

# The Future of Dialect Data

## Representation & Analysis

*Jelena Prokić, Michael Cysouw, Johann-Mattis List*  
Deutscher Sprachatlas  
Philipps-Universität Marburg

# Dealing with Survey Data

1. **Case Study:**

Phonetischer Atlas von Deutschland

2. **Methods:**

Multiple sequence alignment and  
Automatic isogloss drawing

3. **Looking forward:**

The git-way of data management

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

Besprochen von 20.07.1985  
25.07.1985 UStv

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter 178	'vɪntə <sup>h</sup> ɐ	ɐ = kontinuierlich ɔ bereit velarisierter
* 2	fliegen 56	'flaɪ→ə <sup>h</sup> ə	'fliegen die', Sequenzierung unklar - kein geminiertes [t]
3	Blätter 23	'blɛ:də <sup>h</sup> ɐ	
4	Luft 103	lu <sup>h</sup> ft <sup>h</sup>	ɐ = kontinuierlich
5	hört 89	hɪɪt <sup>h</sup>	ɐ = kontinuierlich
6	gleich 78	klɛ <sub>+</sub> ɪk <sub>h</sub>	folgt P
7	schneien 130	ʃnɛɪən	
8	Wetter 174	/	statt demoa 'vɪtəʀɔgə
9	tu 151	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 wɪntɐ	b = kontinuant d = velarisiert
2	fliegen	56 flaɪ̯ə	'fliegen die', Sequenzierung unklar - kein gemischtes [t]
3	Blätter	23 blætɐ	
4	Luft	103 lʊft	= kontinuant
5	hört	89 hœrt	b = kontinuant
6	gleich	130 kleɪ̯ç	
7	schneien	130 ʃneɪ̯ən	
8	Wetter	178 vɛtɐ	statt dem 'vɛtəʀɔgə'
9	tu	151 dɛʊ	

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

Besprochen von 20.07.1985  
25.07.1985 UStv

Lfd. Nr.	Stichwort	Transkription	Bemerkungen
1	Winter 178	'vɪntə <sup>h</sup> ɐ	ɐ = kontinuierlich ɔ bereit velarisierter
* 2	fliegen 56	'flaɪ→ə <sup>h</sup> ə	'fliegen die', Sequenzierung unklar - kein geminiertes [t]
3	Blätter 23	'blɛ:də <sup>h</sup> ɐ	
4	Luft 103	lu <sup>h</sup> ft <sup>h</sup>	ɐ = kontinuierlich
5	hört 89	hɪɪt <sub>~</sub>	ɐ = kontinuierlich
6	gleich 78	klɛ <sub>+</sub> ɪk <sub>~</sub>	folgt P
7	schneien 130	ʃnɛɪən	
8	Wetter 174	/	statt demoa 'vɪtəʀɔgə
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
56	fliegen	'flɛgən	'fliegende', Sequenz
23	Rätter	'blɛ:ɔtə	
103	fliegen	'flɛgən	b = kontinuierlich
89	hö	hɔ	b = kontinuierlich
78	fliegen	'flɛgən	folgt ?
130	fliegen	'flɛgən	
174	wetter	'vɛtə	statt dem 'vɛtə
151	tu	dɛv	'vɛtə

# 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	a:b <sup>h</sup> ə
Ahrbergen	o→ɔphə
Albersloh	a:p <sup>h</sup> ə
Allna	aϕh
Altenberg	ʌfɛ
Altentrüdin	af
Altlandsberg	a'fə'
Altwarp	o:ph
Astfeld	ɒ':p <sup>h</sup> ə
Atzendorf	afɛ
Ballhausen	ʌ'fə
Bardenfleth	ɔ:p̄ϕ
Barssel	ɒ:p <sup>h</sup> ə
Bempflingen	af:
Bennin	ɔp <sup>h</sup>
Billingsbach	af
Bockelwitz	ʌvə
Bonn	a:p'
Borstendorf	ʏf:
Breddin	ɒ:ph
Brelingen	a <sup>h</sup> fβə
Bremscheid	ɒ':phə
...	...

A	FF	E
a:	ph	-
a:	b <sup>h</sup>	ə
o→ɔ	ph	ə
a:	p <sup>h</sup>	ə
a	ϕh	-
ʌ	f	ɛ
a	f	-
a'	f	ə'
o:	ph	-
ɒ':	p <sup>h</sup>	ə
a	f	ɛ
ʌ'	f	ə
ɔ:	p̄ϕ	-
ɒ:	p <sup>h</sup>	ə
a	f:	-
ɔ	p <sup>h</sup>	-
a	f	-
ʌ	v	ə
a:	p'	-
ʏ	f:	-
ɒ:	ph	-
a	f̄β	ə
ɒ':	ph	ə
...	...	...

- **Workflow:**

- ▶ Tokenisation of segments
- ▶ Automatic alignment using LingPy ([github.com/lingpy](https://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



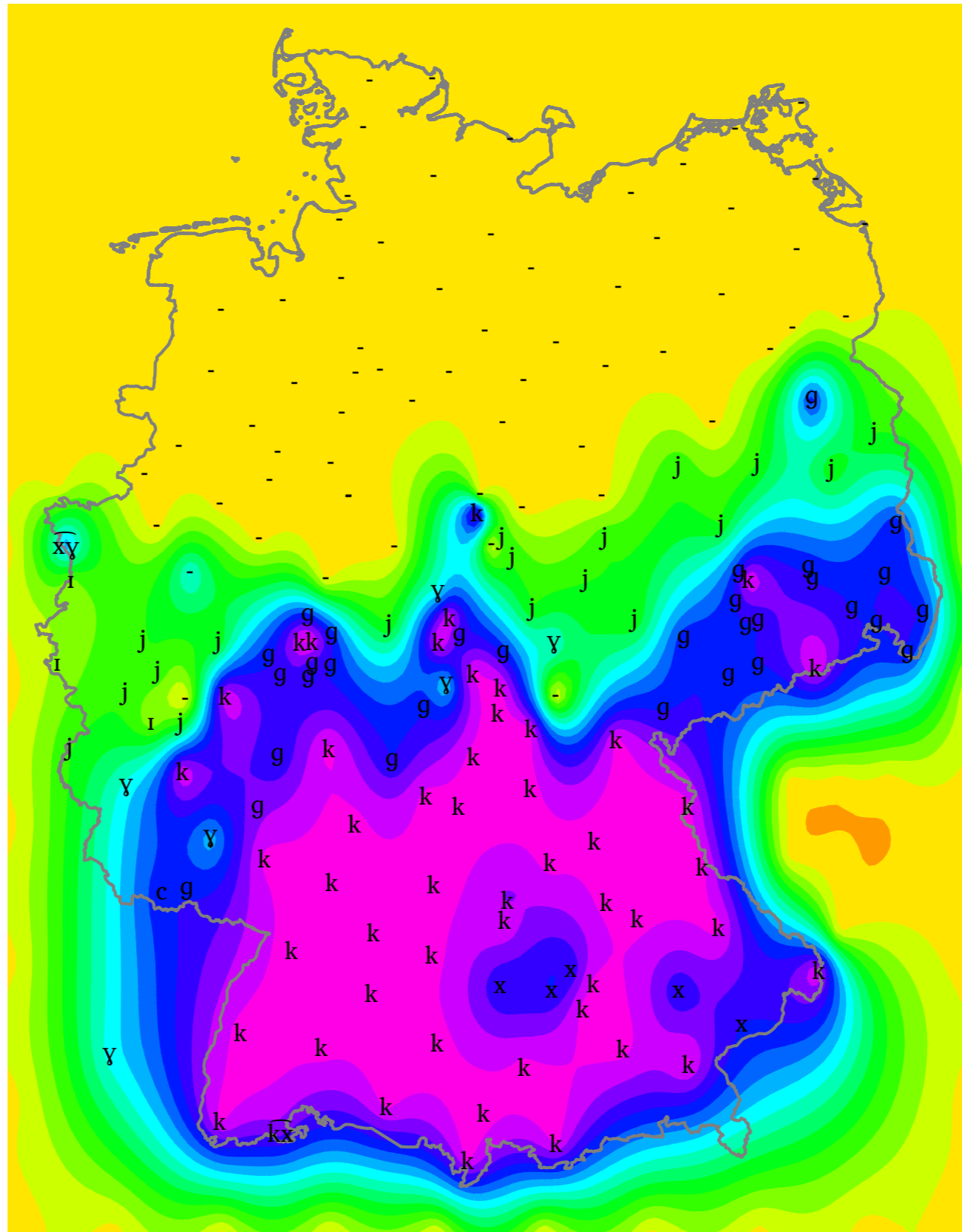
# Automatic Isoglossing

- 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)

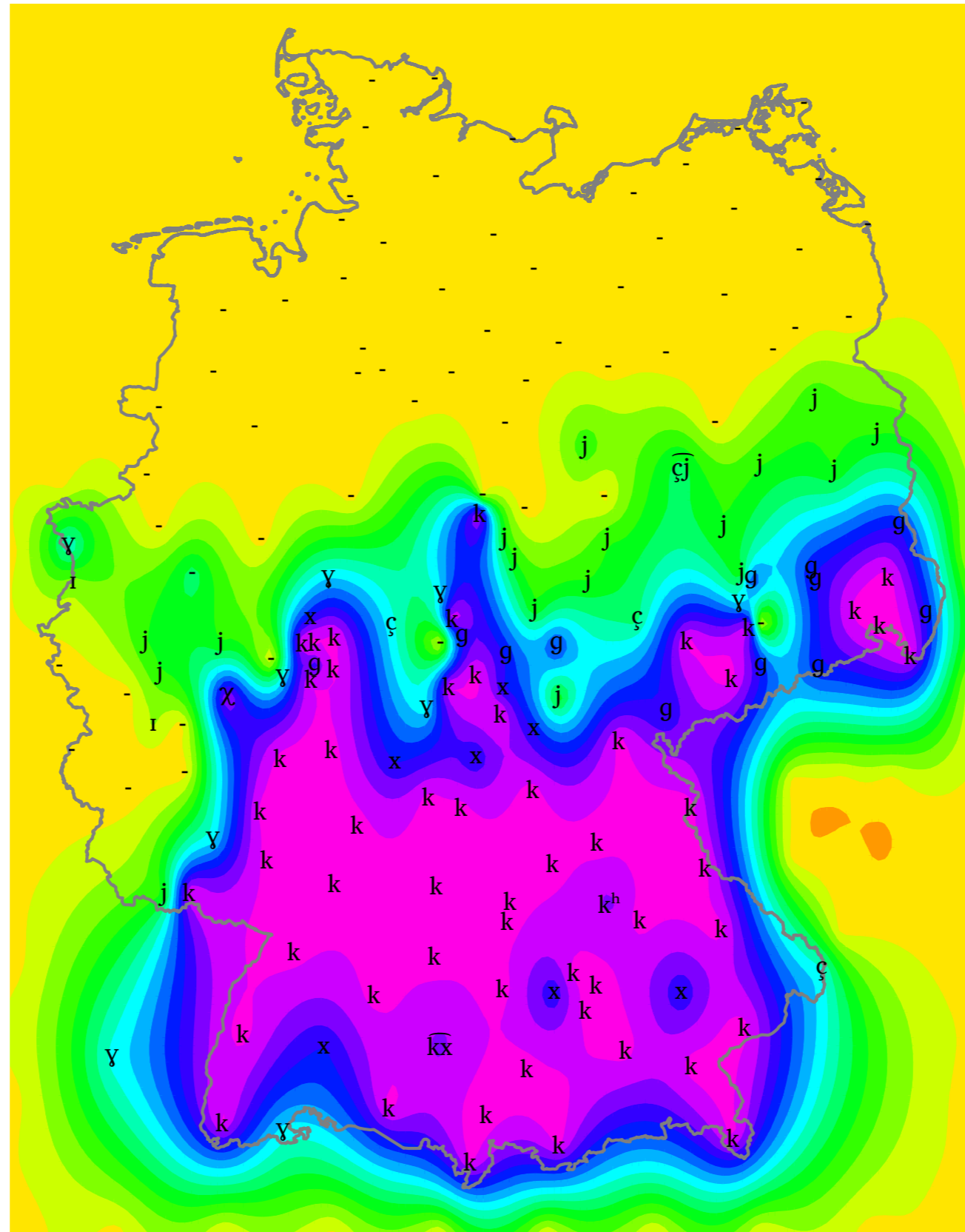
# Example 1

- **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 ʧ̥ ʧ̥ ʧ̥ ʧ̥ ʧ̥ ʧ̥ ʧ̥ ʧ̥ ʧ̥  
 k̥ k̥ ʀ̥ ʧ̥ k̥ k̥ q̥ k̥'

