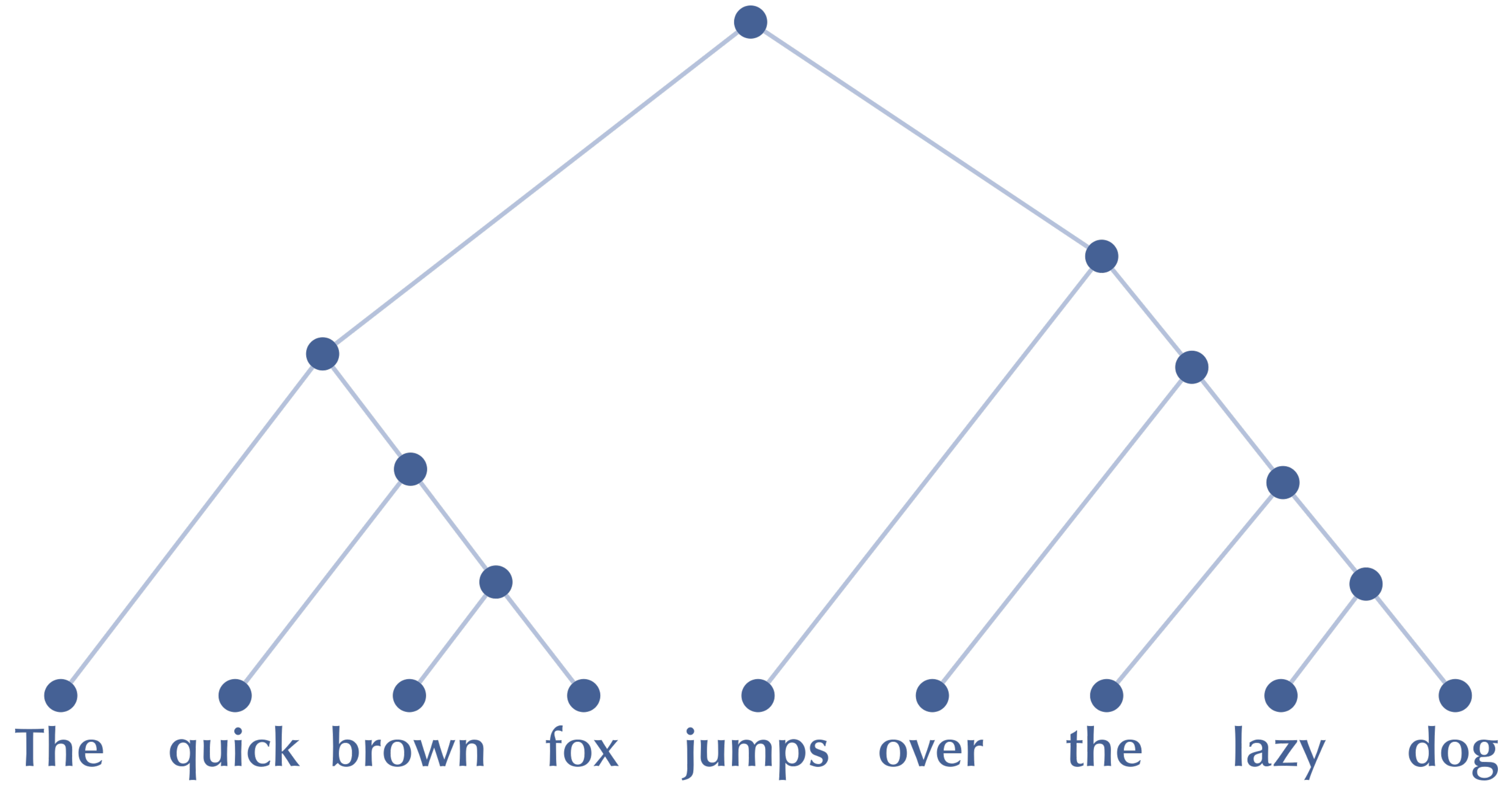


# Topology of Language

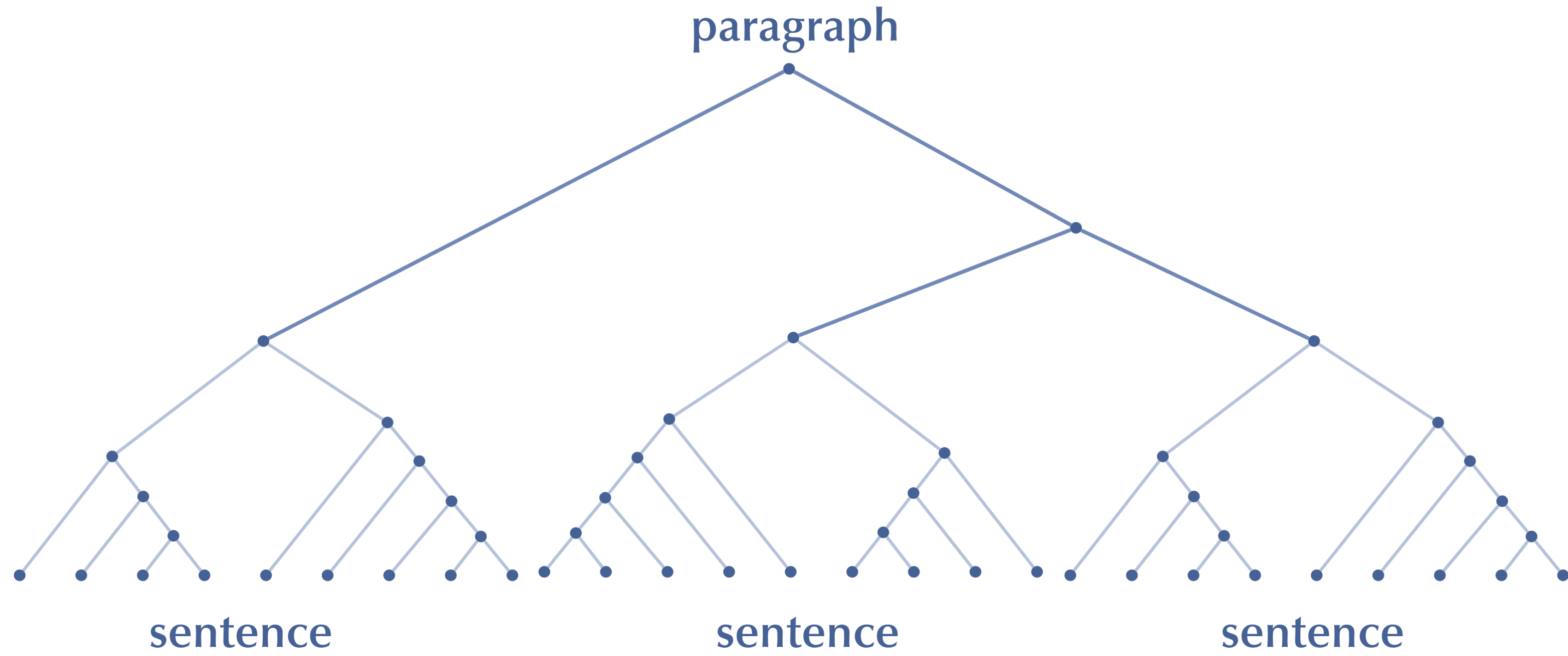
Topological Techniques for Natural Language Processing

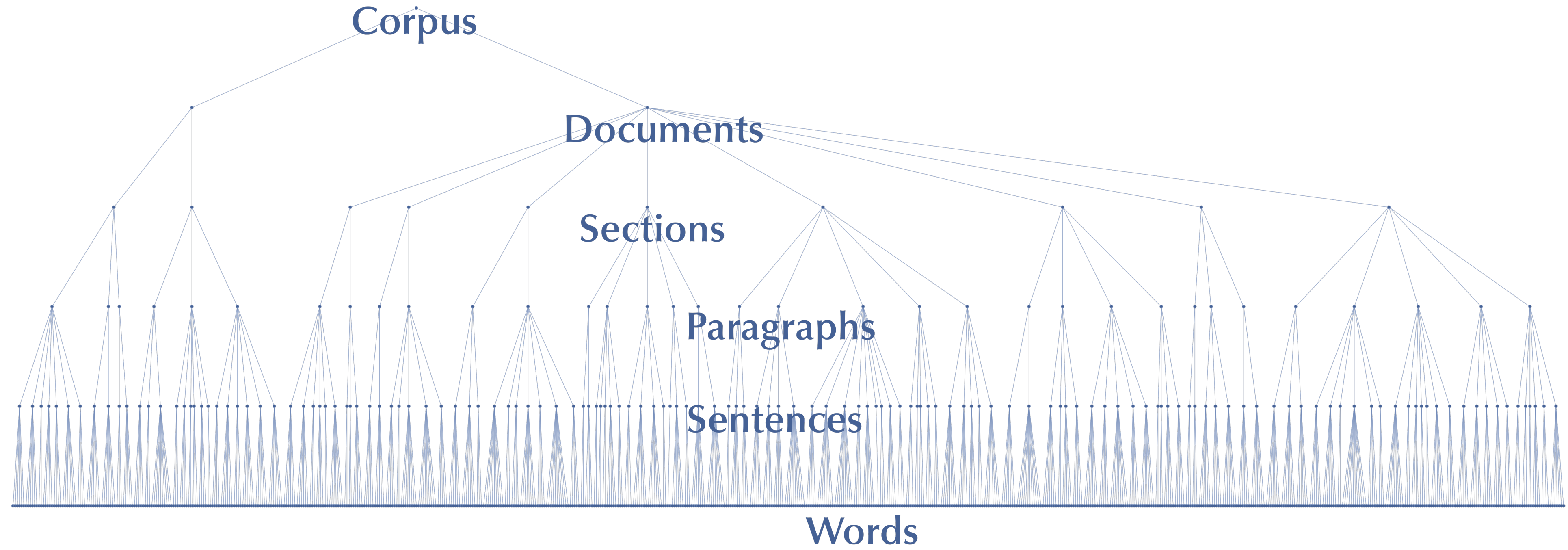
# Vision

Language has structure



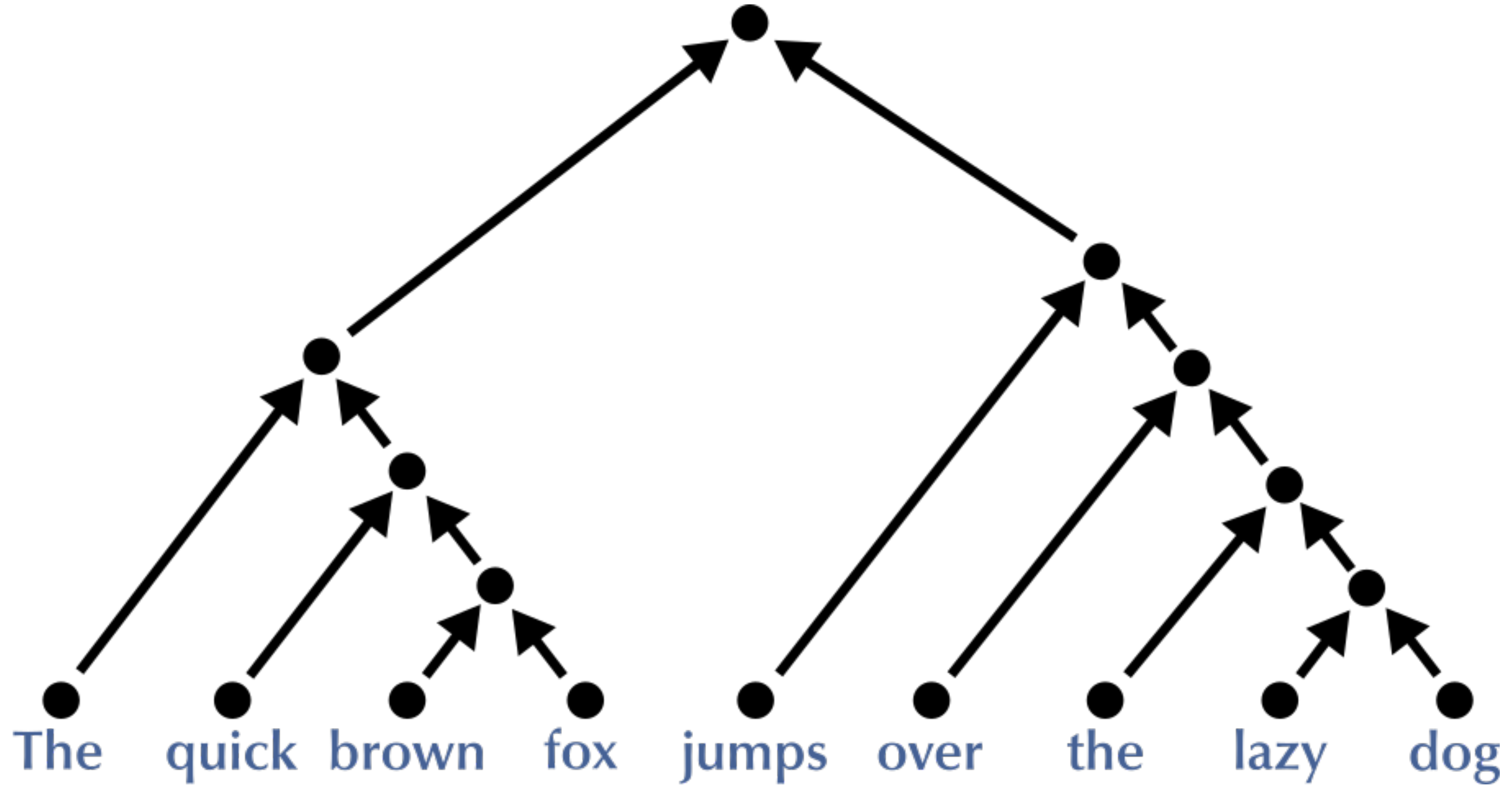
Many layers of  
structure





We could build a  
category from this





Yoneda embedding of a  
document?

All the constituent parts  
of the document

Simplify: all the word usages in a document

This is the classical “bag-of-words” model

What about words?  
In the opposite category...

A word is the set of all  
its usages in a corpus

Simplify: all the sentences  
that use the word

You shall know a word  
by the company it  
keeps  
— John Rupert Firth

# Words as Distributions



# Overview

Represent a word as a  
document of all the  
sentences that use the  
word

Represent such a  
document using a bag-of-  
words model

Represent a word as the  
multinomial distribution of  
words that occur nearby

Information geometry  
tells us that multinomial  
distributions live on a  
manifold

$$d(p, q) = \arccos \left( \sum_i \sqrt{p_i} \sqrt{q_i} \right)$$

We can explore the  
geometry of words  
under this distance  
metric

# Messy Details



## Some problems:

- Word order matters
- Words carry different amounts of information
- Language use is noisy

Words order matters:

Words used *before* a given word  
Words used *after* a given word

Information in words:

Approximate the information  
of a word in a given context  
and weight words accordingly

Language is noisy:

Approximate a background  
noise model of word use and  
decompose the multinomial  
into a mixture of signal and  
noise distributions

Where do we get models  
of entropy and noise?

$$X \approx UV$$

Where

$X$  is an  $N \times D$  matrix

$U$  is an  $N \times d$  matrix

$V$  is an  $d \times D$  matrix

Minimize

$$\sum_{i=1}^N \sum_{j=1}^D \mathbf{Loss} \left( X_{ij}, (UV)_{ij} \right)$$

Subject to constraints...

