

## Zielsetzung

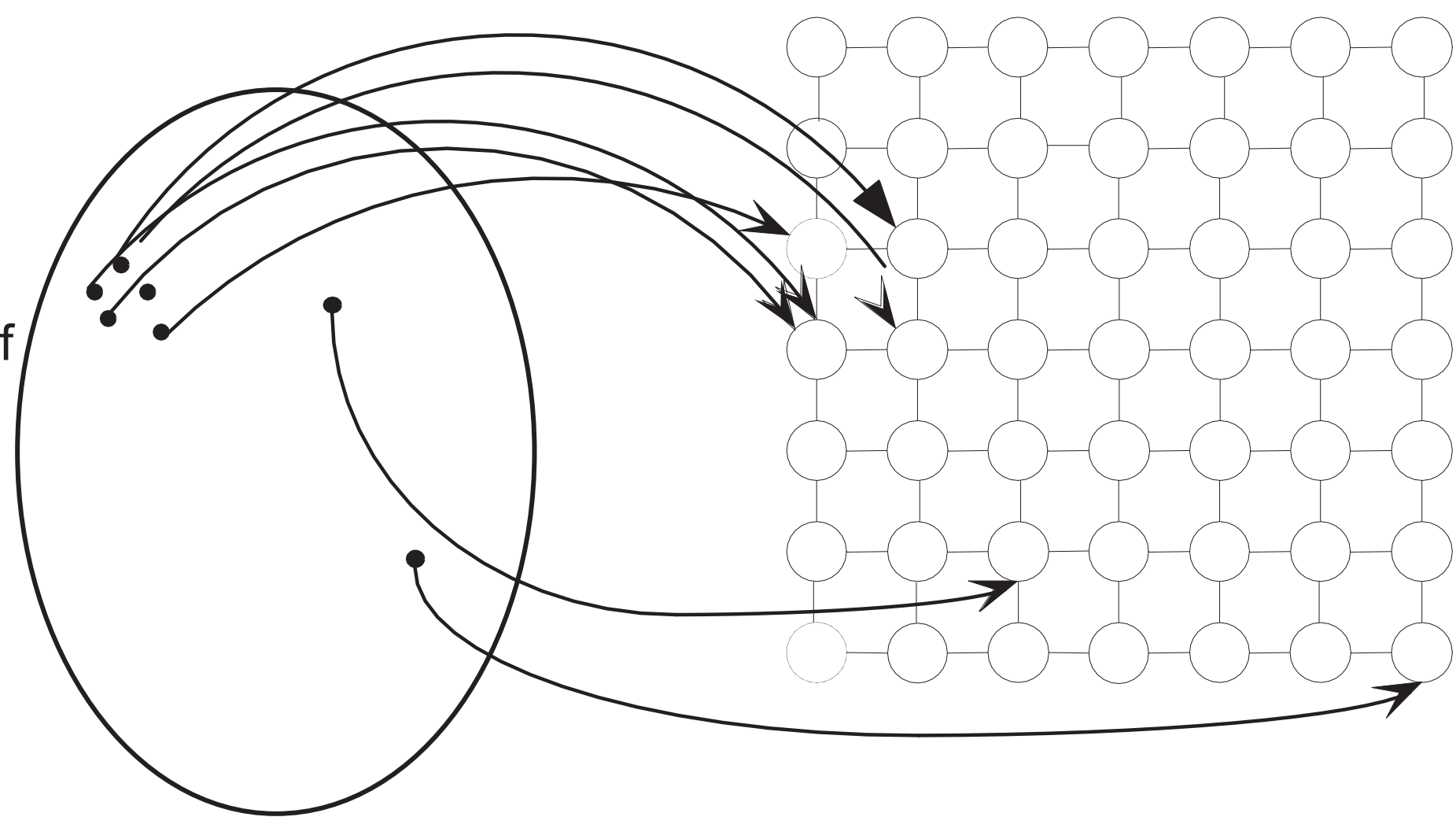
Das Ziel der Diplomarbeit ist die bestehende SOM Toolbox (Software für SOMs) durch Clustering zu erweitern und die entstehenden Cluster grafisch darstellen und zu beschriften. Insbesondere soll das Erkunden der Clusterstrukturen interaktiv möglich sein und die Beschriftungen manuell verbessert werden können.

