

# Probing visualizations of neural word embeddings for lexicographic use

Ágoston TÓTH, Esra ABDELZAHER

University of Debrecen, Faculty of Humanities, Department of English Linguistics

E-mail: [toth.agoston@arts.unideb.hu](mailto:toth.agoston@arts.unideb.hu), [esra.abdelzaher@gmail.com](mailto:esra.abdelzaher@gmail.com)



UNIVERSITY of  
DEBRECEN



**T·U·D·A**  
UNIVERSITY OF DEBRECEN  
FACULTY OF HUMANITIES  
SCHOLARLY FUND



# Agenda, aims

- to visualize the distributional features of occurrences of selected headwords in dictionary examples to see if lexicographical sense delineation is reflected in the distributional data
- to check if visualizations of BERT data are useful for assisting manual sense delineation
- to better understand the distributional information stored in BERT representations

# Motivation | the power of Distributional Semantics

Meaning is a function of distribution (Harris, 1954)

In practice,

- the more similar the context, the more similar the meaning
- morphological, semantic, etc. paradigms can be reproduced using vector arithmetics

## *Phases of Distributional Semantics:*

- 1.** Non-ANN phase (“count methods”)
- 2.** 2013- ANN-based static word embeddings  
e.g. word2vec
- 3.** 2018- contextualized ANN-based token embeddings:  
non-generative: ELMo, **BERT**, etc.  
generative: GPT, T5, etc.

# Methods | Data Collection

Example sentences were collected for 4 words: *full*, *mouth*, *risk* and *sound*:

- all matching examples from the online *Oxford Learner's Dictionaries* (<http://www.oxfordlearnersdictionaries.com>)
- 1000 randomly selected corpus sentences for each word from the *British National Corpus* via <http://www.sketchengine.eu>



# Methods | BERT embeddings

- We produced BERT embeddings for the headword in each example sentence by running the neural network; the neural activations for the target words were extracted and saved for visualization.
- Language Model: the largest pretrained BERT LM from Huggingface, *bert-large-uncased* (<https://huggingface.co/bert-large-uncased>)
- LM size: 336 million pre-trained parameters with 24 layers and 16 attention heads
- word embedding size: 1024 floating point numbers per embedding



# Methods | Dimension reduction 1024D→2D

- t-SNE (van der Maaten & Hinton, 2008) is a non-linear method that constructs a probability distribution over pairs of high-dimensional data points and a similar distribution over pairs of low-dimensional points, and it minimizes the difference between these two distributions using gradient descent in an iterative fashion. t-SNE is considered very effective at preserving the local structure of data at the expense of non-local structure.
- Isomap (Tenenbaum, de Silva & Langford, 2000) uses geodesic distance, which is a path between two points on a surface – rather than along a straight line. A graph is created by connecting neighbouring points and computing the geodesic distance between each pair of points.
- Spectral clustering: the top eigenvectors of the Laplacian matrix are considered to capture the global structure of the data.
- MDS creates a low-dimensional representation by minimizing the difference between distances of data point pairs in the high-dimensional space and pairwise distances in the low-dimensional space.