Suppose

$$X \sim \Pr(\cdot \mid \Theta)$$

where

$$\Theta = UV$$



Let the loss be the  
negative log likelihood  
of observing  $X$  given  $\Theta$

# Multinomial Matrix Factorization

Minimize

$$\sum_{i=1}^N \sum_{j=1}^D - (UV)_{ij} \cdot \log (X_{ij})$$

Subject to

$$(UV)\mathbf{1} = \mathbf{1} \text{ and } (UV)_{ij} \geq 0$$

# Probabilistic Latent Semantic Analysis

Minimize

$$\sum_{i=1}^N \sum_{j=1}^D - (UV)_{ij} \cdot \log (X_{ij})$$

Subject to

$$U\mathbf{1} = \mathbf{1}, V\mathbf{1} = \mathbf{1} \text{ and } U_{ij} \geq 0, V_{ij} \geq 0$$

$\Theta$  provides a low rank  
multinomial distribution  
model

$\Theta_{ij}$  is the “background” probability of word  $j$  occurring in the context of word  $i$

$-\log(\Theta_{ij})$  is an approximation of the information carried by word  $j$  in the context of word  $i$

We can weight each entry  
by the (approximate)  
information it carries

A rank 1 model gives  
classical TF-IDF from NLP

$\Theta_i$  provides a model of the  
“background” multinomial  
of word occurrence in the  
context of word  $i$



We can use an EM  
algorithm to decompose  
a given multinomial into  
noise and signal

We are back to where  
our overview left off...

# Embedding Words

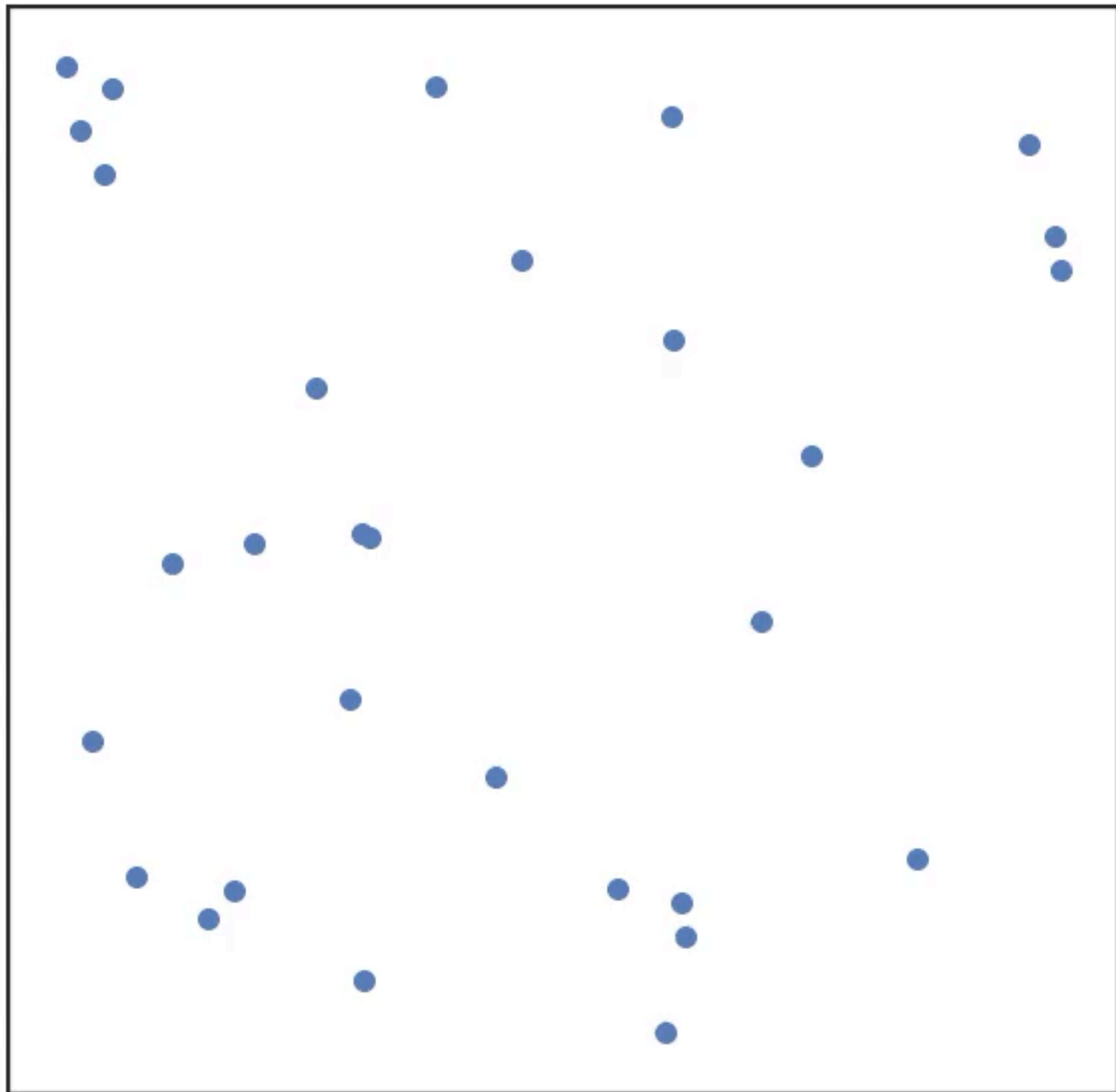
We have high quality and  
high dimensional  
representations of words  
with a natural ambient  
geometry ...

We need a low  
dimensional representation  
that preserves topology



UMAP is a topology based  
dimension reduction  
algorithm

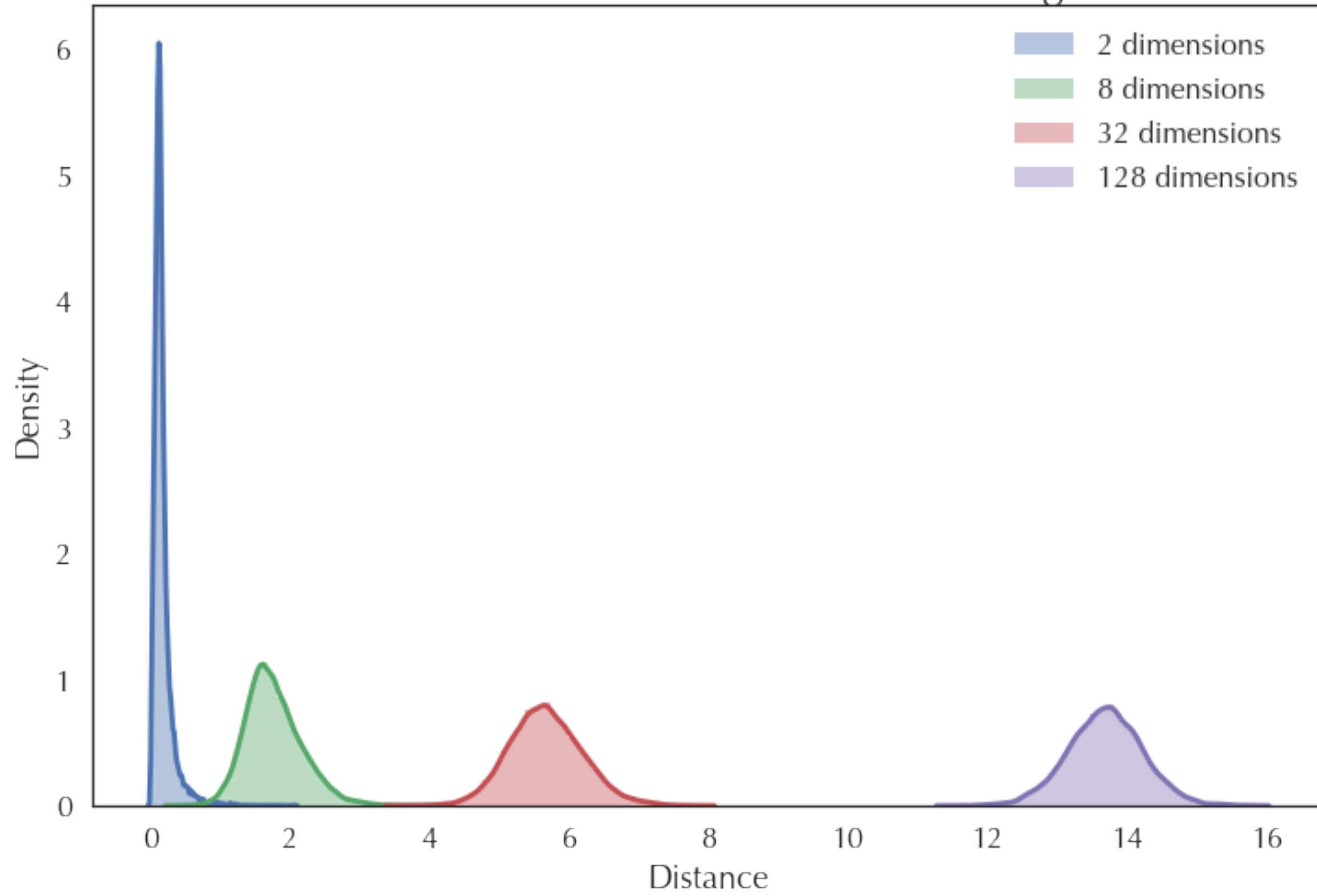
Assumptions:  
**Uniform distribution**  
Locally connected



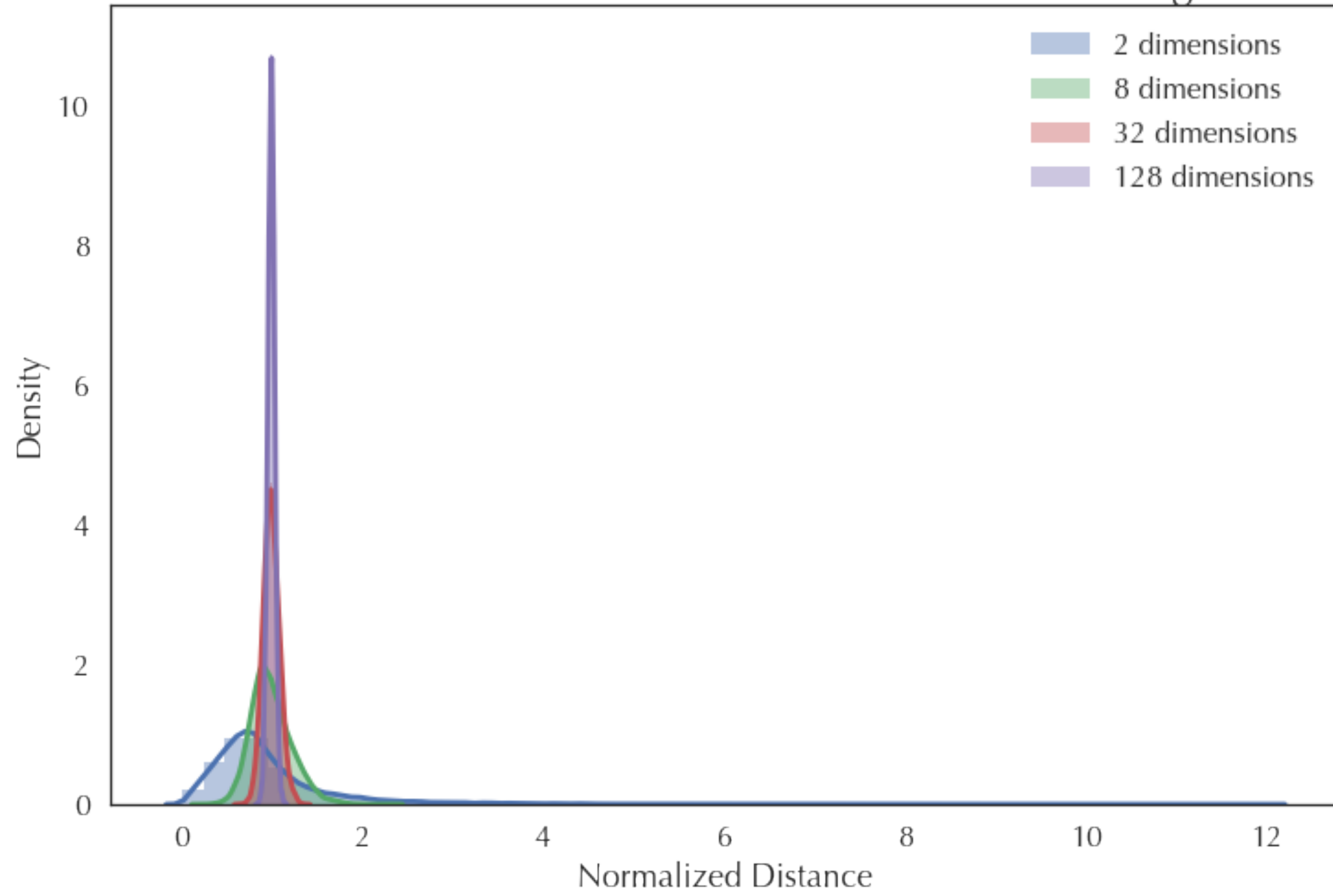


Assumptions:  
Uniform distribution  
**Locally connected**

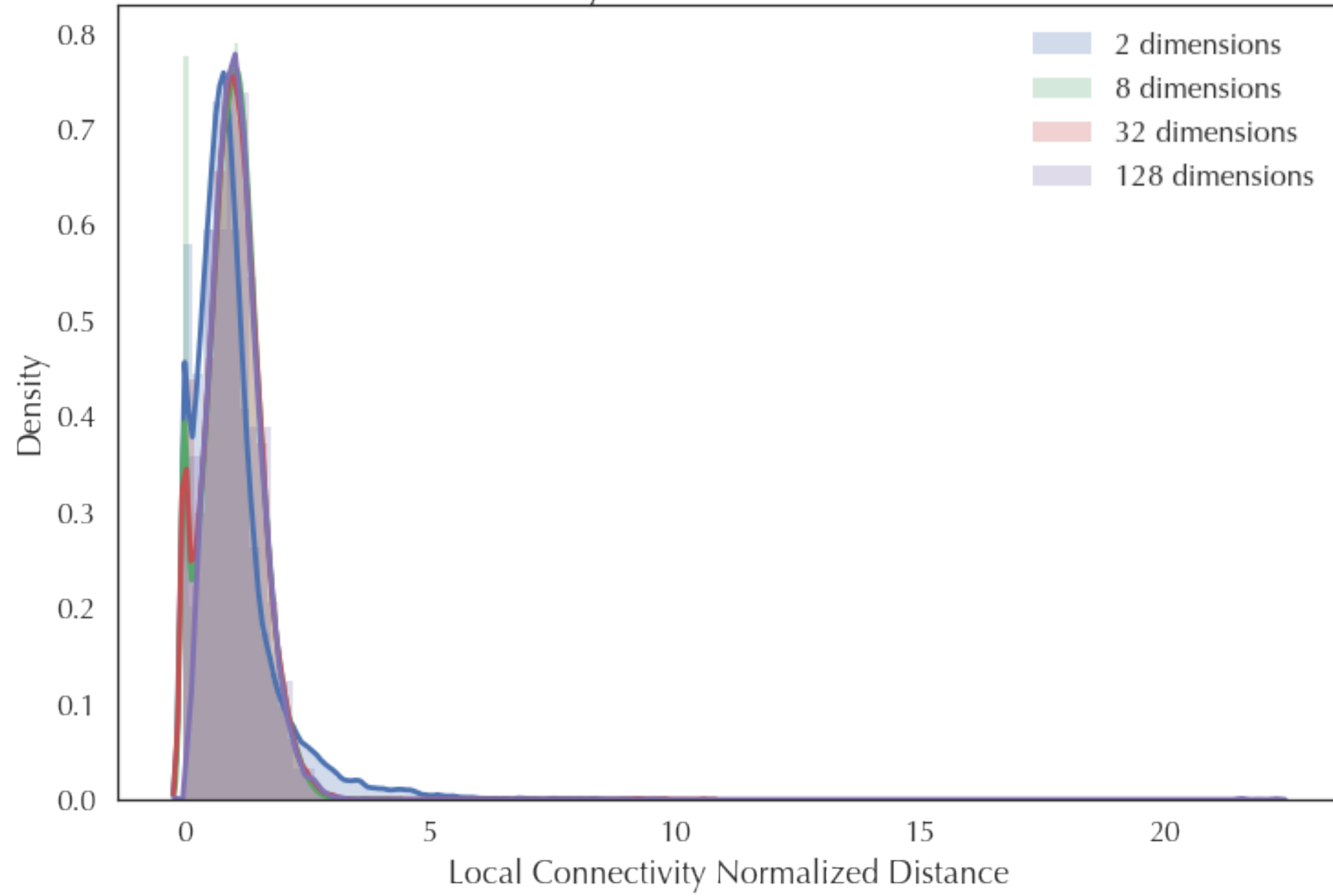
Distribution of distances to 20 nearest neighbors