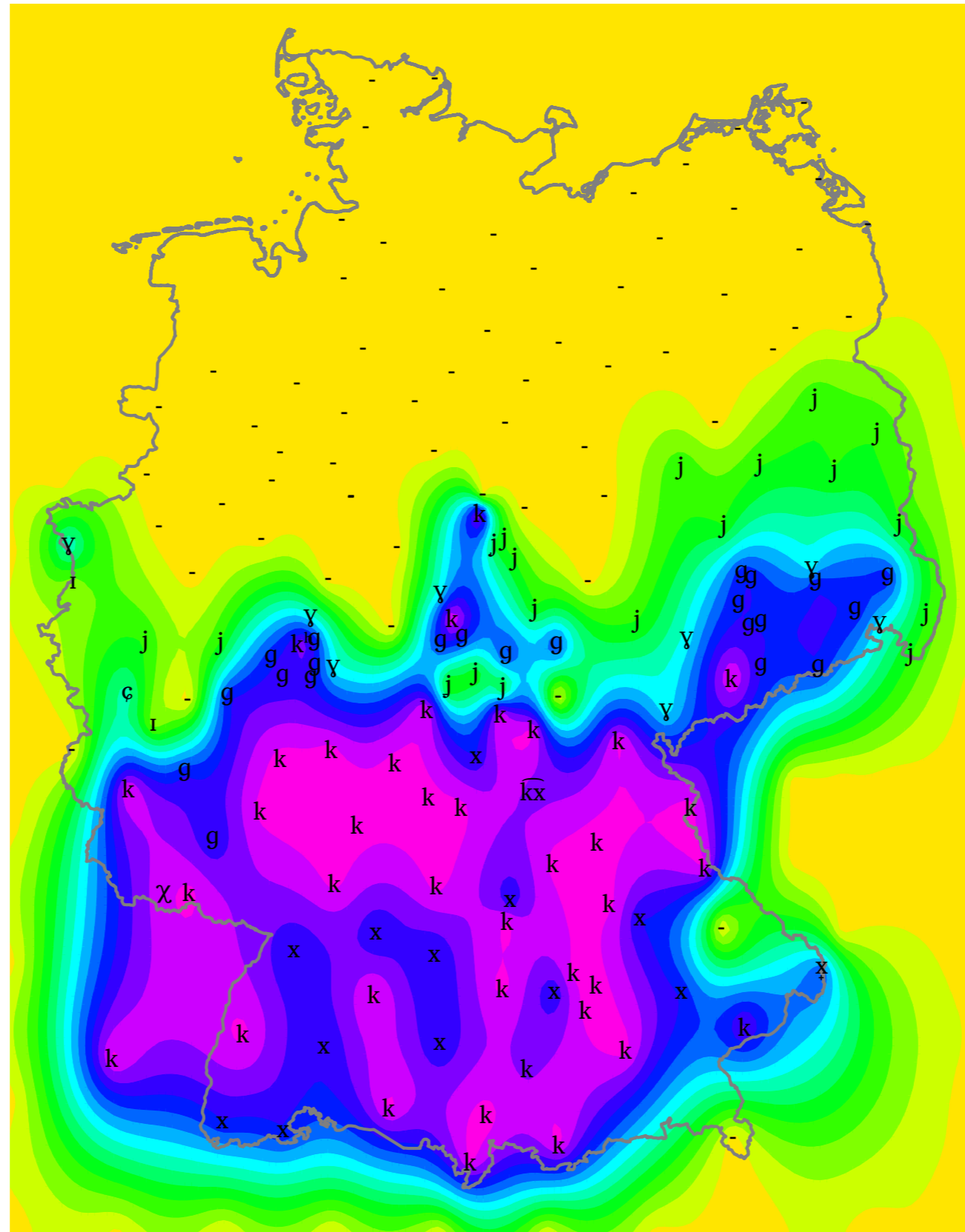
# gefahren



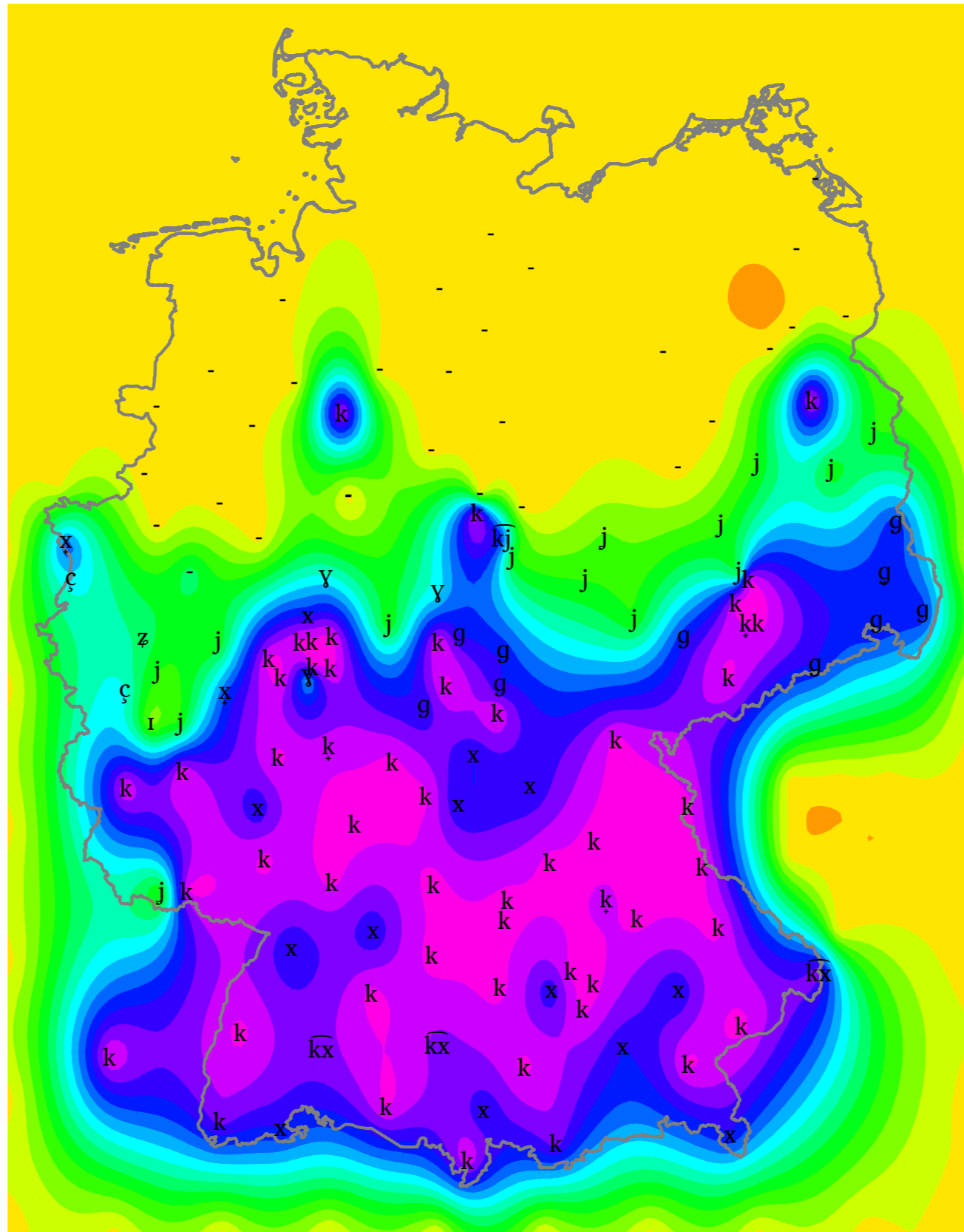
# gefunden



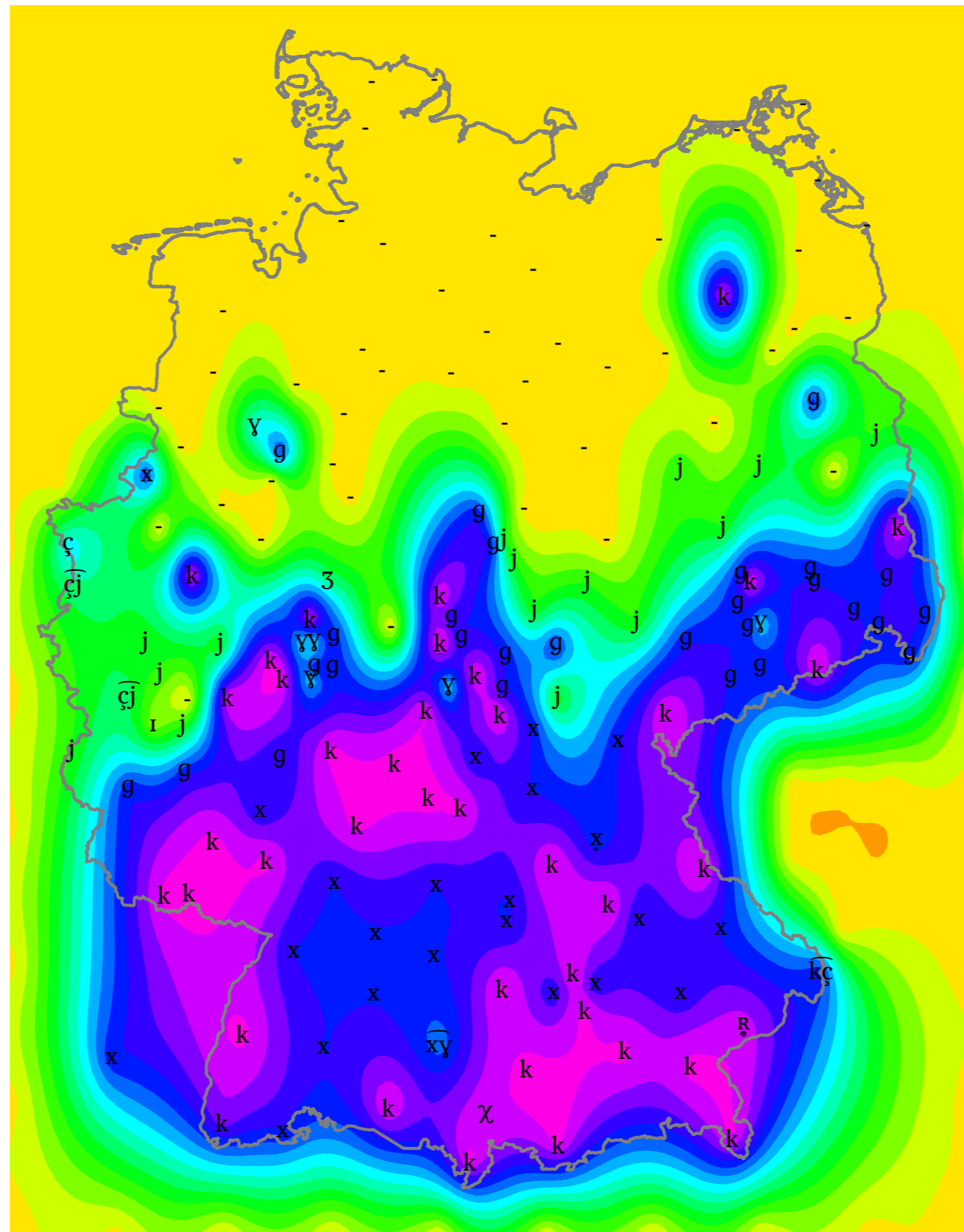
# eingeschlafen



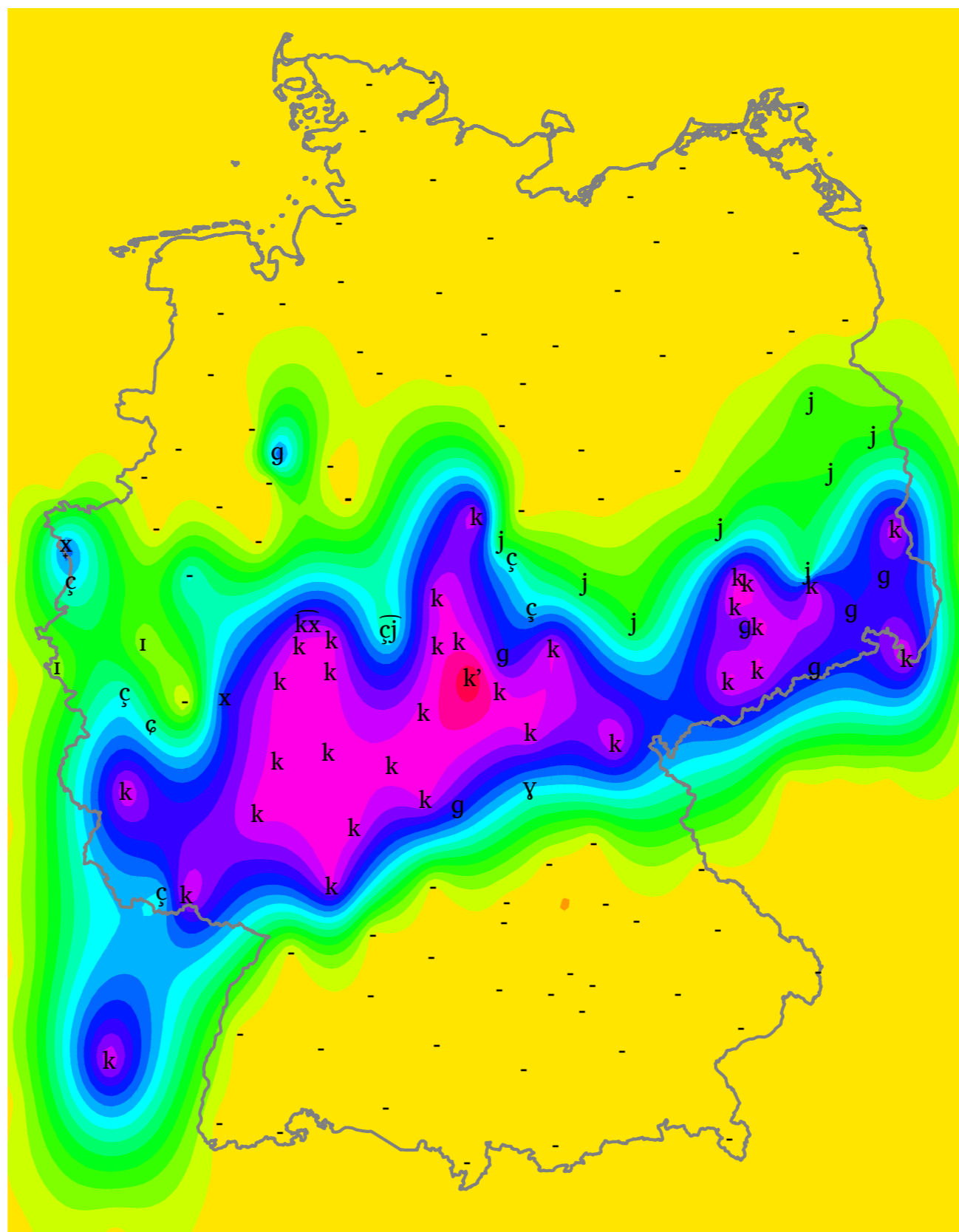
# gestohlen



# gestorben

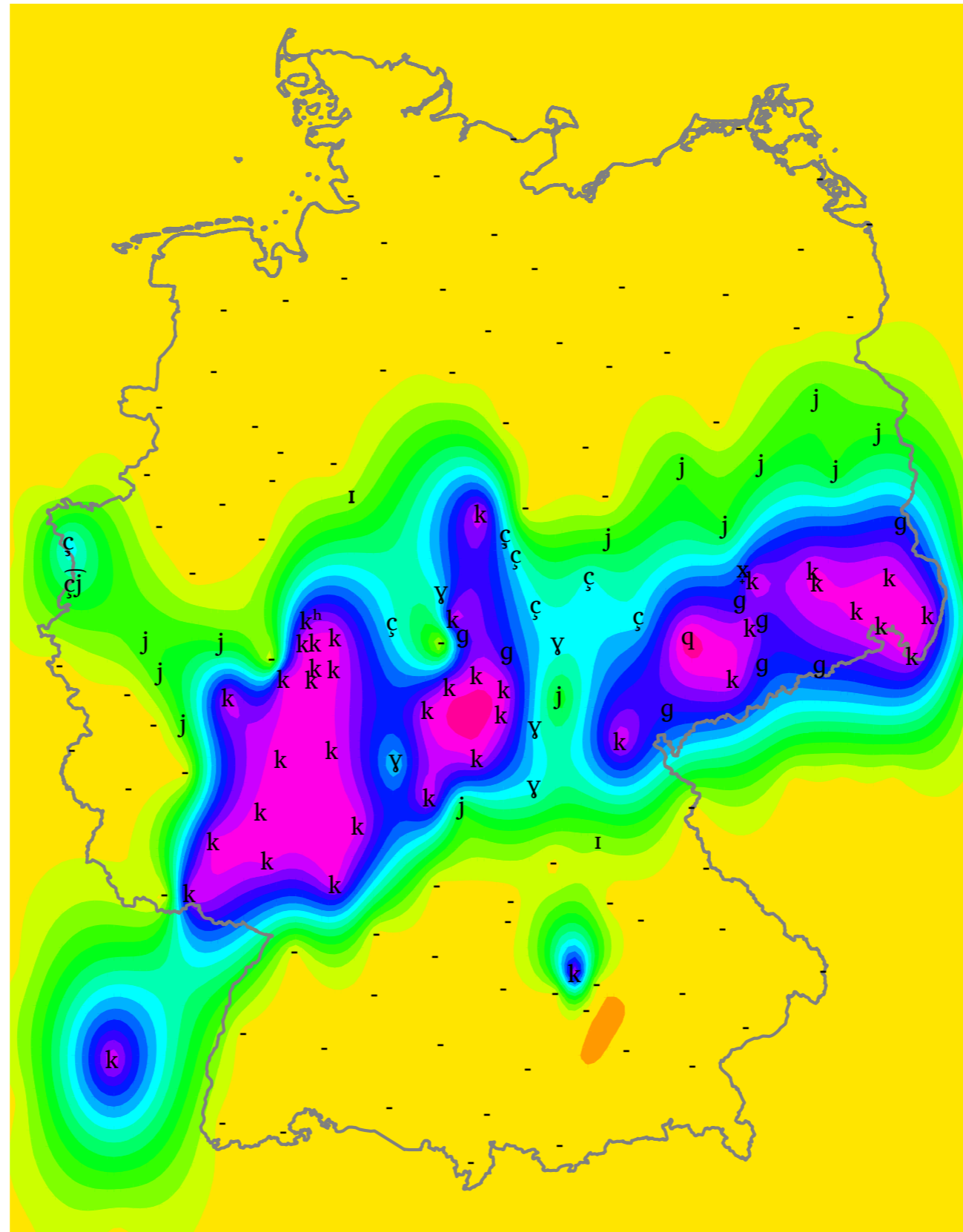


