figure 3 demonstrates, this measure is dependent on the number of linguistic variants in a data distribution, making it difficult to compare CohG across maps with different numbers of variants. For example, if a map shows two linguistic variants a complete random distribution results in  $0.5 \leq \text{CohG} \leq 1$  and  $0.33 \leq \text{CohG} \leq 1$  for three variants etc.

In order to solve this problem, we perform a CohG\* correction in which CohG is divided by the number of variants and scaled  $0 < \text{CohG}^* \leq 1$ . As becomes evident by Figure 3, CohG\* is robust against the number of variants, while CohG, in contrast, is sensitive to it and converges to CohG\* as the number of variants increases. Similar holds for the number of neighbors against which CohG\* is robust while CohG is sensitive to it (not reported).

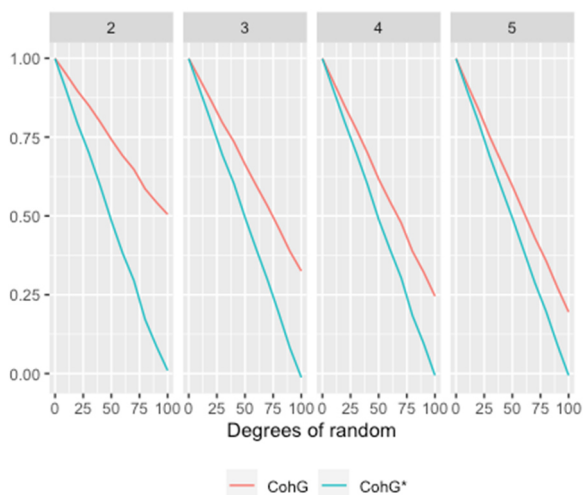


Figure 3: Comparison of CohG and CohG\* based on simulated degrees of both spatial coherence and random data filling (0-100%) for a data distribution with 2 to 5 linguistic variants.

Another view on CohG\* is provided in Figure 4 and Figure 5. In these figures, data simulations are performed for the locations of the corpus, generating different degrees of random data distributions. Starting from a uniform distribution, 20% of the data of each map are successively overwritten with a random distribution.

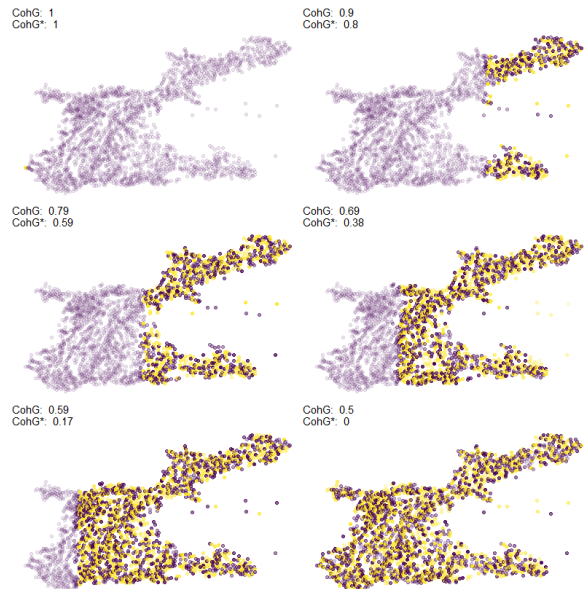


Figure 4: Simulation of different degrees of spatial heterogeneity (0%, 20%, 40%, 60%, 80%, 100%) for a map with two linguistic variables. Variant 1 = purple, variant 2 = yellow,  $\alpha = 1 - \text{Coh}$ .