Distribution of normalized distances to 20 nearest neighbors



Distribution of local connectivity normalized distances to 20 nearest neighbors



We now have a filtered  
simplicial complex

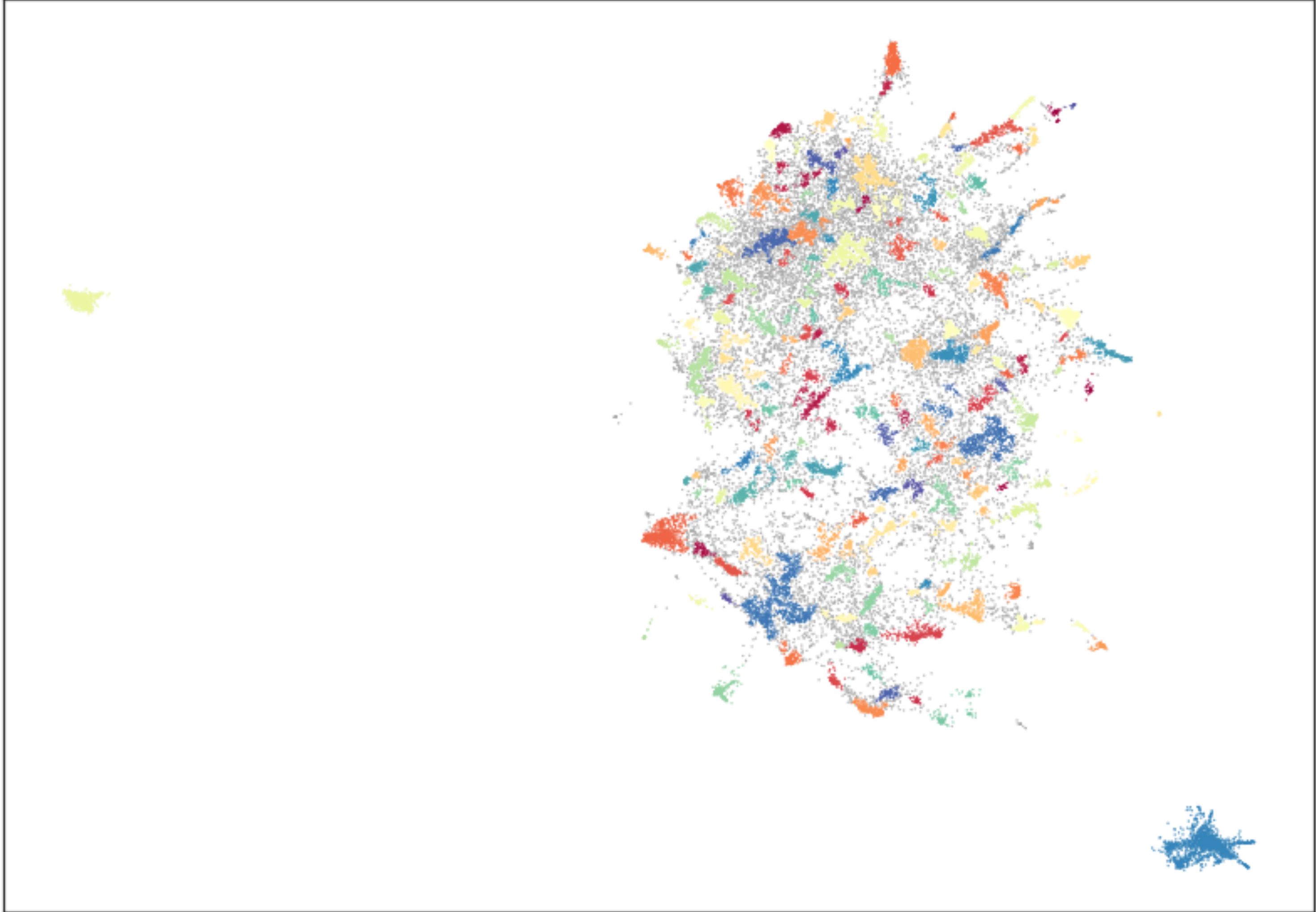
# Optimize a low dimensional representation via cross-entropy

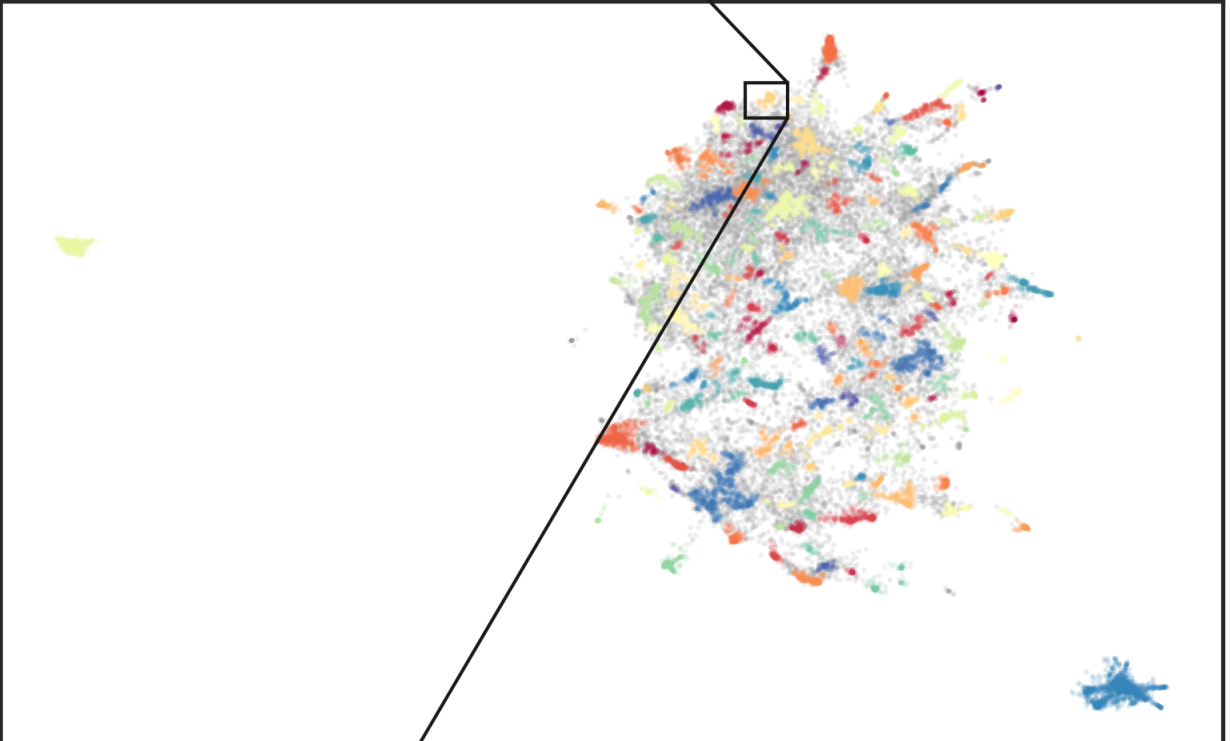
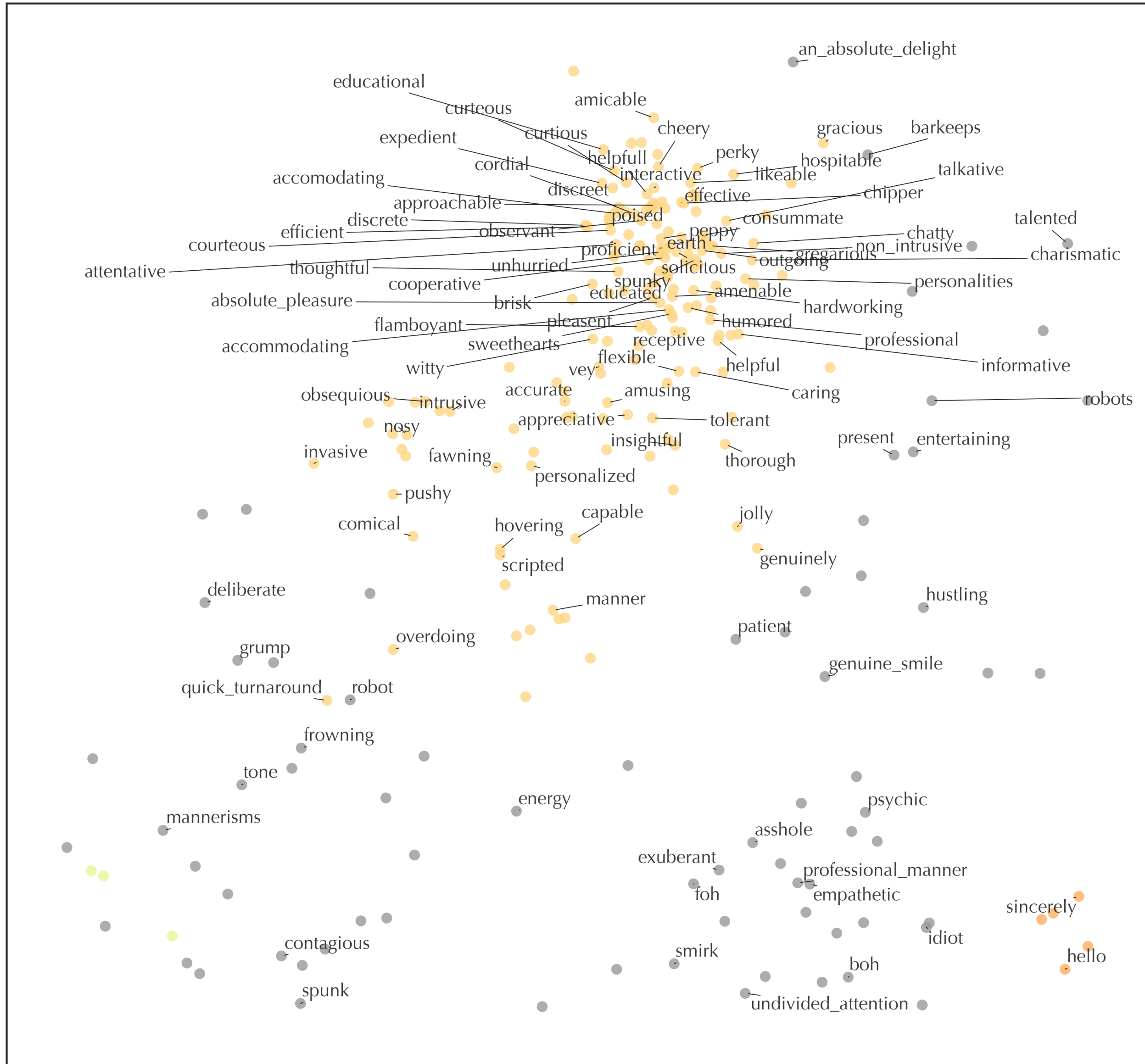
$$\sum_a \mu(a) \log \left( \frac{\mu(a)}{\nu(a)} \right) + (1 - \mu(a)) \log \left( \frac{1 - \mu(a)}{1 - \nu(a)} \right)$$

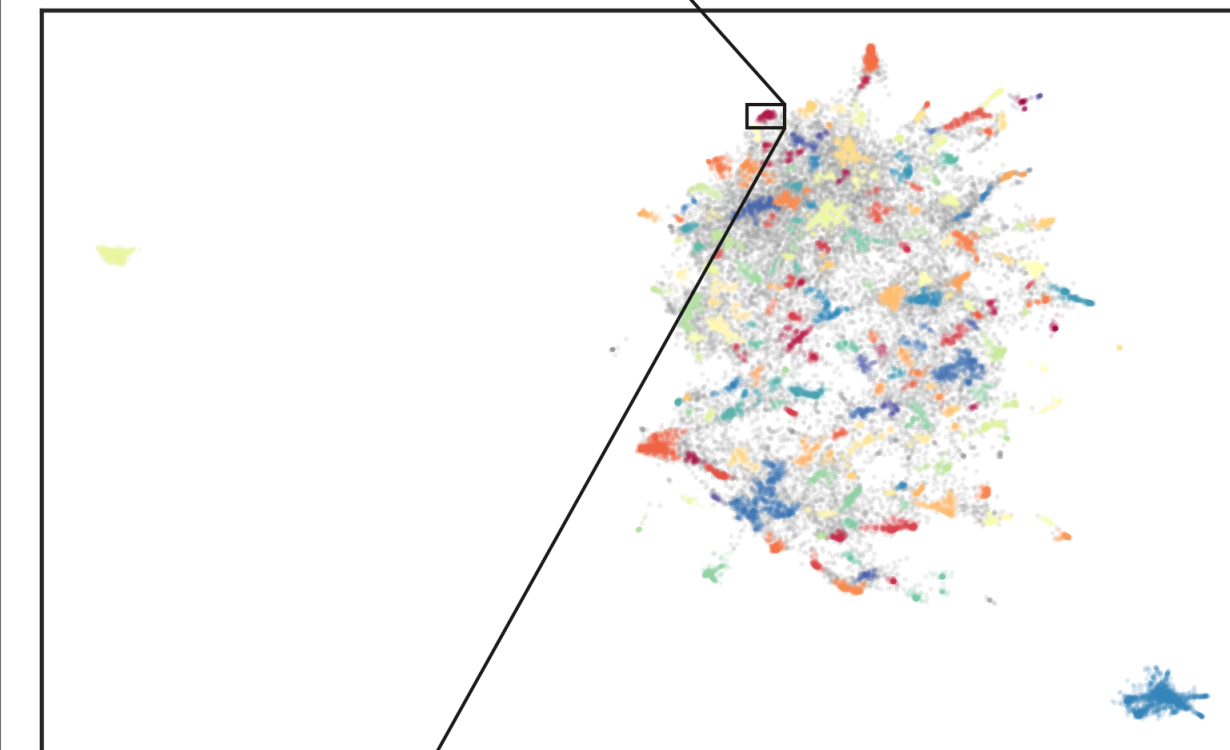
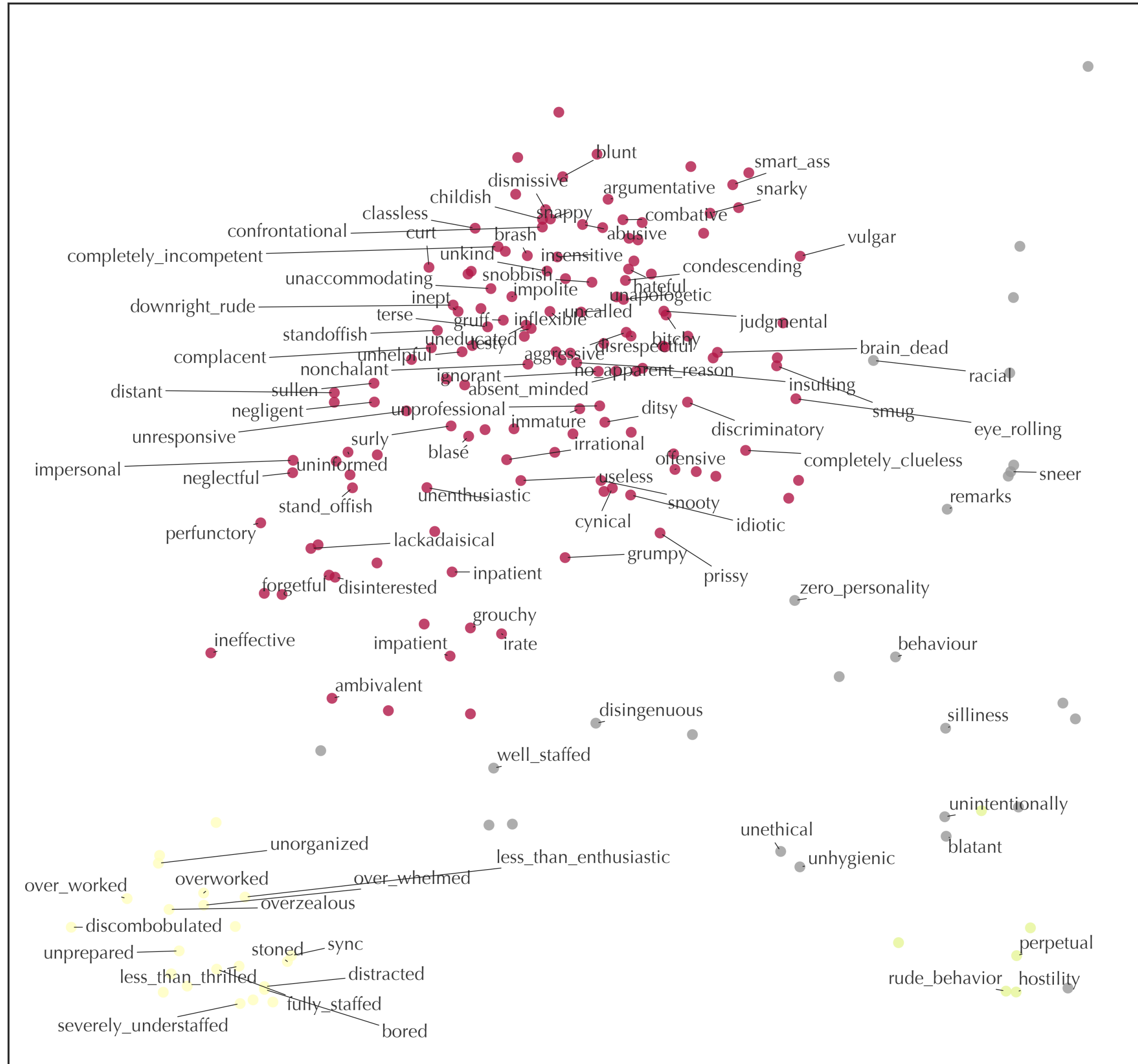
$$d(p, q) = \arccos \left( \sum_i \sqrt{p_i} \sqrt{q_i} \right)$$