## Self-Organizing Map (SOM)

2-dimensionales neuronales Netz

Mapping von mehrdimensionalen Daten auf 2-dimensionale Fläche

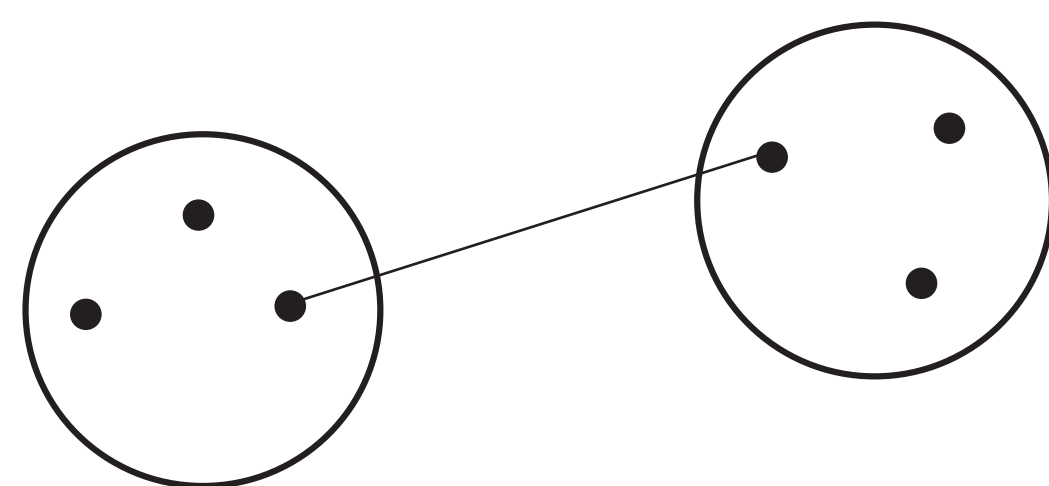
Topologieerhaltend



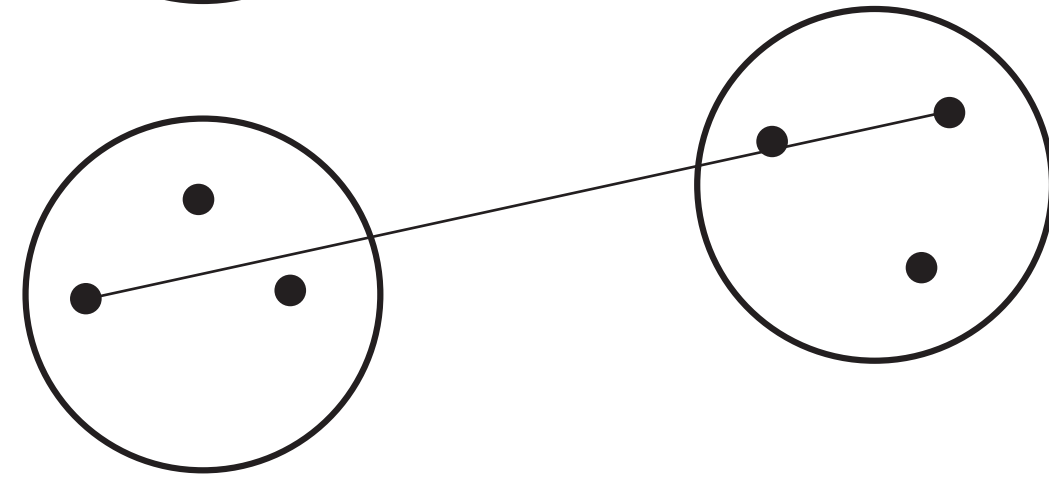
## Clustering Methoden

Bottom-up: Iteratives Zusammenfügen der jeweils nächsten Cluster. Drei Methoden zur Bestimmung der "nächsten Cluster" implementiert:

Single linkage (kürzest mögliche Distanz zweier Clusterelemente)



Complete Linkage (längste mögliche Distanz zweier Clusterelemente)



Ward's Linkage (Anstieg der Varianz, wenn die beiden Cluster zu einem vereint werden)