# gebrannt

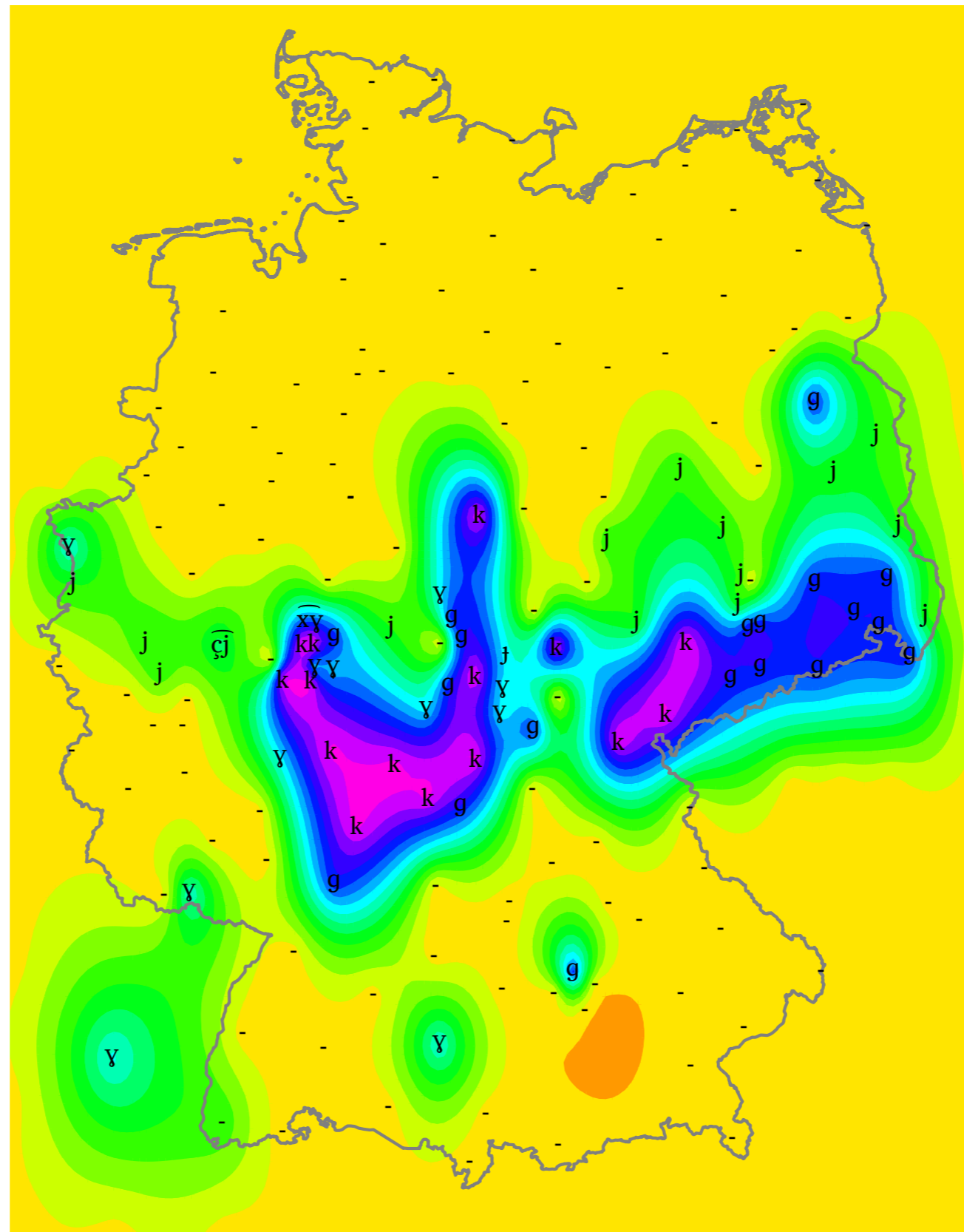




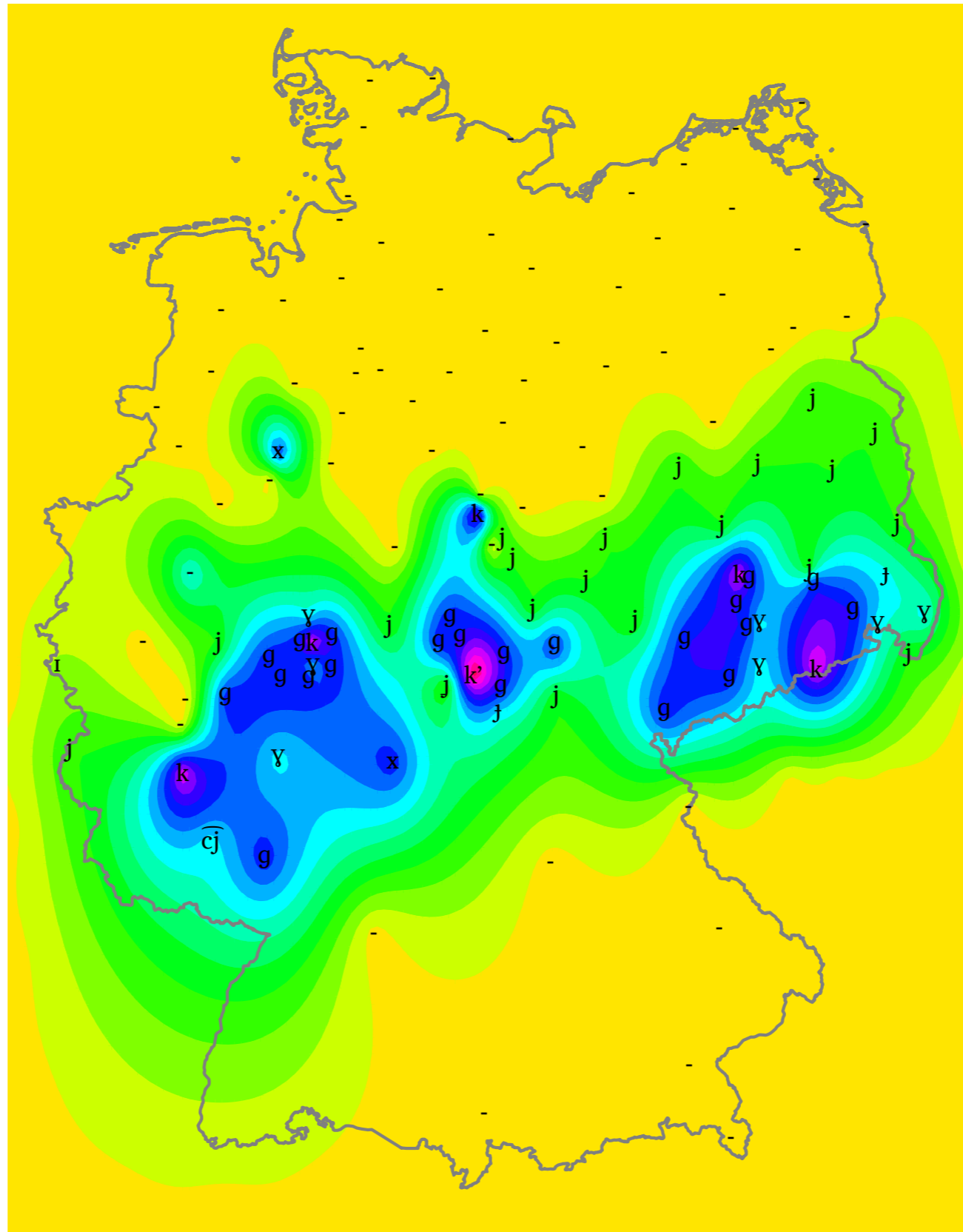
# gebracht



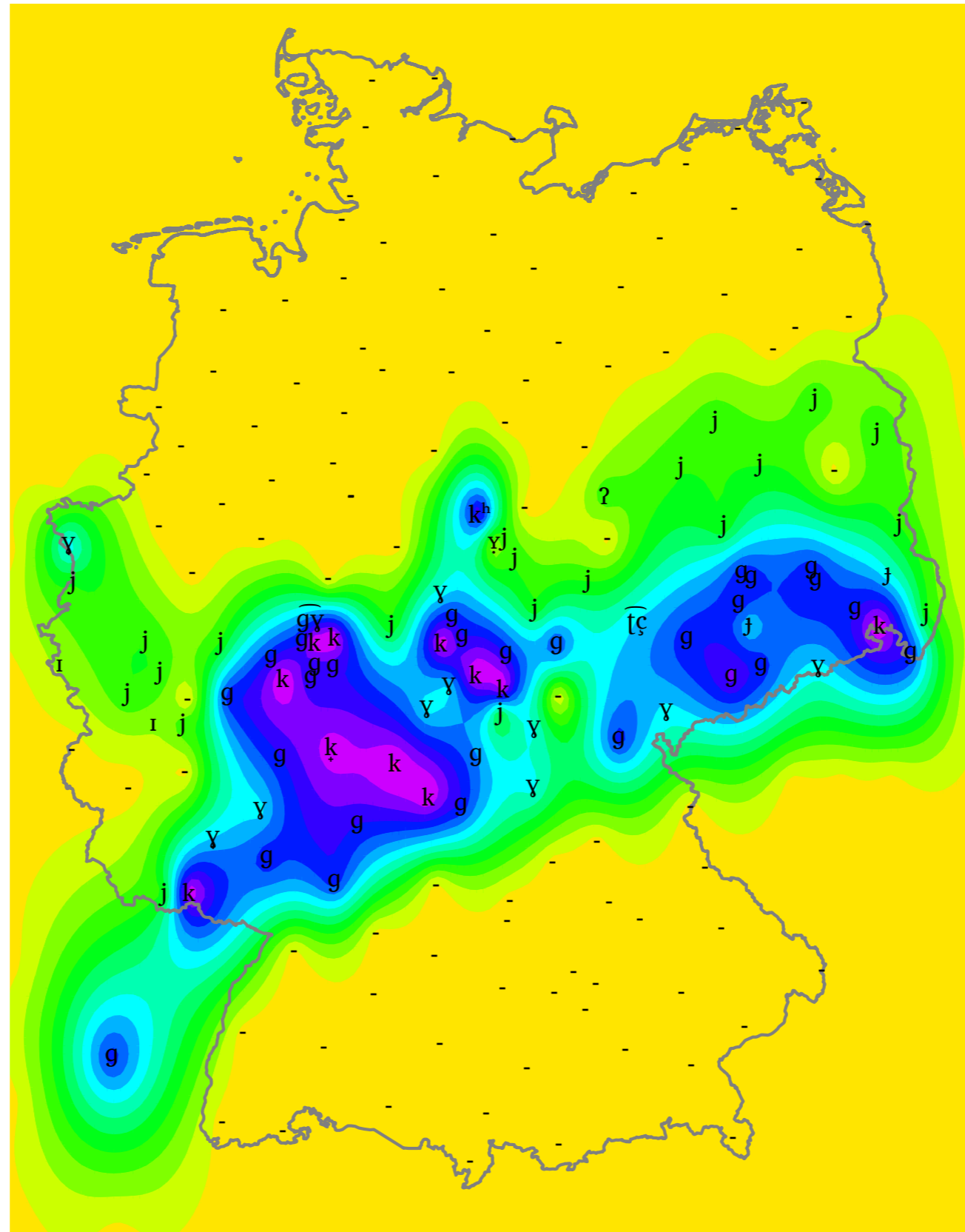
# geblieben



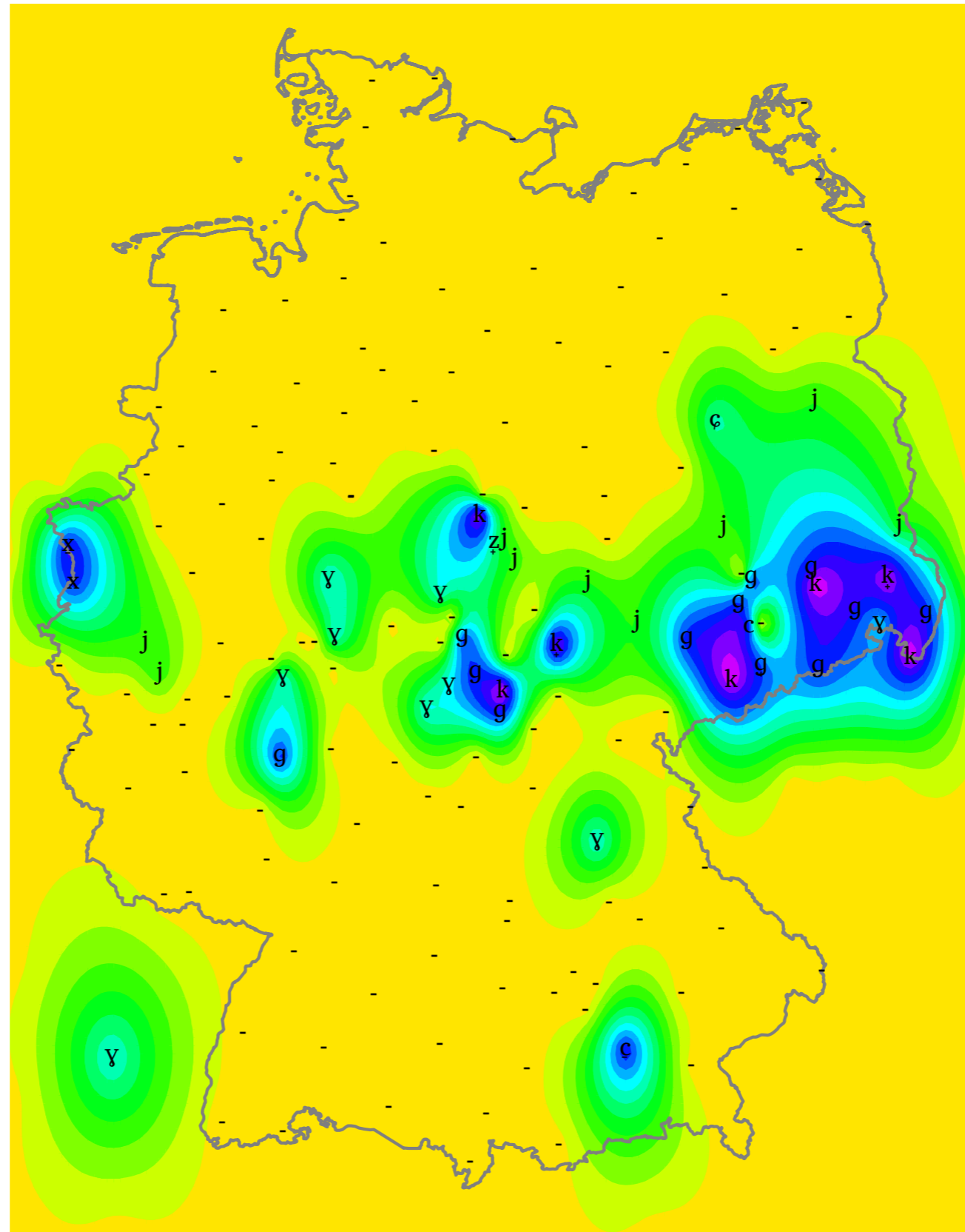
# eingebrochen



# gekannt



# gekommen

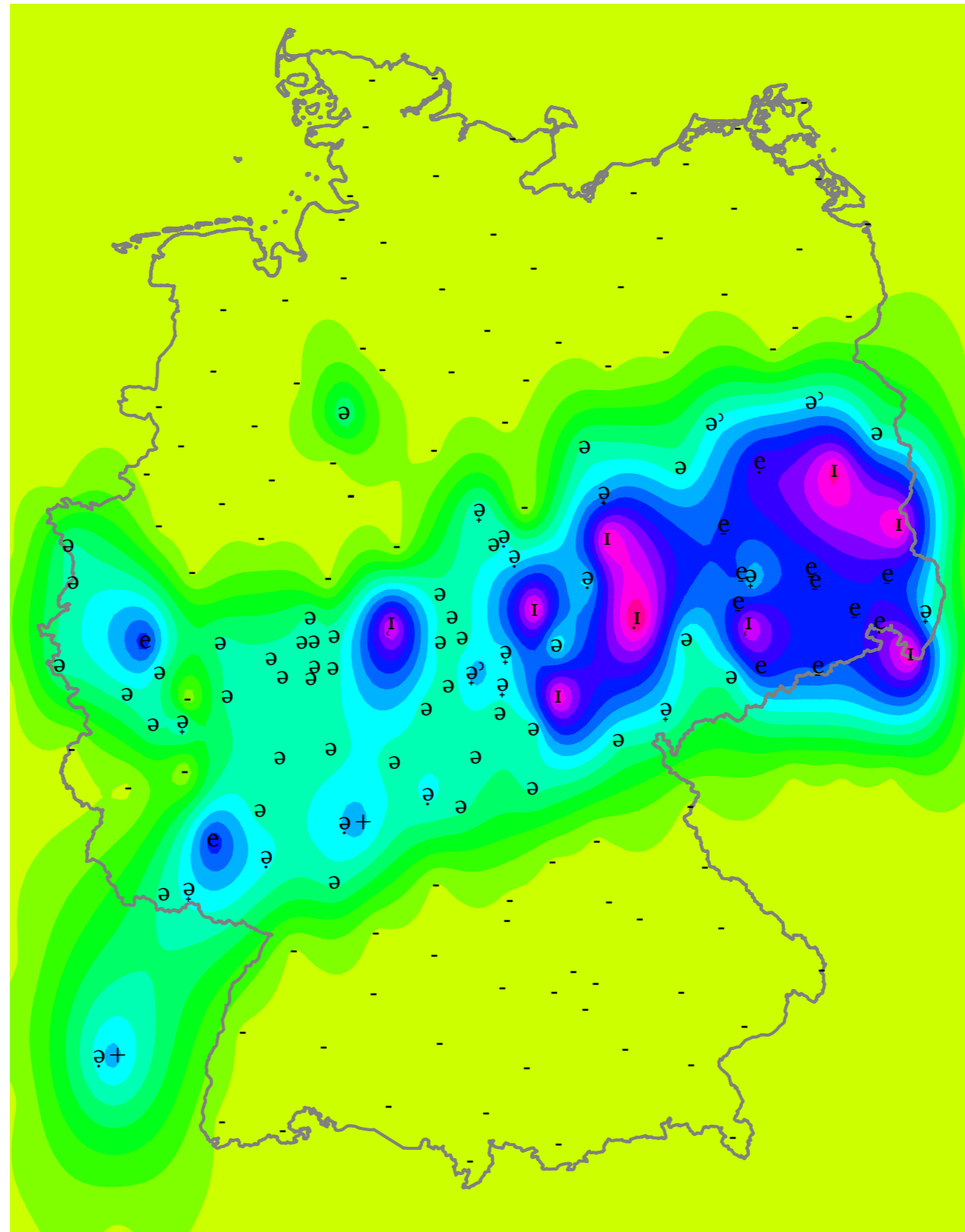