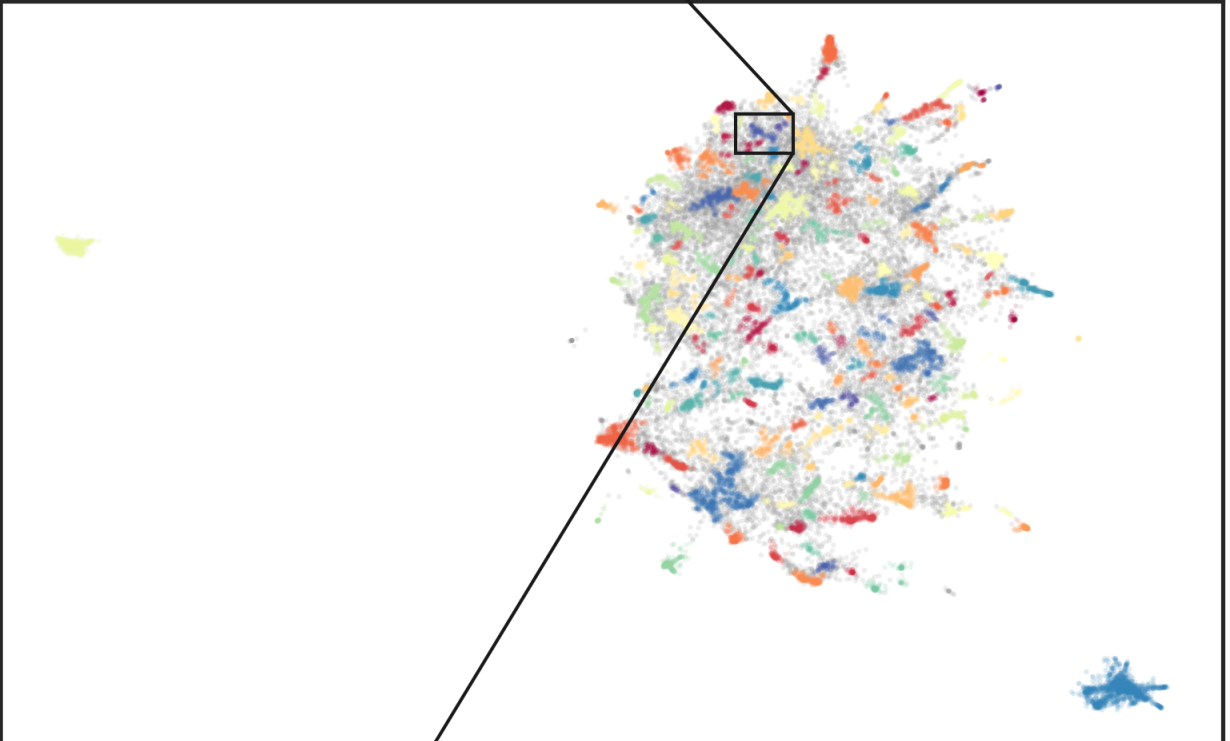
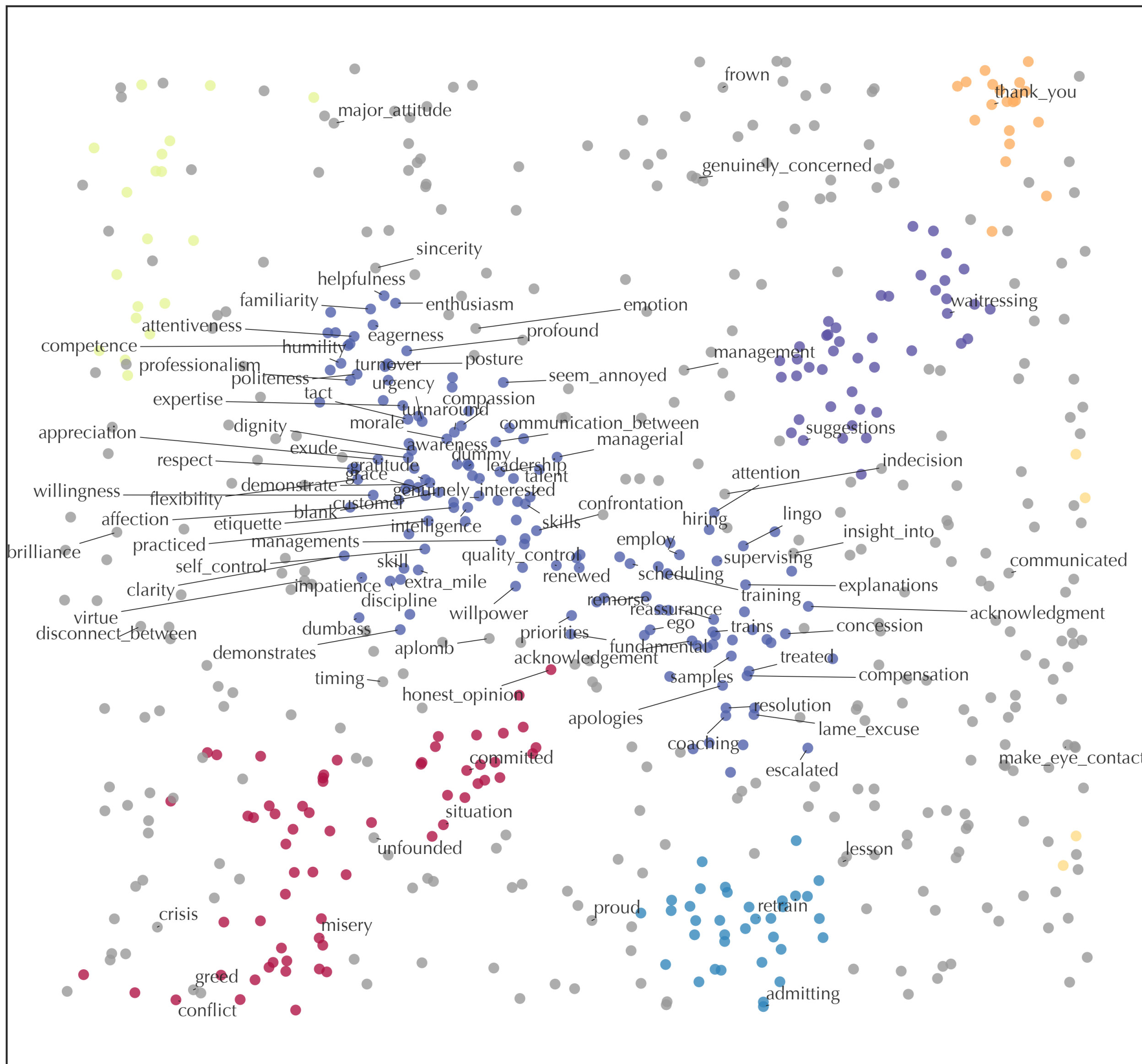


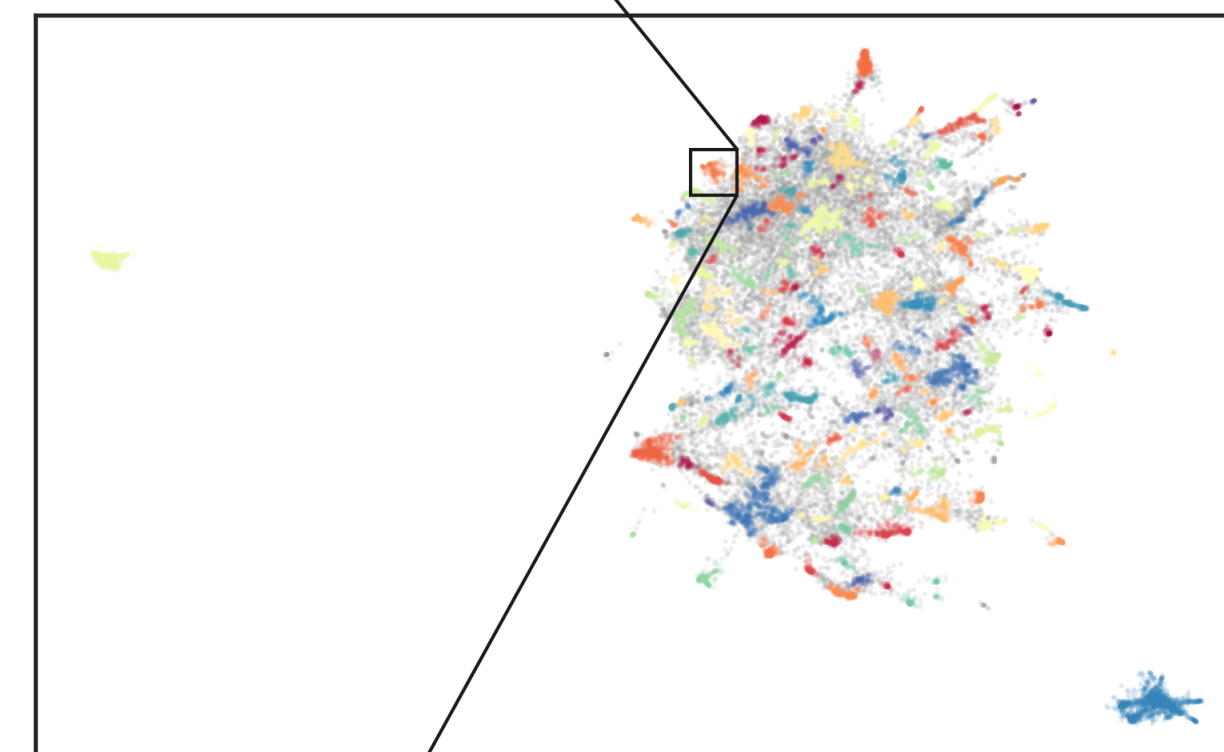
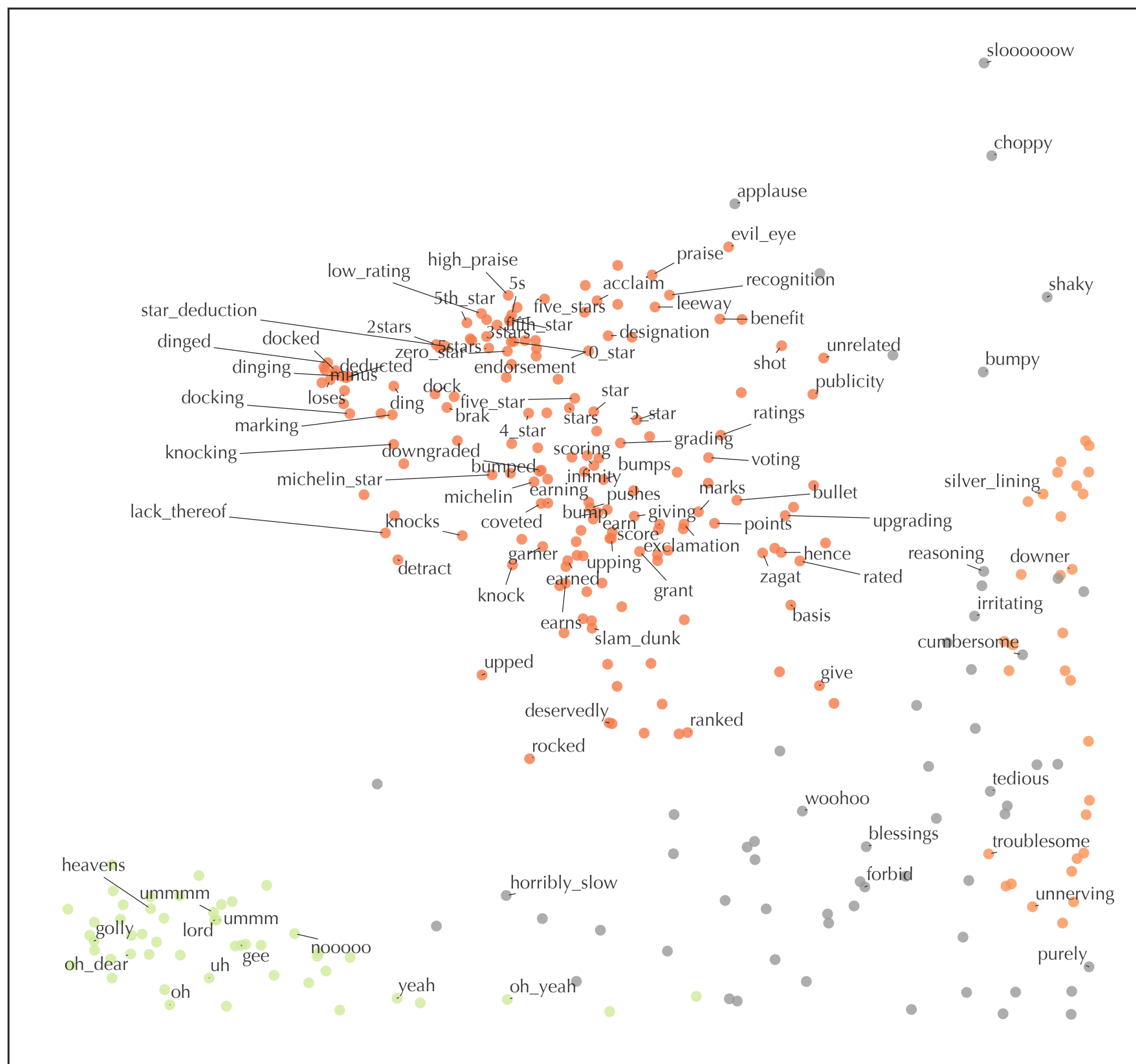
# Embedding Yelp Reviews

