

**Interpolation, and the probabilistic
view on linguistic diversity in general
– how to convince the linguists? –**

Michael Cysouw
Deutscher Sprachatlas
Philipps-Universität Marburg

Beyond the symbol-map

- Visualisation of geographic distributions in linguistics is mostly symbol-based with added lines to signal boundaries
- But: the data are empirical probabilistic distributions without real boundaries
- Why do we linguists like boundaries so much?

Händler, Harald & Carl Ludwig Naumann. 1976.
Automatisierung der Isoglossenfindung.
Germanistische Linguistik 76(3-4). 123-159.

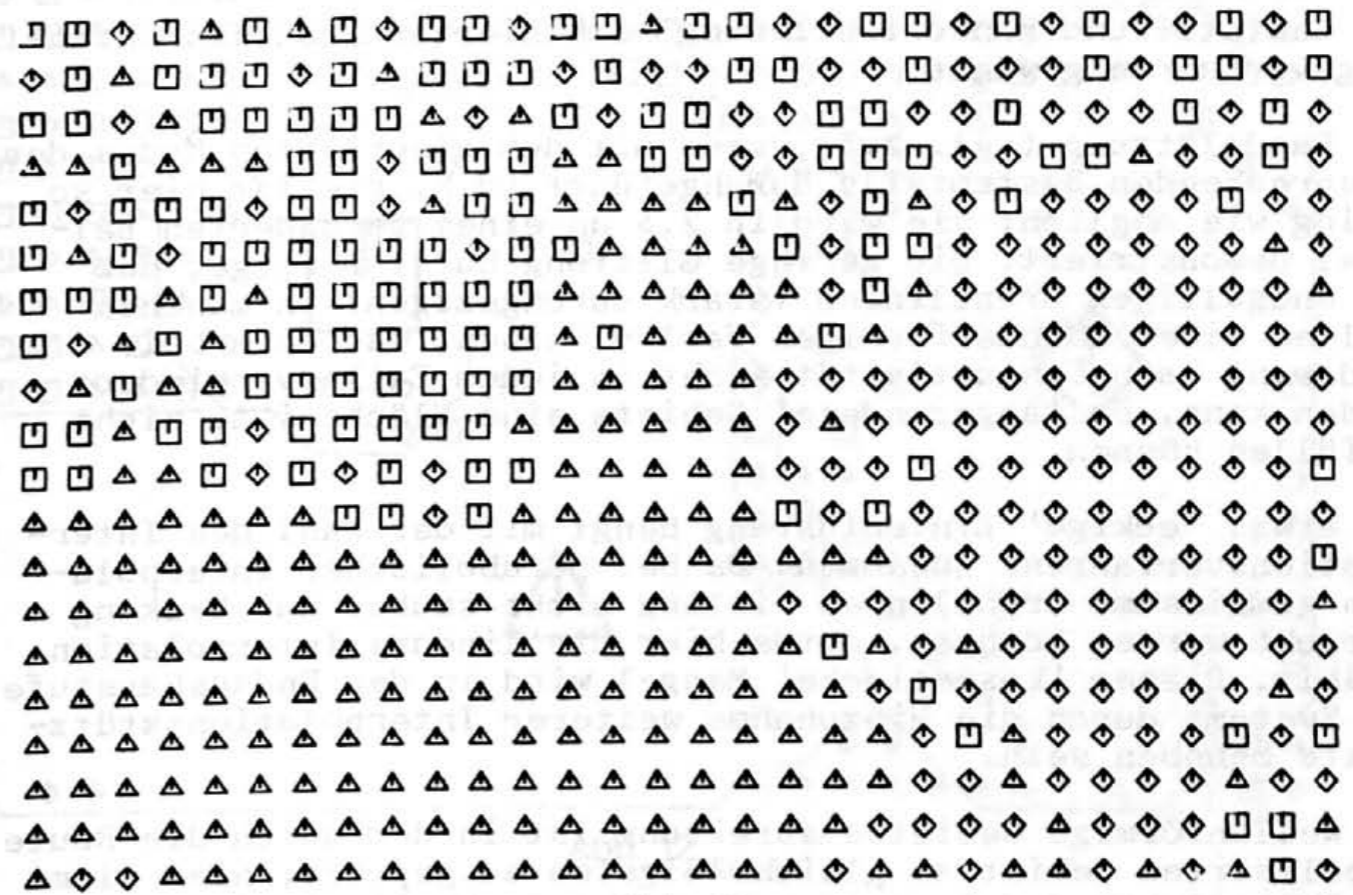


Abb. 10

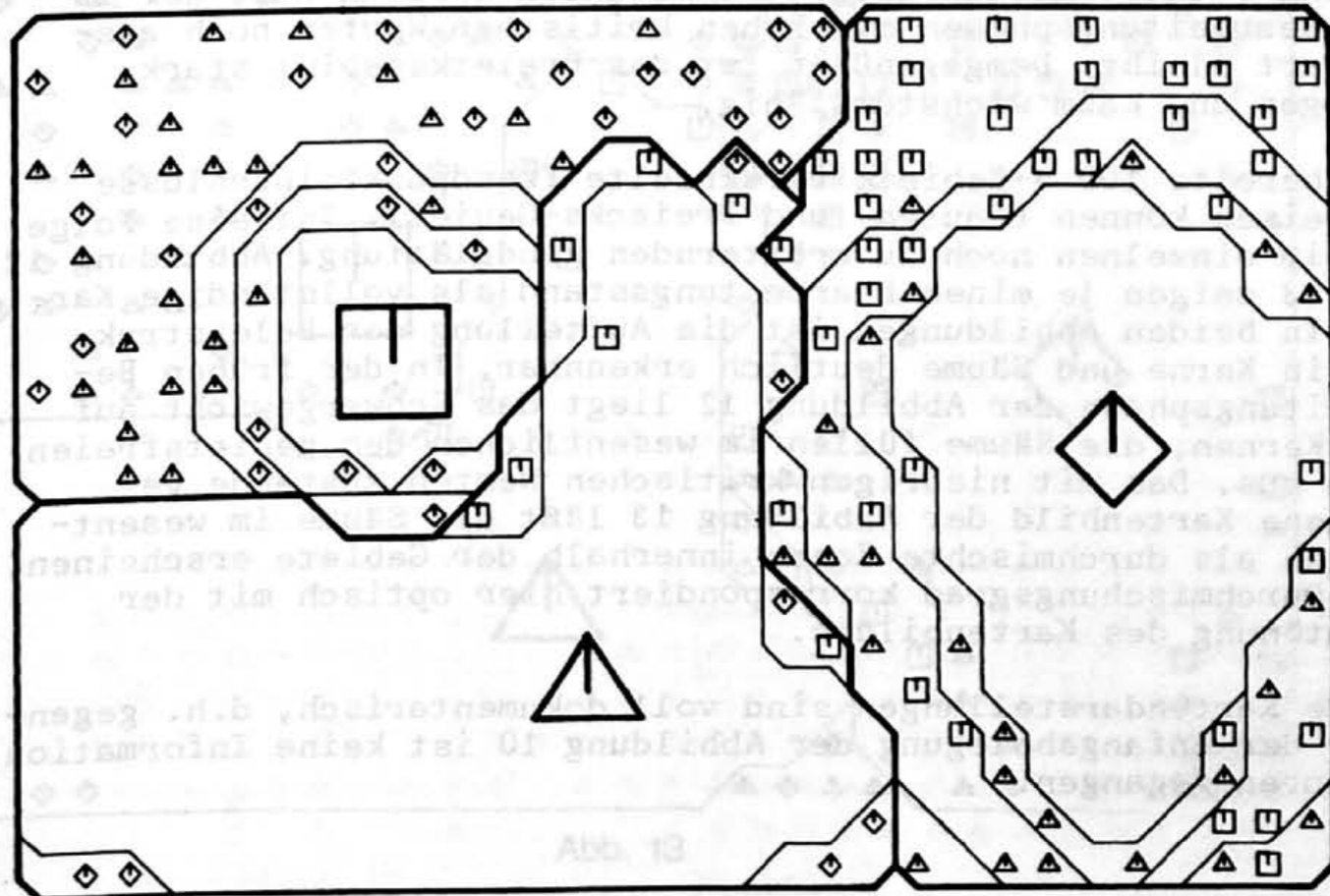


Abb. 11

Grad der Randglättung