# Example 2

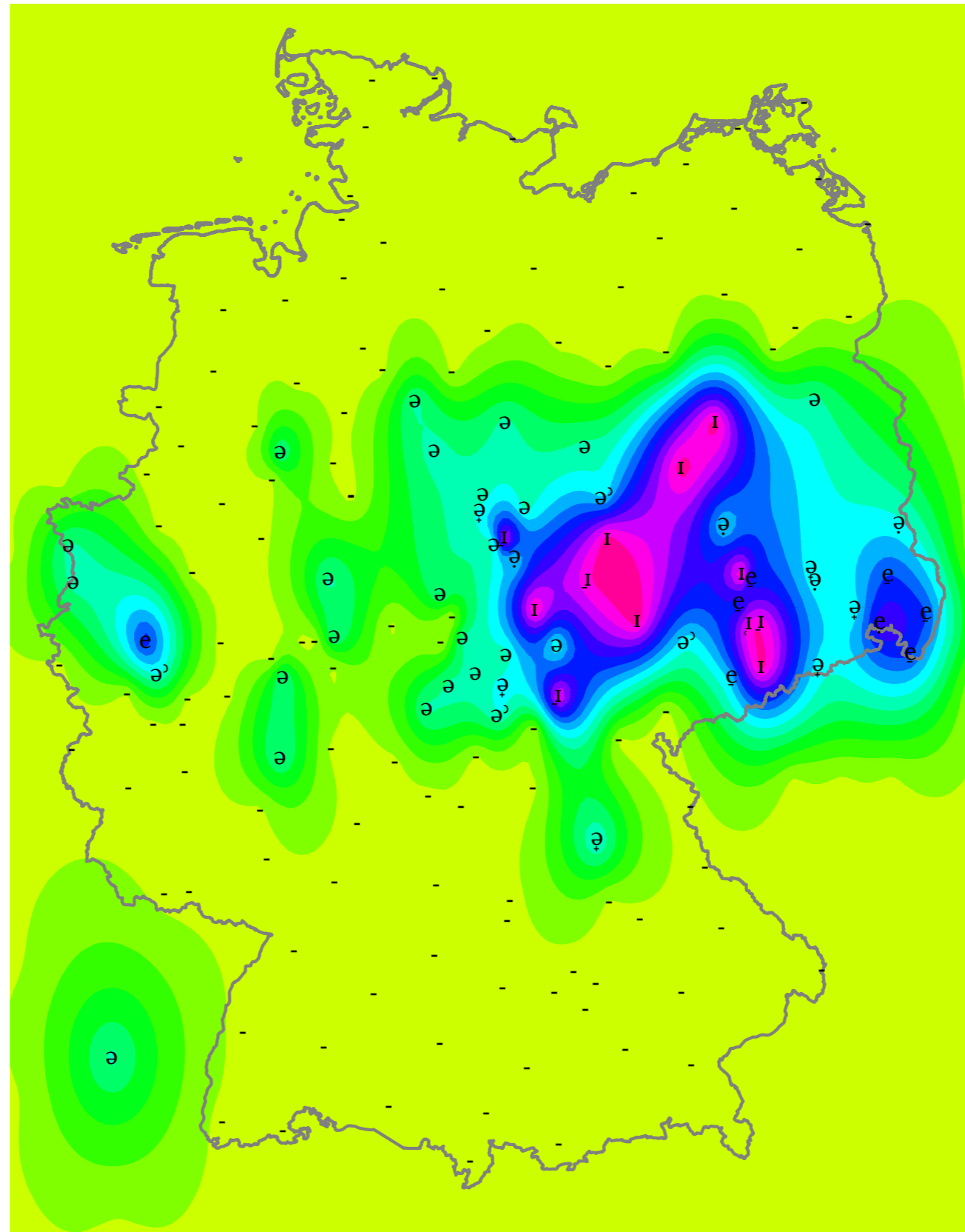
- **Vowel** in German Perfect-prefix ge-
- Ordered to 'strength', visualised as 'height'

- (CLUSTER) (NULL) ə ə ə ə<sup>c</sup> ə ə ə<sup>ʰ</sup> ə ə<sub>+</sub> ə<sub>+</sub>  
ə + ə + ə<sup>ʰ</sup> ə<sup>ʰ</sup> e e e<sup>ʰ</sup> e œ<sup>c</sup> œ ɛ ɛ ɪ ɪ ɪ ɪ ɪ