# Methods | Observation of 1024D & 2D diagrams

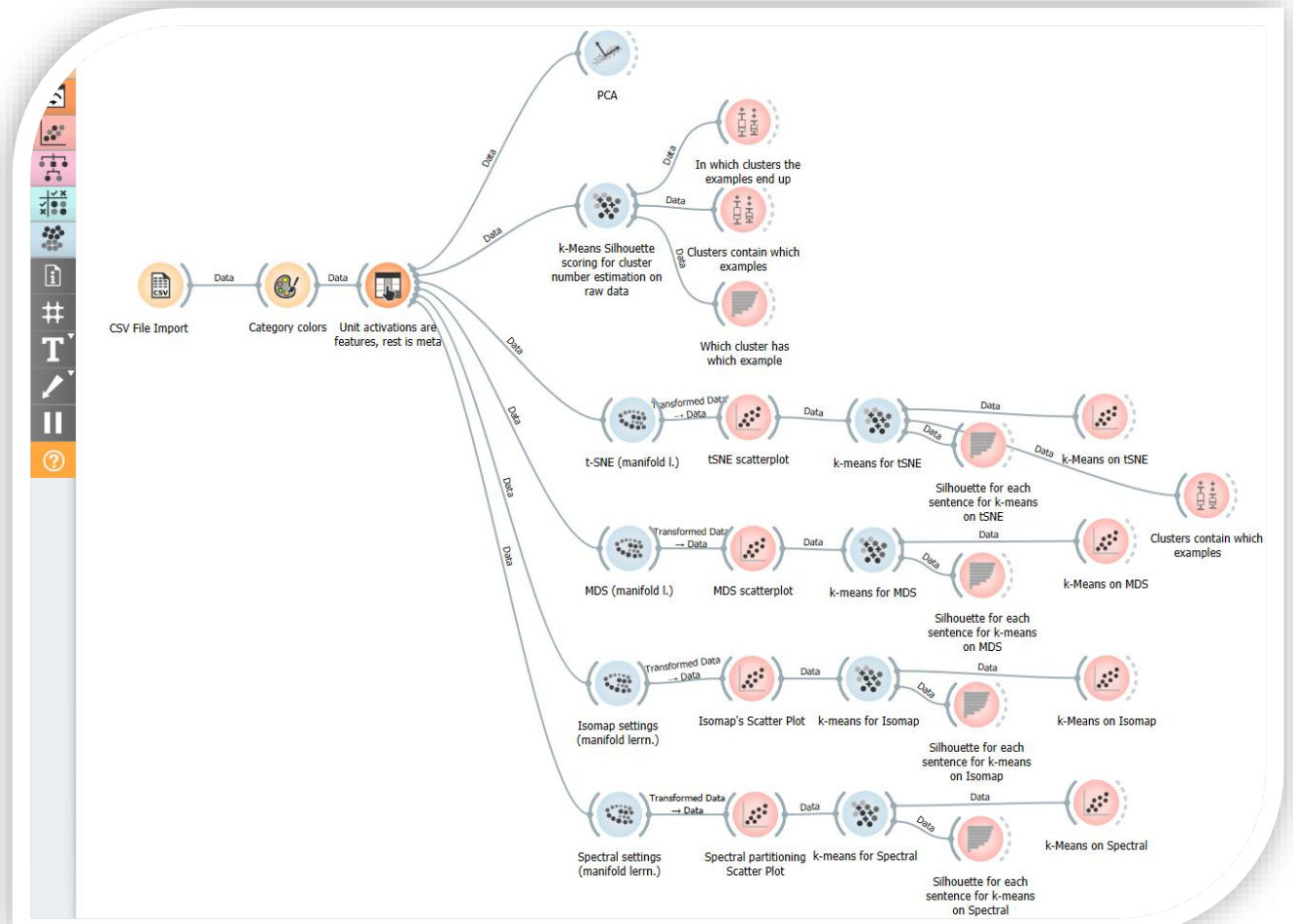
- visual observation, looking for clusters, patterns
- $k$ -means clustering (5000 iterations, 20 reruns)

For selecting  $k$ : Silhouette scoring (Rousseeuw, 1987); a measure of how well data points fit into their clusters, and it “shows which objects lie well within their cluster, and which ones are merely somewhere in between clusters” (ibid.). A higher score indicates better clustering.