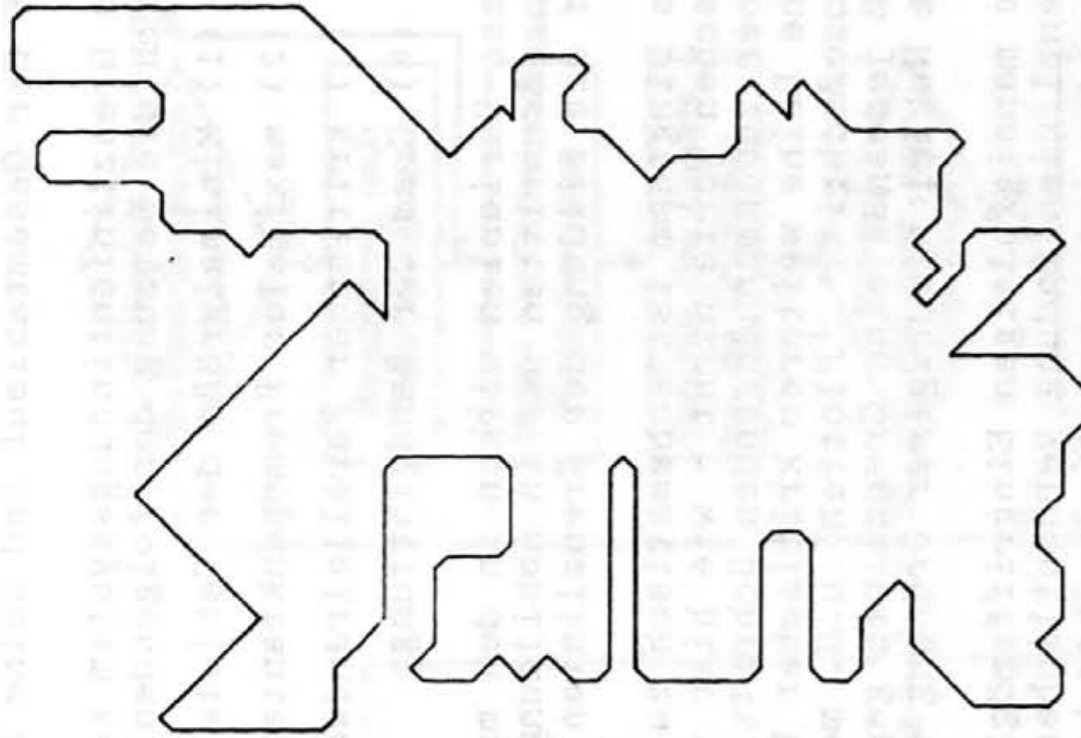


Abb. 18 a

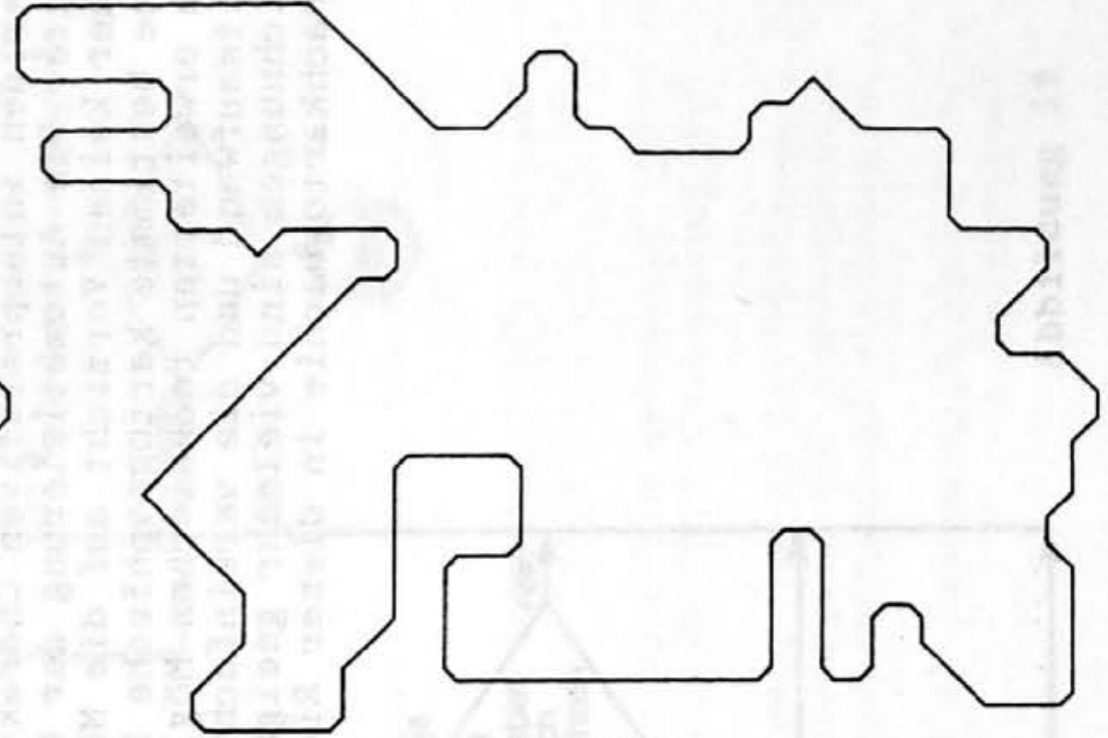


Abb. 18 b

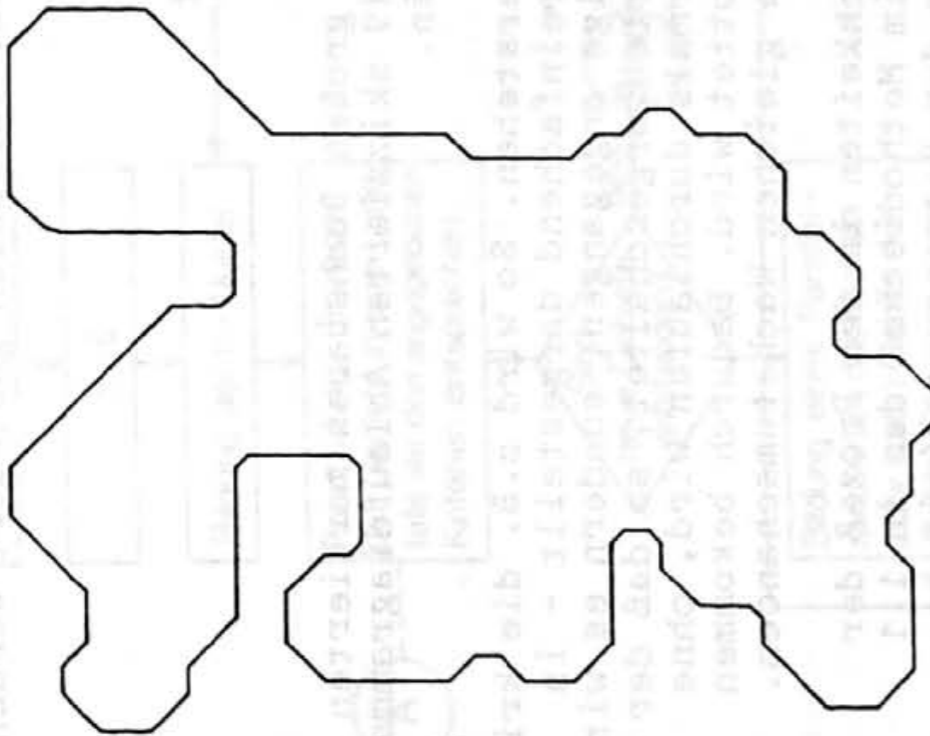


Abb. 18 c

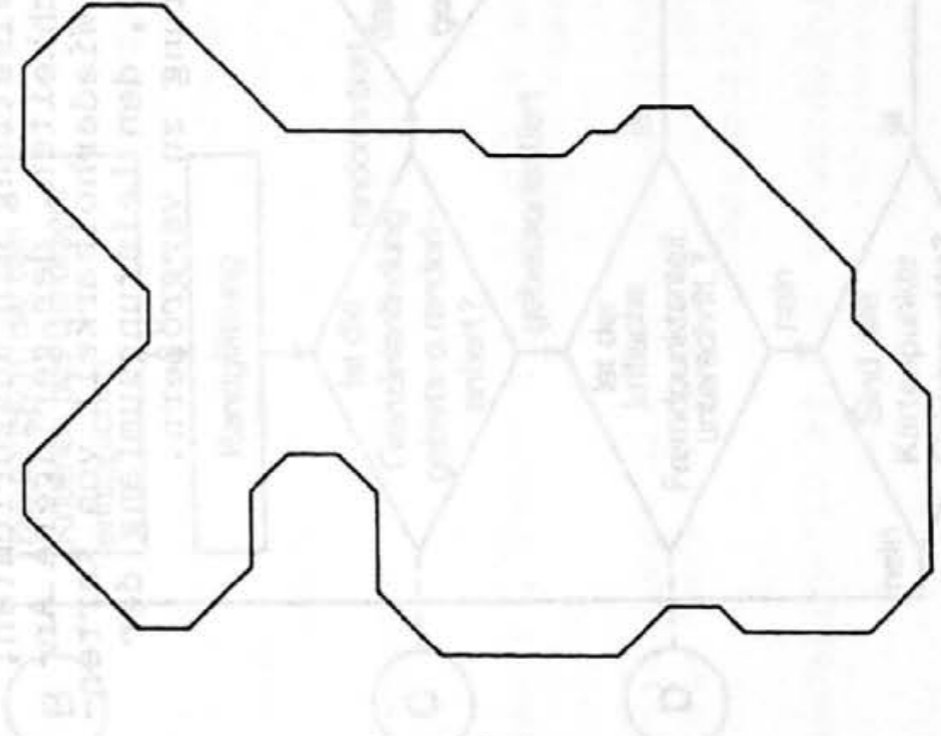
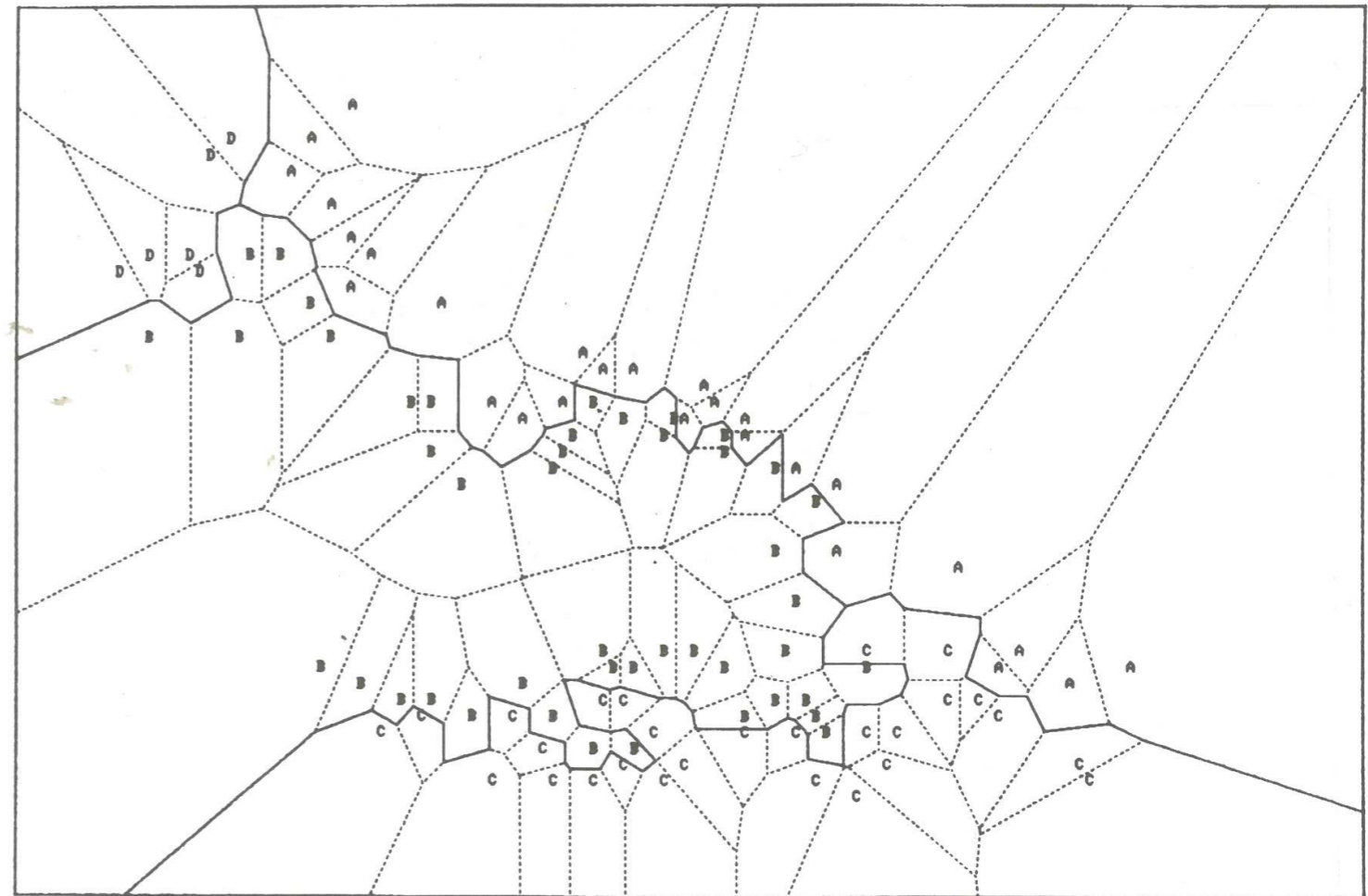
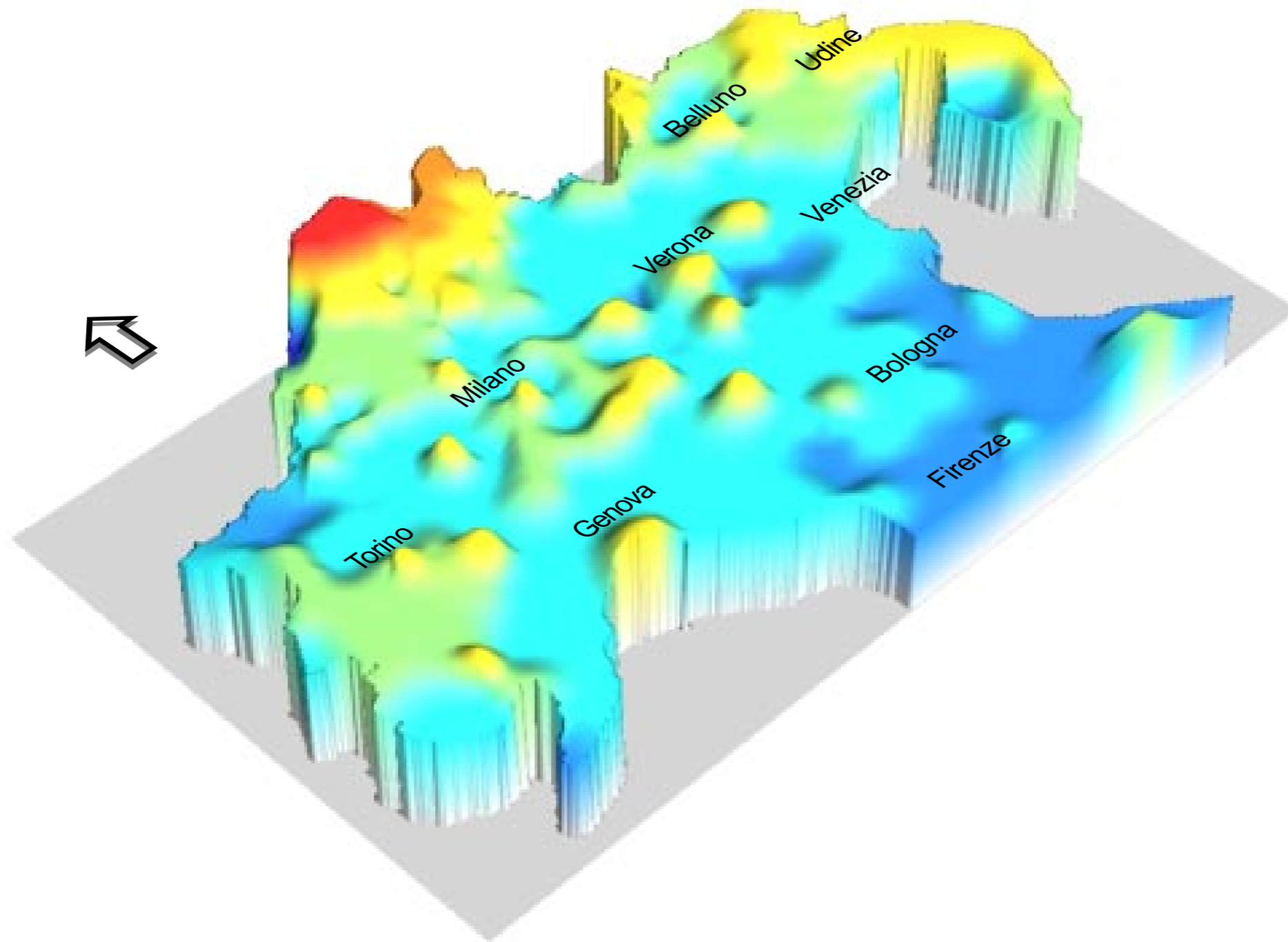


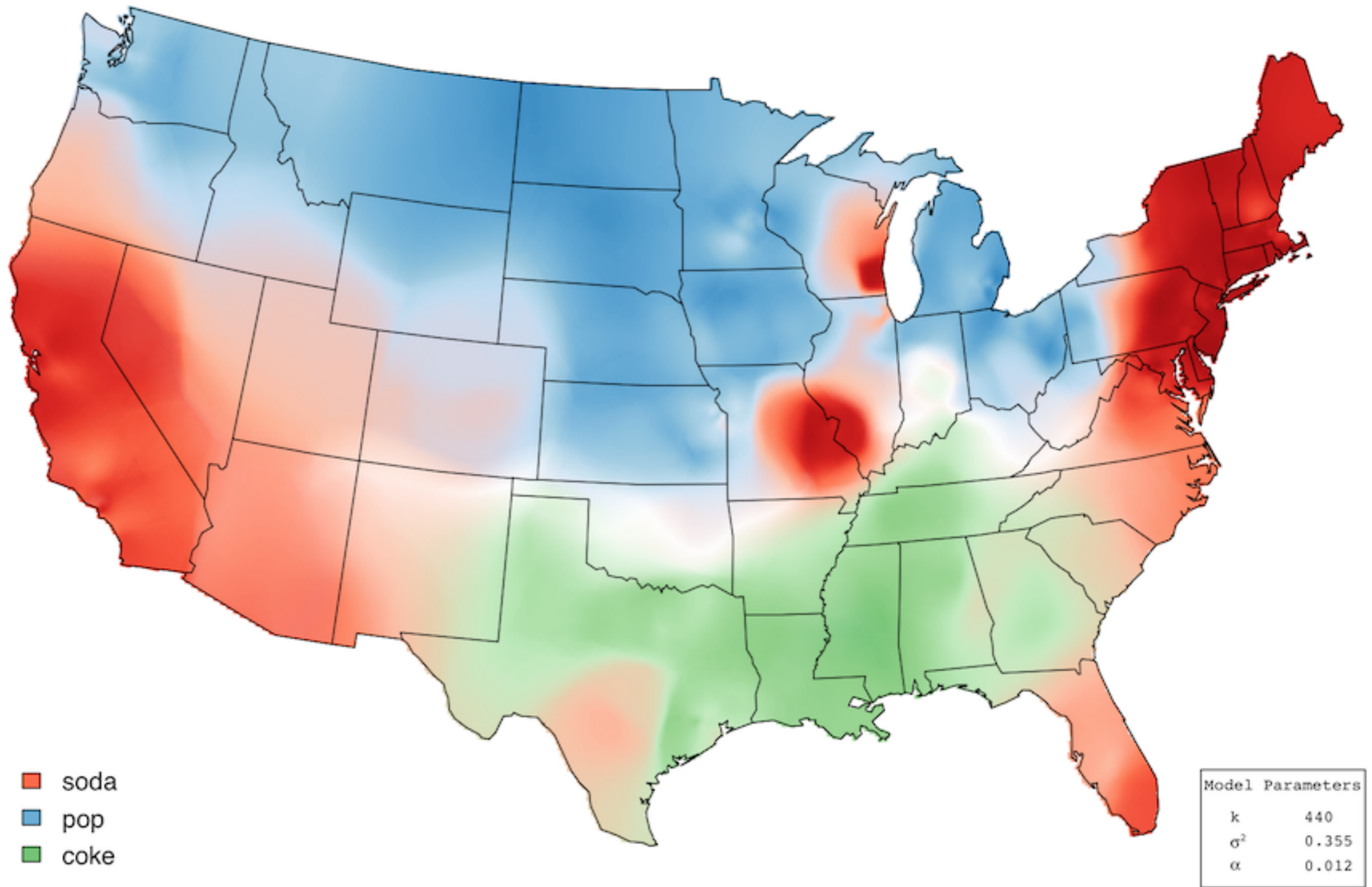
Abb. 18 d

Pudlatz, Hilmar. 1977. Automatische Erzeugung von Isoglossen auf dem Plotter mit Hilfe von Thiessen Polygone.
Germanistische Linguistik 77(3-4), 245-258.





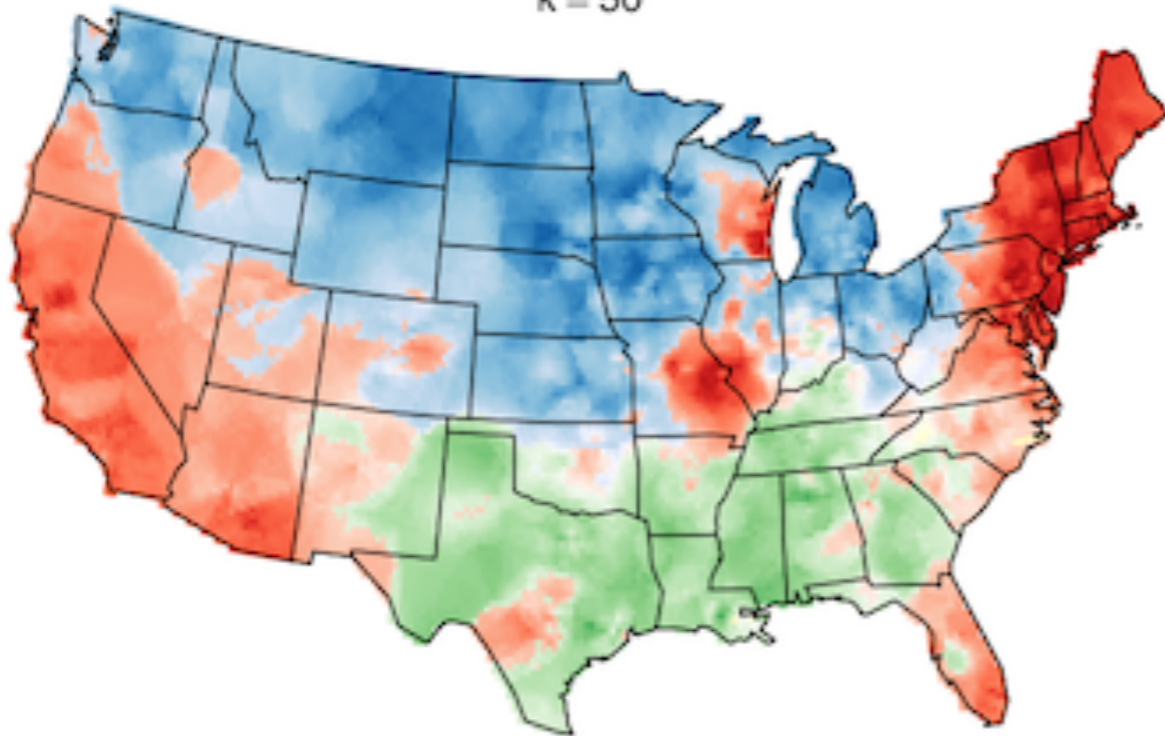
Hans Goebel (1982/2000/2012)



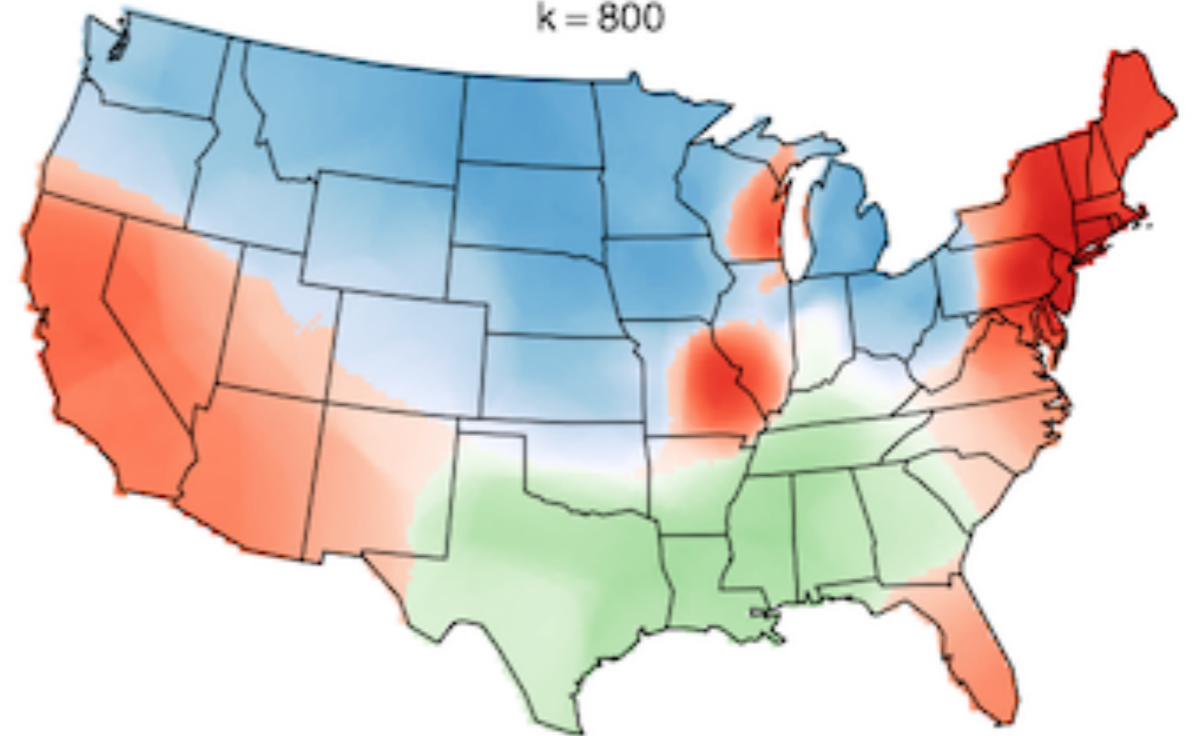
Katz, Joshua. 2013. Beyond “Soda, Pop or Coke”.
 Online at <<http://www4.ncsu.edu/~jakatz2/files/dialectposter.png>>