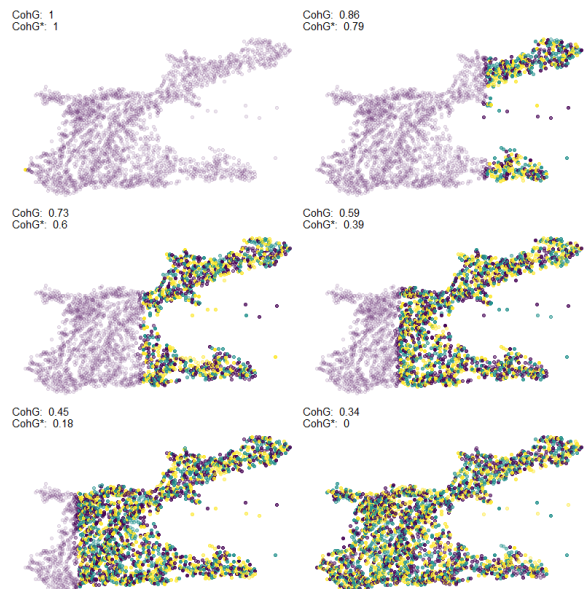


Figure 5: Simulation of different degrees of spatial heterogeneity (0%, 20%, 40%, 60%, 80%, 100%) for a map with three linguistic variables. Variant 1 = purple, variant 2 = yellow, variant 3 = green,  $\alpha = 1 - \text{Coh}$ .

While Figure 4 illustrates data simulation with two linguistic variants, Figure 5 illustrates the same procedure based on three linguistic variants. The figures show that while the CohG is related to the amount of variants, the CohG\* values describe the same amount of coherence/homogeneity unattached to the number of variants.

Against this background, the Coh measure, and also the CohG\* measure, yields plausible results as far as different degrees of coherence or heterogeneity are concerned. However, it is still an open question how the values turn out in concrete use cases and what more detailed conclusions can be drawn from them.

## 4 Use Cases

### 4.1 Lambdacism in *Kirche* ‘Church’

As a first example we focus on a rather simple spatial pattern provided by the distribution of *-r-* and *-l-* sounds in the word *Kirche* ‘church’ (*Kirche* vs. *Kilche*) in the southern part of our study area (Figure 6). The phonological process behind this is the so-called lambdacism, which is typical for some regions of the German-speaking area (cf. Lameli 2015).

Figure 6 illustrates the distribution of the variants in the southern part of the study area. At each site one variable is documented, where *Kirche* (blue) occurs 1008 times, *Kilche* (red) 222 times (1230 sites in total). Hence, 81.94 % of the sites in the study area show *-r-*.

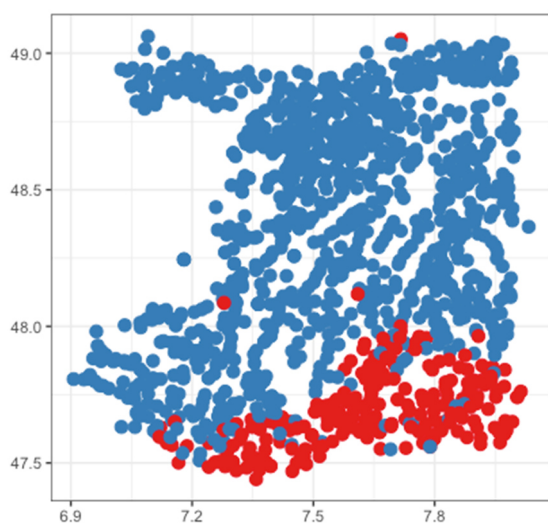


Figure 6: Example of a spatial distribution of linguistic variants *-r-* (blue) and *-l-* (red) in the word *Kirche* ‘church’.

In a random distribution the expected probability that a particular site’s neighbor shares the same variant is  $EV = (1008-1) / (1230-1) = 81.94\%$ . For the same distribution we reveal under the consideration of 5 nearest neighbors  $CohG^* = .94$  ( $Coh = .9$ ) indicating that, on average, 94 % of the neighboring 5 sites share the same variant *-r-* as the site under observation. However, the question remains open as to how high CohG\* turns out to be in a random distribution when 5 nearest neighbors are considered, as in the present case. For this purpose, 1000 data simulations were performed in which the existing occurrences of *-r-* and *-l-* sounds were randomly distributed among the study sites. The resulting mean of  $CohG^* = .41$  indicates that, given a random distribution of data, statistically 41 % of the neighboring five locations share the same variant as a particular site under observation with a range of  $CohG^* = .37-.44$ .

By CohG\* being higher than both the random distribution and the expected value EV, (1) spatial clustering of *-r-* and *-l-* is indicated and, as a consequence, (2) a clear separation of the variants. Indeed, very few locations aside, all variants cluster in contiguous areas as already becomes clear by visual inspection.

Testing the distribution of local Coh values against a normal distribution using a Wilcoxon rank sum test reveals a statistical difference between the expected value EV and the empirically found Coh measure ( $z = -4.21, p < .001, r = .94$ ). What these measures refer to becomes evident when plotting  $1-Coh (= Div)$  on a map (Figure 7).