$$D = V(AB) - [V(A) + V(B)]$$

## Cluster darstellen

Grenzlinien einzeichnen

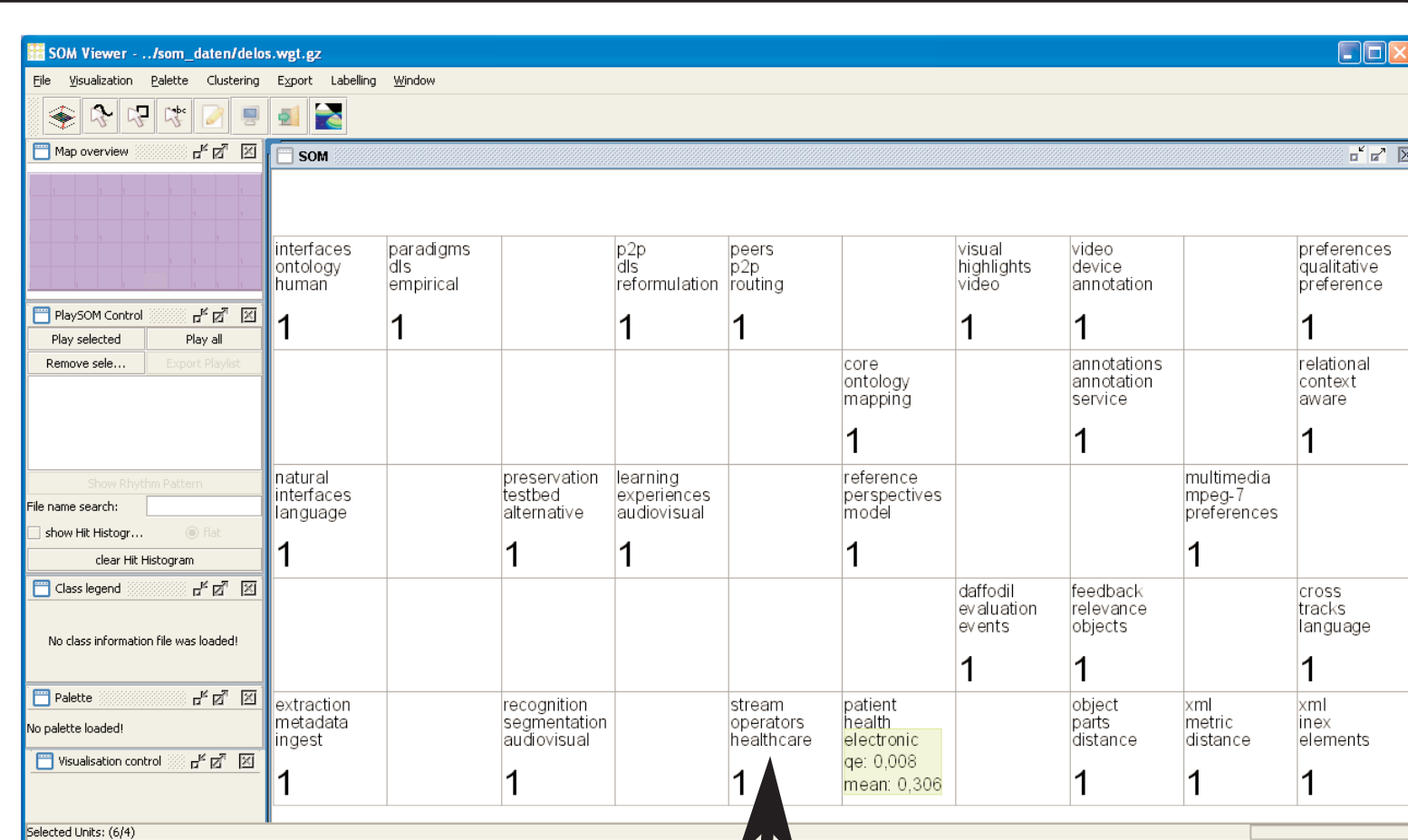
Cluster einfärben

## Cluster beschriften

Text (anhand bestehender Unit Labels)

Größe (abhängig von der Kartenbreite)

Position (Schwerpunkt oder Mittelpunkt des umgebenden Rechtecks)



## Unit Labels

Textdaten: Wörter, die in den enthaltenen Dokumenten am häufigsten vorkommen  
Bei anderen Daten: die Bezeichnungen der Merkmale welche alle Daten auf einer Unit am ähnlichsten haben.

## Bearbeiten der Beschriftungen

Text (auch mehrzeilig)