smoothing

k = 50



k = 800



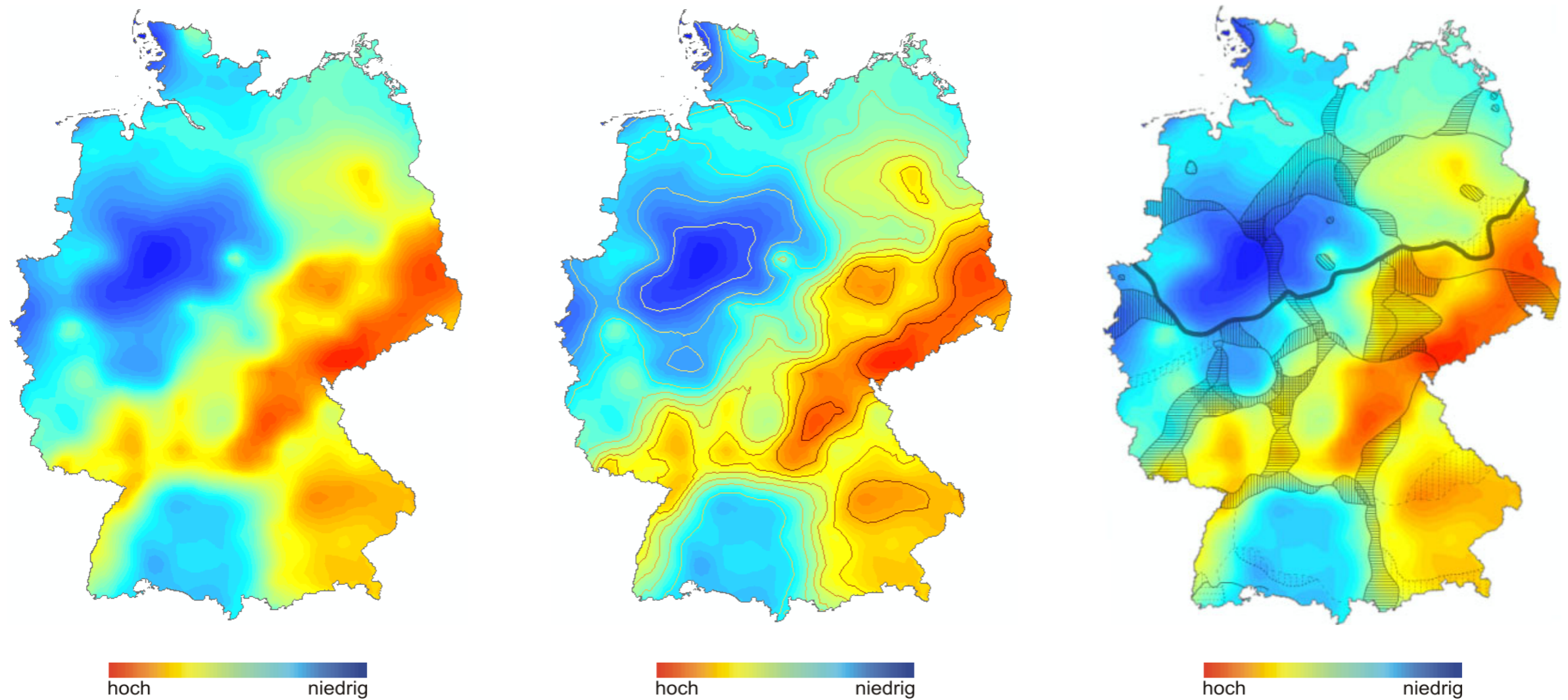
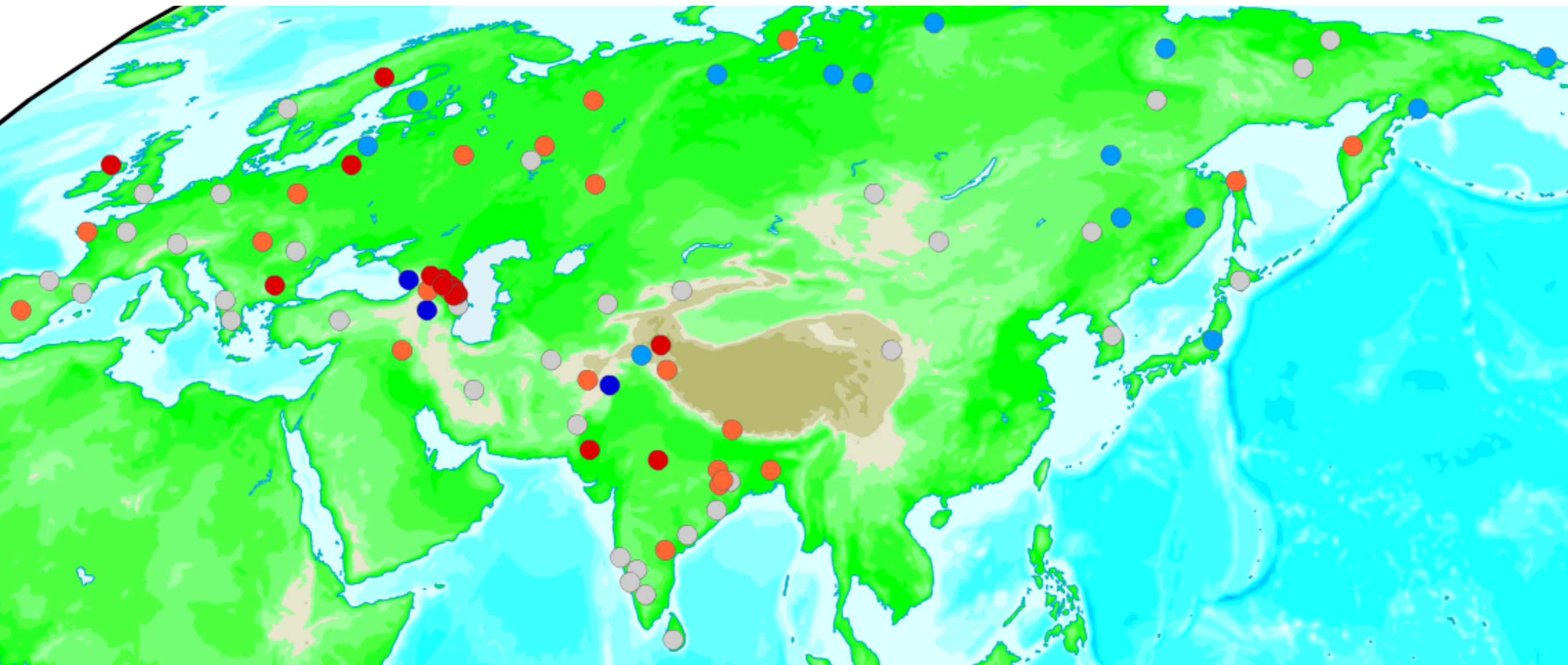


Abbildung 5-5: Geostatistische Interpolation der Similaritätswerte (Universal Kriging); links: Interpolationsergebnis; Mitte: definierte Konturlinien; rechts: Überblendung mit der Dialekteinteilungskarte von Wiesinger (1983a)

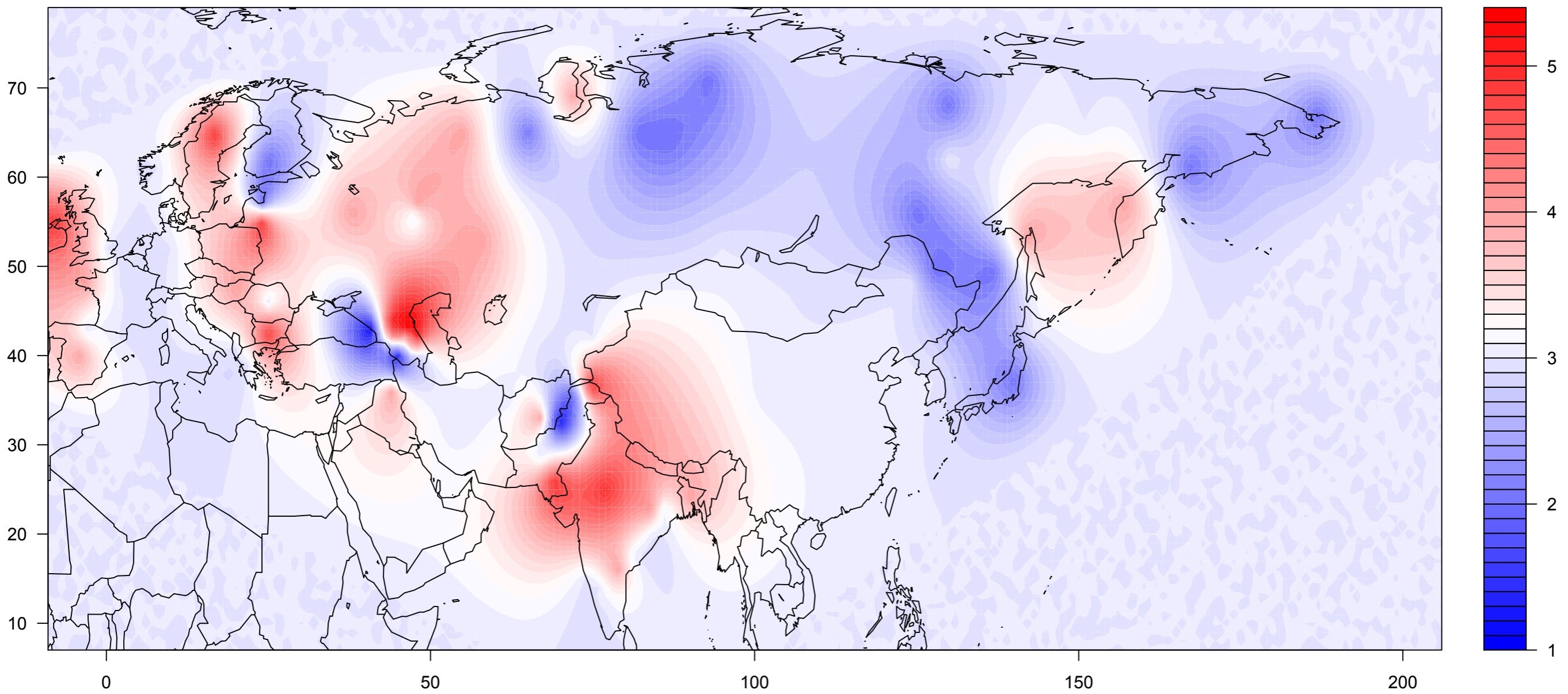
Consonant Inventories



few consonants

many consonants

Maddieson, Ian (2005) 'Consonant Inventories' in: Martin Haspelmath, Matthew S. Dryer, David Gil, & Bernard Comrie (eds.) *World Atlas of Language Structures*. Oxford: Oxford University Press.



- Interpolation (Zwischenräume ausfüllen)
 - ▶ Glättung ist notwendig
 - ▶ Implizite “Vorhersage” über die Zwischenräume

Ort der Mundart/Kreis Astfeld/ Gandersheim	Aufnahme-Nr. I/62	Transkribent Angelika Trauer	Listentyp A
Planrechteck x 29	Aufnahmedatum 20.11.1965	Transkribiert von 14.6.1985 bis 24.7.1985	
besprochen von 24.07.1985 25.07.1985 UStv			

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter	'vɪntə ^h ə	ɸ = kontinuant ə bereit velarisier
* 2	fliegen	'flaɪ→ə ^h ə	'fliegen die', Sequenzie- rung unklar - kein geminiertes [t]
3	Blätter	'blɛ:də ^h ə	
4	Luft	lu ^h ft ^h	ɸ = kontinuant
5	hört	hɪ ^h t ^h	ɸ = kontinuant
6	gleich	klɛ ^h →ɪk ^h	folgt P
7	schneien	'ʃnɛ→ɪən	
8	Wetter	/	statt demoa 'vɪtə ^h ʀɔŋə
9	tu	dɛ→u	

Ort der Mundart/Kreis Astfeld/ Günthersheim	Aufnahme-Nr. I/62	Transkribent Angelika Braun	Listentyp A
---	----------------------	--------------------------------	----------------

Phonetischer Atlas von Deutschland

- Wenker-sentences recorded in the 1960s (with additions in the 1970s)
- Selected words from the recordings were transcribed on paper in the 1980s
- A joint project between Marburg and Groningen digitised the data in the 2000s
- In total 29530 words distributed over 183 locations and 186 cognate sets

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter	178 vɪntɐ	k = kontinuant b = kontinuant
2	fliegen	56 flaɪ̯ən	'fliegen die', Sequenzierung unklar - kein genügendes [t]
3	Blätter	23 blɛtɐ	
4	Luft	103 lʊft	= kontinuant
5	hört	89 hœrt	= kontinuant
6	gleich	118 glɛɪ̯ç	
7	schneien	130 ʃnɛi̯ən	
8	Wetter	171 vɛtɐ	statt dem 'vɛtəʀɔŋə
9	tu	151 tu	

Ort der Mundart/Kreis Astfeld/ Gandersheim	Aufnahme-Nr. I/62	Transkribent Angelika Trauer	Listentyp A
Planrechteck X 29	Aufnahmedatum 20.11.1965	Transkribiert von 14.6.1985 bis 24.7.1985	

Besprochen von 24.07.1985
25.07.1985 UStv

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter 178	'vɪntə ^h ə	ɸ = kontinuant ə bereit velarisier
* 2	fliegen 56	'flaɪ→ə ^h ə	'fliegen die', Sequenzie- rung unklar - kein geminiertes [t]
3	Blätter 23	'blɛ:də ^h ə	
4	Luft 103	lu ^h ft ^h	ɸ = kontinuant
5	hört 89	hɪ ^h t ^h	ɸ = kontinuant
6	gleich 78	klɛ ^h →ɪk ^h	folgt P
7	schneien 130	ʃnɛ→ɪən	
8	Wetter 174	/	statt demoa 'vɪtə ^h ʀɔŋə
9	tu 151	dɛ→u	

Ort der Mundart/Kreis Astfeld/ Günthersheim	Aufnahme-Nr. I/62	Transkribent Angelika Braun	Listentyp A
---	----------------------	--------------------------------	----------------

Phonetischer Atlas von Deutschland

- Digitised in X-SAMPA, converted back to match original transcriptions, minor corrections for consistency of encoding
- The data is transcribed in high phonetic detail (3786 different phonetic segments)
- We will make the complete data available
 - ▶ electronically, separated by phonetic segments
 - ▶ as close as possible to the original source
 - ▶ including all idiosyncrasies

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter	vɪntə	b = kontinuierlich Plosiv, velarisches
56	Fliegen	'flɛgən	'flɛgən die', Sequenz kən
23	Bätter	'blɛ:dət	
103	Wust	'vʊst	b = kontinuierlich
89	hört	'hɔ:t	b = kontinuierlich
78	Klappert	'klɛpɛrt	folgt P
130	Glober	'glɔbɛr	
174	Wetter	'vɛtɛr	statt dem 'vɛtɛrɔgɔ
151	tu	dɛv	

Multiple Sequence Alignment

- Just a fancy name for sound correspondences
- Each sound correspondence is “aligned” in a column, possibly adding empty cells
- It is a useful and consistent way to represent comparative data (both between languages or dialects)

LOCATION	WORD
Aachen	a:ph
Adorf	ɑ:b ^h ə
Ahrbergen	o̘→ɔ̘phə
Albersloh	ɑ:p ^h ə
Allna	ɑϕh
Altenberg	ʌfɛ
Altentrüdin	af
Altlandsberg	ɑ'fə'
Altwarp	o:ph
Astfeld	ɒ':p ^h ə
Atzendorf	afɛ
Ballhausen	ʌ'fə
Bardenfleth	ɔ̘:p̄ϕ
Barssel	ɒ:p ^h ə
Bempflingen	af:
Bennin	ɔ̘p ^h
Billingsbach	af
Bockelwitz	ʌvə
Bonn	ɑ:p'
Borstendorf	ʏf:
Breddin	ɒ:ph
Brelingen	ɑfβ̄ə
Bremscheid	ɒ':p̄h̄ə
...	...

A	FF	E
a:	ph	-
ɑ:	b ^h	ə
o̘→ɔ̘	ph	ə
ɑ:	p ^h	ə
ɑ	ϕh	-
ʌ	f	ɛ
a	f	-
ɑ'	f	ə'
o:	ph	-
ɒ':	p ^h	ə
a	f	ɛ
ʌ'	f	ə
ɔ̘:	p̄ϕ	-
ɒ:	p ^h	ə
a	f:	-
ɔ̘	p ^h	-
a	f	-
ʌ	v	ə
ɑ:	p'	-
ʏ	f:	-
ɒ:	ph	-
ɑ	f̄β̄	ə
ɒ':	p̄h̄	ə
...

- **Workflow:**

- ▶ Tokenisation of segments
- ▶ Automatic alignment using LingPy (github.com/lingpy)
- ▶ Manual correction
- ▶ Separation of cognates (e.g. *Samstag* vs. *Sonnabend*)
- ▶ Annotation of columns (e.g. many-to-one alignments, metathesis)
- ▶ Merging of complex columns and removing boundaries

Automatic Isoglossing

- Each column represents an isogloss
- Plotting one column gives a dialect map (after some interpretation of the data)
- Areas can be added by using 3D-interpolation (or any other of the many methods being developed currently)

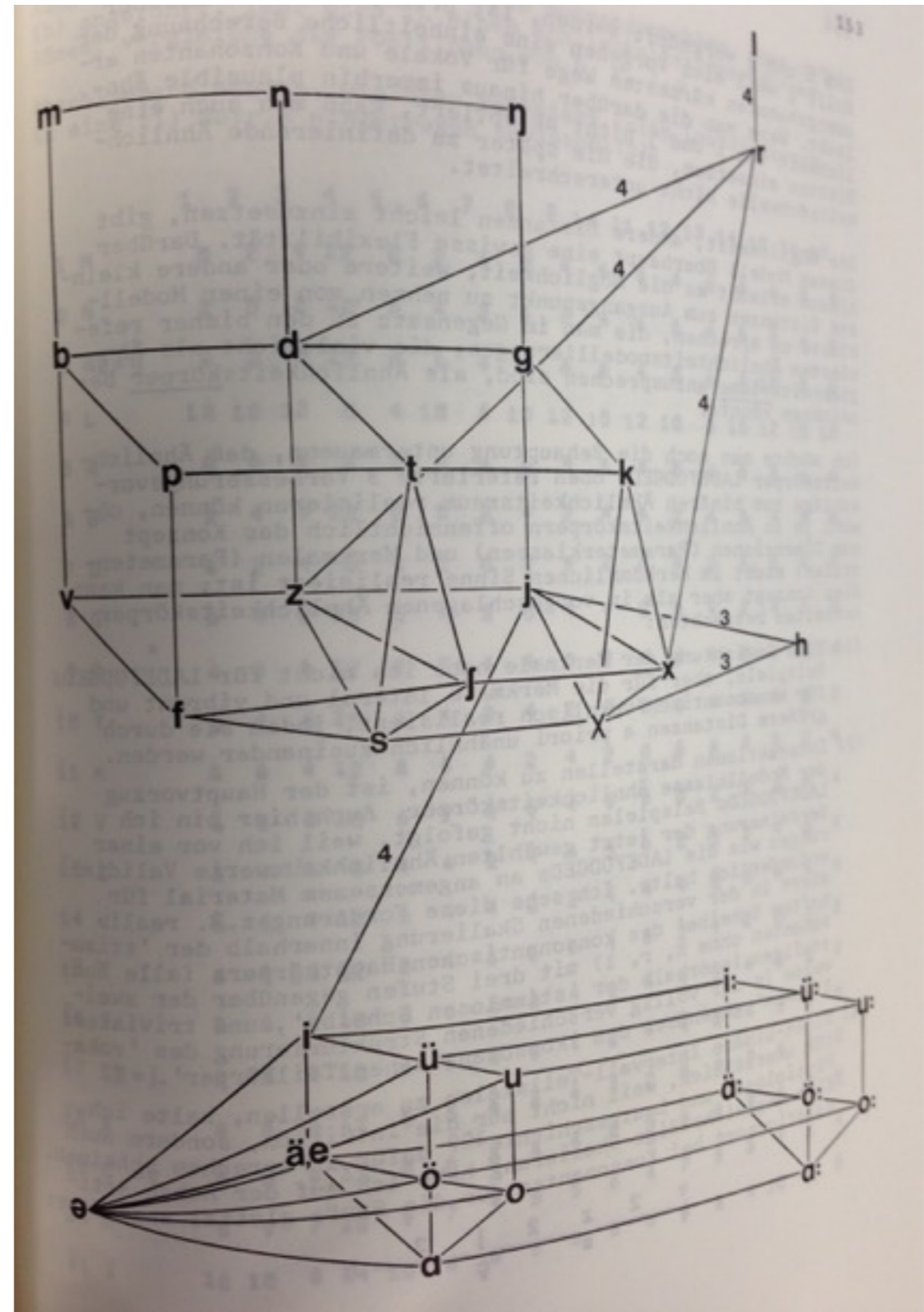
Example 1

- **Consonant** in German Perfect-prefix ge-
- Ordered to ‘strength’, visualised as ‘height’
- (NULL) ʏ ɪ j̥ j̥ j̥ ʔ ʧ̥ ʒ̥ ʒ̥ ʒ̥ ʧ̥ ʏ̥ ʧ̥ ʧ̥ ʏ̥ ʧ̥
 ʧ̥ ʏ̥ t̥ ɡ̥ x̥ ɡ̥ ʒ̥ ʒ̥ ʒ̥ ʧ̥ ʧ̥ ʧ̥ ʧ̥ ʧ̥ ʧ̥
 k̥ k̥ ʀ̥ χ̥ k̥ k̥ q̥ k̥'

Character Model

- ▶ Needleman, SB & CD Wunsch. 1970. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology* 48(3). 443-453.
- ▶ Widely used in bioinformatics ever since
- ▶ Multiple time reinvented in linguistics
- ▶ I will use a 'linear' character model, which is easier to visualise. For 'nominal' character models, multiple maps have to be combined