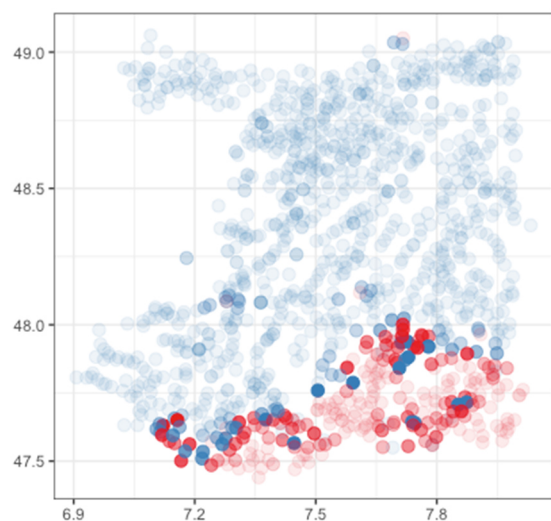


Figure 7: Local measure of linguistic coherence ( $Div = 1-Coh$ ) applied to the data of Figure 6.

As expected, the highest Div values are at the border zone between the variants. Most interestingly, there are differences depending on the spatial alternation of the variants. For example, on the left, where we find a mix of variants, Div values are high. In contrast, in the center, where we find a separation of *Kirche* and *Kilche*, Div values are low. The spots illustrated by Figure 7 thus allow conclusions to be drawn about zones of increased linguistic dynamics: around the sites with high values (intense colors) there is a high degree of variation, around the sites with low values (pale colors) there is a lower degree of variation. While the former can be expected to be more sensitive to language change regarding the variable under discussion, the latter can be expected to be more robust to language change.

Methodologically, it should be emphasized that, due to the nearest neighbor approach, the described procedure always computes a gradient-like result. Even if there is a sharp separation between variants

(Figure 6) a gradient would be computed (Figure 7).

The intensity of this gradient-like effect depends on the number of nearest neighbors. Using the minimum of two nearest neighbors will result in exactly three index values and the resulting map would set a focus on areas which differ from their surroundings (Figure 8/A). This may be useful to detect islands of variation in rather coherent areas. With increasing numbers of nearest neighbors, the amount of possible index values will increase and return much more smoother transitions. This is helpful for the detection of areas with variation in a cluster-like way. Areas with variation in close distances would be smoothed to clusters which would be differentiated from surrounding homogeneous areas (e.g., Figure 8/D). This way of proceeding captures, for example, border regions in a more schematic way and those regions which are most likely unaffected by these border regions.

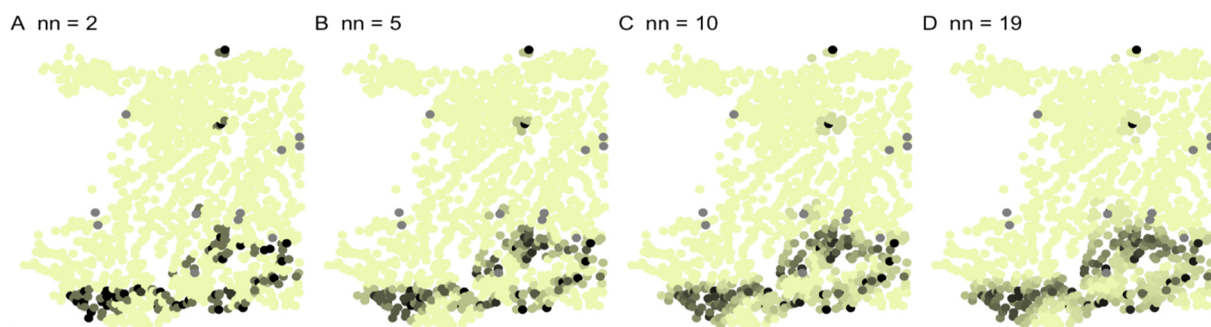


Figure 8: Local measure of linguistic coherence applied to the data of Figure 6 with different number of nearest neighbors and without information on linguistic variants.