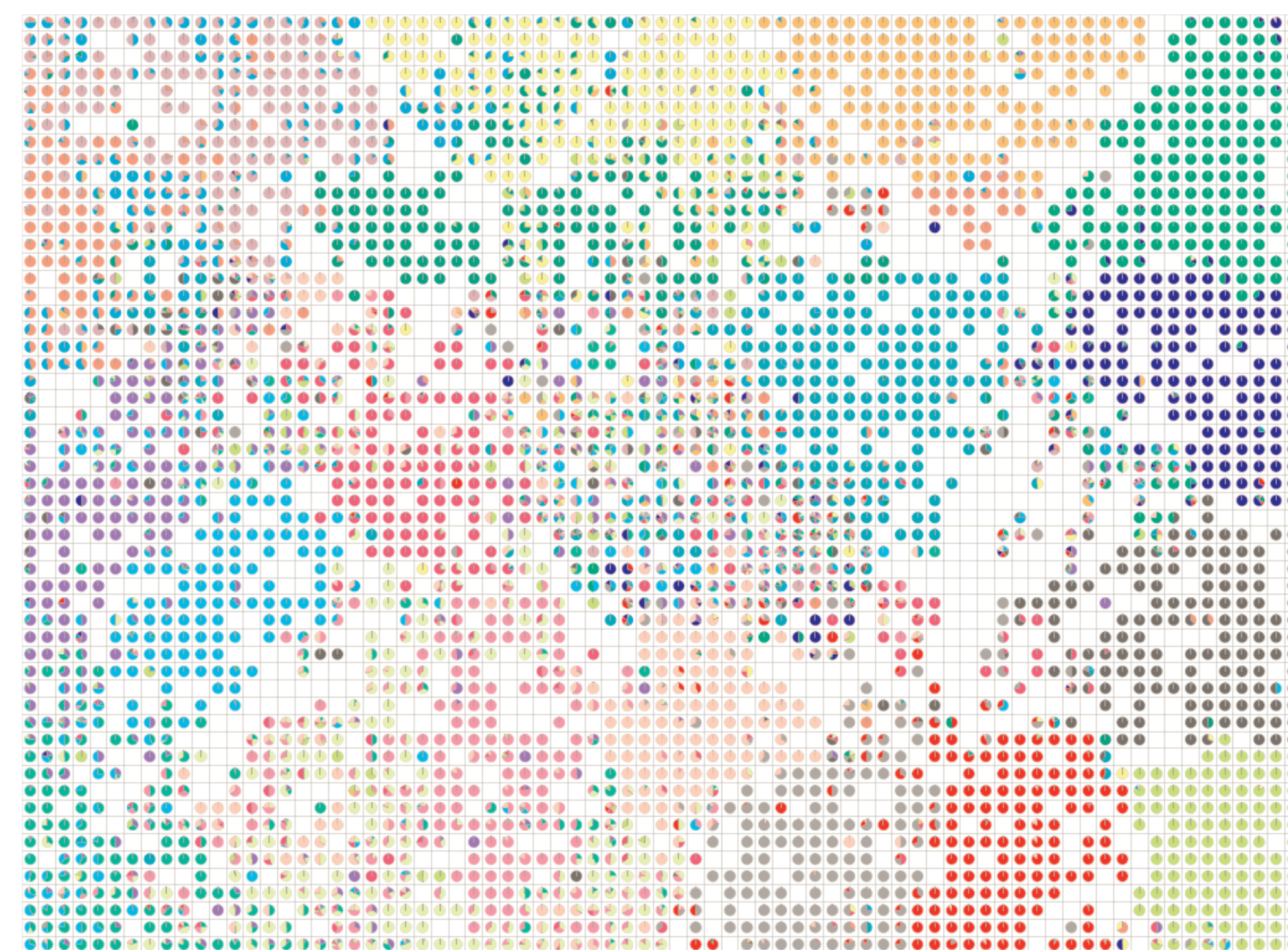
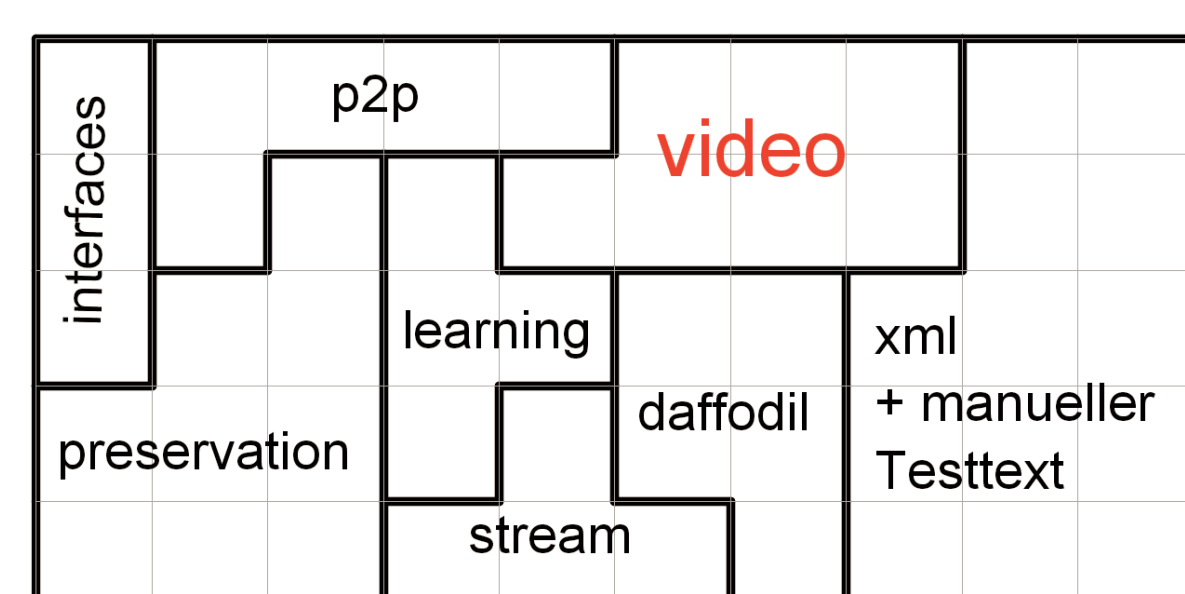
Schriftgröße