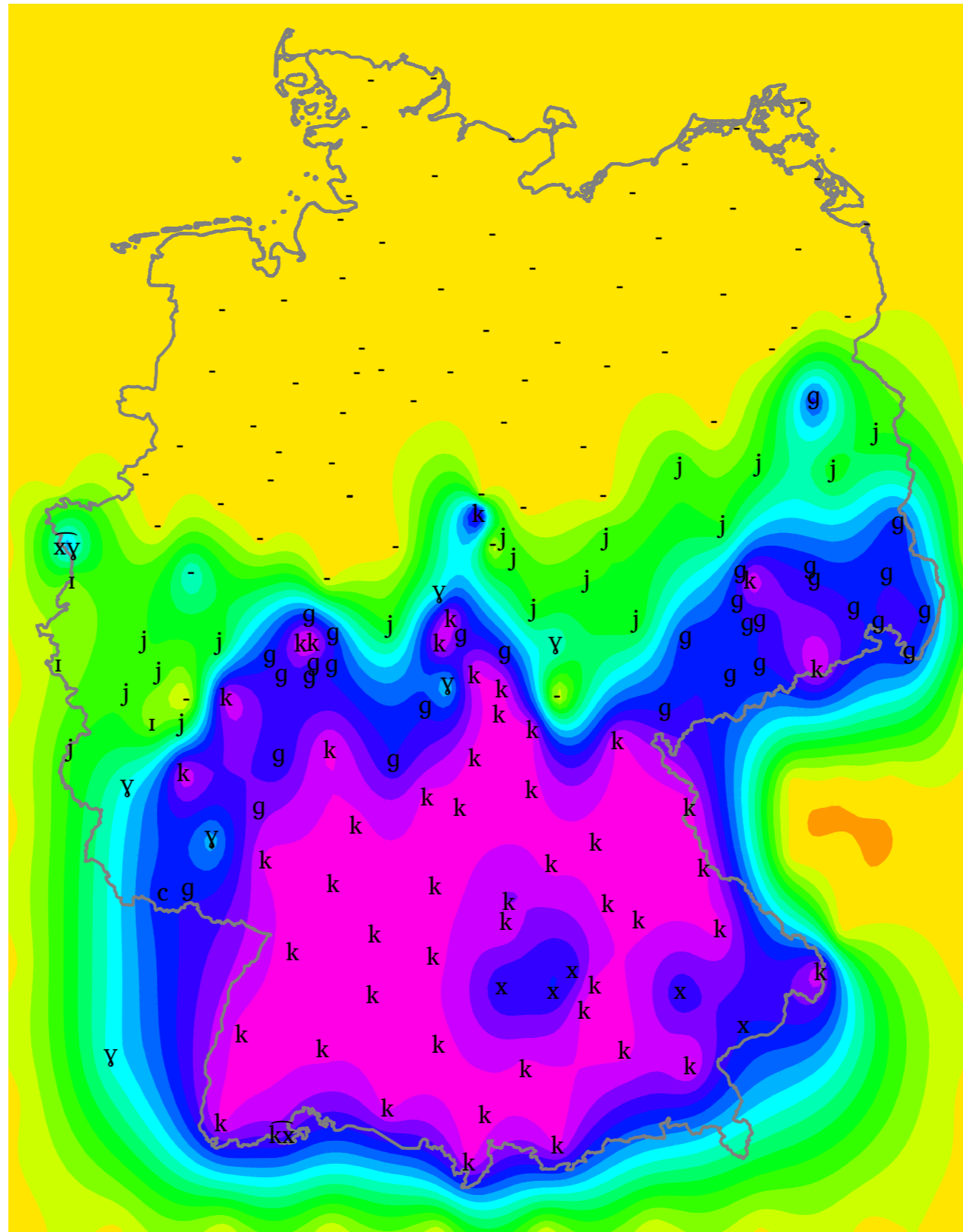
Naumann, Carl Ludwig. 1976. Grundzüge der Sprachkartographie und ihrer Automatisierung. *Germanistische Linguistik* 76(1-2). 1-285.



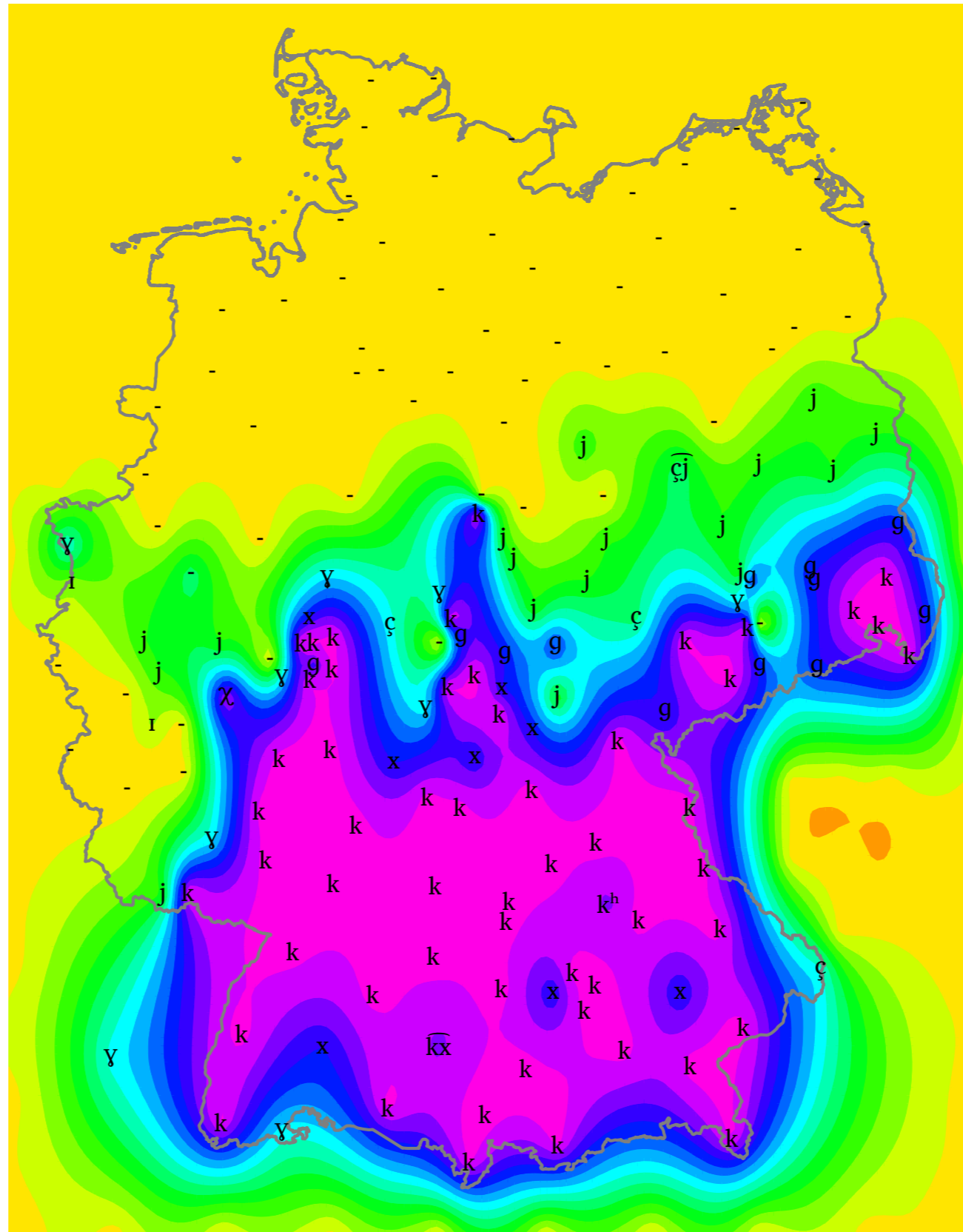
Example 1

- **Consonant** in German Perfect-prefix ge-
- Ordered to 'strength', visualised as 'height'
- (NULL) ʏ ɪ j̥ j̄ j̄ ʔ ʧ̣ ʒ̣ ʒ̣ ʒ̣ ʧ̣ ʏ̣ ʧ̣ ʧ̣ ʏ̣ ʧ̣
 ʧ̣ ʏ̣ ṭ̣ ɡ̣̣ x̣̣ ɡ̣̣ x̣̣ ʧ̣̣ ɡ̣̣ ḳ̣ ʧ̣̣ ḳ̣ x̣̣ ʧ̣̣ x̣̣ k^h
 ḳ̣ x̣̣ ḳ̣ ʀ̣̣ ʧ̣̣ ḳ̣ ḳ̣ q̣̣ ḳ̣'