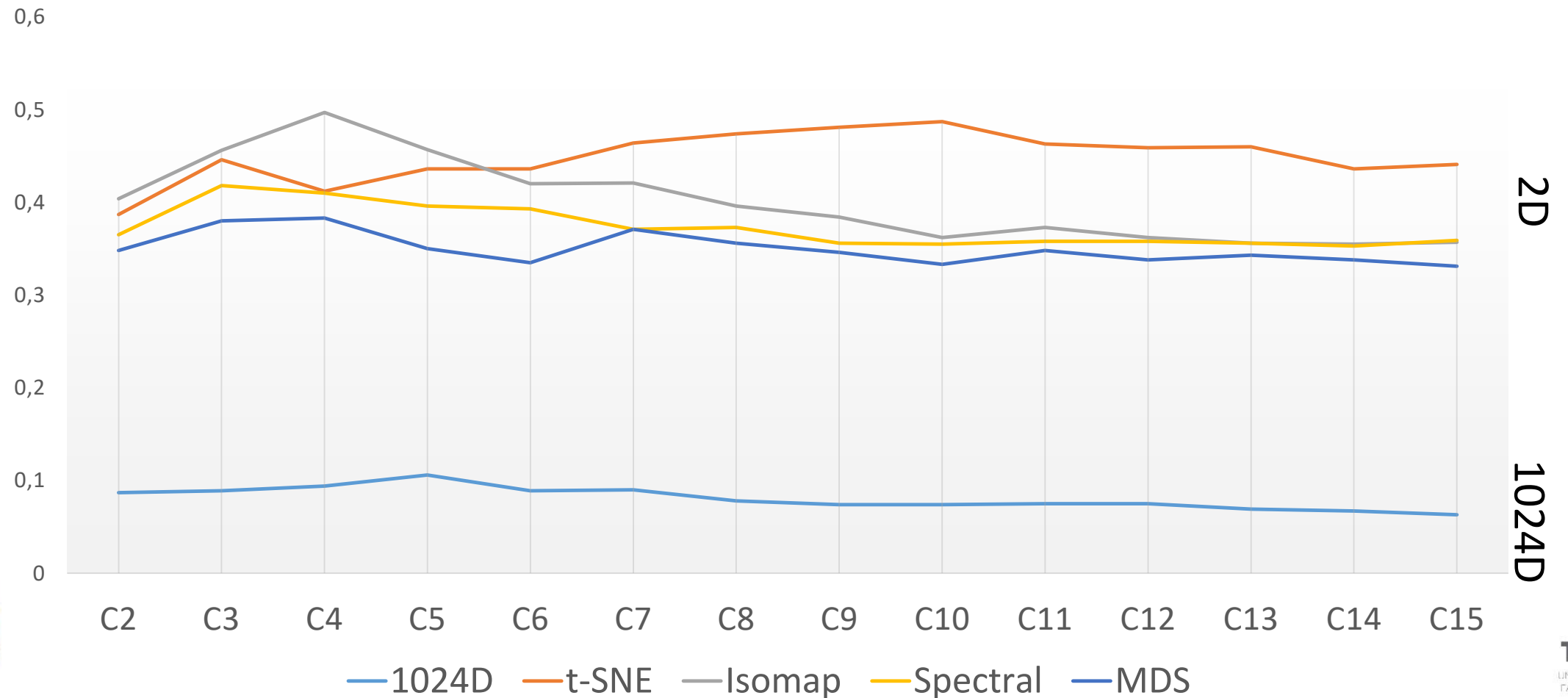
# Methods | The software tool we used

The software used for dimension reduction and  $k$ -means clustering, also the source of our illustrations: *Orange Data Mining toolkit* (Demsar et al., 2013; <https://orangedatamining.com> )

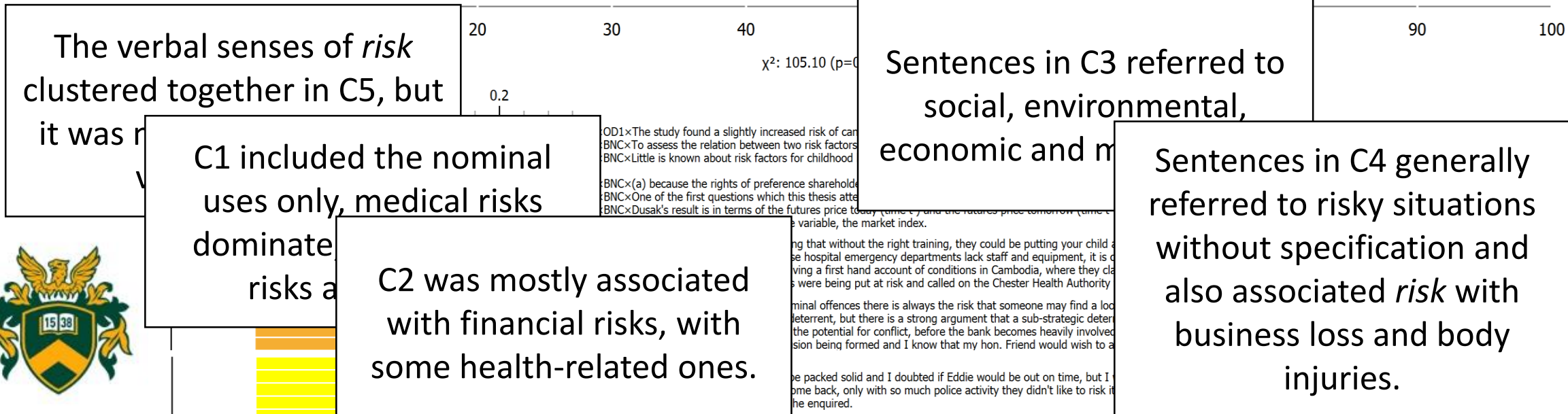
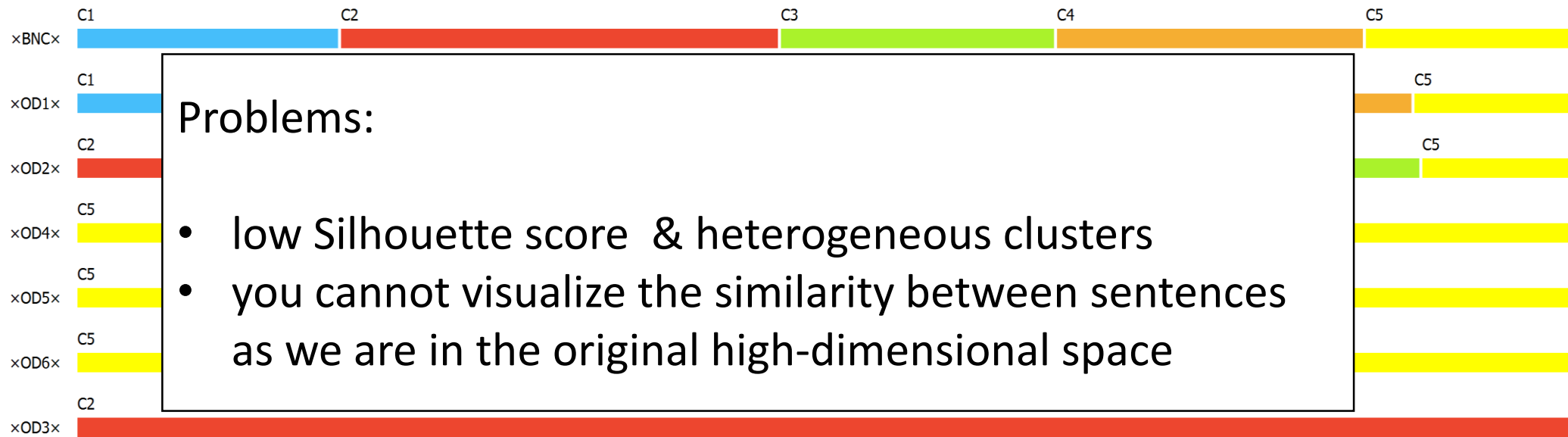