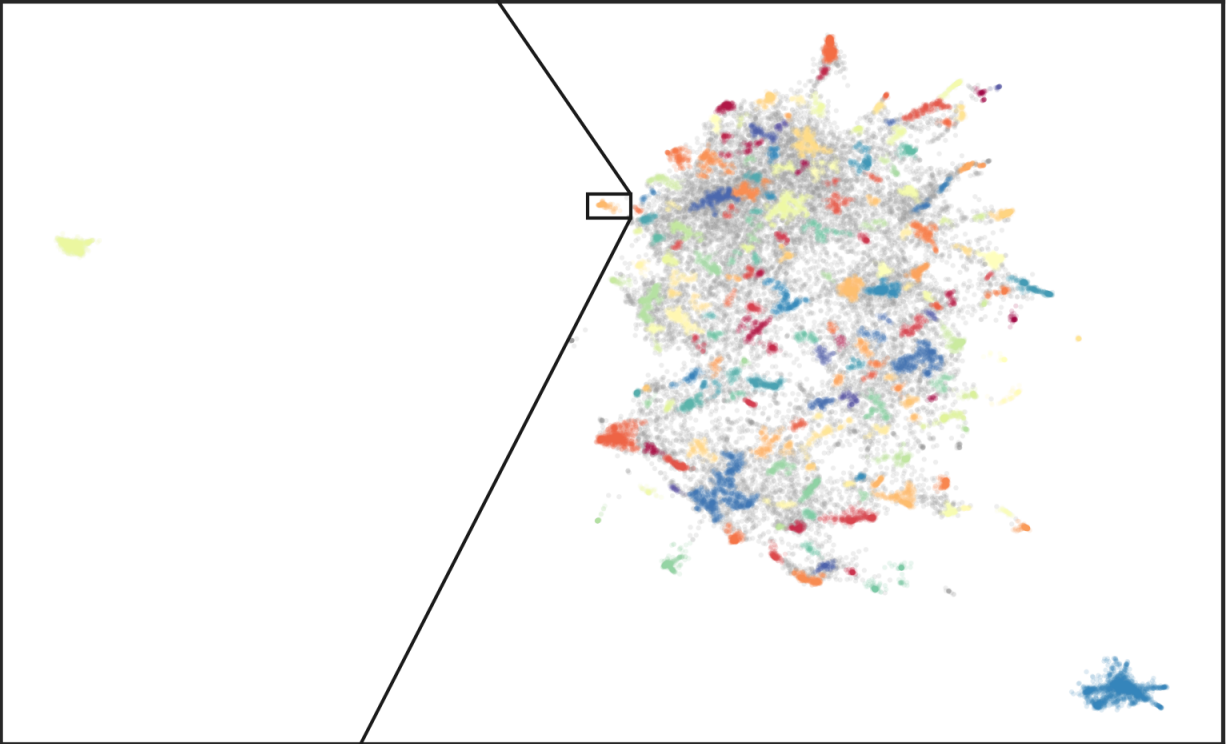
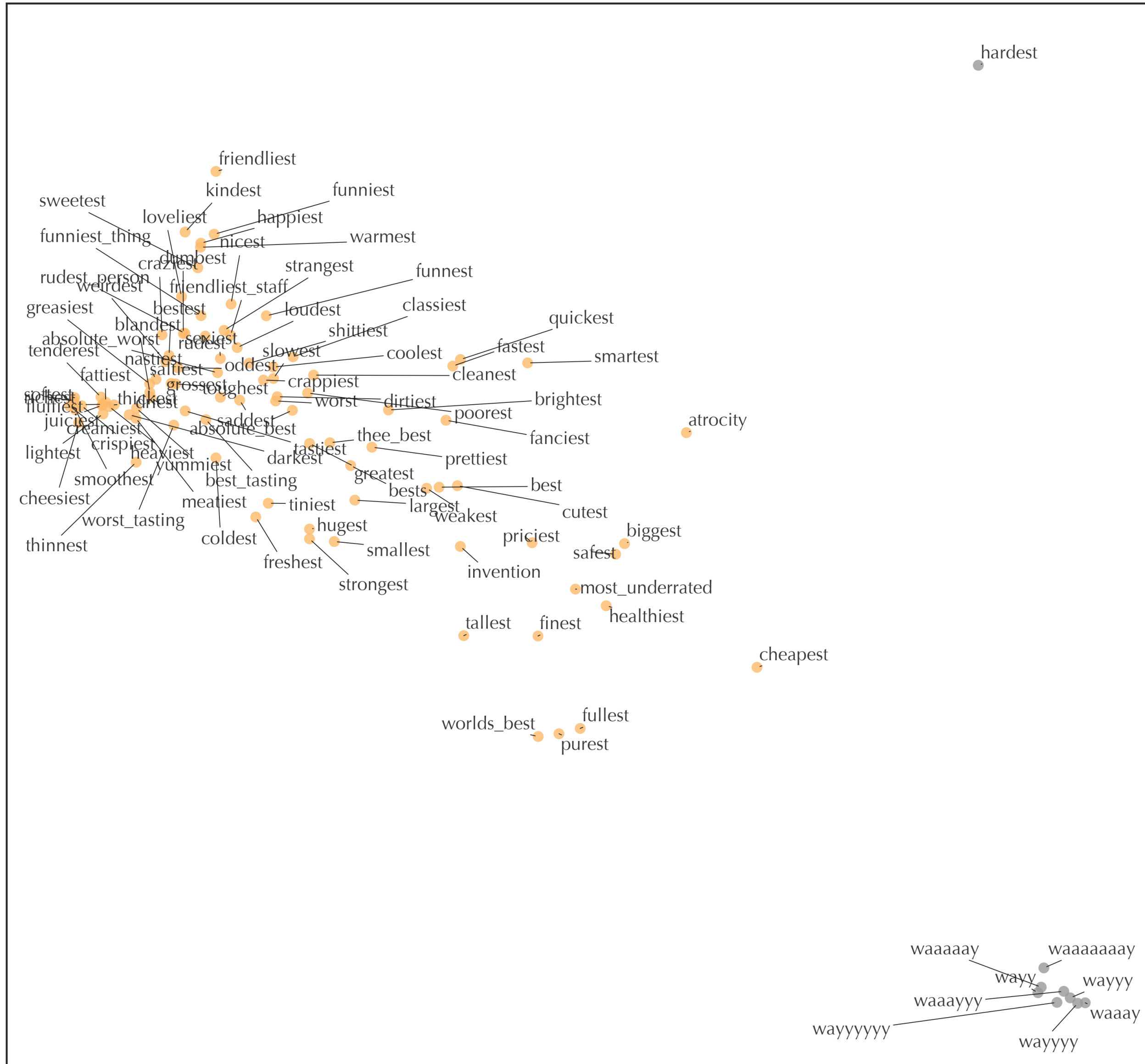


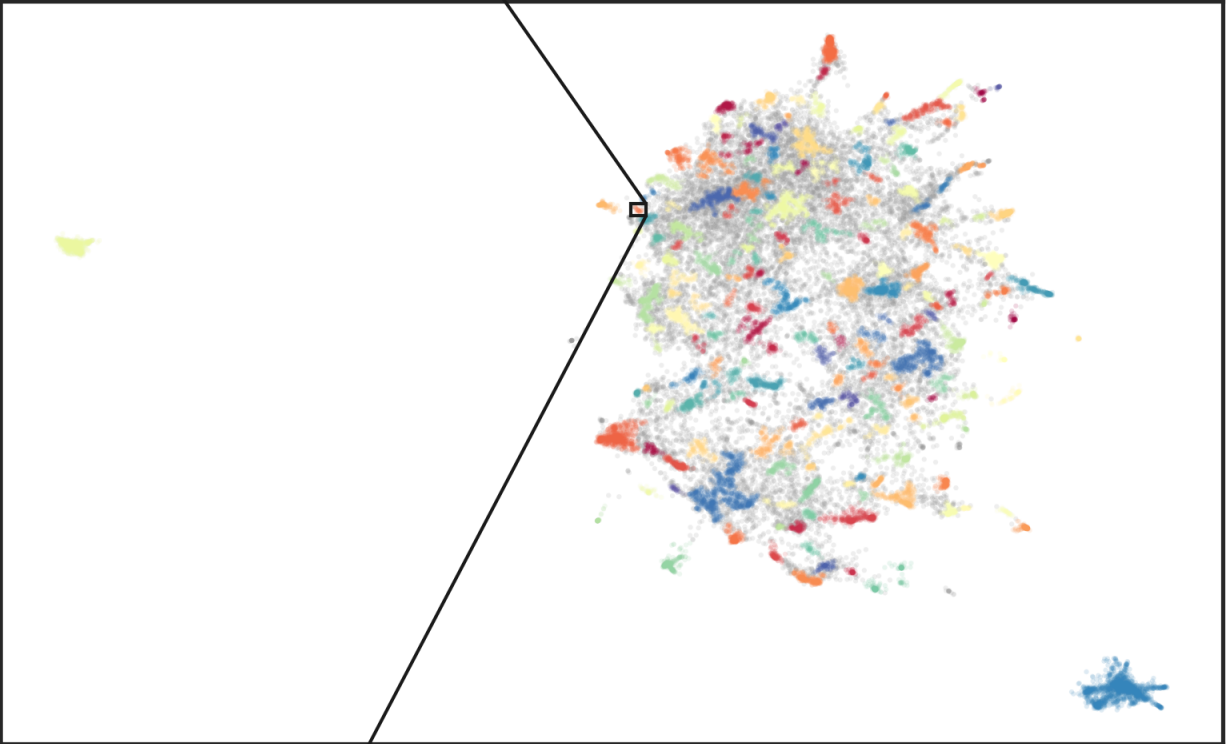
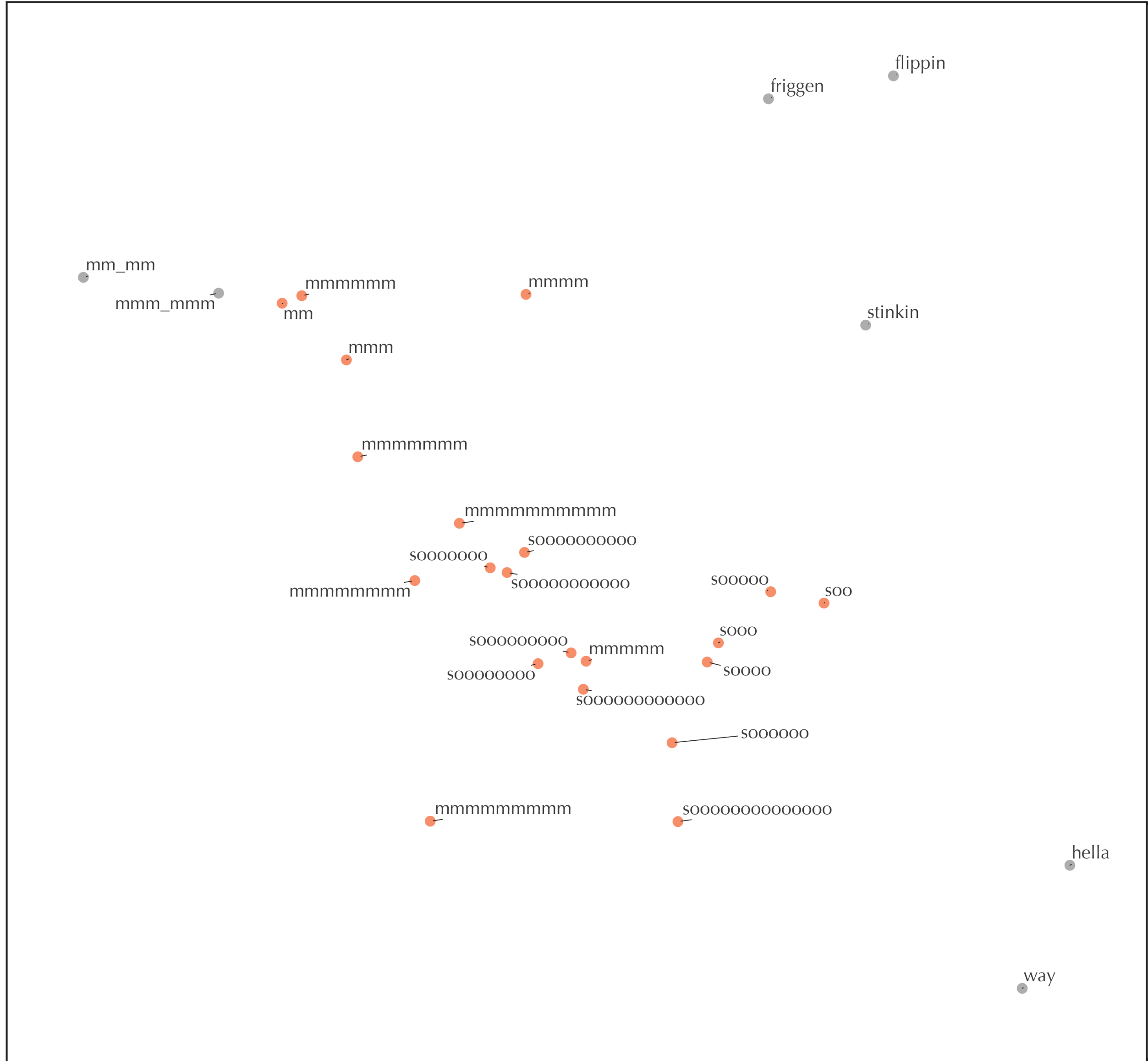