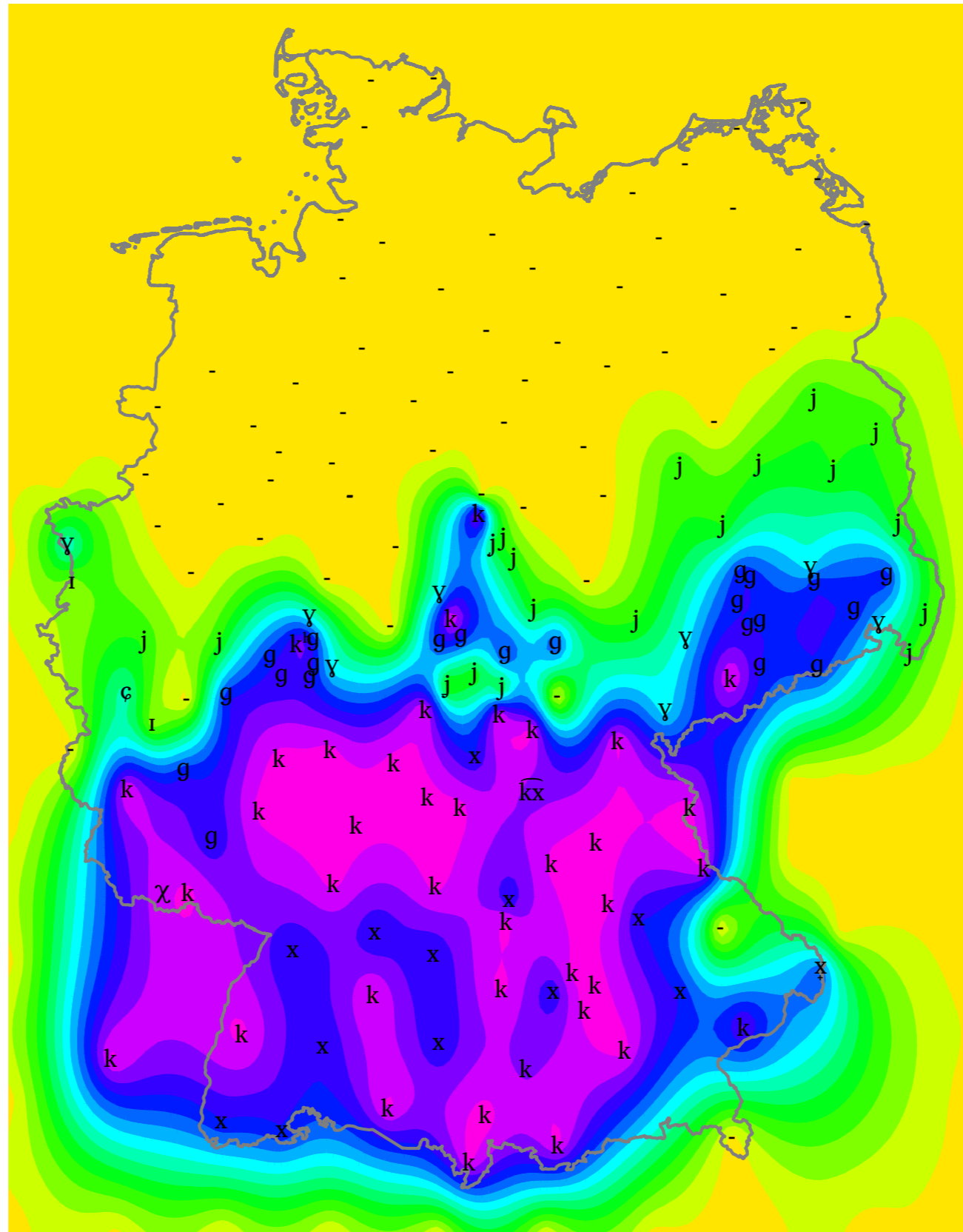
gefahren



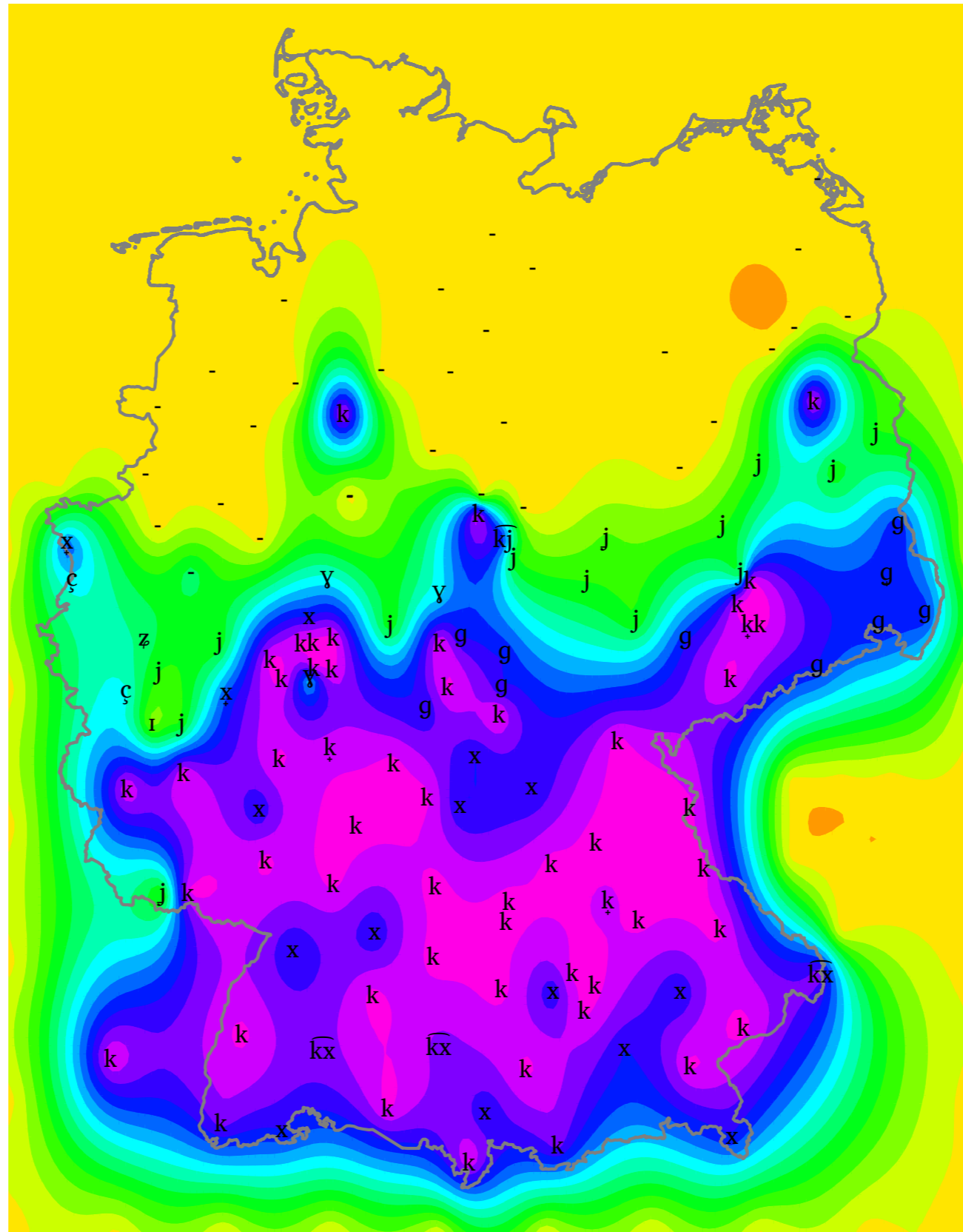
gefunden



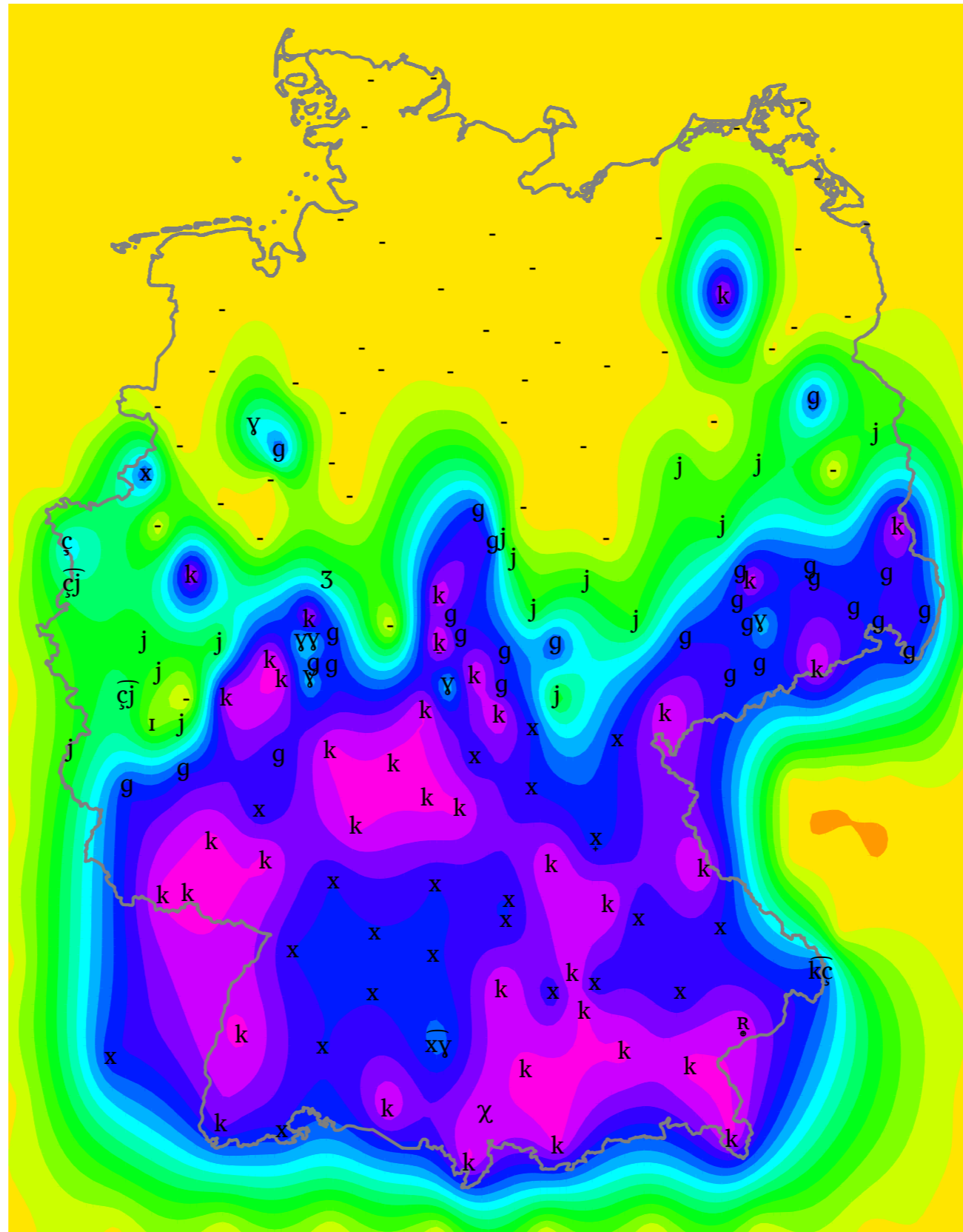
eingeschlafen



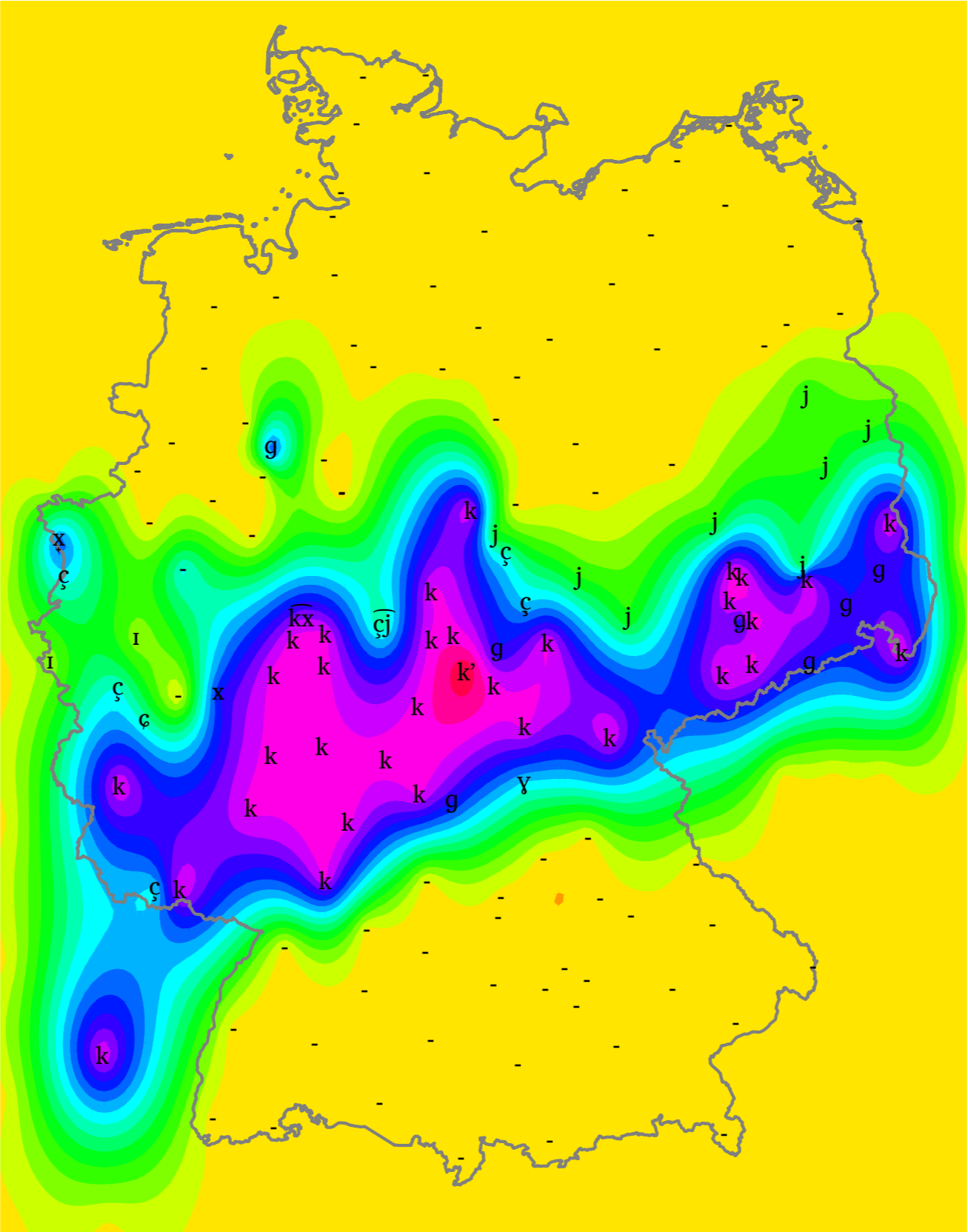
gestohlen



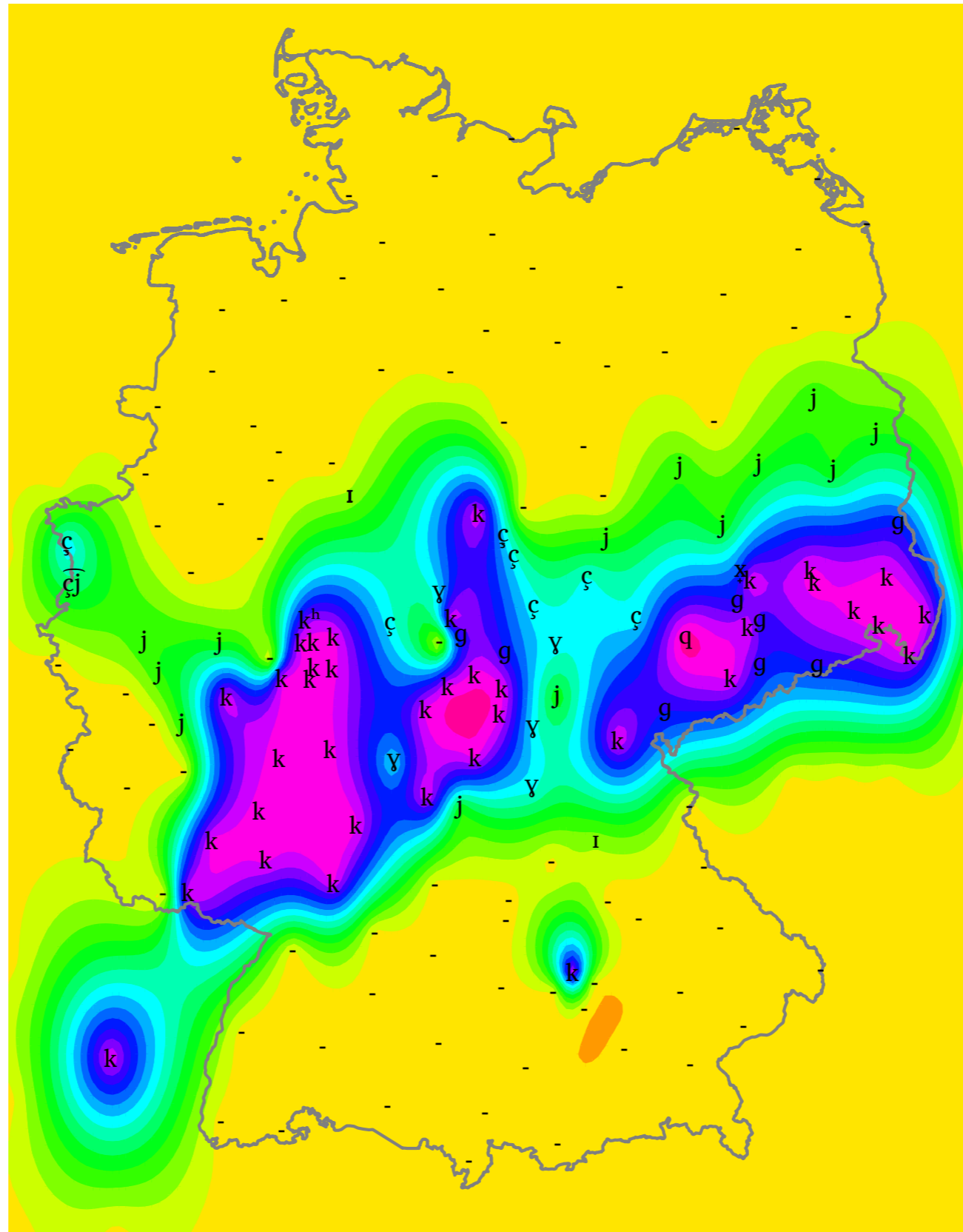
gestorben



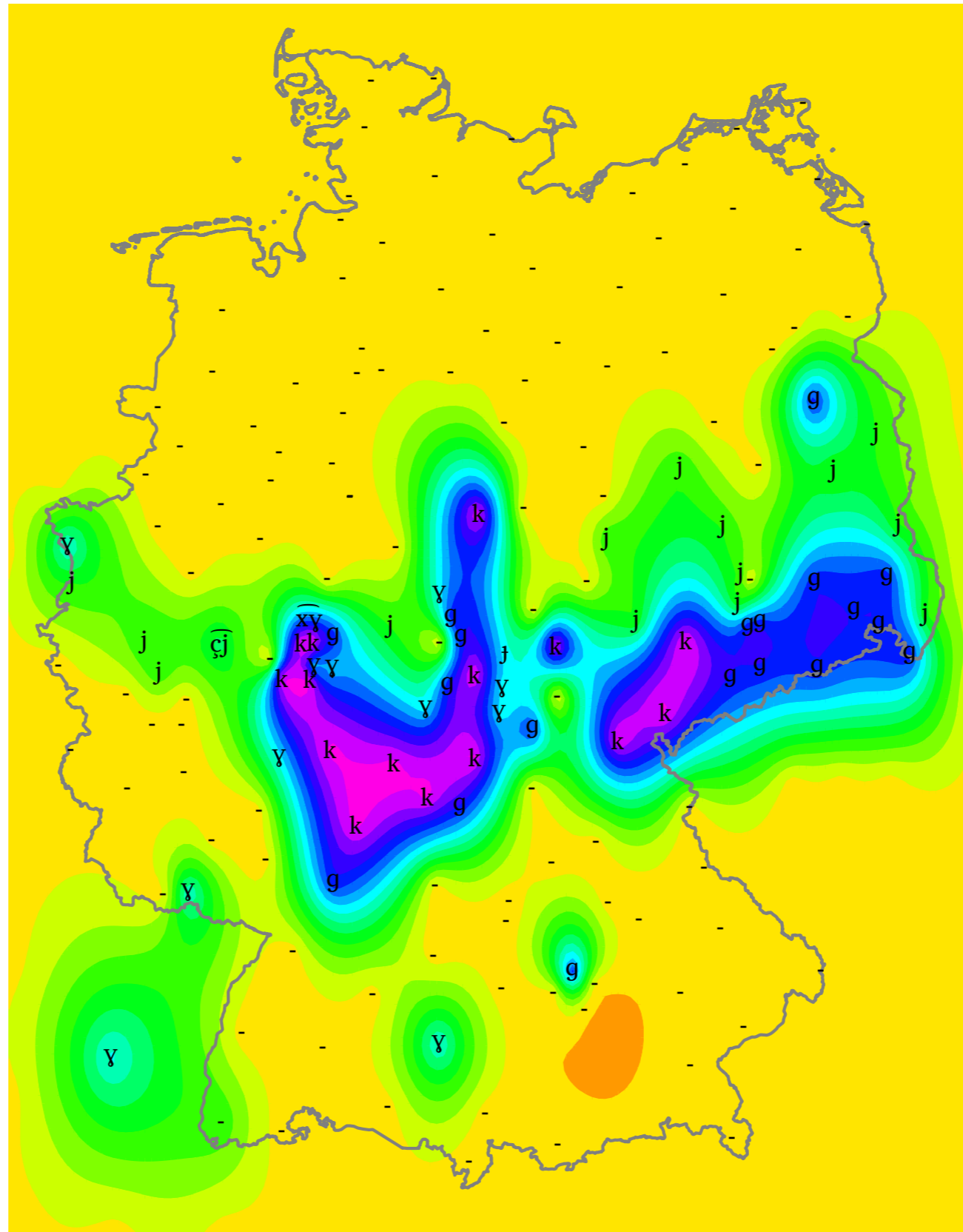
gebrannt



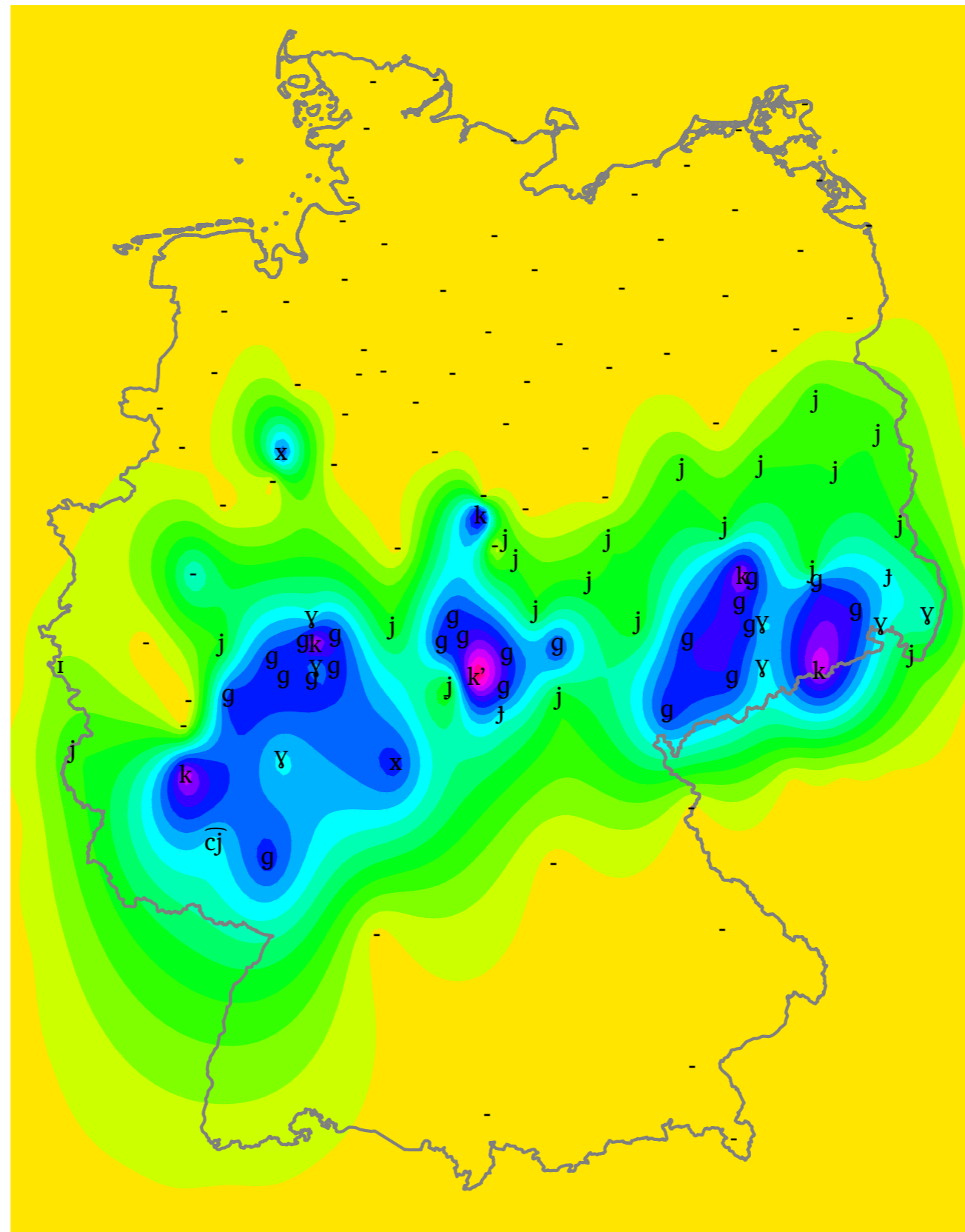