## 4.2 Subtractive Plural in *Hunde* ‘Dog-PL’

Another example is provided by Figure 9, which focuses on the whole language area of the Maurer data. The map illustrates the variation of the word ending in *Hunde* (‘dog-PL’; CohG\* = .87) considering three graphemic variants (<nd>, <ng>, <nn>), of which <nn> (phonologically /n/) and <ng> (phonologically /ŋ/) have been considered as subtractive plurals (Birkenes 2014). While the *Kirche* example considers only two linguistic variants, Figure 9 refers to three linguistic variants. The figure combines three different views. On the left side is the distribution of variants without any preparation, in the middle the representation of the coherence measure (expressed in Div) including

information on the variants and on the right side the representation of coherence (Div) without information on the linguistic variants.

Obviously, the coherence map in the middle clearly highlights the spots of linguistic variation. Among them are areas where only two variants interact (e.g., <nd> and <nn> in the South, <nd> and <ng> in the North), but also areas where all three variants meet (in the center). Similar to the previous example the coverage of individual variants is mapped.

The map on the right, on the other hand, emphasizes where generally such patterns of variation are encountered. This map consequently emphasizes the contrast between homogeneous and heterogeneous moments of the spatial data distribution. In this case, too, conclusions can be

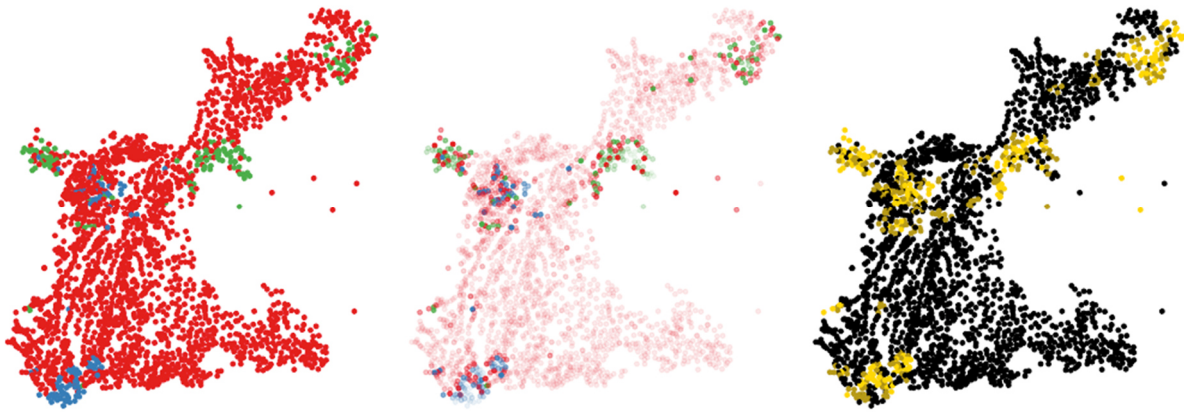


Figure 9: Local measure of linguistic coherence ( $Div = 1 - Coh$ ) for a linguistic variable with three variants (*Hunde* ‘dog-PL’); green = <ng>, red = <nd>; blue = <nn>; left: distribution of variants; middle: Div measure with information on linguistic variants; right: Div measure without information on linguistic variants.

drawn (as in the previous example) about the extent of regional variation and possible language change events; it is in the yellow zones where variation is highest and possible language change is most likely.

From a methodological perspective, the following is worth mentioning. By integrating the nearest neighbors, a smoothing effect is created, which shows linguistic variation in places where actually no variation is documented by data collection. The idea behind this is that variation is probably more widespread than what is captured by data collection. For example, if only one person is asked about a particular linguistic variant at each of two surveyed locations (which is very often the case in dialectological studies), it would possibly be wrong to take different answers per se as evidence of strict linguistic differences between those locations. Instead, it must be expected that both variants would be encountered in both localities and would be appropriately documented with other participants if data were repeatedly collected. However, the probability of this decreases with increasing geographical distance. The measure thus provides a prediction for the communicative reach of language variants.

## 5 Discussion

The Coh measure, as well as the Div measure respectively, reveals spots of local variation, which indicate horizontal (i.e. geographical) or vertical (i.e. social, pragmatic) heterogeneity. As Labov (2004) points out, these spots of increased language variation might be possible starting points of language change. In this regard, Bellmann (1983) considers the model in Figure 10.

Starting from a situation where variant A is the only available realization of a particular linguistic variable, at a certain time variant B becomes an alternative. This is the situation illustrated by Figure 10 for both scenarios (above and below). However, the Coh measure goes beyond local variation by modeling the closest relative area of influence of that alternative.