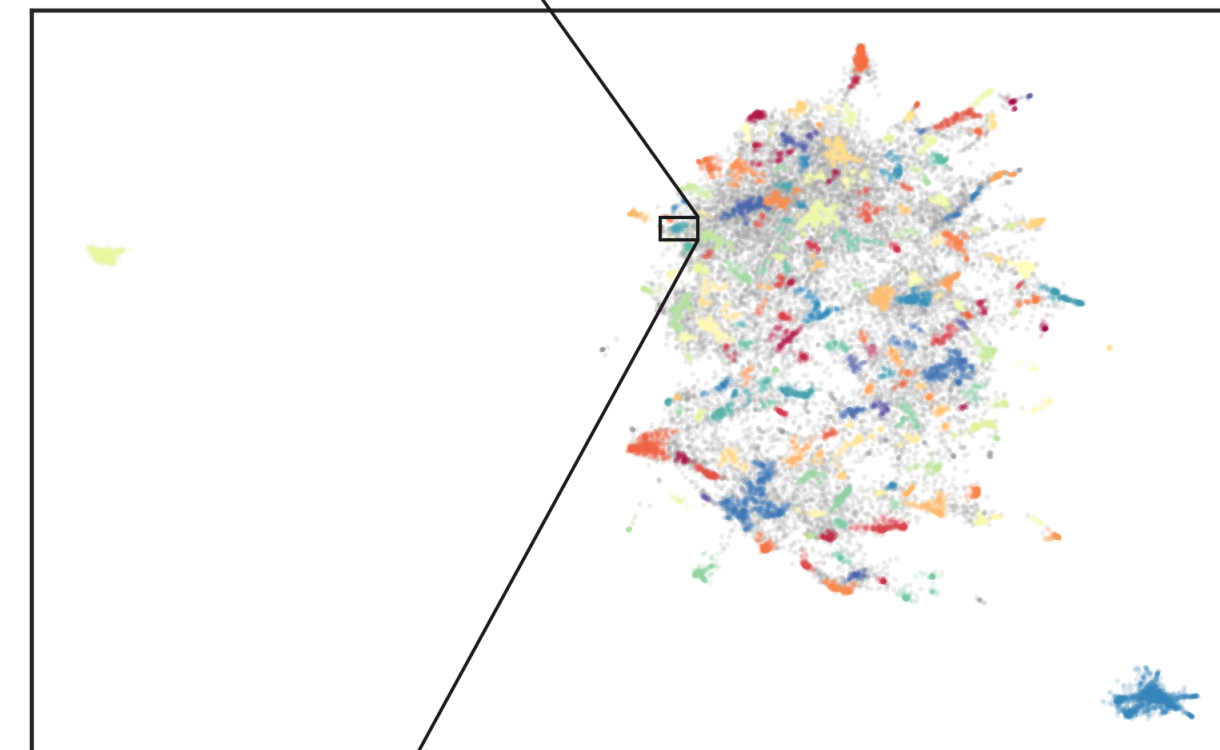
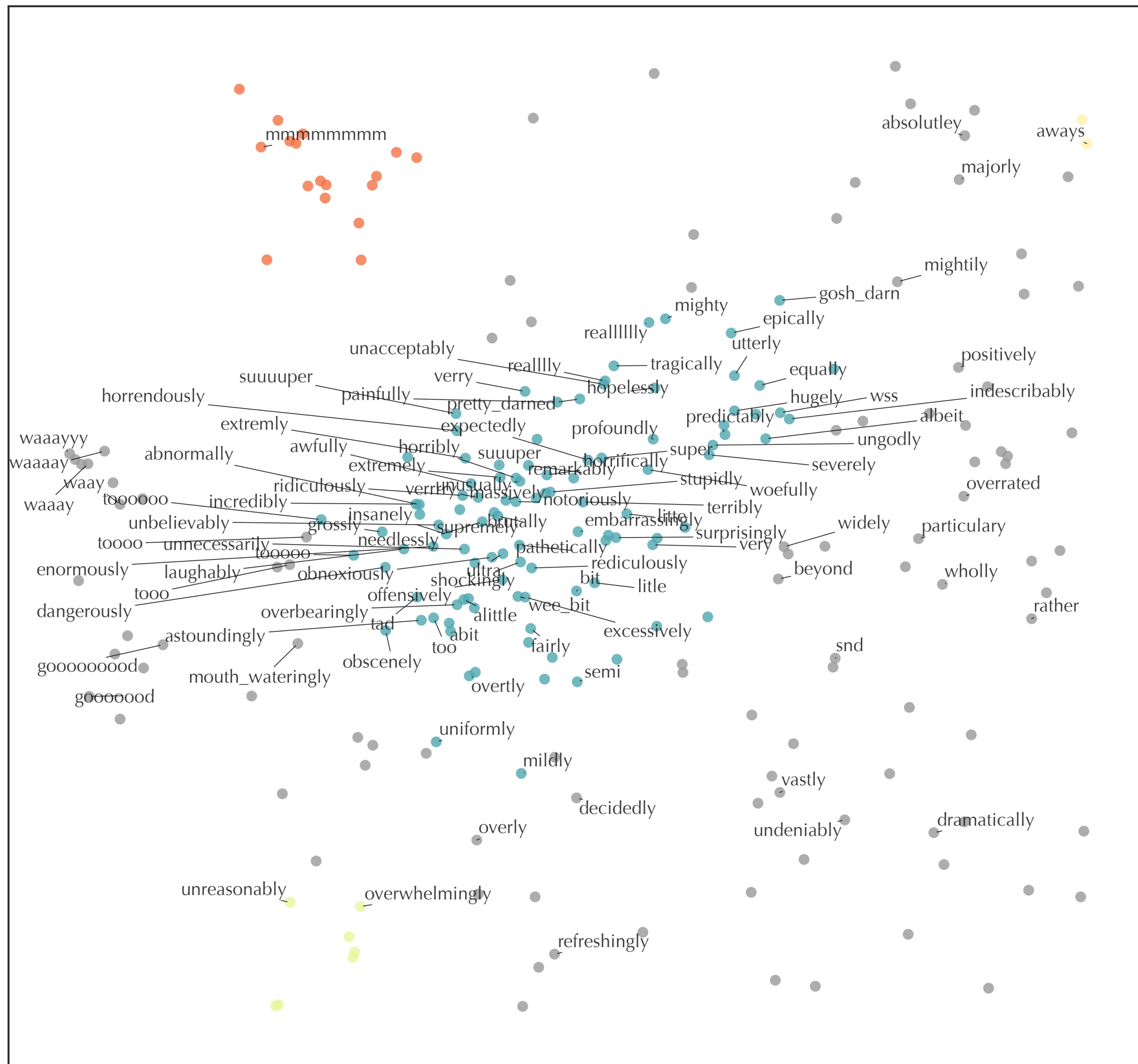


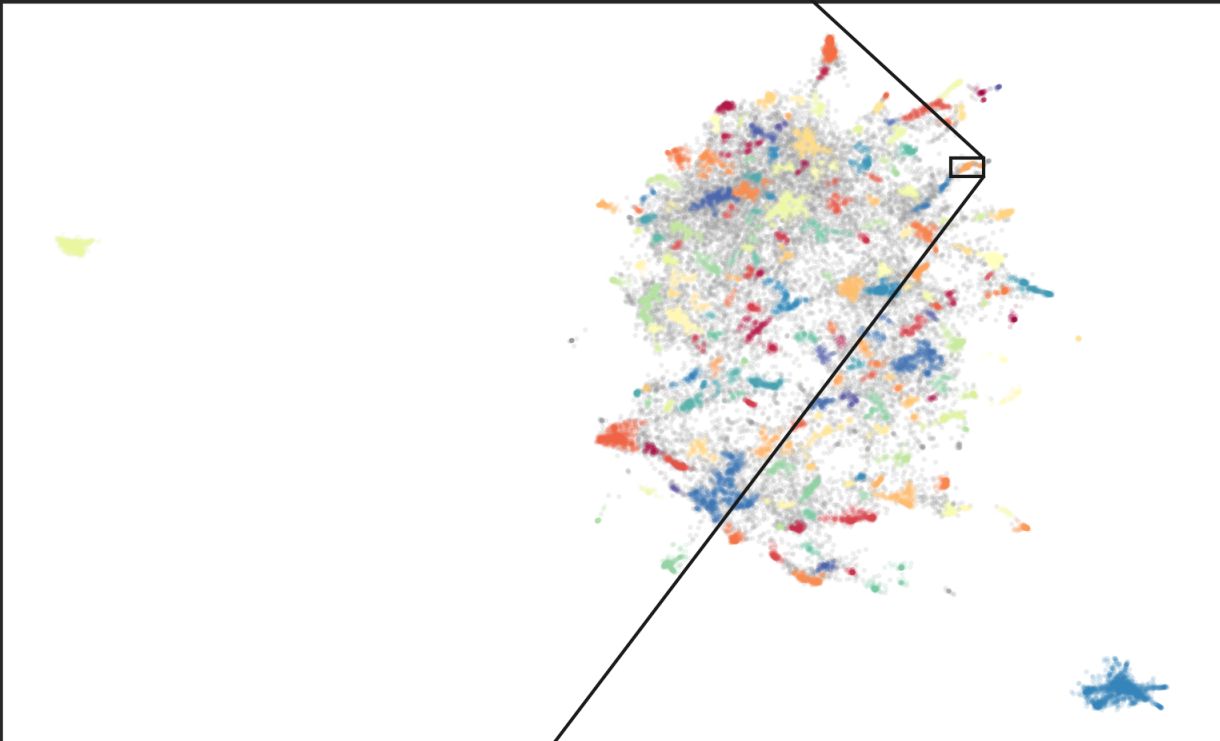
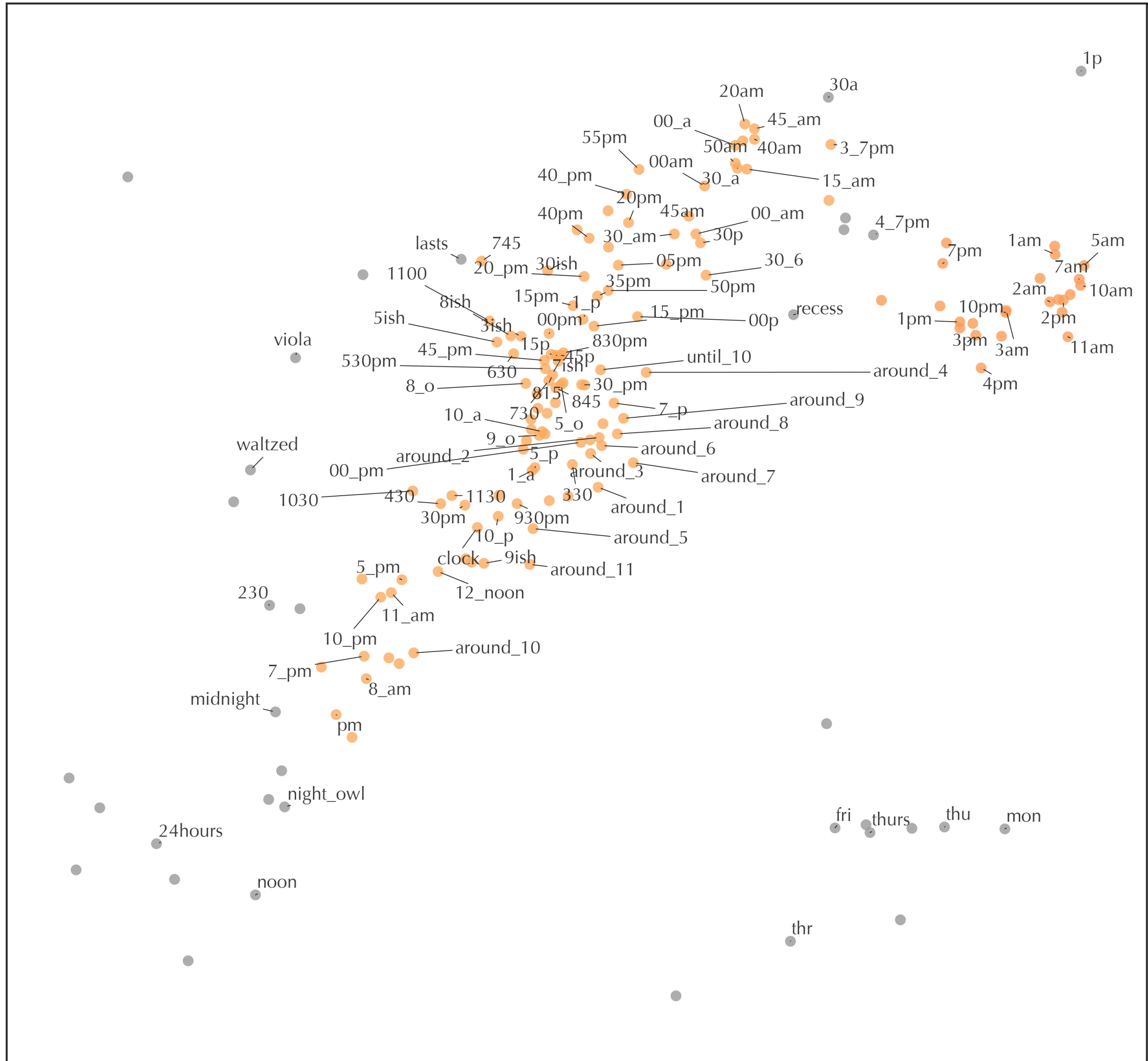


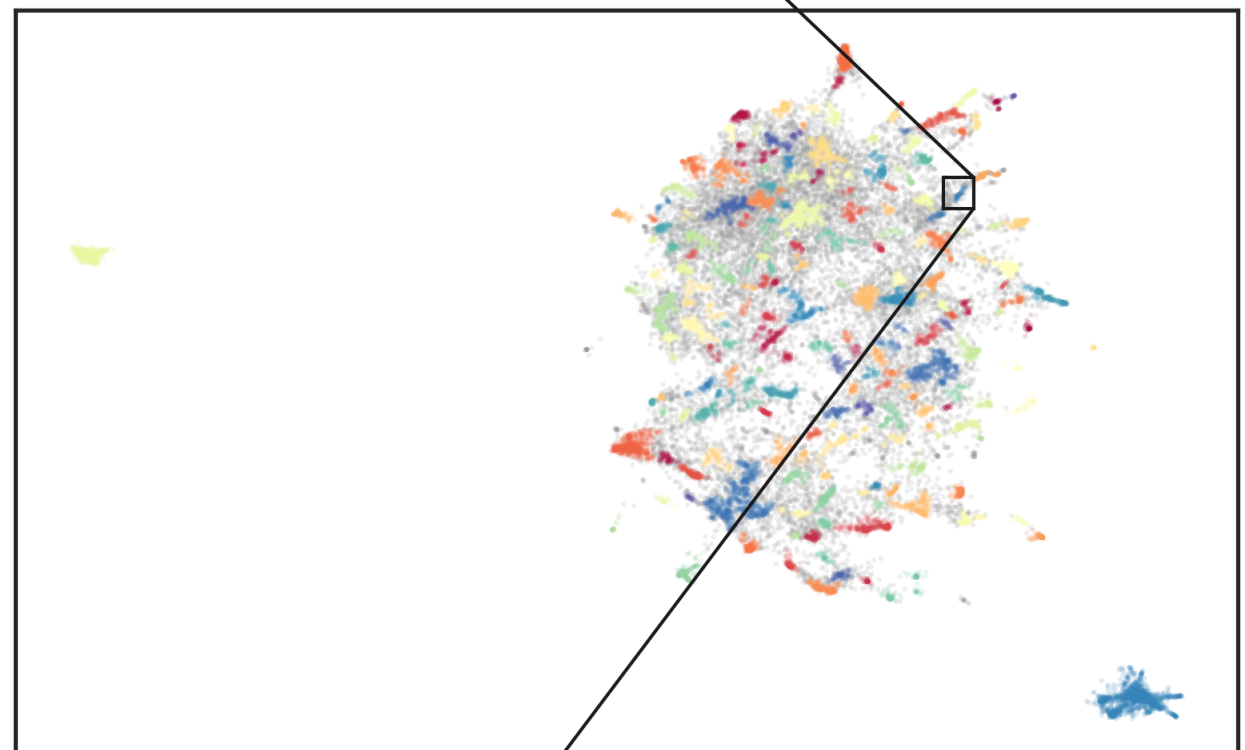
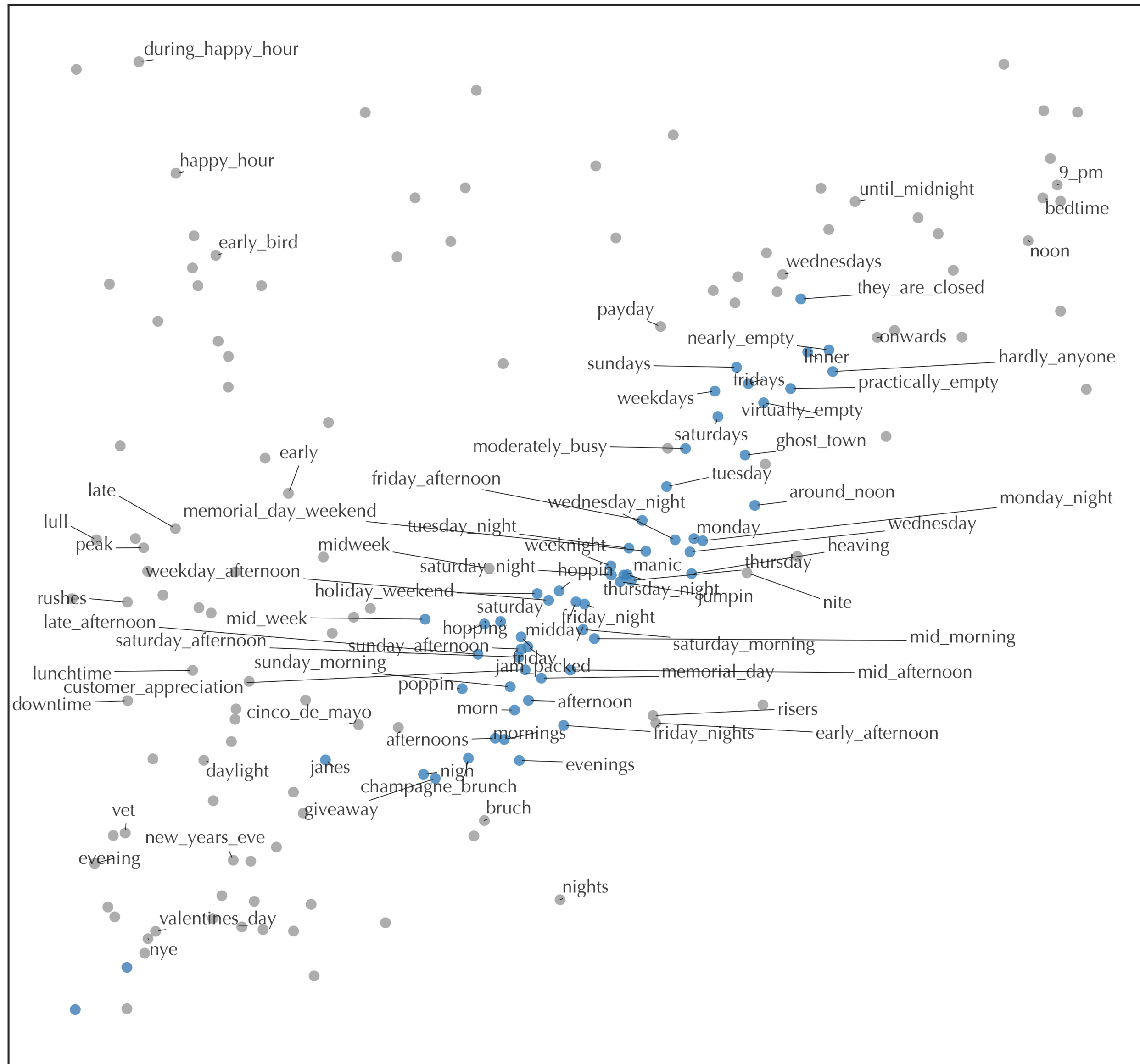