gebracht



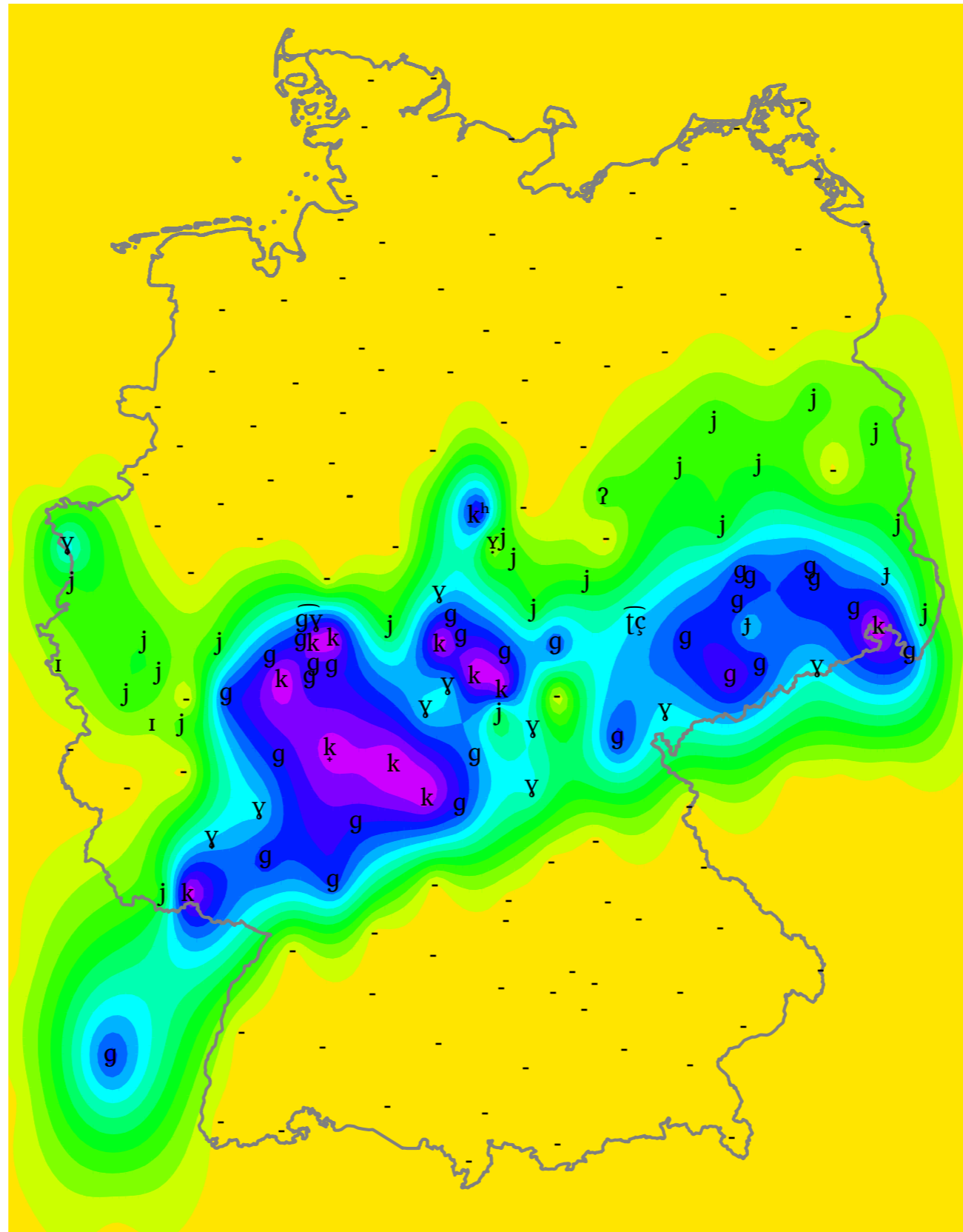
geblieben



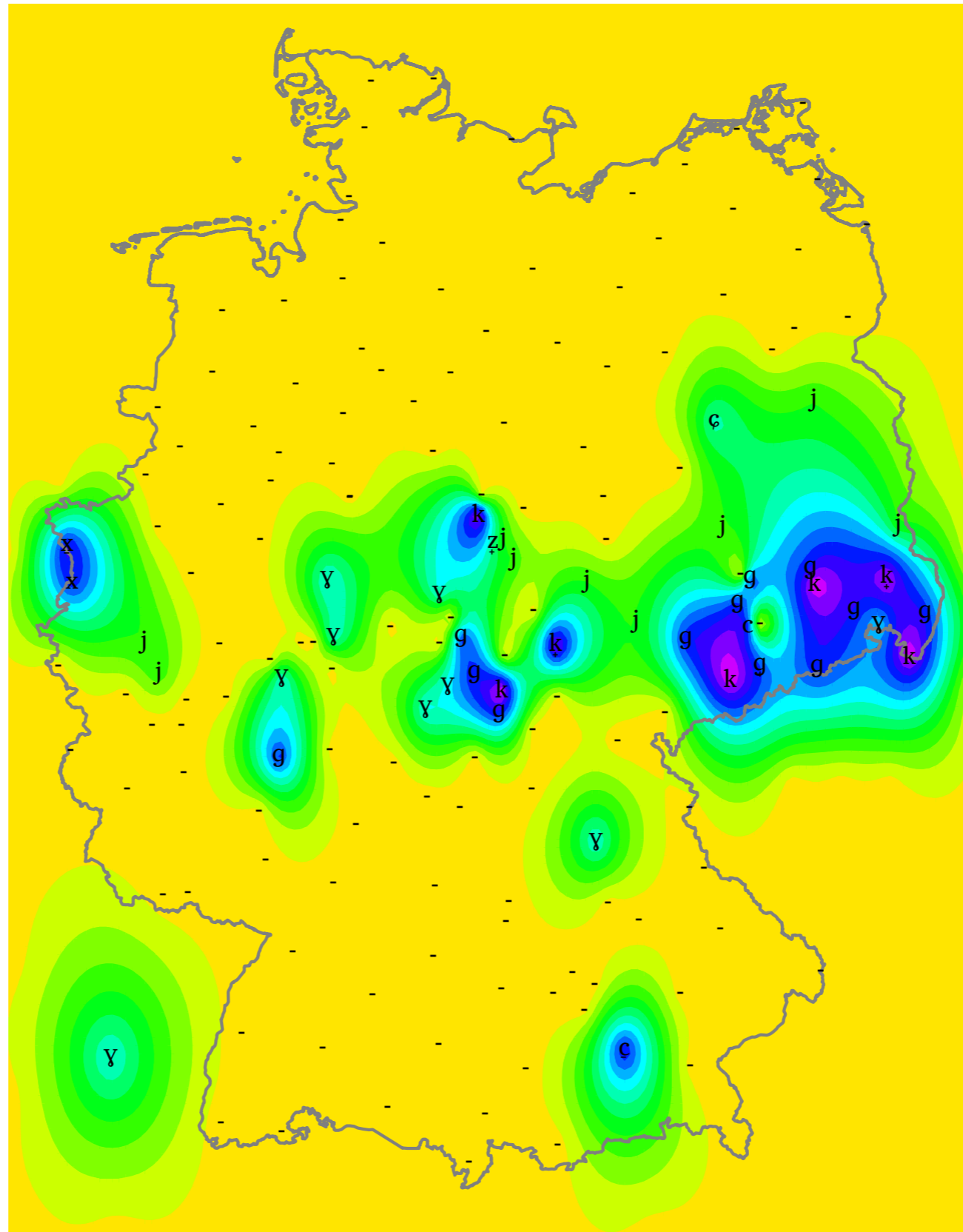
eingebrochen



gekannt



gekommen



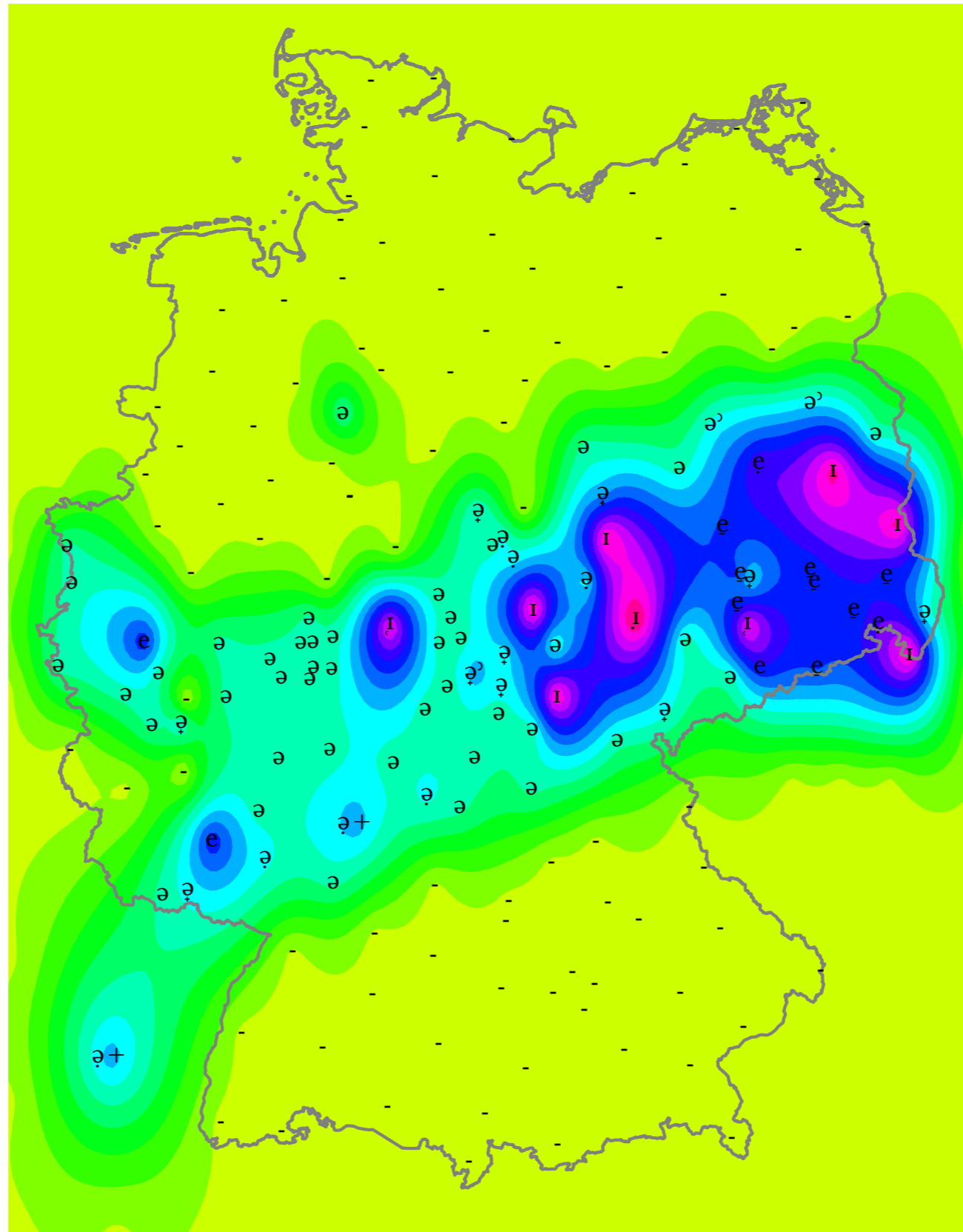
Example 2

- **Vowel** in German Perfect-prefix ge-
- Ordered to 'strength', visualised as 'height'

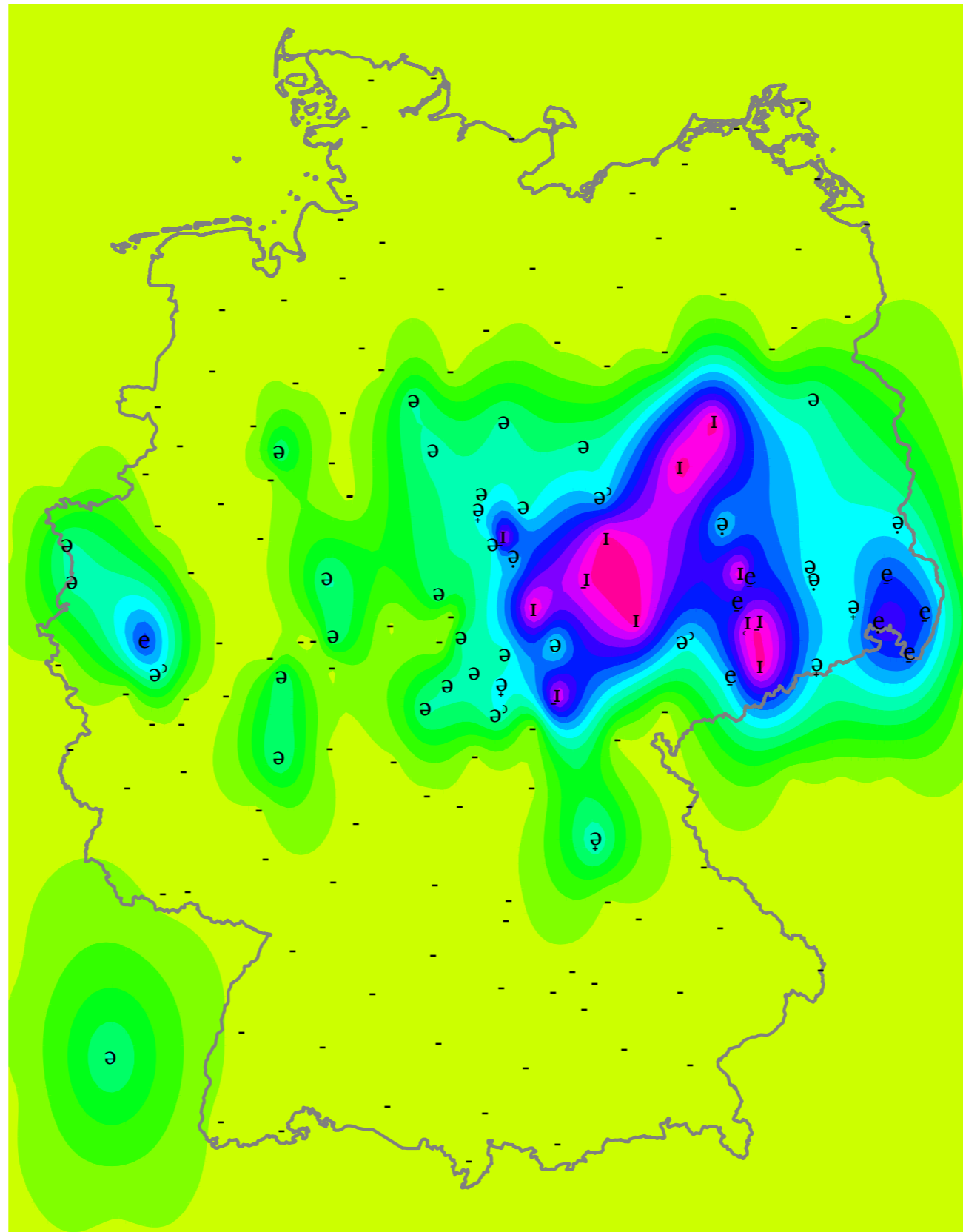
- (CLUSTER) (NULL) ə ə ə ə^c ə ə ə^c ə ə₊ ə₊