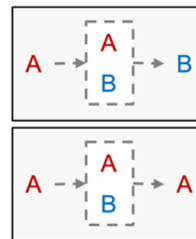


Figure 10: Possible stages in the formation of language variation and/or language change on the example of two variants A and B; above: scenario 1 (language change); below: scenario 2 (temporary language variation).

Obviously, analysis using Coh (like Figure 7) does not specify how long the variative phase will persist. Furthermore, it could be that variant B disappears again (Figure 10 below), and it could just as well be that variant B prevails (Figure 10 above) while A disappears. Consequently, Coh does not allow for a clear prediction of the process of language change, but it does illustrate that, if language change does occur, it is likely to occur at the spots with high Div ( $= 1 - Coh$ ). Against this background, the relevance of the Coh measure is to indicate spots of particular linguistic dynamics. Identifying these spots enables both prediction and explanation of ongoing and/or completed language change.

On the other hand, with  $\text{CohG}^* \rightarrow 1$  it can also be shown directly whether a language region has proto-typical variants, which can then be easily identified in the data distribution.

Furthermore, applying the coherence measure to a collection of multiple linguistic phenomena, as shown in Figure 11, leads to a new perspective on the structuring of linguistic space. Instead of highlighting the clusters of linguistic similarity,

rather the zones of particular linguistic dynamics are identified. From looking at the coherence values, even without mapping, a first impression is given whether the lemmas in question show a strong spatial clustering or not. This is useful for huge datasets with lots of linguistic variables. At the same time, it becomes evident that the measure is sensitive for outliers (i.e., isolated sites), which are evident by individual points.