# gekannt

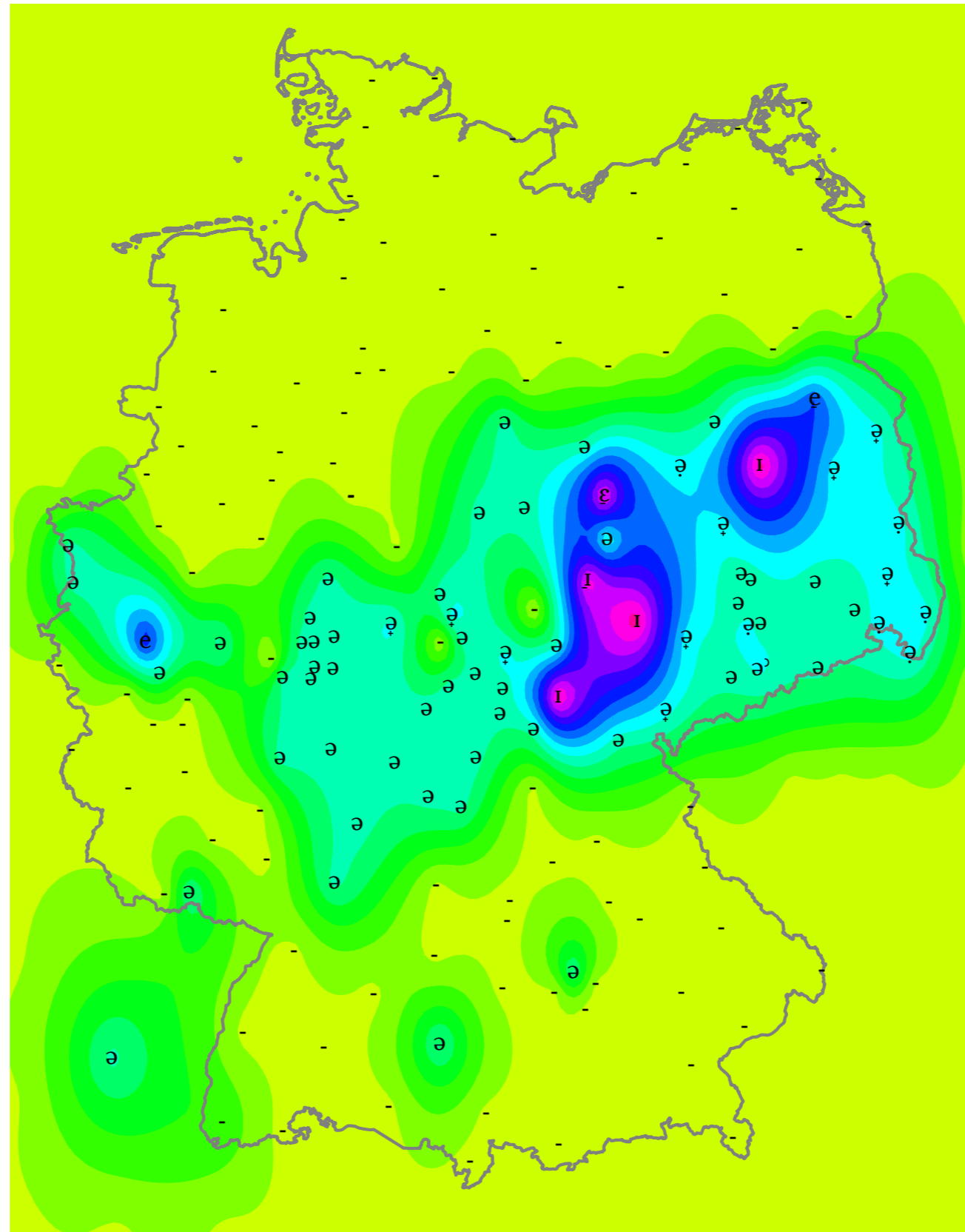


gekommen

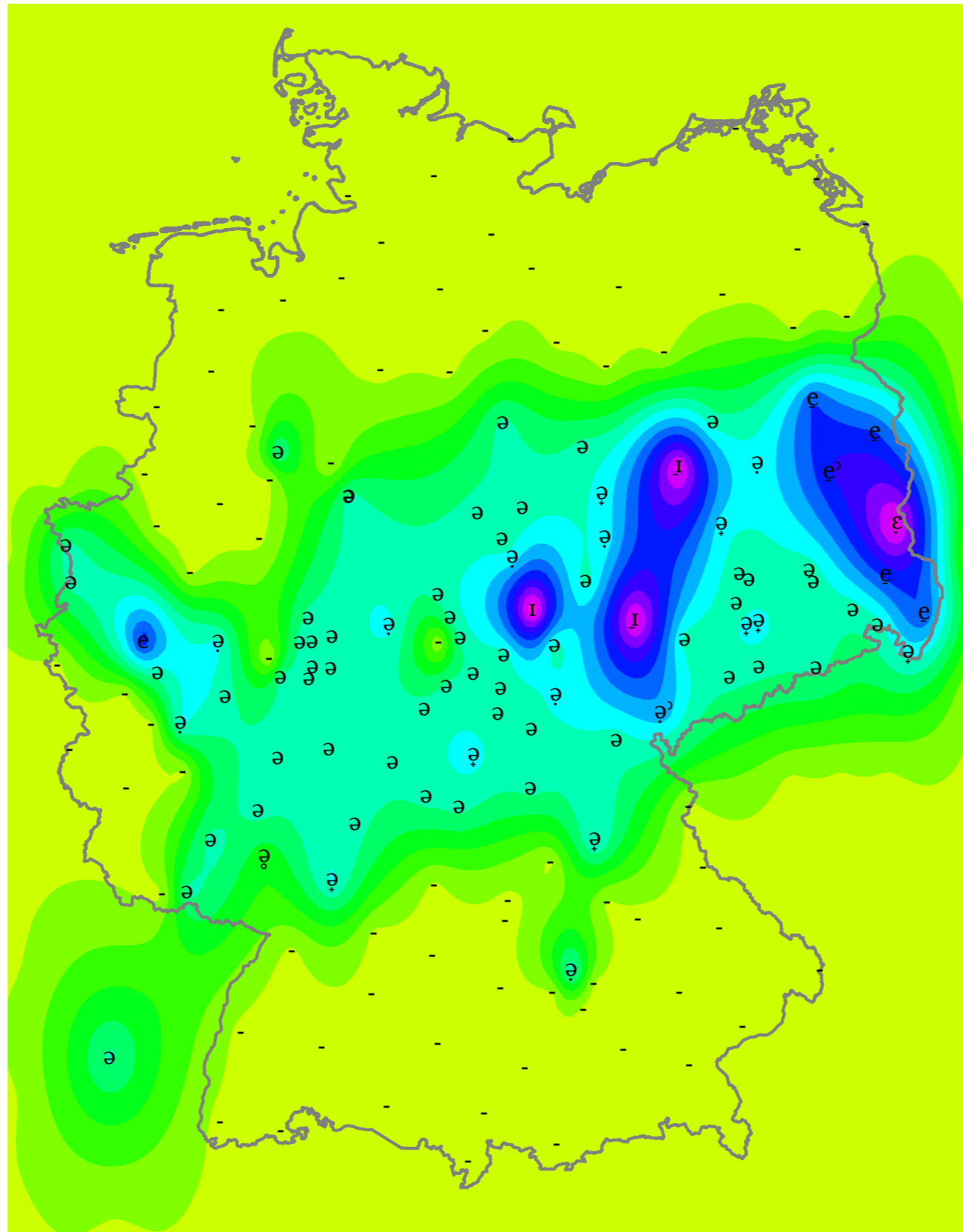




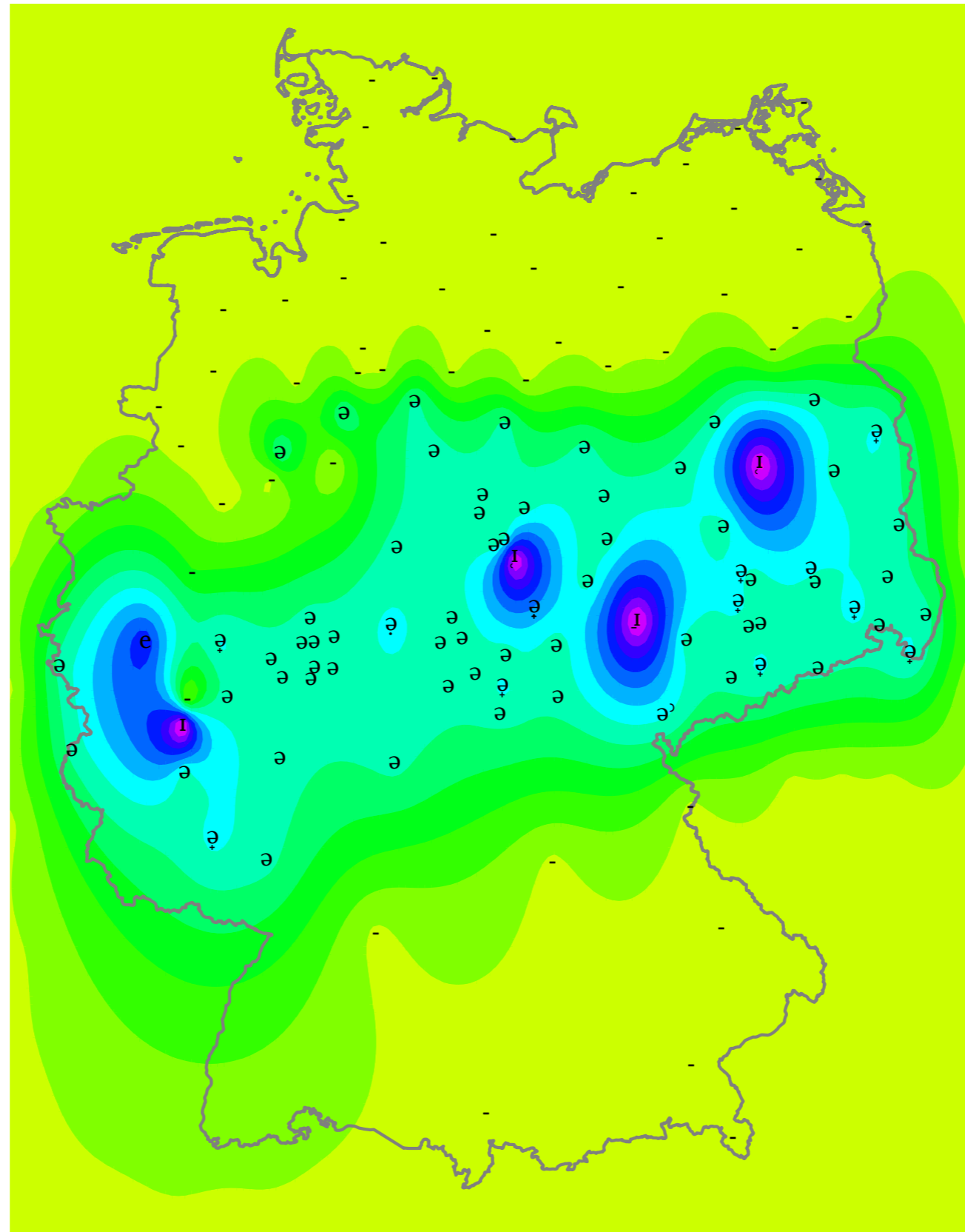
# geblieben



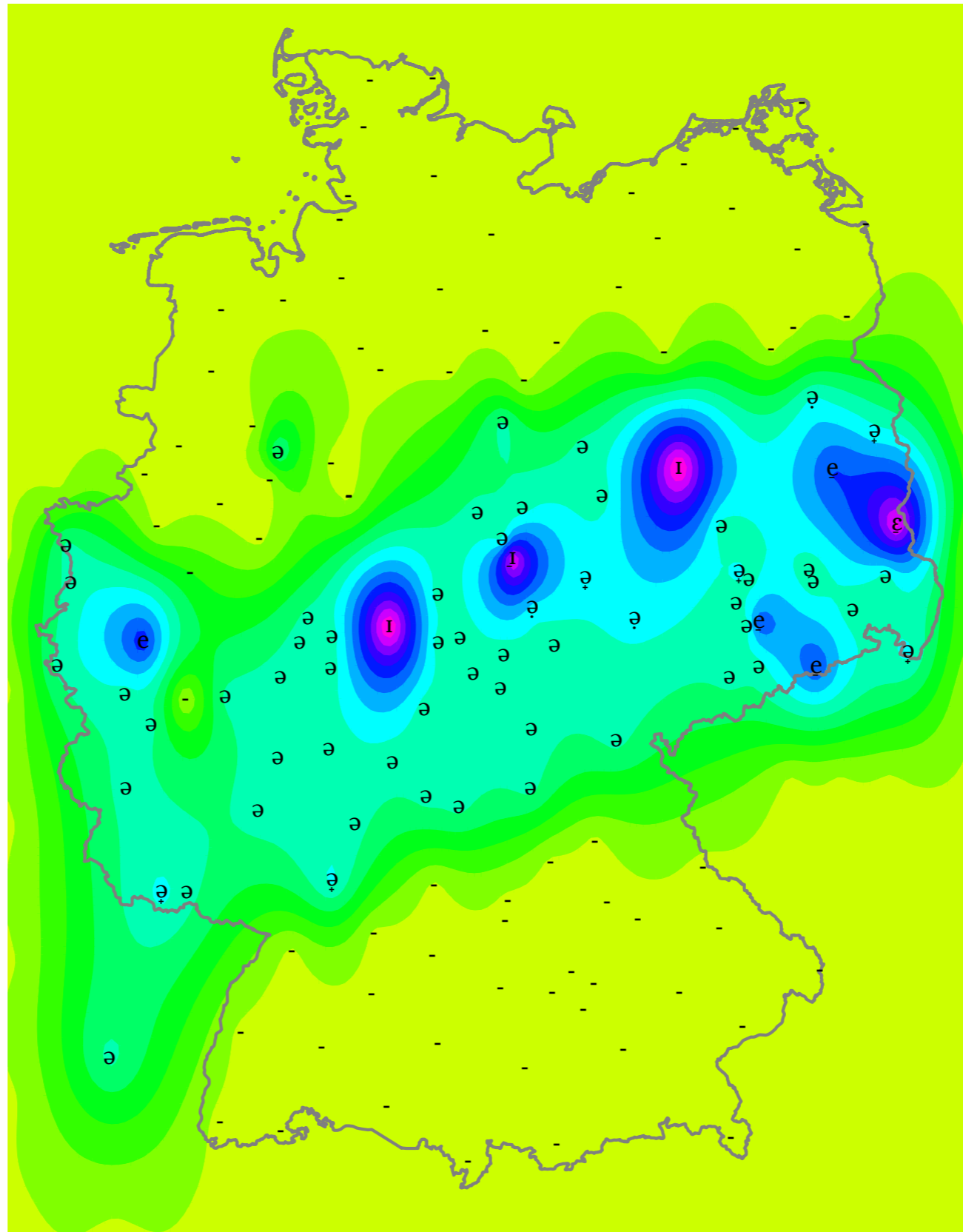
# gebracht



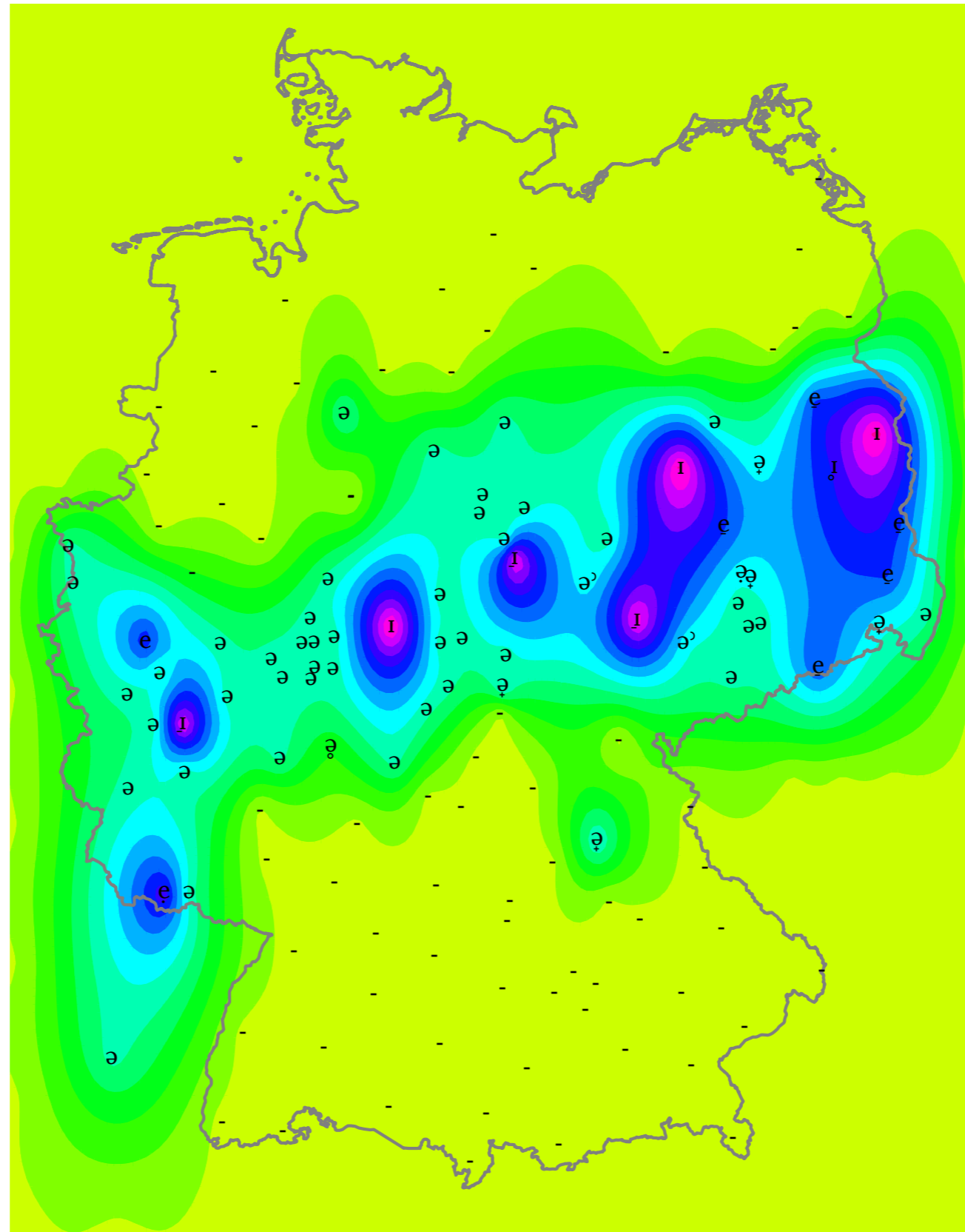
# eingebrochen



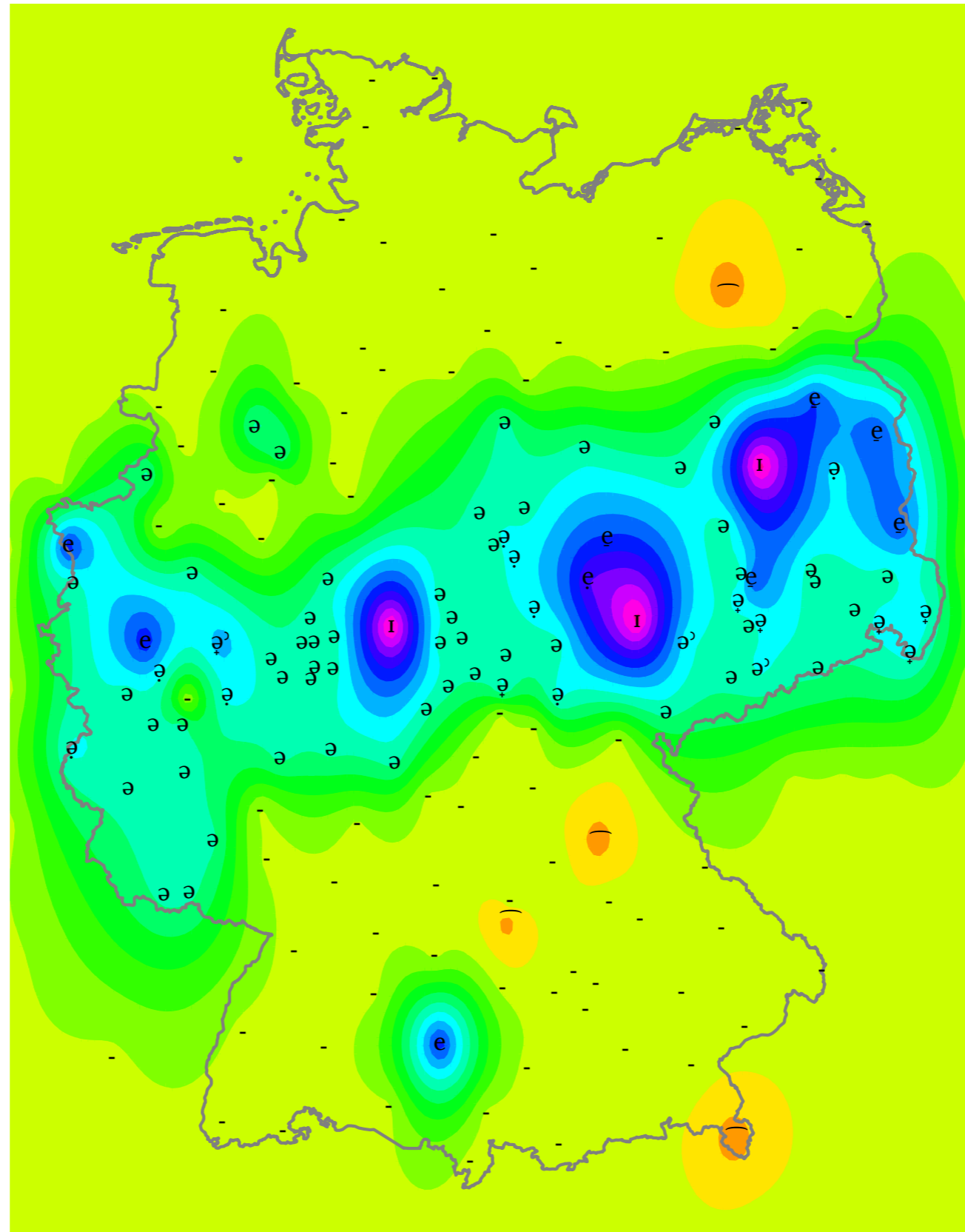
# gebrannt



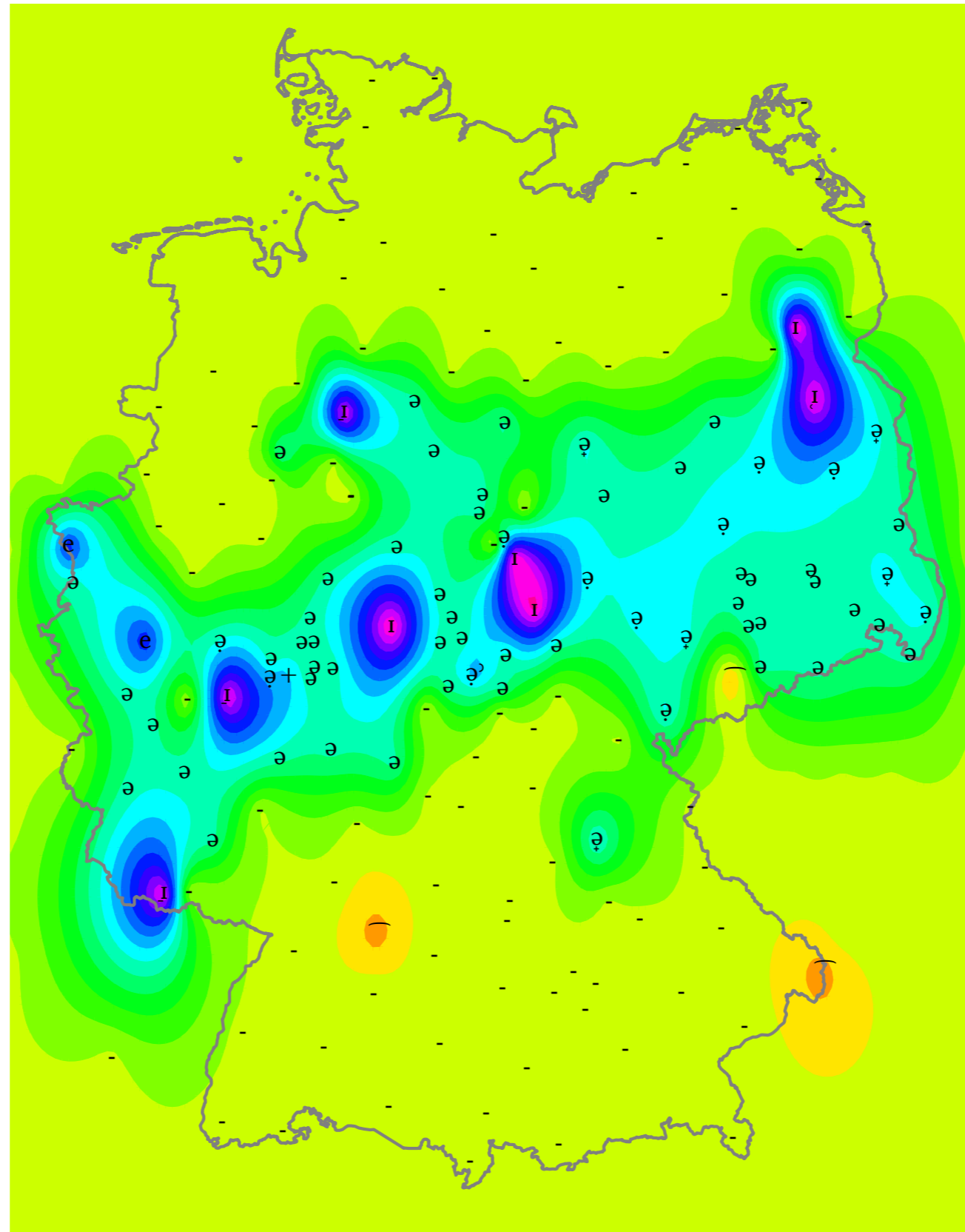
# gestohlen



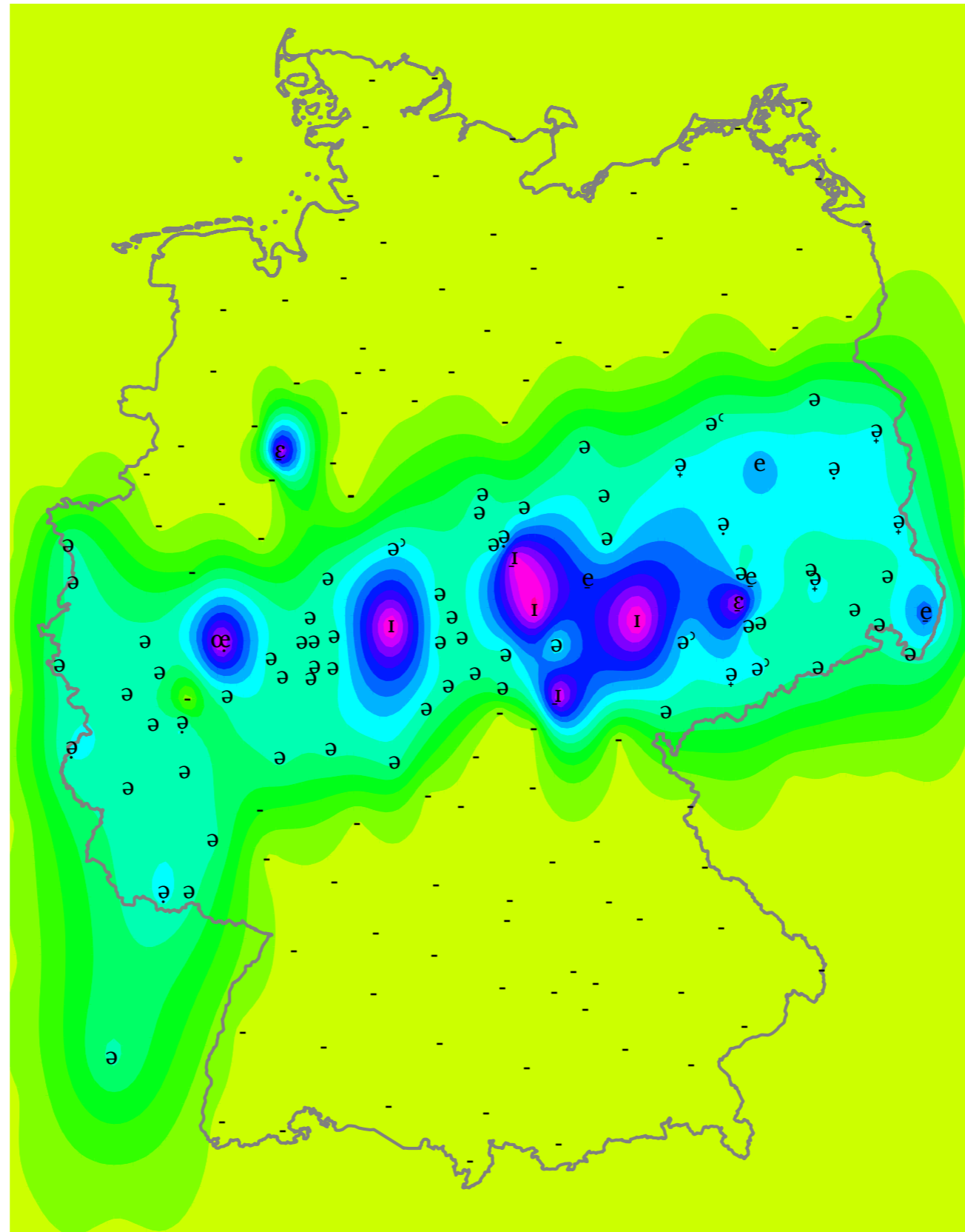