Platzierung (mittels drag and drop)

Farbe (hervorheben / abheben vom Hintergrund)

Rotation

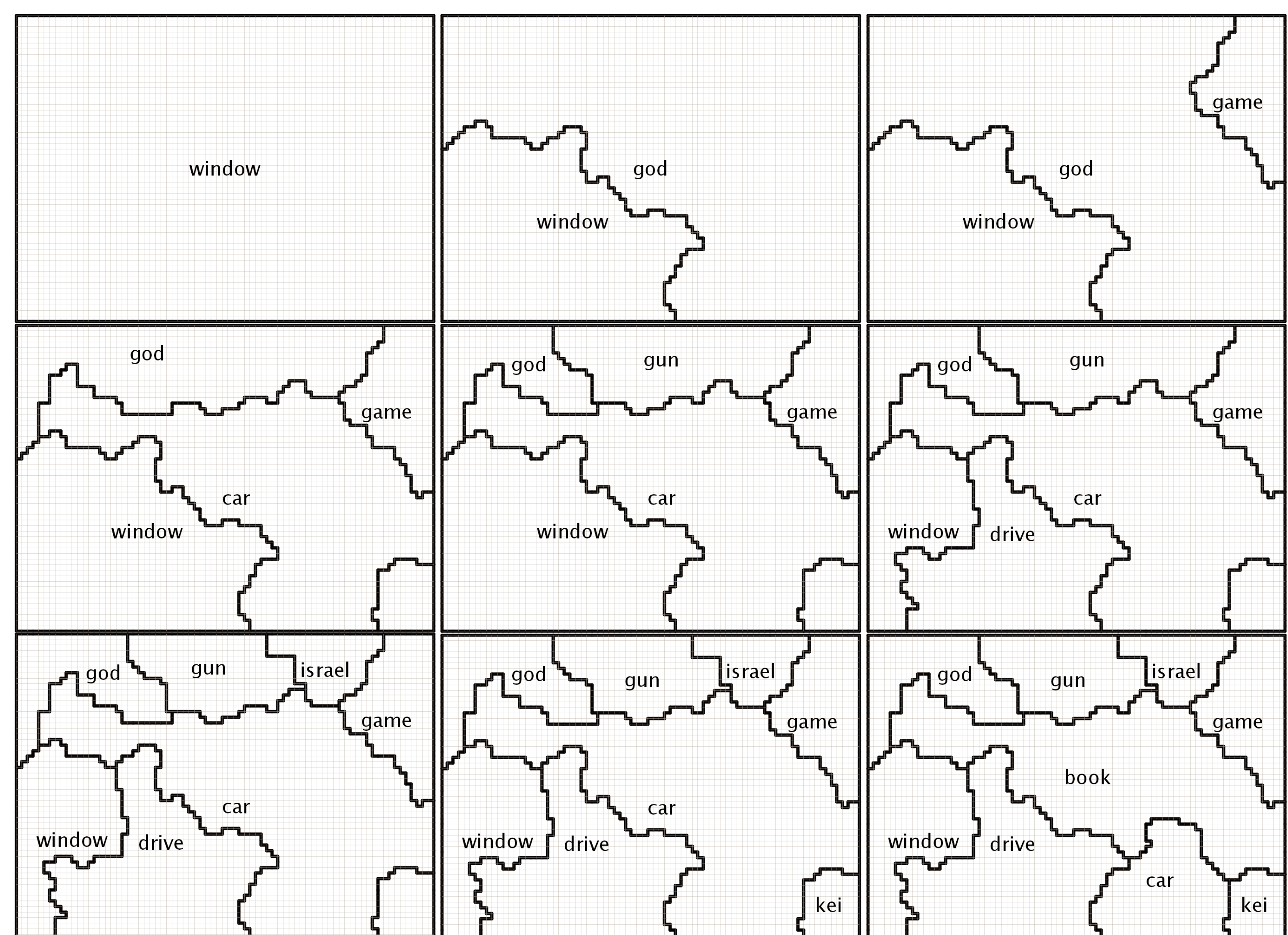
nicht anzeigen



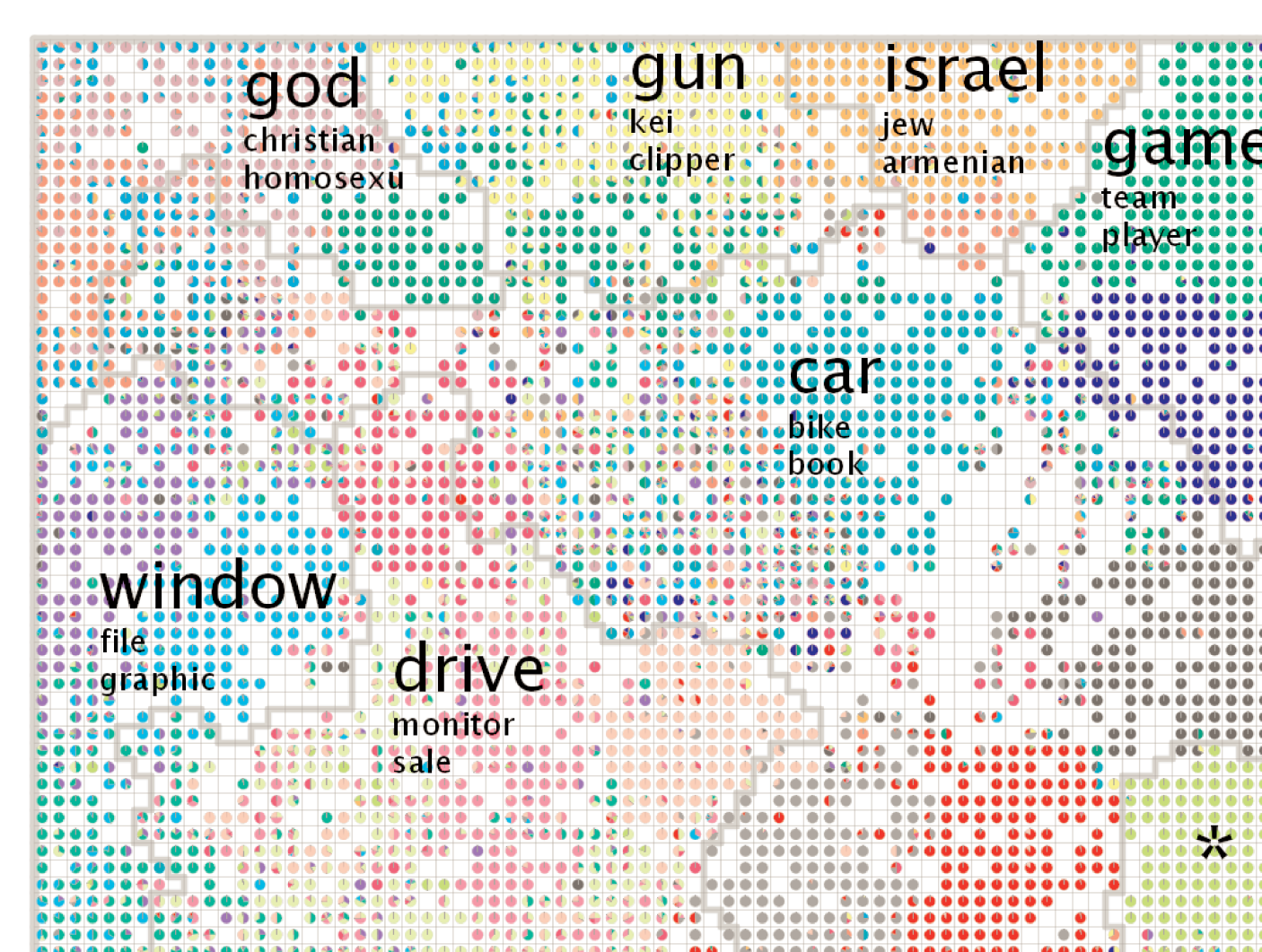
alt.atheism	Orange
comp.graphics	Purple
comp.os.ms-windows.misc	Green
comp.sys.ibm.pc.hardware	Yellow
comp.sys.mac.hardware	Pink
comp.windows.x	Cyan
misc.forsale	Light Blue
rec.autos	Grey
rec.motorcycles	Red
rec.sport.baseball	Dark Blue
rec.sport.hockey	Dark Green
sci.crypt	Light Green
sci.electronics	Red
sci.med	Dark Blue
sci.space	Dark Green
soc.religion.christian	Yellow
talk.politics.guns	Orange
talk.politics.mideast	Light Blue
talk.politics.misc	Dark Green
talk.religion.misc	Cyan