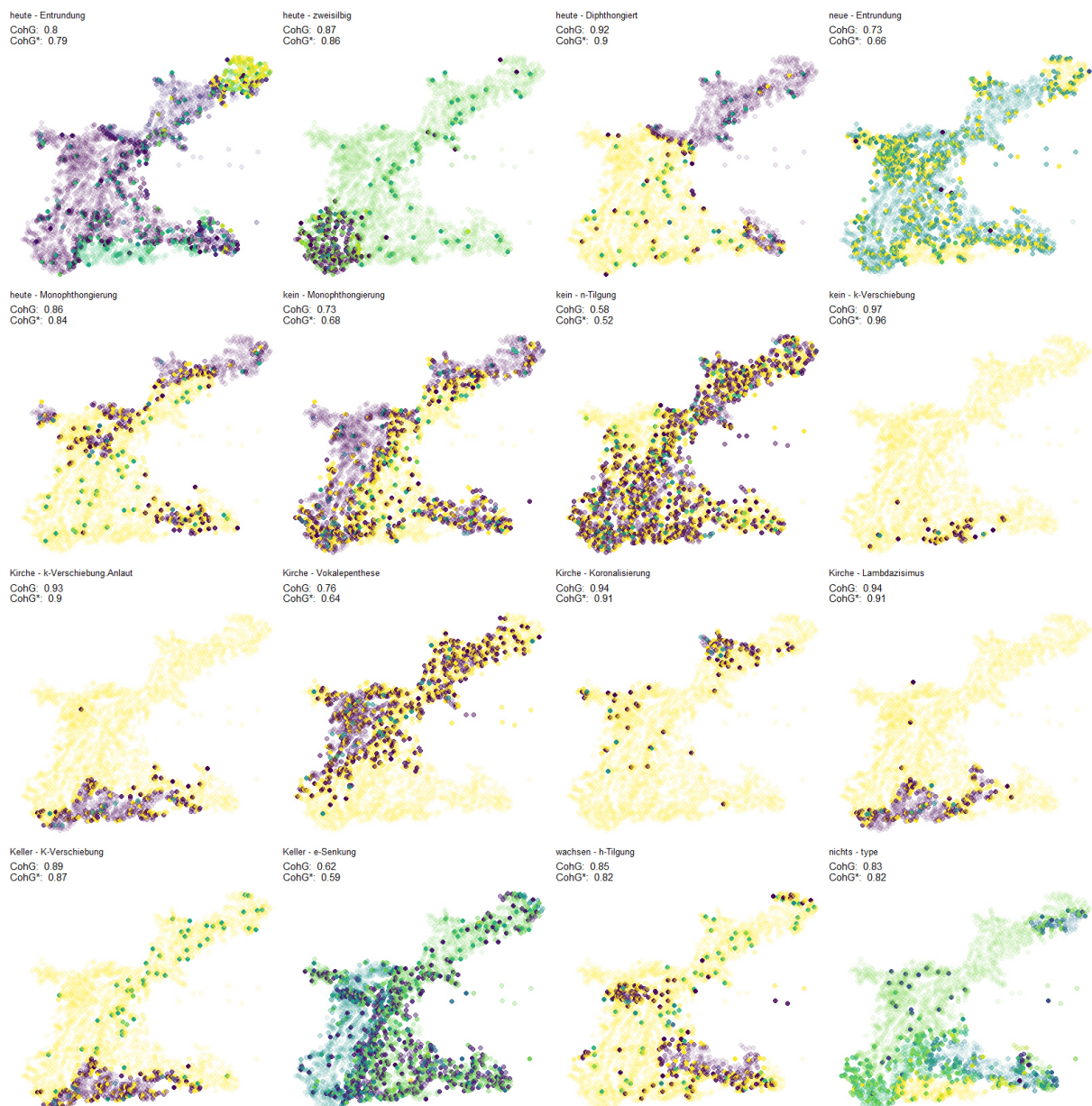


Figure 11: Local measure of linguistic coherence ( $\text{Div} = 1 - \text{Coh}$ ) for different linguistic variables.

Among the existing dialectometric literature, our coherence measure is comparable to the technique introduced by Rumpf et al. (2009) using Kernel Density Estimation (KDE). Our measure explicitly considers geographical neighborhood, but, in

contrast to the KDE approach, it is more focused on local variation. Instead of calculating an adequate bandwidth, we choose a certain number of neighbors in order to test for the integration of an individual site into the linguistic area. In this

respect, the underlying concept is that linguistic space develops in small-scale communication zones, not in large-scale continua. From a technical perspective, a difference to the KDE approach is that we do not rely on the definition of individual variant-occurrence maps as an intermediate step of analysis, but process the variation given in the data set directly.

Notwithstanding this, there are other studies that work with the notion of coherence or focus on transitional spaces. Nerbonne & Kleiweg (2007), for example, introduce a local measure of incoherence, which, however, focuses on linguistic rather than geographic distances. Our measure thus provides an alternative view of the relationship between spatial and linguistic proximity based on individual maps and not on aggregated data. Goebel (2010), nonetheless, illustrates the importance of skewness as a global statistical measure of the linguistic integration of individual sites into the linguistic area and the assessment of transitional zones. Similar to Nerbonne & Kleiweg (2007), the basis of linguistic measurement is in Goebel's approach not the individual map, but a set of aggregated data. Unlike Goebel (2010), we focus exclusively on concrete geographic neighbors of an individual site with both the local and global measures, which makes our approach, in the case of the local measure, independent from the overall statistical distribution, which is in dialectometric studies typically shaped by linguistic distance or similarity.

## 6 Conclusion

This paper introduces a nearest neighbor approach as a diagnostic tool in order to find regions which are more sensitive to language variation and change than others. For this purpose, a local measure of coherence is used (Coh). In addition, a global coherence measure (CohG) as well as a corrected global measure (CohG\*) was used to quantitatively assess the spatial coherence of more comprehensive data distributions (e.g., on maps) and to automatically identify linguistic items with higher/lesser language variation. Two case studies illustrate the application of the method and the informative quality of the measures.

## Limitations

The method works reliably, even if a map contains multiple variants. However, if there are more than,

say, 10 or 15 variants, it can happen that no clear spots can be identified on the maps. For this matter, a more probabilistic approach would be desirable, which is currently not implemented.

Another limitation is the distance measure used for the identification of nearest neighbors. Currently, nearest neighbors are defined using Euclidean distance. This is not a problem if the analysis takes place in flat terrain (e.g., the Upper Rhine Plain). In mountainous terrain, however, this can lead to slight biases. To solve this problem, we will implement more realistic distance measures such as travel time in the future.

From a linguistic perspective, a limitation of the method is that even if it informs about the variation spots, it does not provide any information about the direction in which a possible language change could develop. However, such a statement is difficult to make without concrete comparative language data (e.g., diachronic data) or social interpretation. Since the Maurer data allow an analysis in apparent-time, further approaches for investigation will be possible in the future.