# Results | Silhouette scores for clustering *risk*

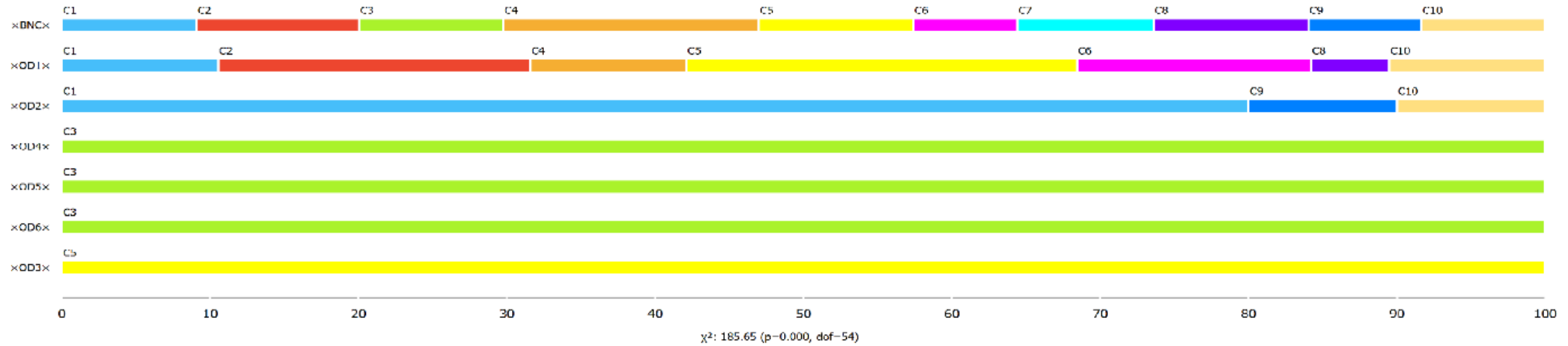
Silhouette scores and k-means clusters for risk's example sentences



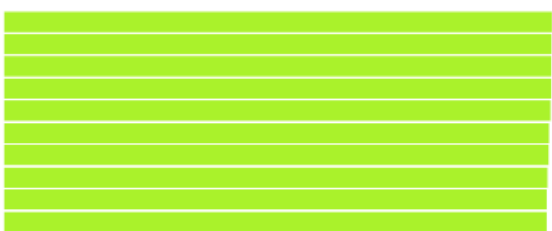
# Results | *risk* in BERT's original **1024D** vector space



# Results | *k*-means clusters after t-SNE, *k*=10, *risk*



0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8



×BNC×'Generally that I had to be convinced any person posed a risk - and to give a warning before I fired.'  
 ×BNC×Is the Minister satisfied that those humanitarian needs have been met and are being met, or that they can be met so long as there is a risk to the Kurdish popu  
 ×BNC×Because of the increased risk due to the speeds of which these vessels are capable the clause restricts the cover granted substantially.  
 ×BNC×It is the risk to public order inherent in the defendant's words or conduct that represents the harm struck at by the section.  
 ×BNC×"Stone & Dobinson's test of"" obvious risk to health or welfare"" would broaden manslaughter perhaps unacceptably, and the position is now that accepted by t  
 ×BNC×These are important areas, but the writers usually concentrate on the technical and physical aspects of securing computer systems against external threats whi  
 ×BNC×Th  
 ×BNC×Ye