ə + ə + ə^c ə^c e e e^c e œ^c œ ε ε ɪ ɪ ɪ ɪ ɪ

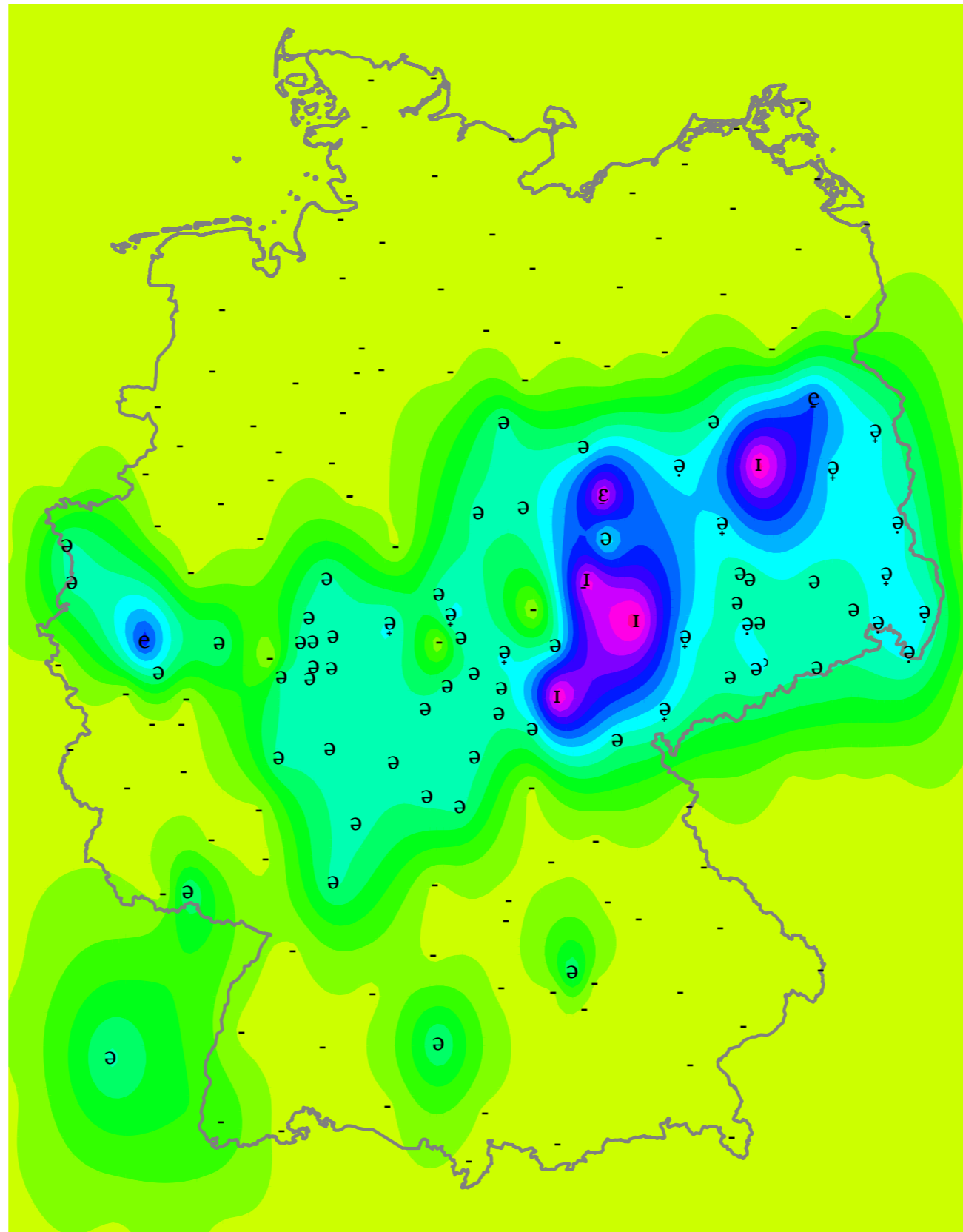
gekannt



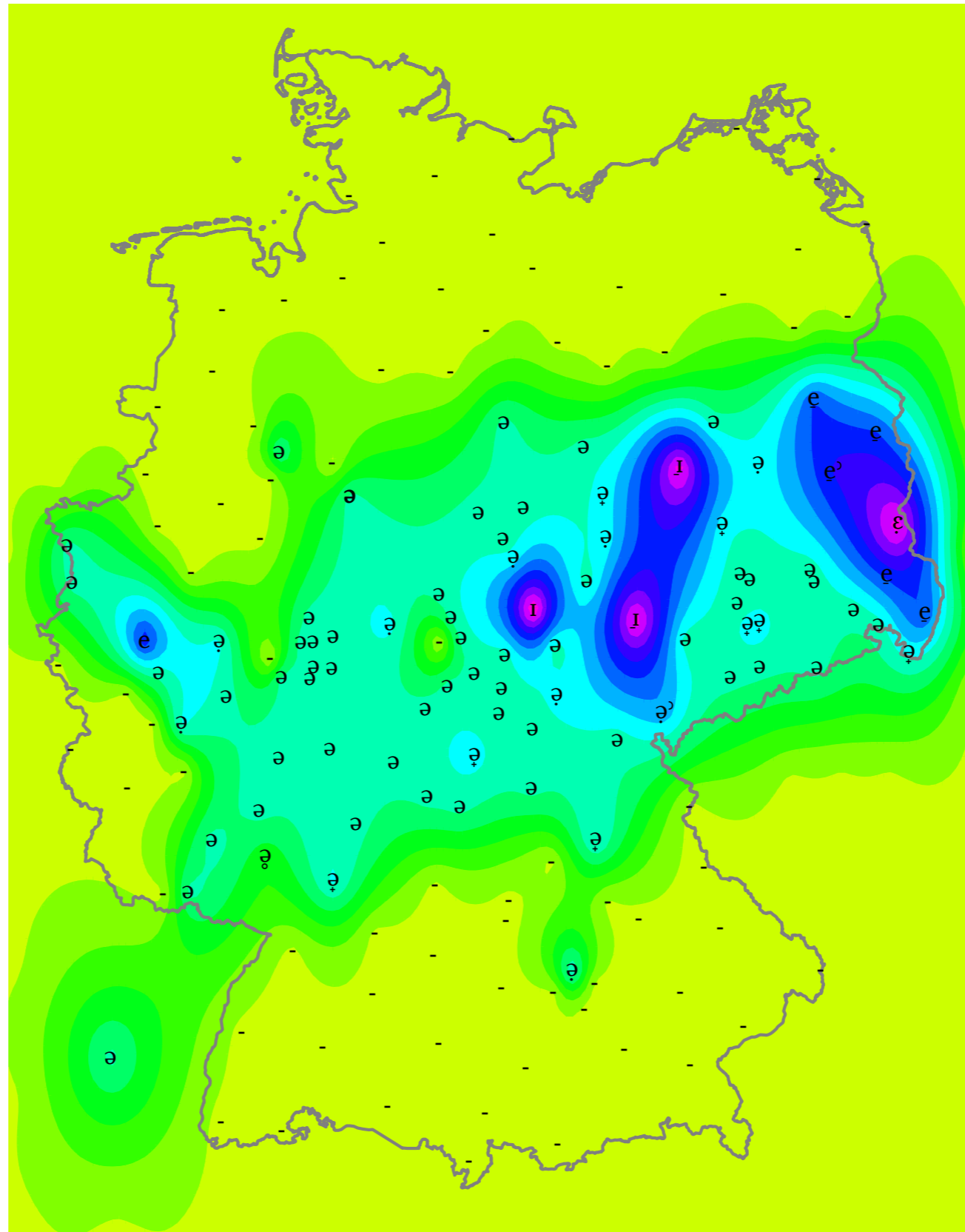
gekommen



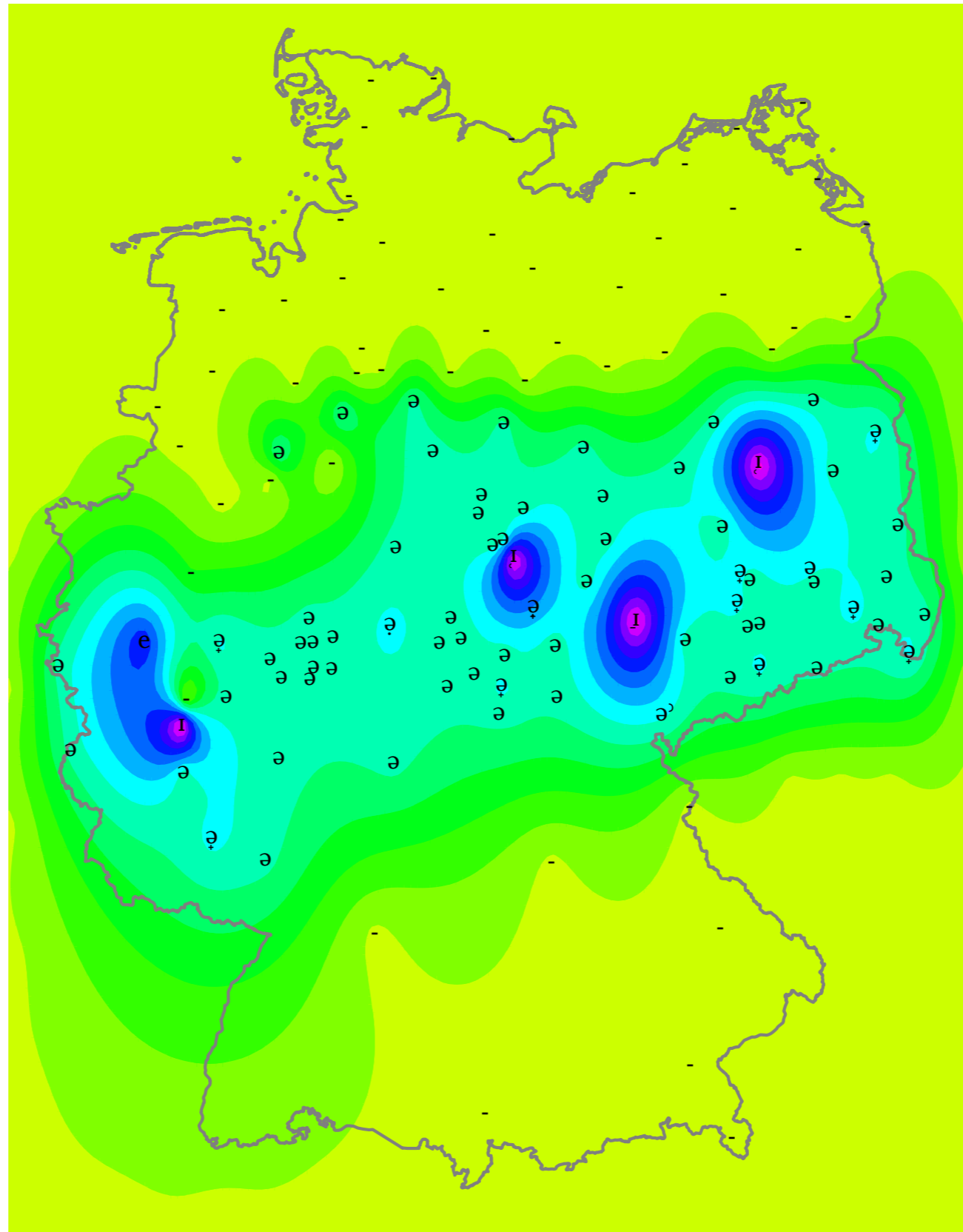
geblieben



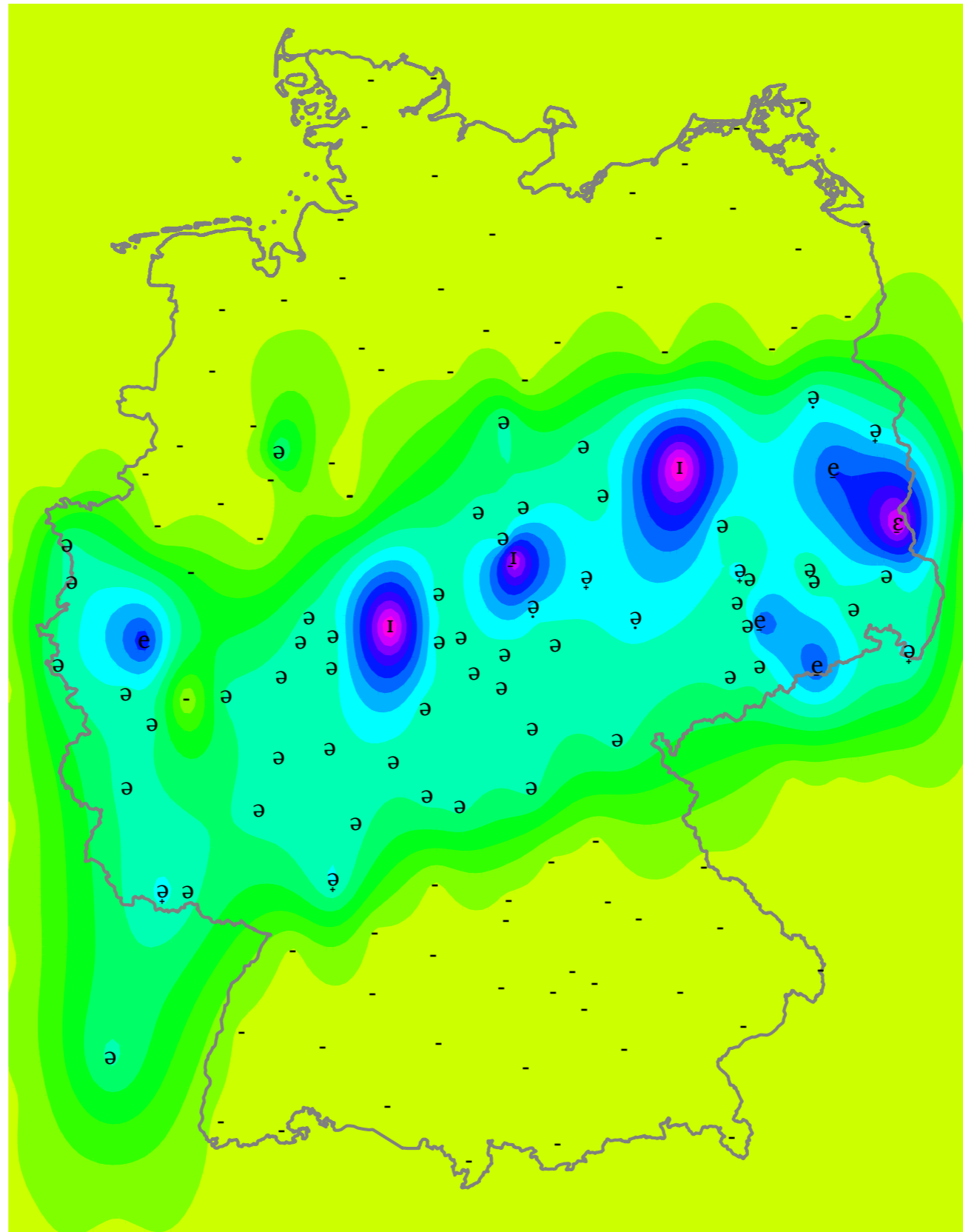
gebracht



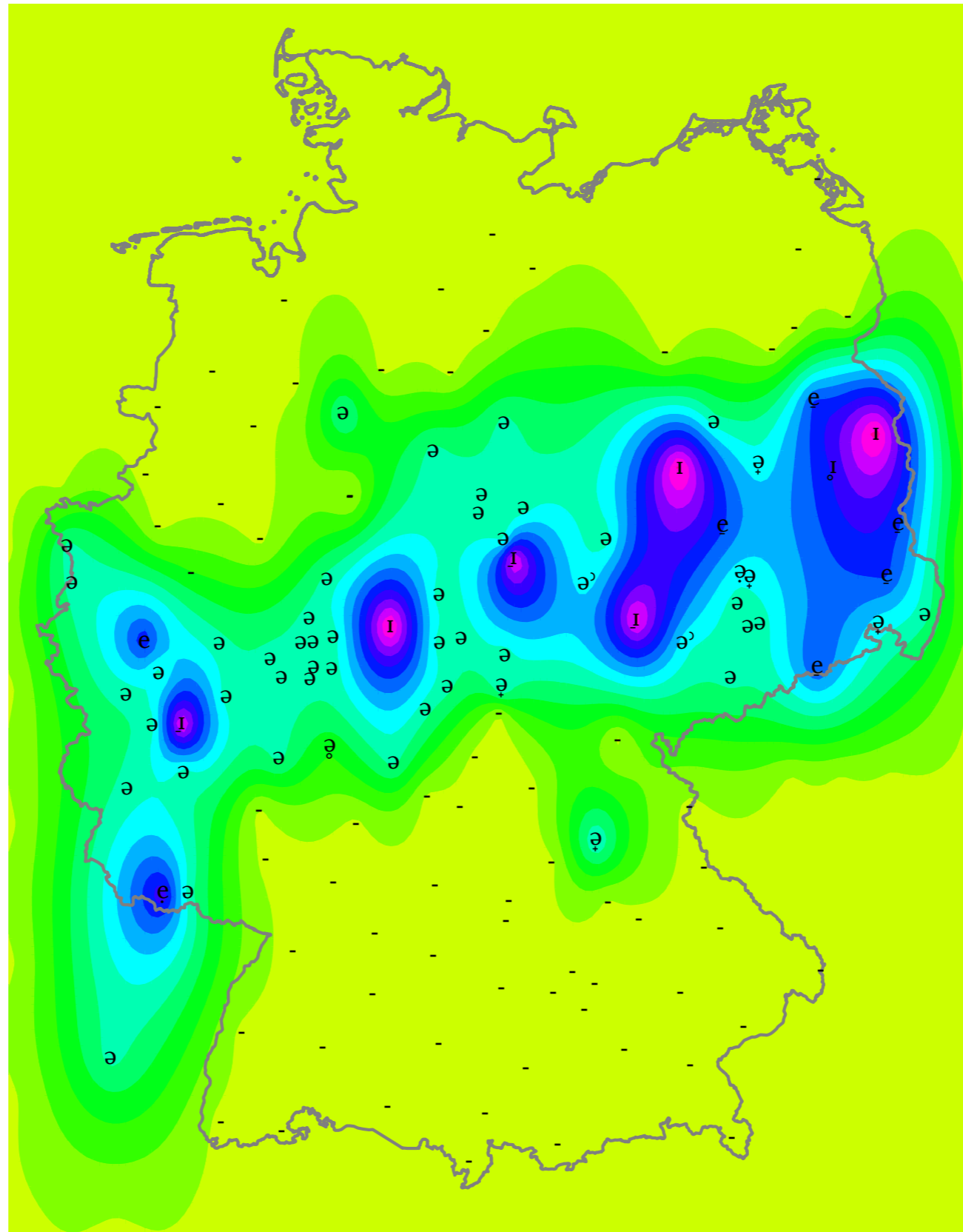
eingebrochen



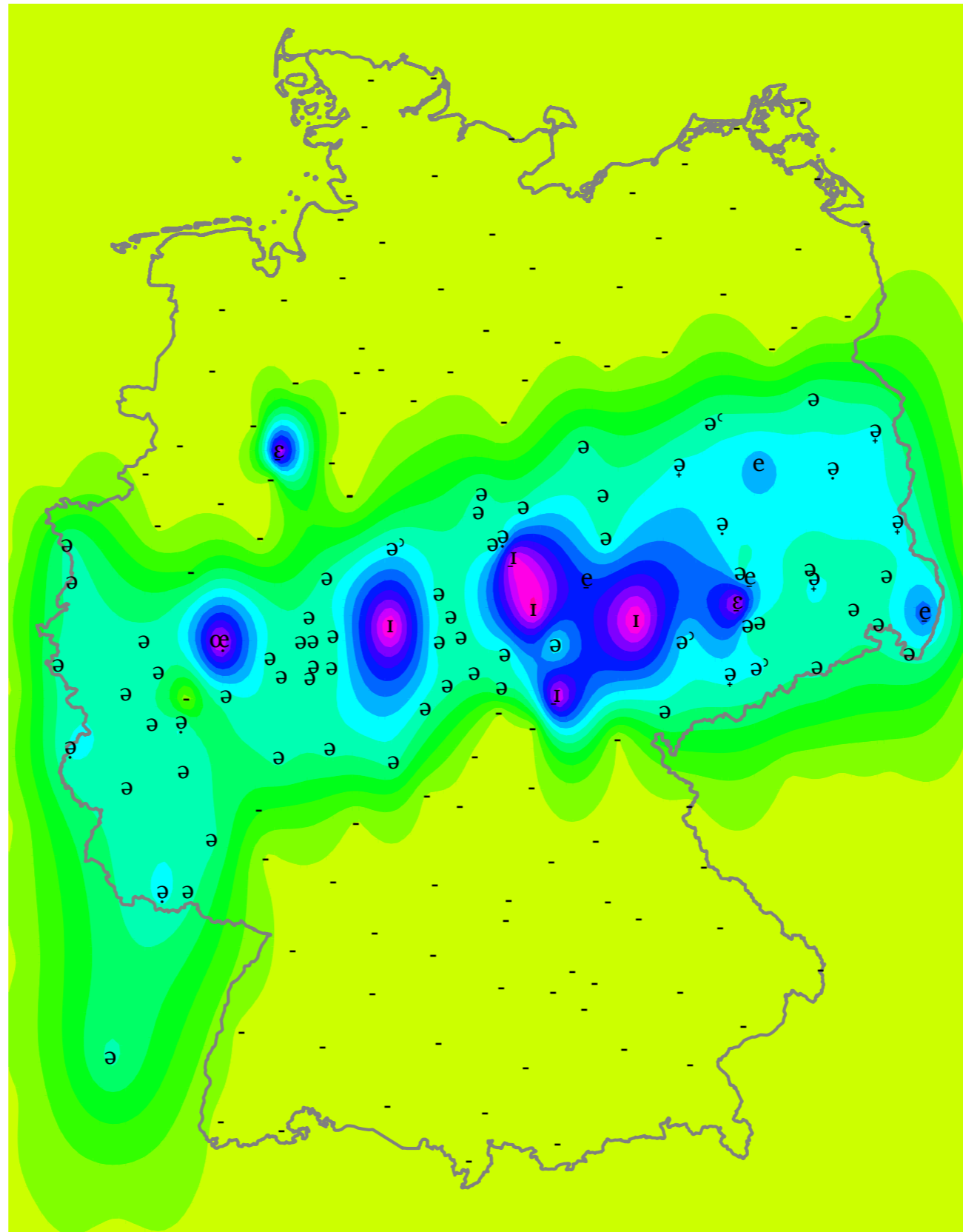
gebrannt



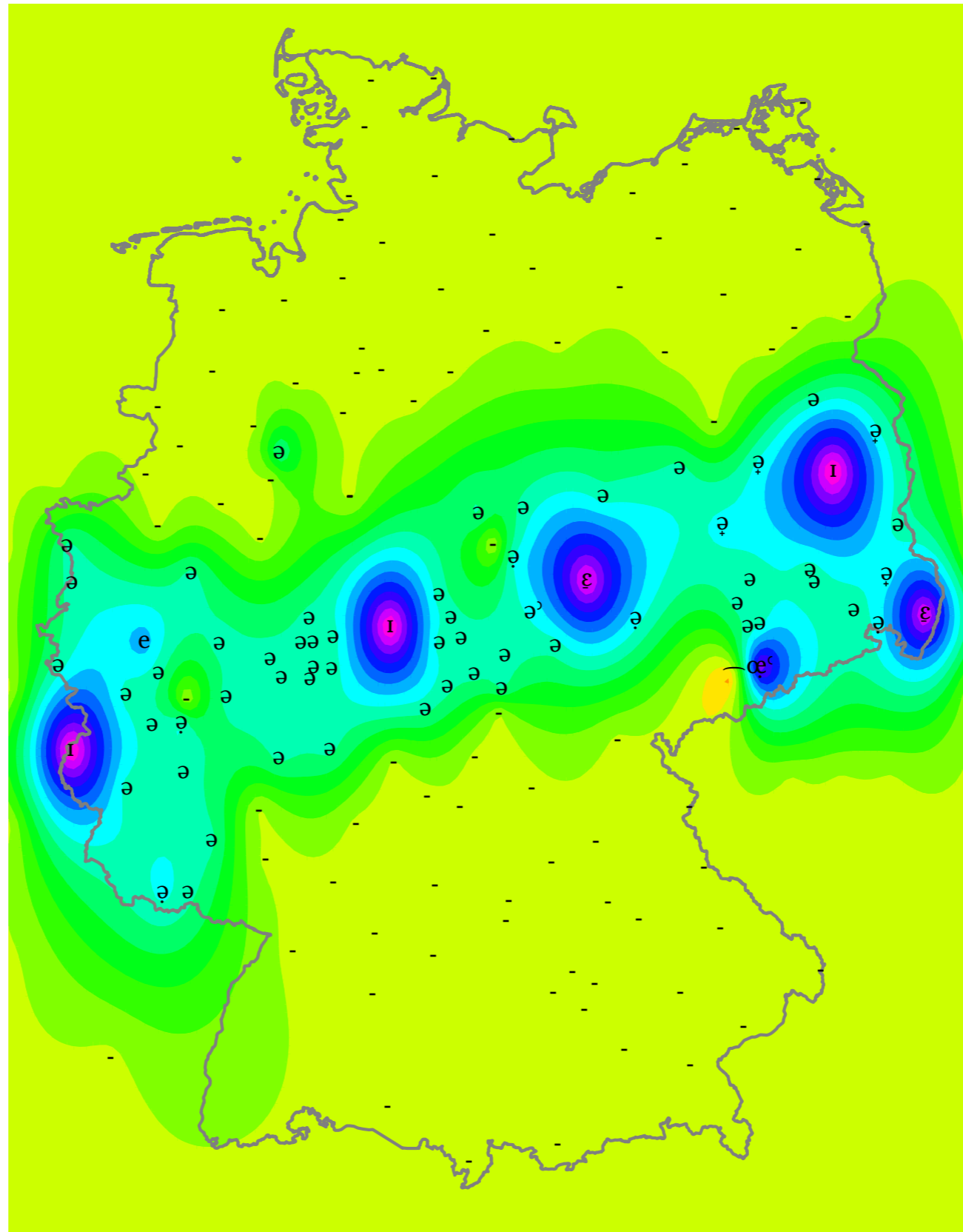
gestohlen



gefahren



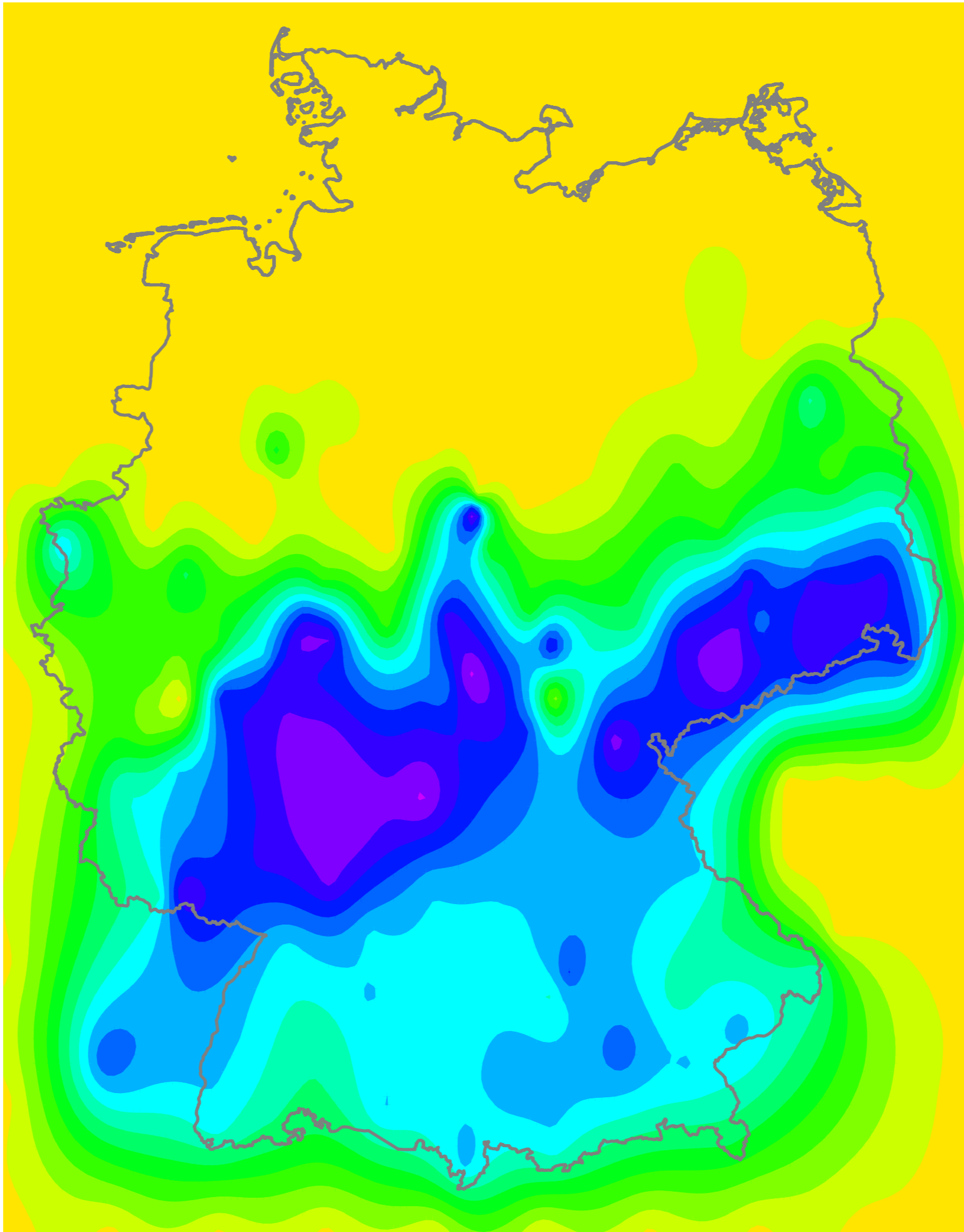
gefallen



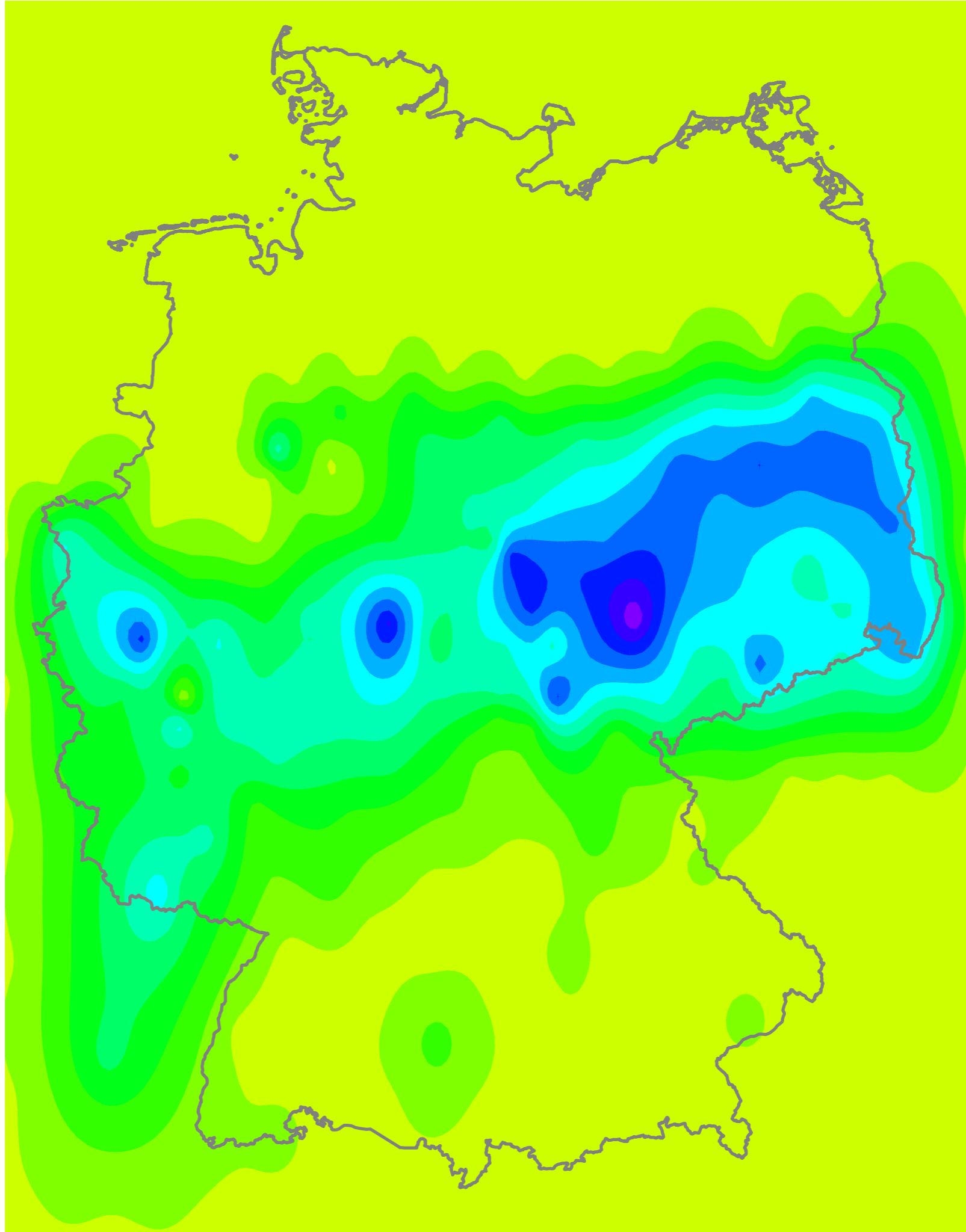
Averaging

- Combining all consonants and all vowels into a single map

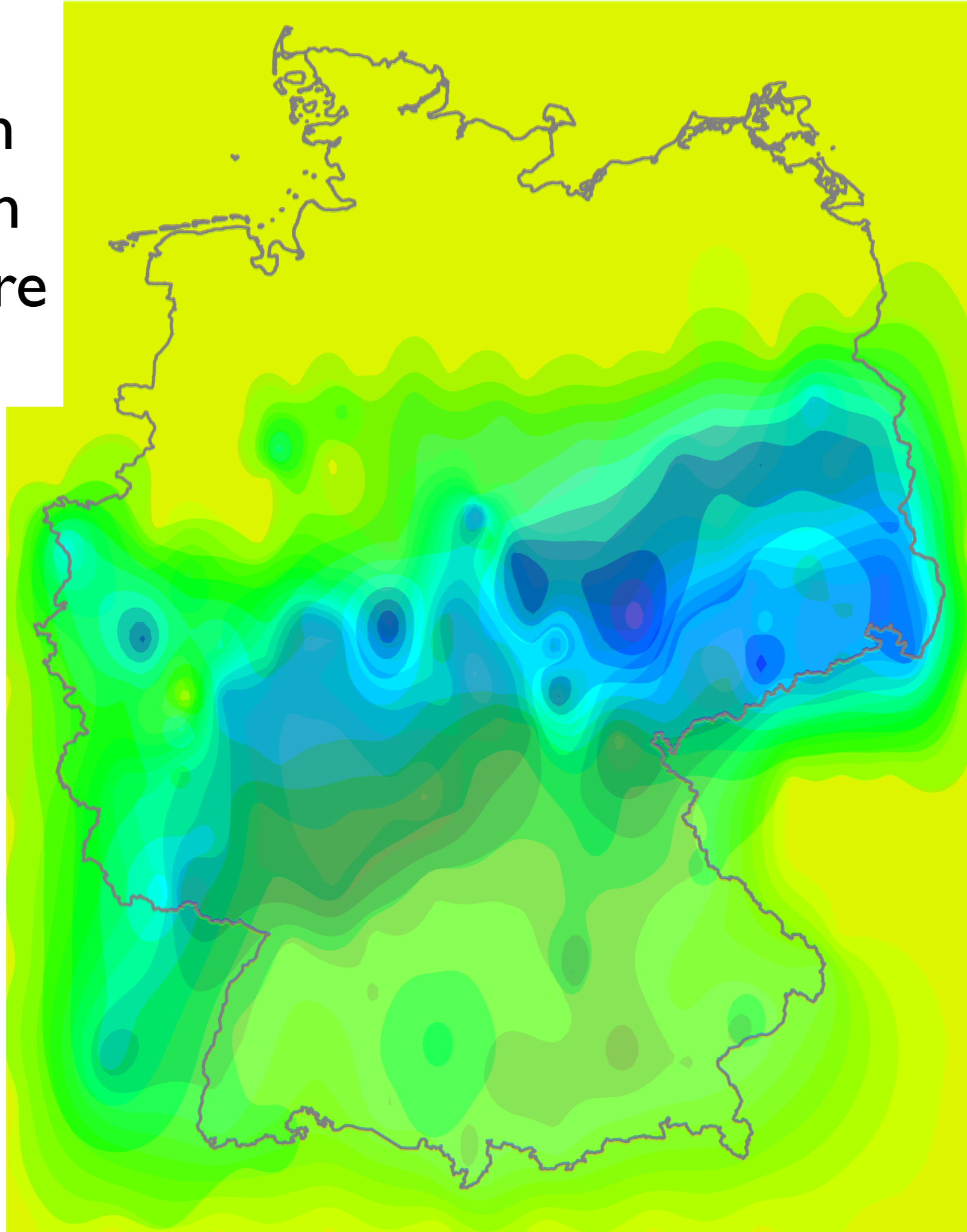
Consonants
of ge- prefix
(averages over
14 lexemes)