# gestorben



# eingeschlafen

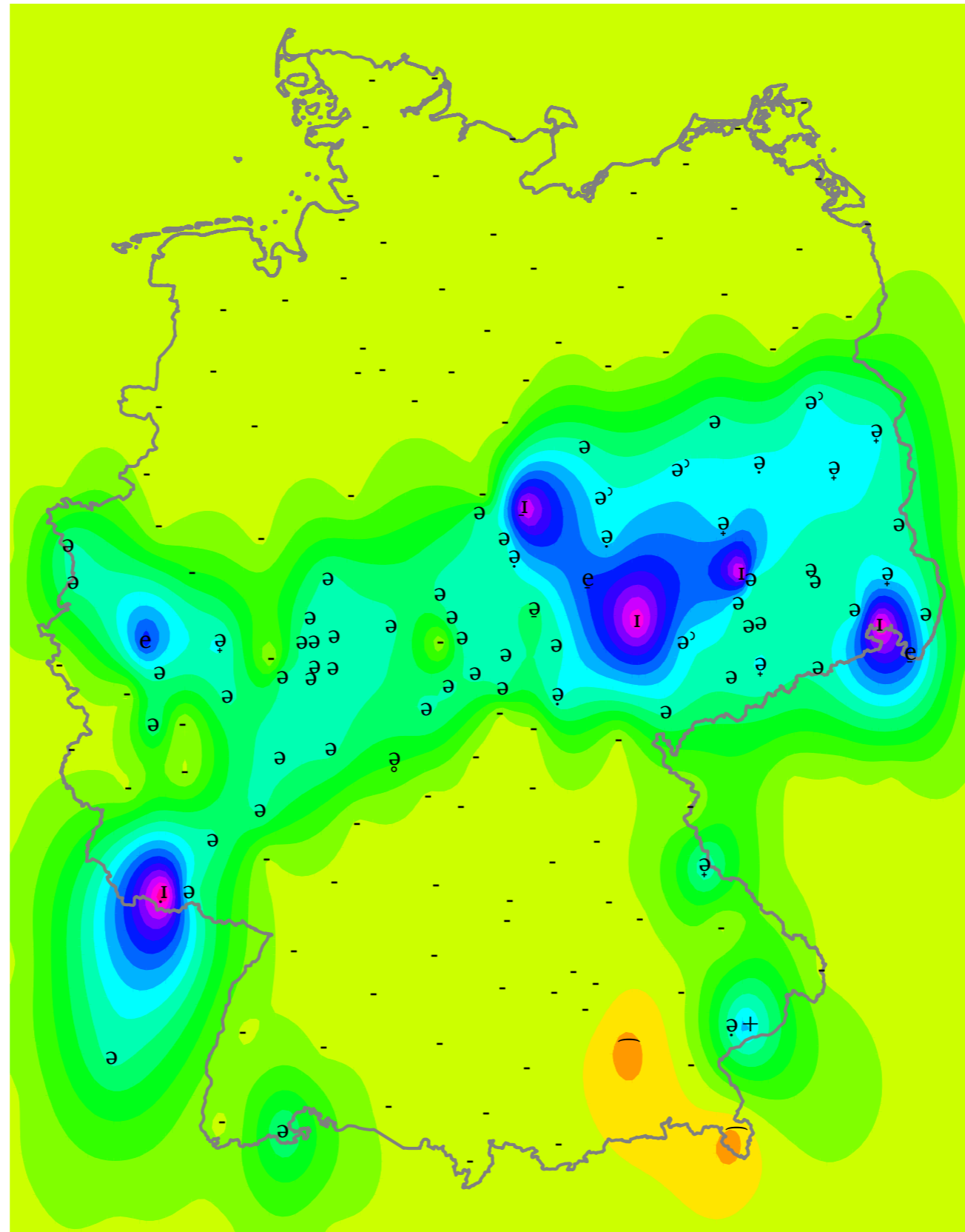


# gefahren

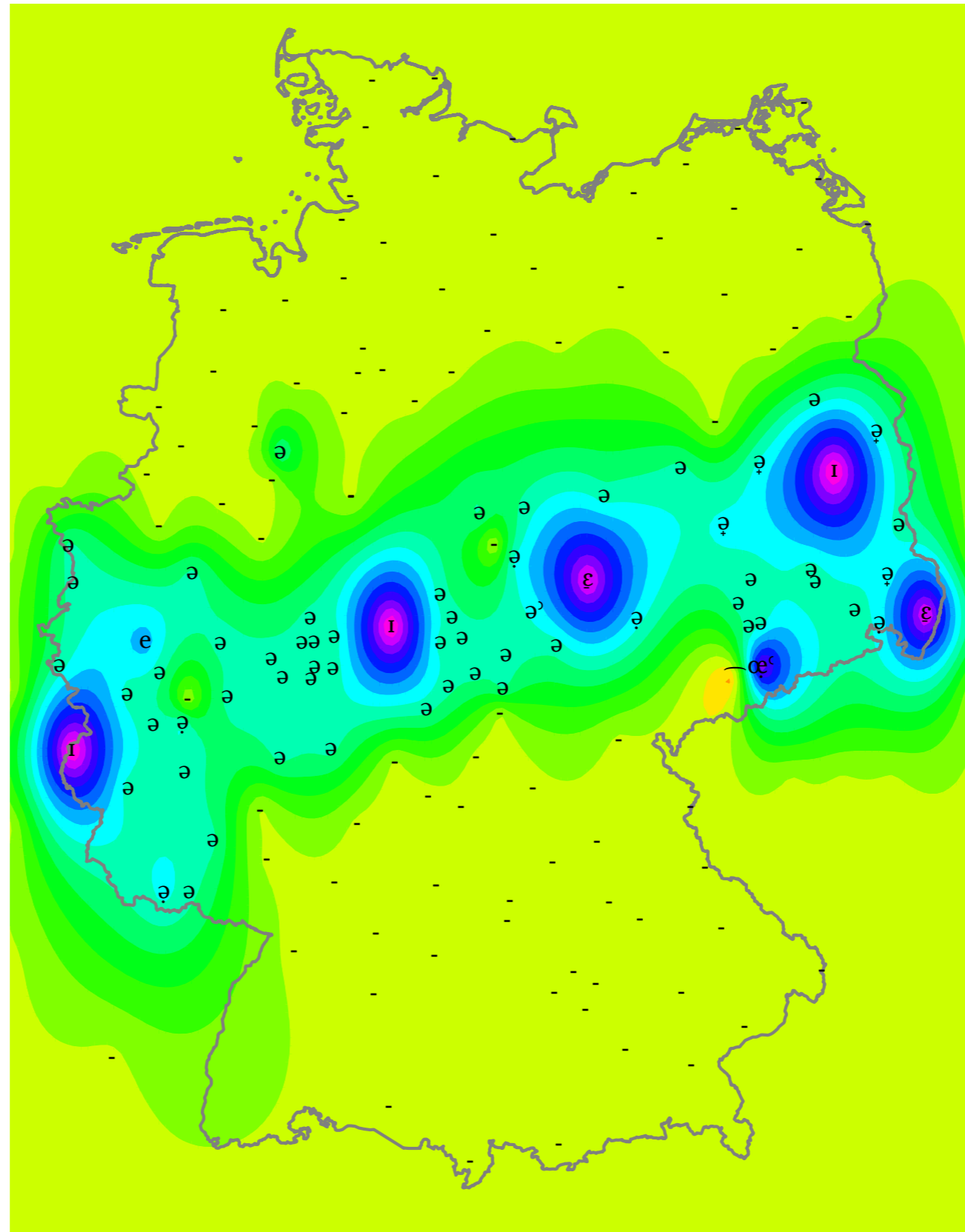




# gefunden



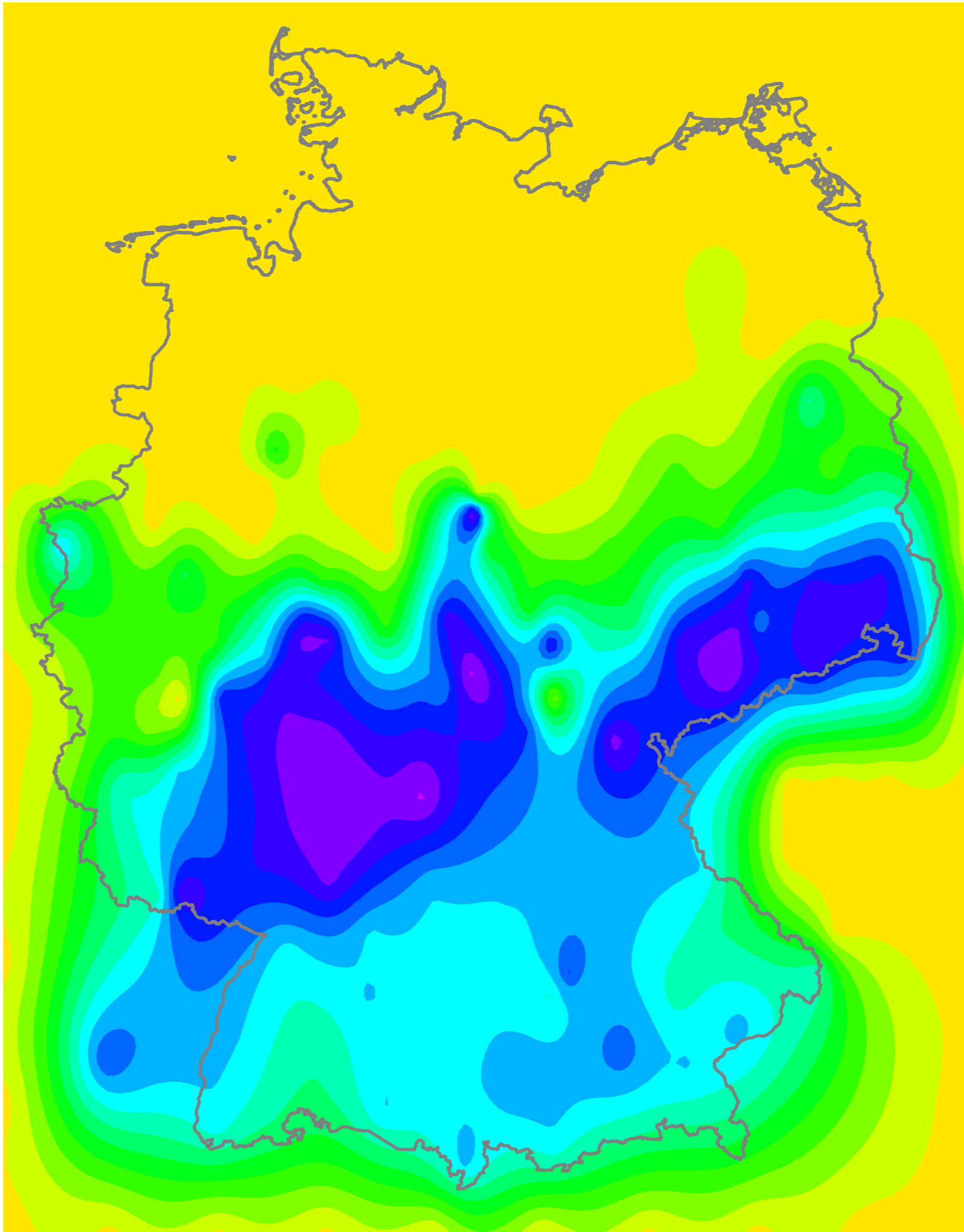
# gefallen



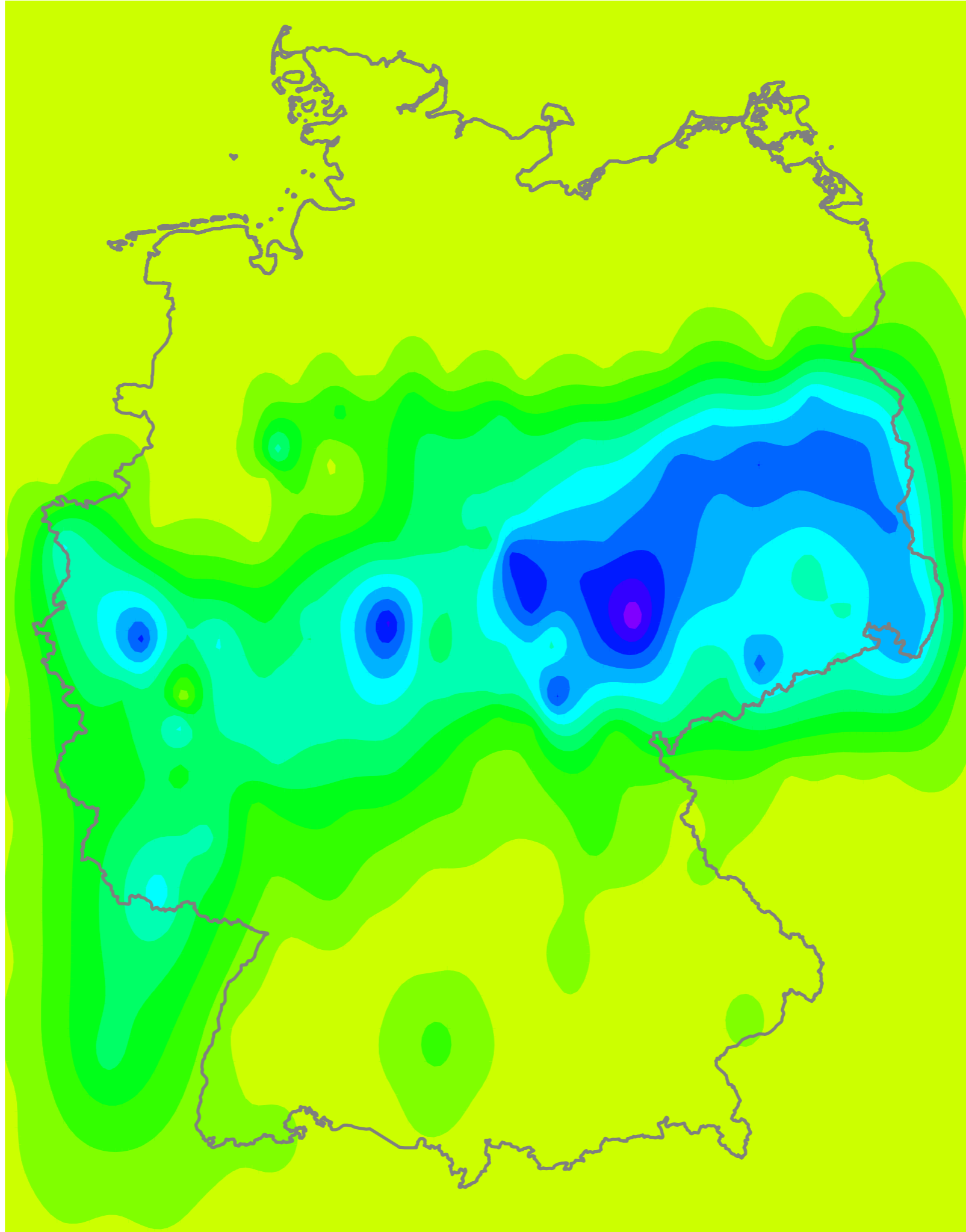
# Averaging

- Combining all consonants and all vowels into a single map

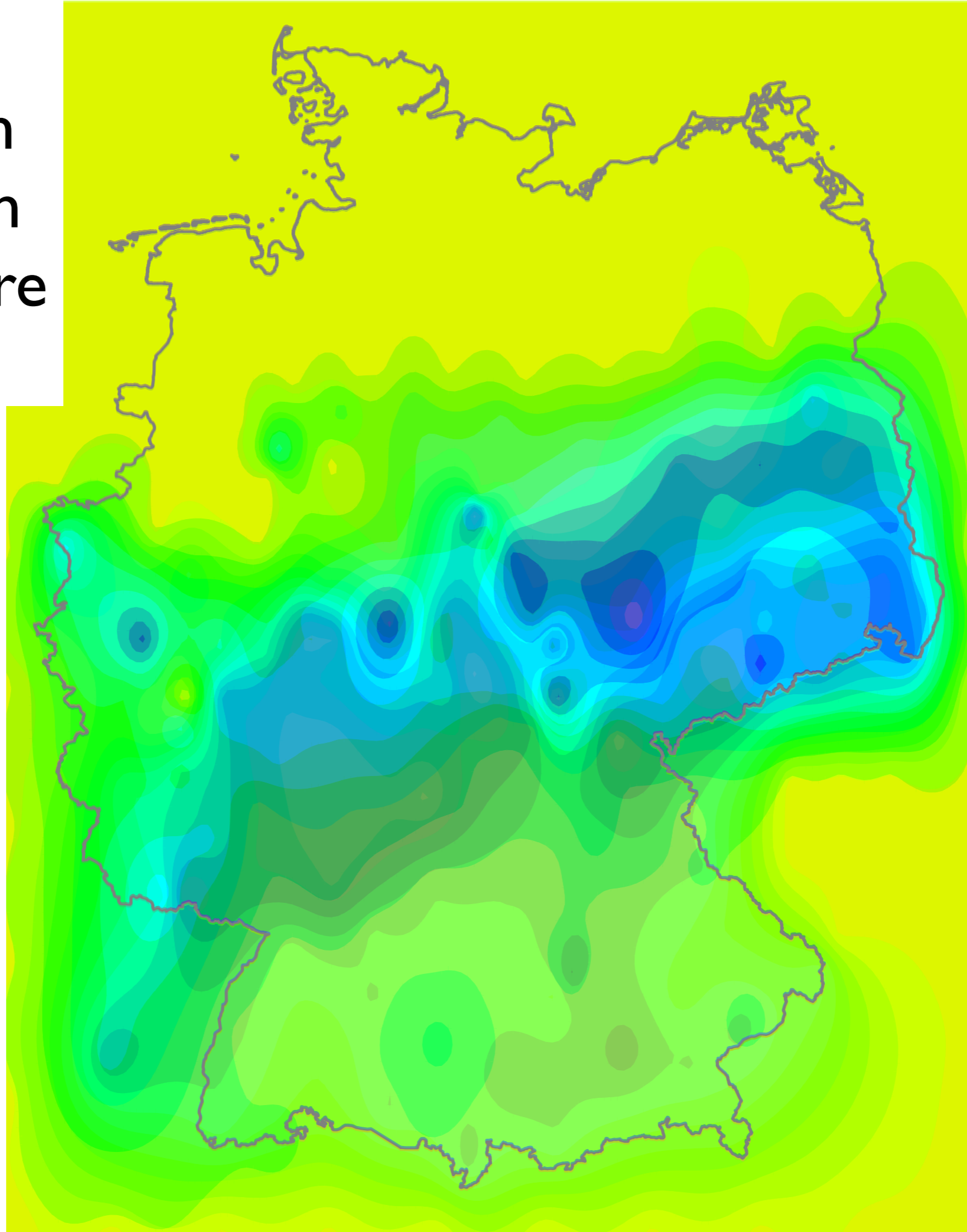
**Consonants**  
of ge- prefix  
(averages over  
14 lexemes)



**Vowels**  
of ge- prefix  
(averages over  