Farbzuordnung zu Klassen

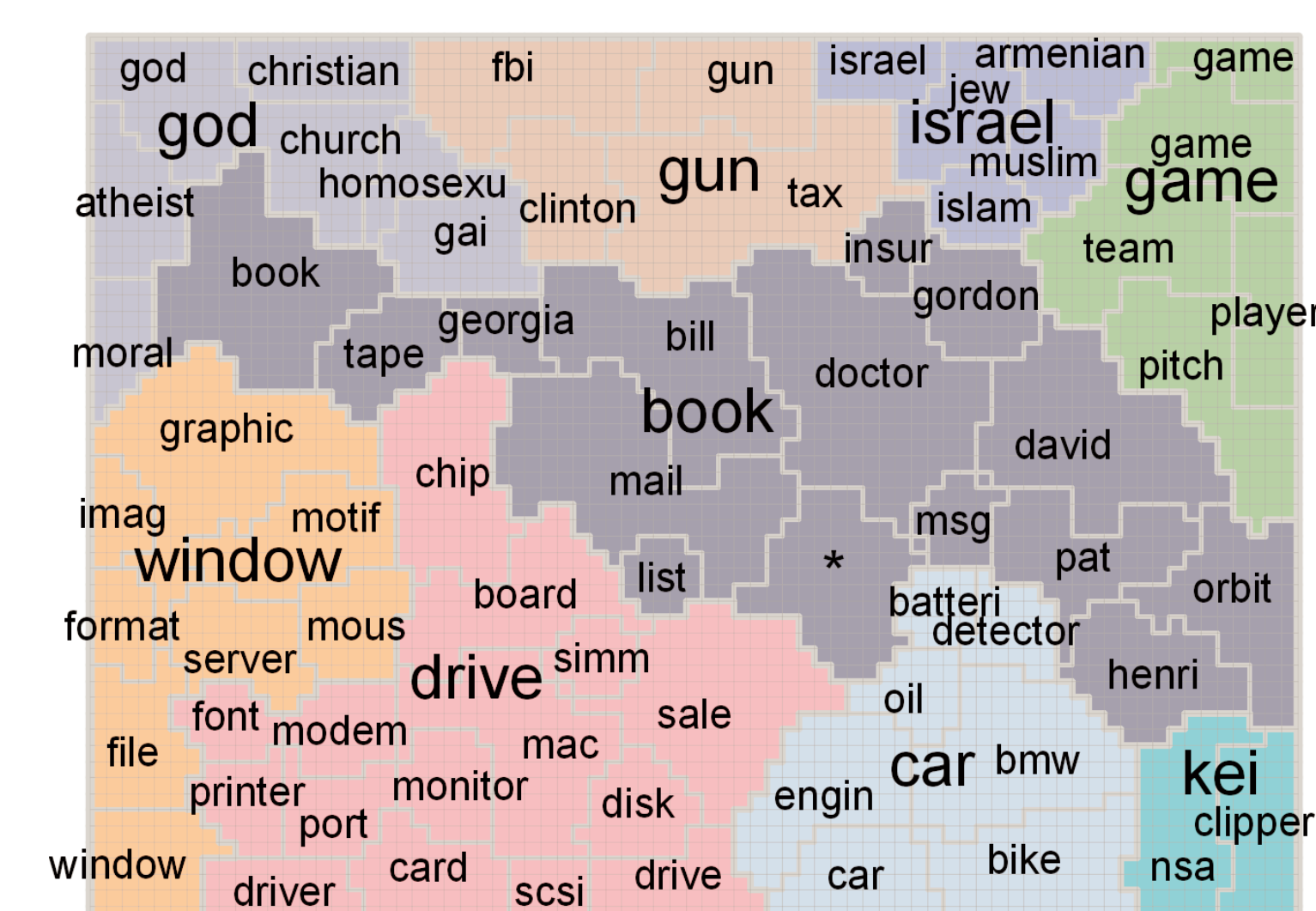
SOM von je 1.000 Postings aus 20 Newsgroups (Klasseninformation)