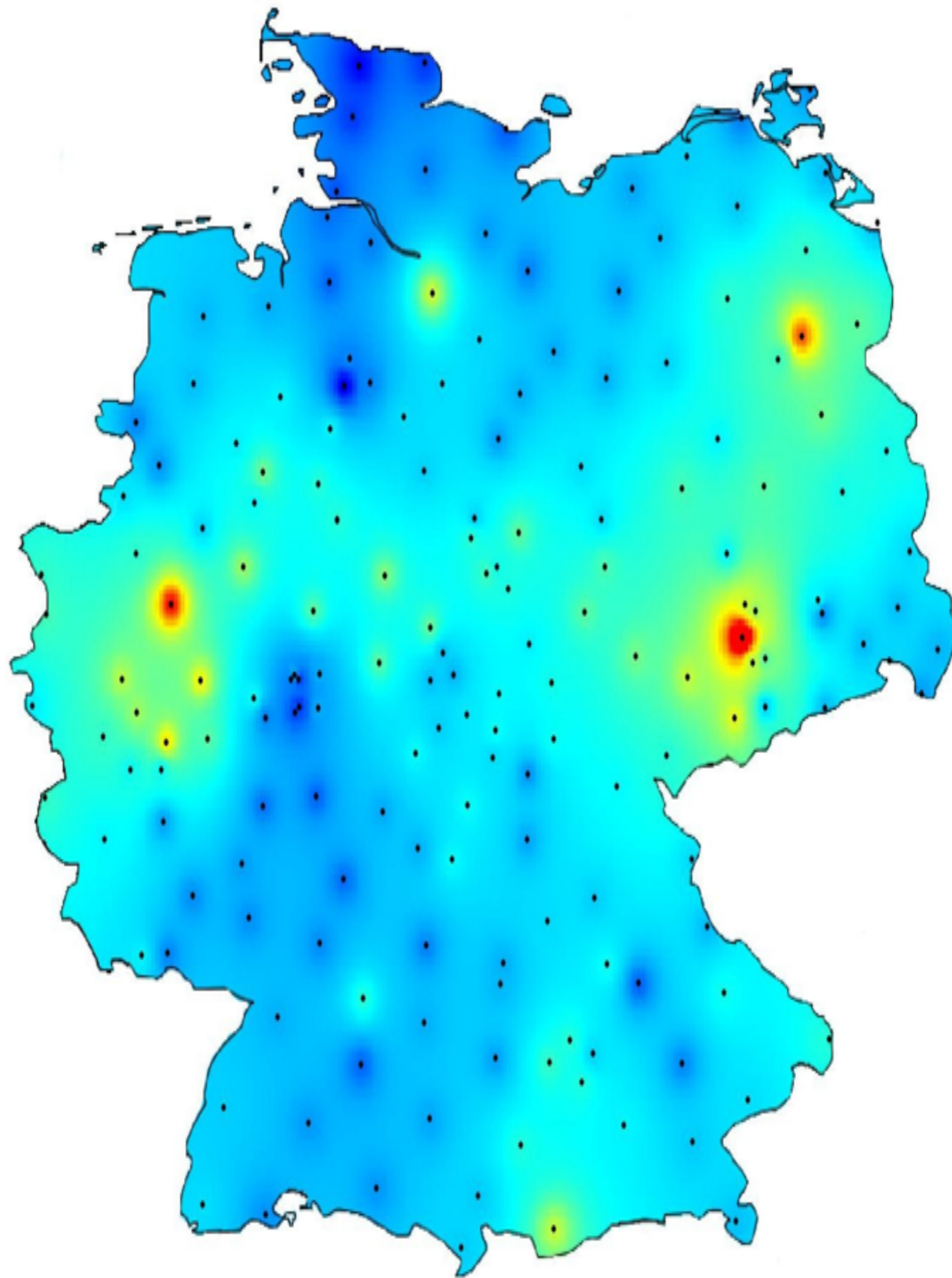


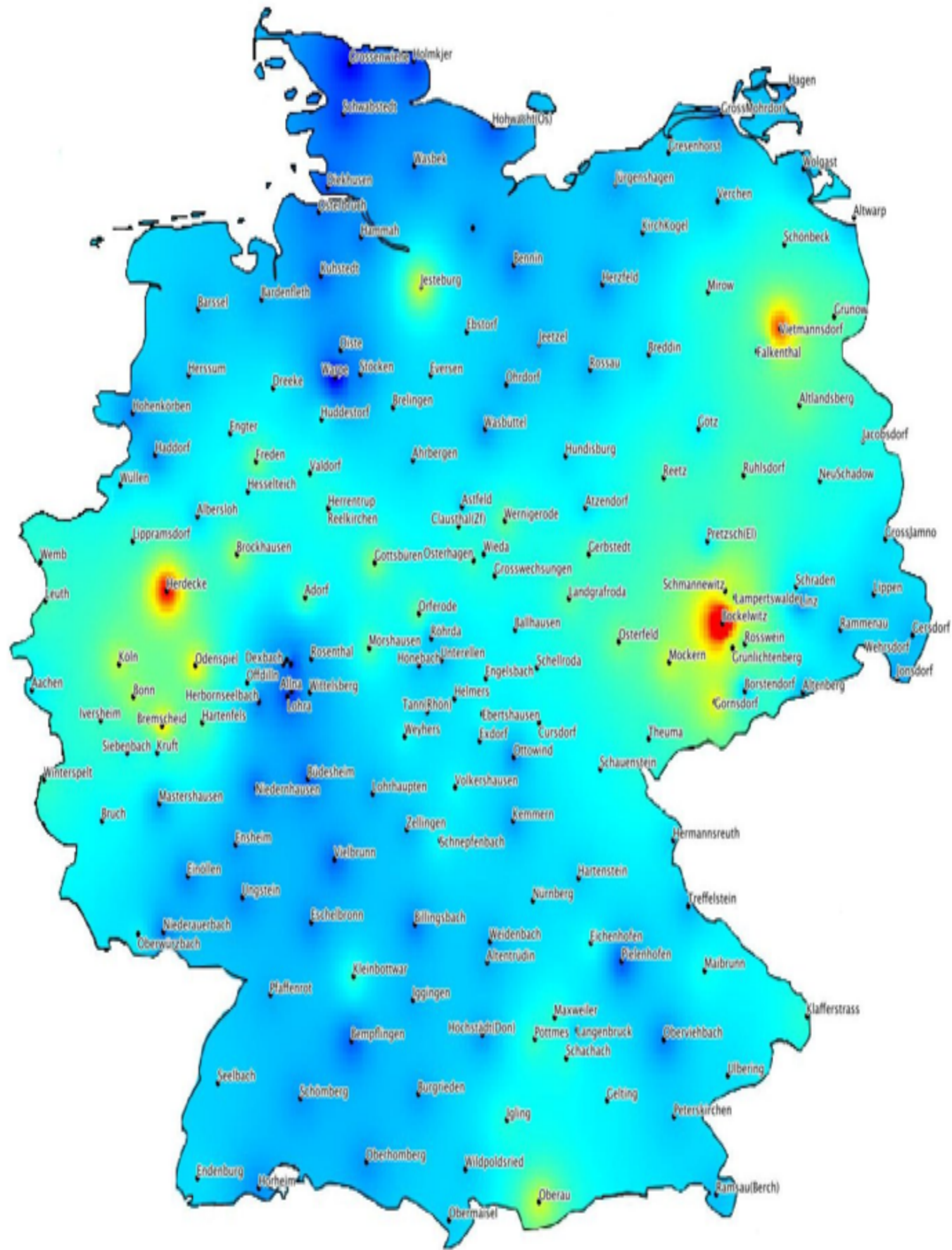
Vowels
of ge- prefix
(averages over
14 lexemes)



Vowels are
'strong' in an
area in which
consonants are
'weakening'







The git-approach

- git is a versioning control system: keeping track of who changed what
- You can run it privately on your own computer, or use services like github or bitbucket
- Text-based: so use simple formats!
(CSV in UTF-8, NFC, LF)

Conclusions

- Interpolation results in nice pictures
- Trends can be made visible
 - ▶ watch out with sparse data
 - ▶ be aware of the many possible interpolations
- Boundary-lines can be drawn (“Isohypeses”), but should not be interpreted as traditional boundaries