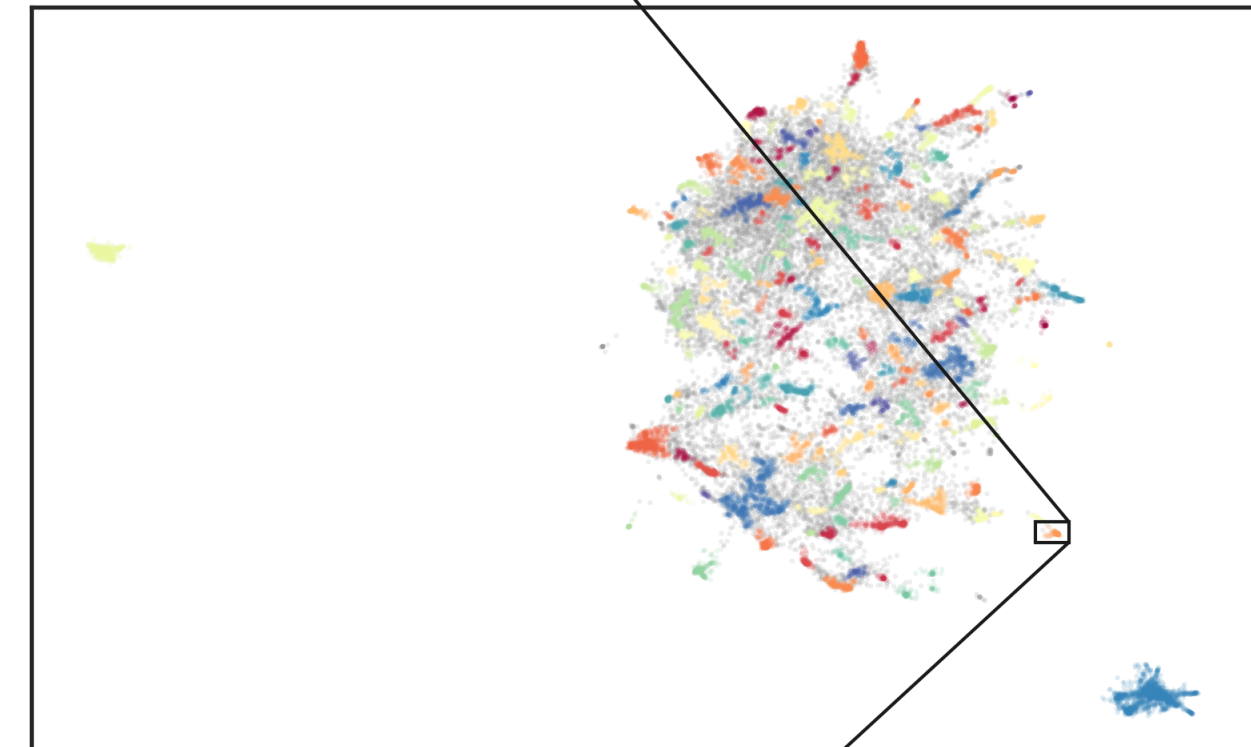
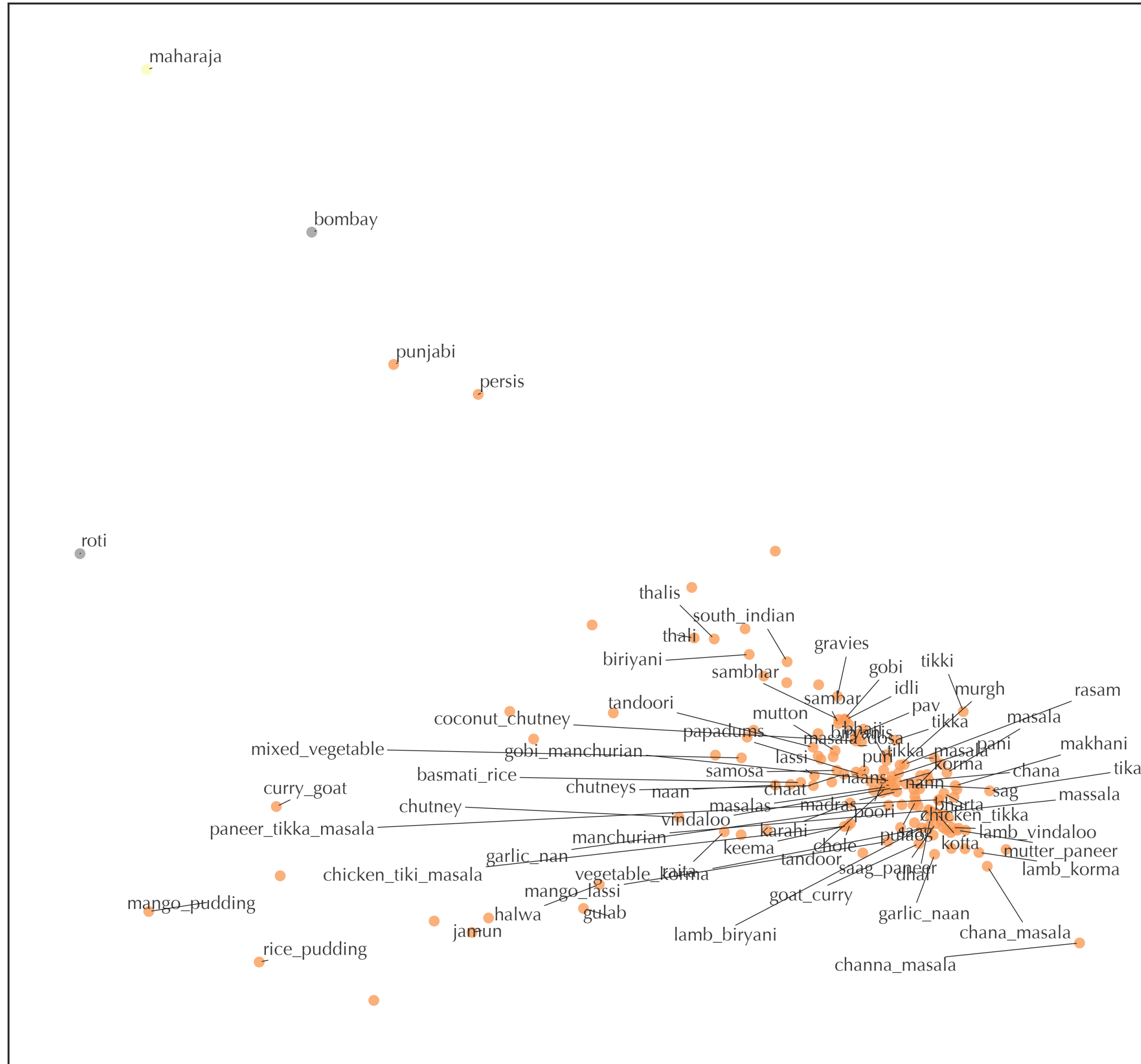




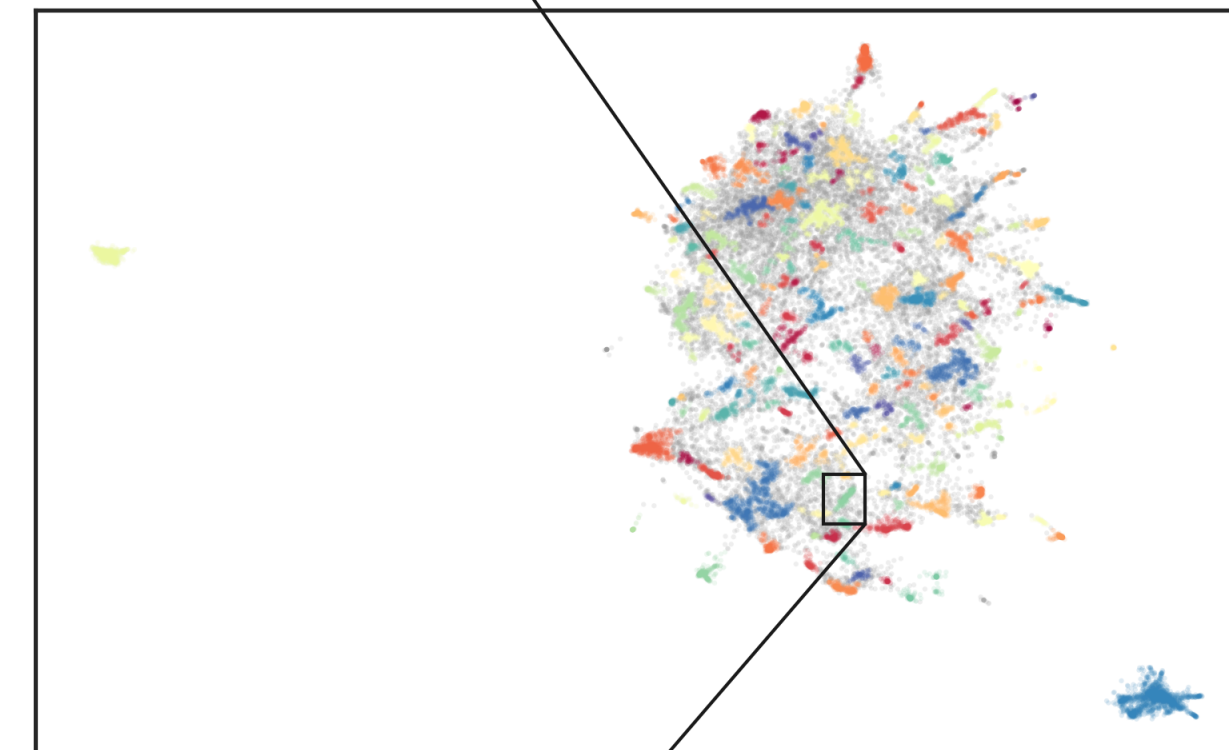
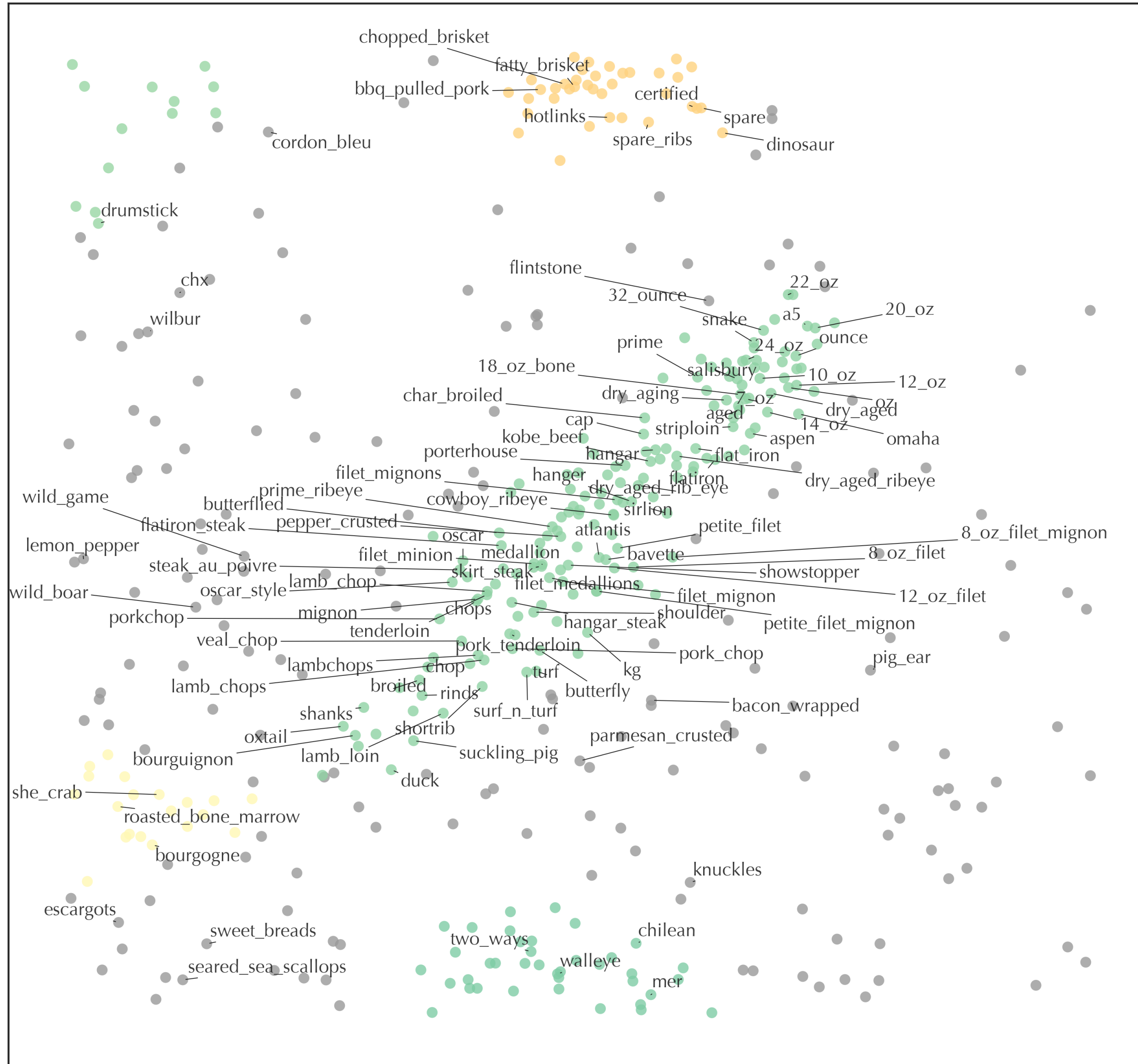