## Ethics Statement

This work complies with the ACL Ethics Policy.

## Acknowledgments

We are grateful to six reviewers for their valuable comments as well as Peter Auer, Michael Cysouw, Alexandra Lieb and Maj-Brit Strobel for discussion. Maj-Brit Strobel was kind enough to provide us with data from her work. This research is funded by the German Research Foundation under the project "Alemannisch variativ" (DFG, grant number 452440801).

## References

- Adrian Baddeley, Rolf Turner. 2005. spatstat: An R Package for Analyzing Spatial Point Patterns. *Journal of Statistical Software*, 12(6):1-42. URL <https://www.jstatsoft.org/v12/i06/>.
- Günter Bellmann. 1983. Probleme des Substandards im Deutschen. In Klaus J. Mattheier (ed.) *Aspekte der Dialekttheorie*. Niemeyer: Tübingen:105-130.
- Magnus Breder Birkenes. 2014. *Subtraktive Nominalmorphologie in den Dialekten des Deutschen. Ein Beitrag zur Interaktion von Phonologie und Morphologie*. Steiner: Stuttgart.
- Simon Garnier, Noam Ross, Robert Rudis, Antônio P. Camargo, Marco Sciaini, and Cédric Scherer. 2021. Rvision - Colorblind-Friendly Color Maps for R. R



- package version 0.6.2. URL <https://sjmgarnier.github.io/viridis/>.
- Hans Goebel. 1984. *Dialektometrische Studien. Anhand italoromanischer und galloromanischer Sprachmaterialien aus AIS und ALF*. Niemeyer: Tübingen:191-193.
- Hans Goebel. 2010. Dialectometry and quantitative mapping. In Alfred Lameli et al. (eds.) *Language and Space. An International Handbook of Linguistic Variation. Language Mapping*. Mouton de Gruyter: Berlin, Boston:433-457.
- Wilbert Heeringa. 2003. *Measuring Dialect Pronunciation Differences using Levenshtein Distance*. University Press: Groningen.
- William Labov. 1994. *Principles of linguistic change. Vol. 1: Internal factors*. Blackwell: Oxford.
- Alfred Lameli. 2015. Zur Konzeptualisierung des Sprachraums als Handlungsraum. In Michael Elmentaler et al. (eds.) *Deutsche Dialekte. Konzepte, Probleme, Handlungsfelder*. Steiner: Stuttgart:59-83.
- John Nerbonne, Peter Kleiweg. 2008. Toward a dialectological yardstick. *Journal of Quantitative Linguistics*, 14:148-166.
- Ferjan Ormeling. 2010. Visualizing geographic space. The nature of maps. In Alfred Lameli et al. (eds.) *Language and Space. An International Handbook of Linguistic Variation. Language Mapping*. Mouton de Gruyter: Berlin, Boston:21-40.
- Jelena Prokić, John Nerbonne, Vladimir Zhobov, Petya Osenova, Kiril Simov, Thomas Zastrow and Erhard Hinrichs. 2009. The computational analysis of Bulgarian dialect pronunciation. *Serdica. Journal of Computing*, 3:269-298.
- R Core Team. 2021. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Jonas Rumpf, Simon Pickl, Stephan Elspaß, Werner König and Volker Schmidt. 2009. Structural analysis of dialect maps using methods from spatial statistics. *Zeitschrift für Dialektologie und Linguistik*, 76(3):280-308.
- Maj-Brit, Strobel. 2021. Die Verschriftungen in der Dialekterhebung Friedrich Maurers in Baden und im Elsass als Evidenz für die Verbreitung der Standardlautung. *Zeitschrift für Germanistische Linguistik*, 49(1):155-188.
- Georg Wenker. 2013. *Schriften zum „Sprachatlas des Deutschen Reichs“*. Gesamtausgabe. Olms: Hildesheim, New York, Zürich.
- Martijn Wieling and John Nerbonne. 2015. Advances in Dialectometry. *Annual Review in Linguistics*, 1(1):243-264.
- Ferdinand Wrede, Walther Mitzka and Bernhard Martin. 1926–1956. *Deutscher Sprachatlas auf Grund des von Georg Wenker begründeten Sprachatlas des Deutschen Reichs und mit Einschluß von Luxemburg in vereinfachter Form*. Elwert: Marburg.