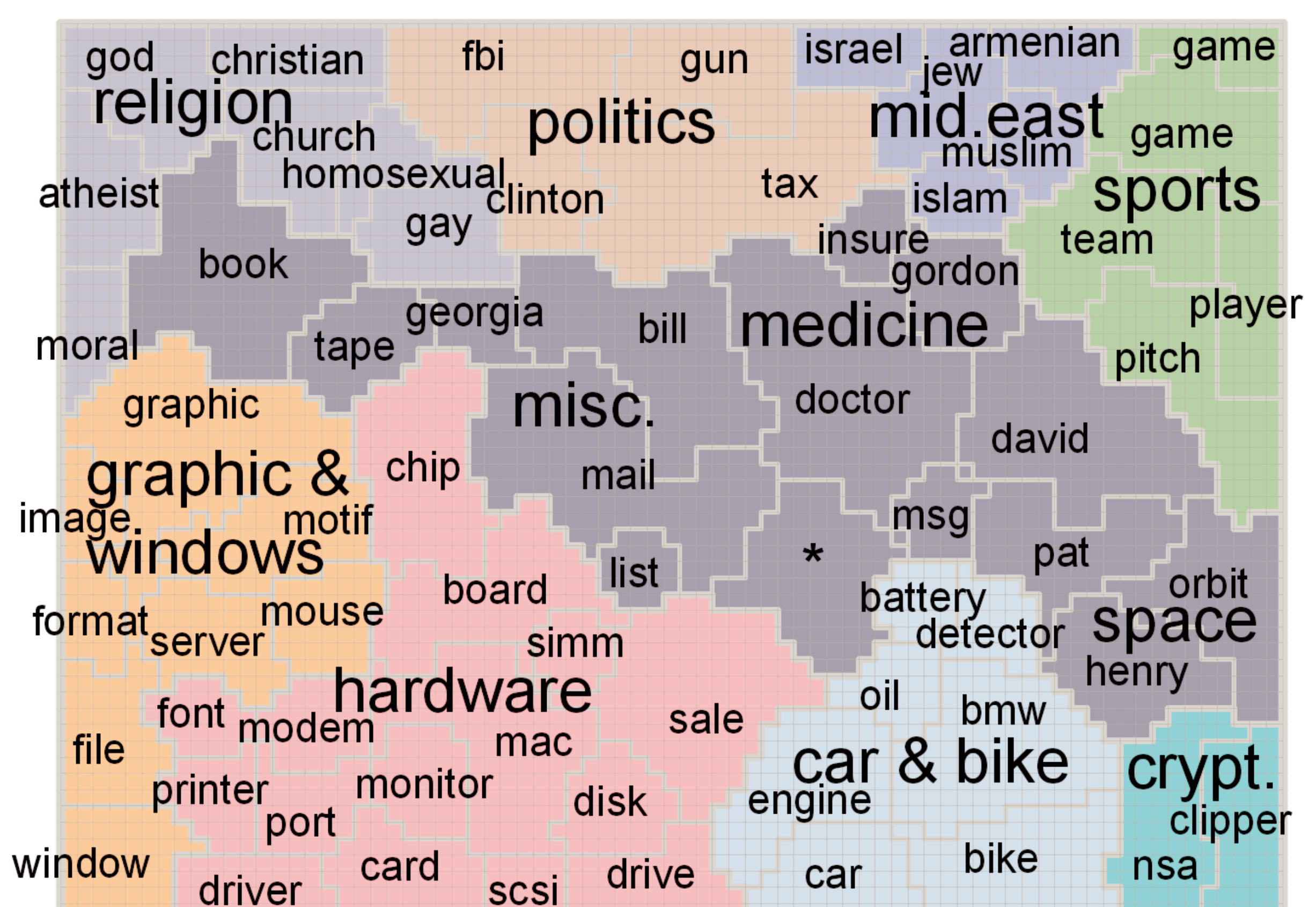
9 Clustering Schritte (in umgekehrter Reihenfolge) mit Beschriftung



7 Cluster mit je 3 Beschriftungen + Klasseninformation



9 große und 67 kleine Cluster



Editierte Karte