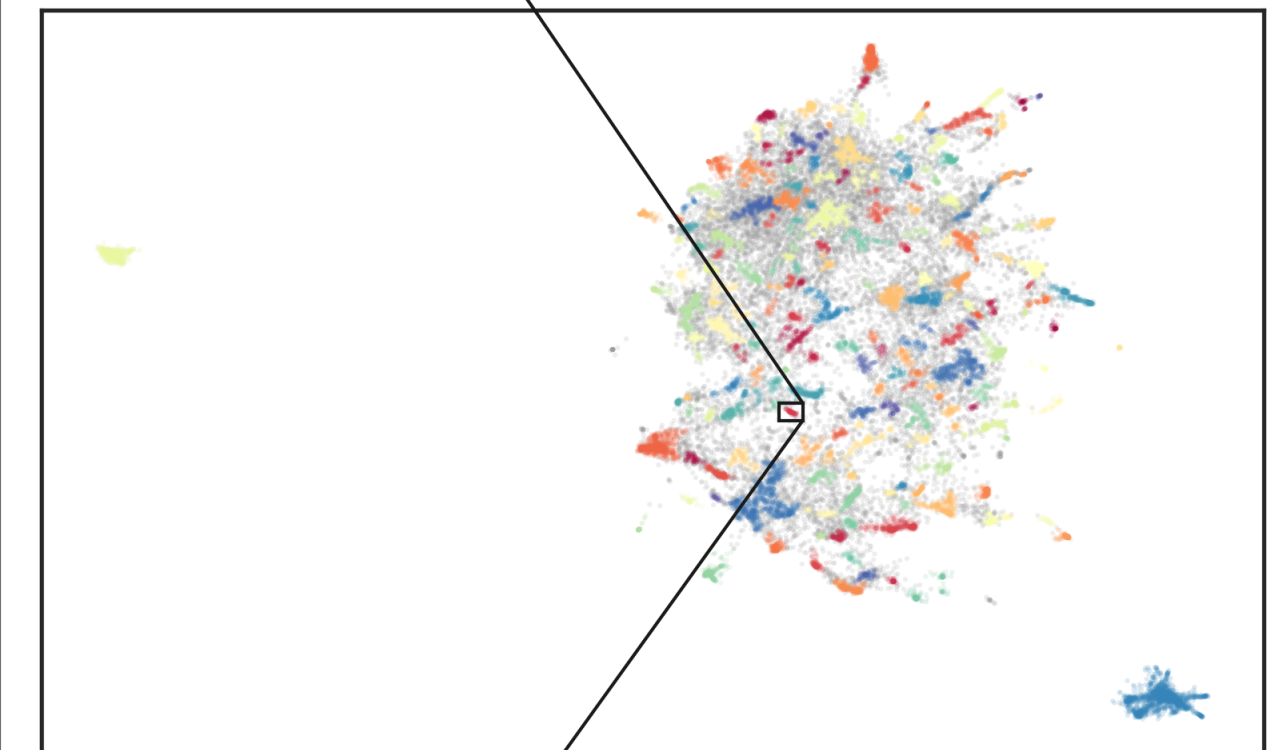
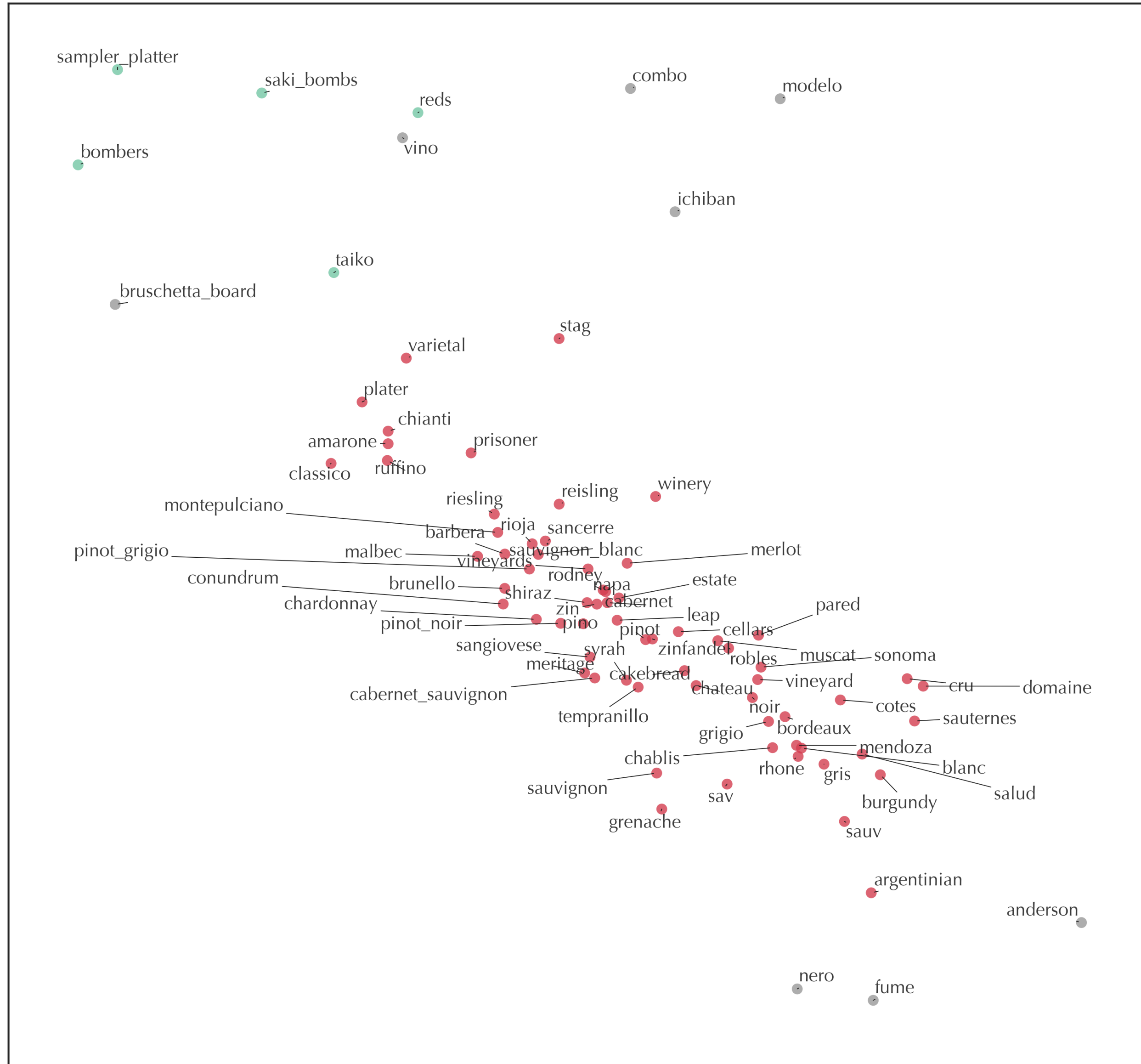














A lot of structure to explore!

What is the persistent  
homology of UMAP's  
simplicial sheaf  
representation?

# Documents as Distributions

Represent a document  
using a “bag-of-words”  
model

Represent a document as  
the multinomial  
distribution of words  
occurring in the document

Information geometry  
tells us that multinomial  
distributions live on a  
manifold



$$d(p, q) = \arccos \left( \sum_i \sqrt{p_i} \sqrt{q_i} \right)$$

## Some problems:

- Word order matters
- Words carry different amounts of information
- Language use is noisy
- Words are not independent / orthogonal

We can't fix word order  
issues this time  
(yet)

Information in words:

Approximate the information  
of a word in a given context  
and weight words accordingly

Language is noisy:

For larger documents we  
can hope to have the noise  
wash out in the background

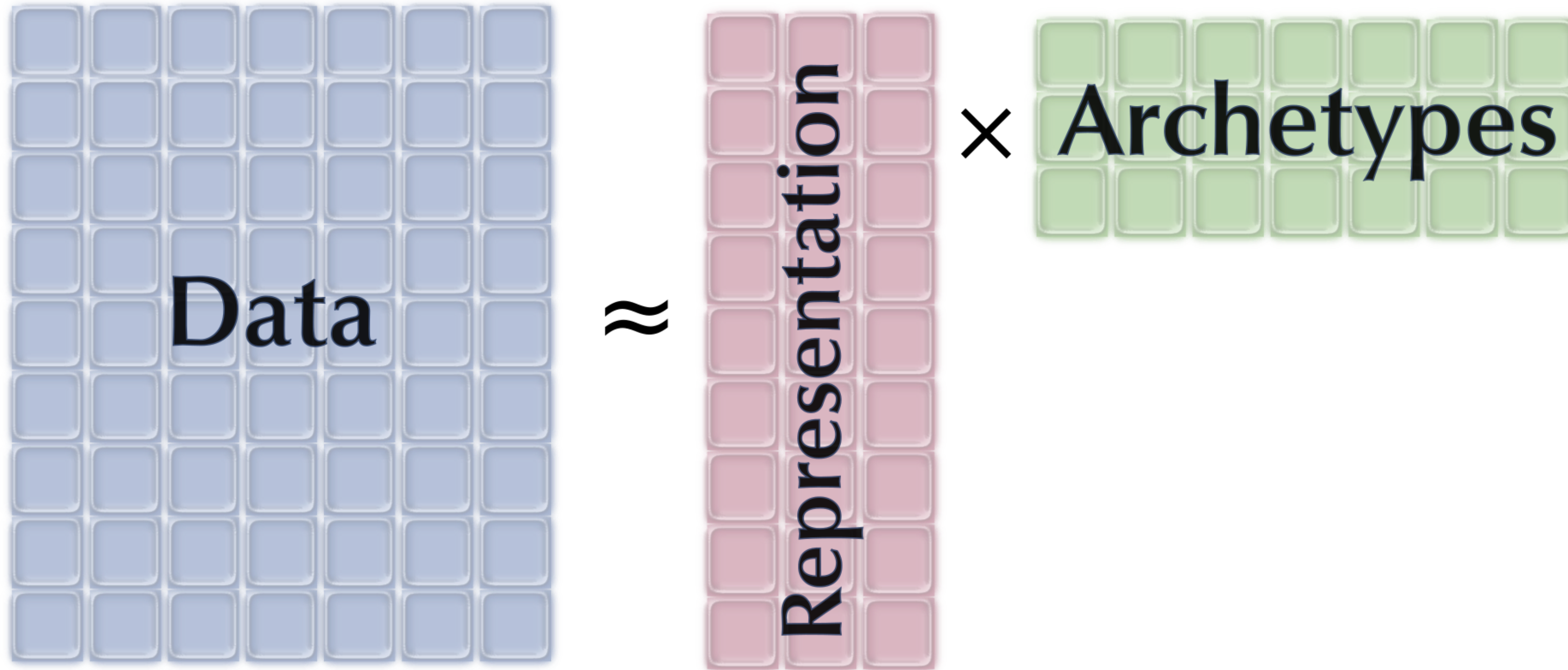
This red rose smells  
sweet

That scarlet flower's  
fragrance is delightful

A document as a  
multinomial is averaging  
one-hot encodings of words

We can perform a change  
of basis to word vectors





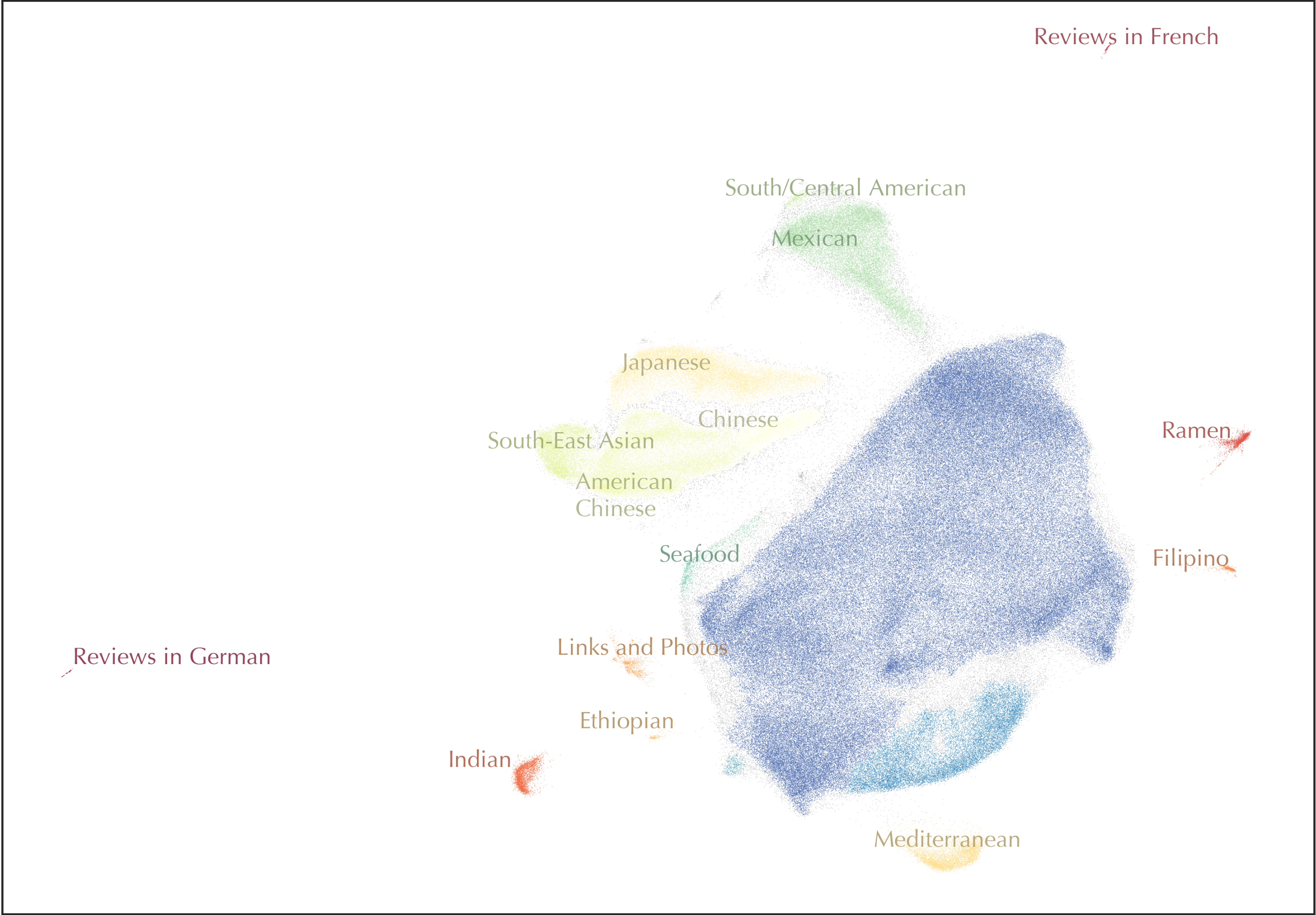
$$D' = D \times W_{pLSA}$$

We are taking averages,  
so we should remove  
noise again ...

# Embedding Documents



UMAP is a topology based  
dimension reduction  
algorithm



Since words and documents now live in the same combined space we can embed both together

This can provide powerful  
tools for topic modelling



# Conclusions

We can use topological  
techniques to explore and  
represent language

Many of these ideas  
provide a mathematical  
basis for variations on  
existing techniques

This remains a rich field  
for further exploration