14 lexemes)

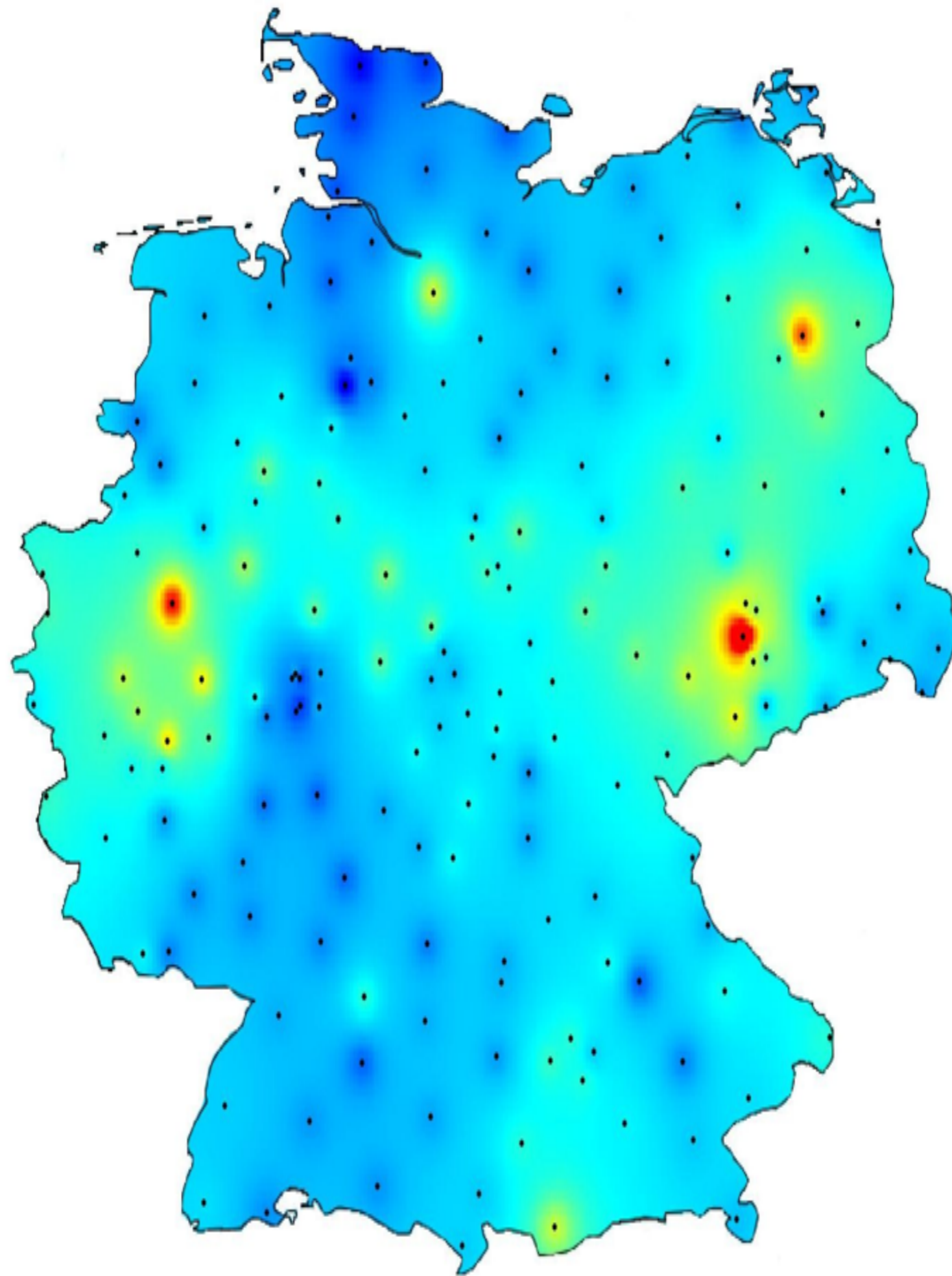


Vowels are  
'strong' in an  
area in which  
consonants are  
'weakening'

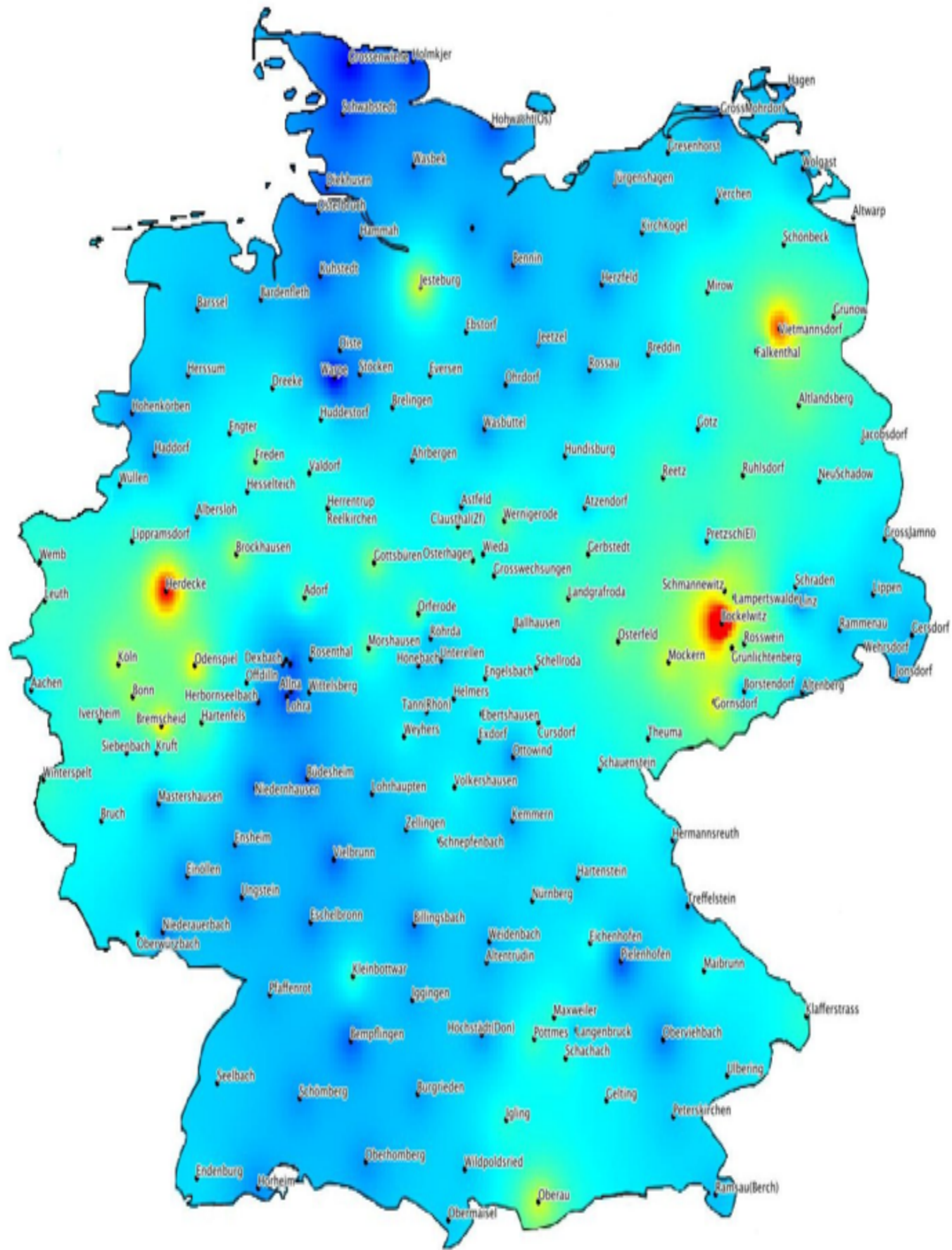


# Looking at Regularity

- Idea: do not look at the actual forms, but at the regularity of the changes between locations
- Practice: compare all neighbouring locations, and establish how regular the correspondences are
- Aggregate all irregularity per location:
  - ▶ blue: regular correspondences to neighbours
  - ▶ yellow/red: irregular correspondences to neighbours







# Working with Survey Data

- Data is never perfectly organised (errors, disagreements about coding)
- Instead of complaining: correct errors, or make your own version of the coding !
- The git-approach: share the data, let others make a fork/clone, possible merging the variants later

# 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)