C1 included *be risk to NP*  
 C2: *increased/reduced/high/low risk of NP*  
 C4: *risk of -ing, risk+that+clause*  
 C6: *Adj+risk+N*  
 C7: *at risk*  
 Health-related risks: C2, financial risks: C5

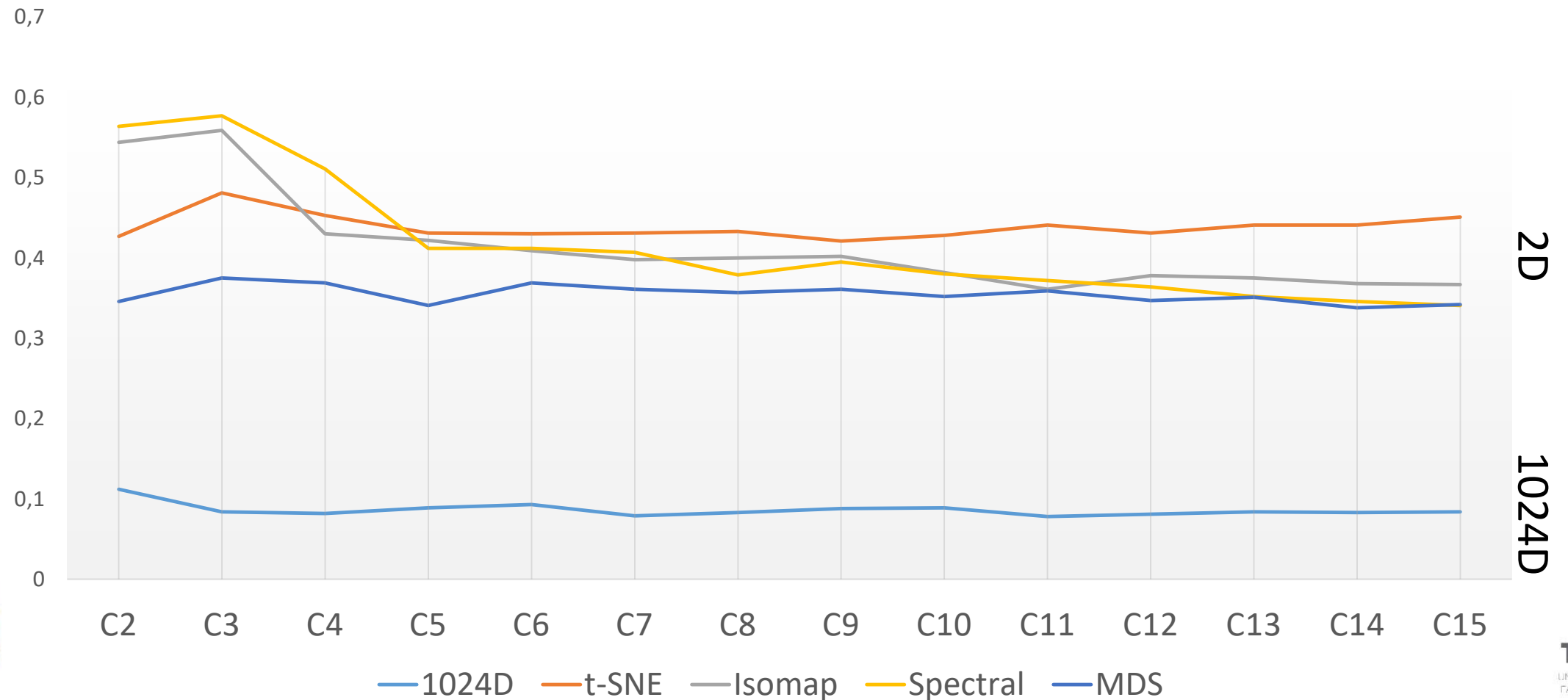
×BNC×Can Bob Halton hear the protest and risk giving up some of his 'I love me'?  
 ×BNC×In the current climate, few executives were prepared to risk airing their anxieties publicly, but privately they express a wide range of fears.

×BNC×Such a system calls for kites with similar stability characteristics, otherwise there is a risk of tangling the two kitelines, as the shared weight draws them togeth  
 ×BNC×There is a real risk of it being regarded not as a mate but a meal.



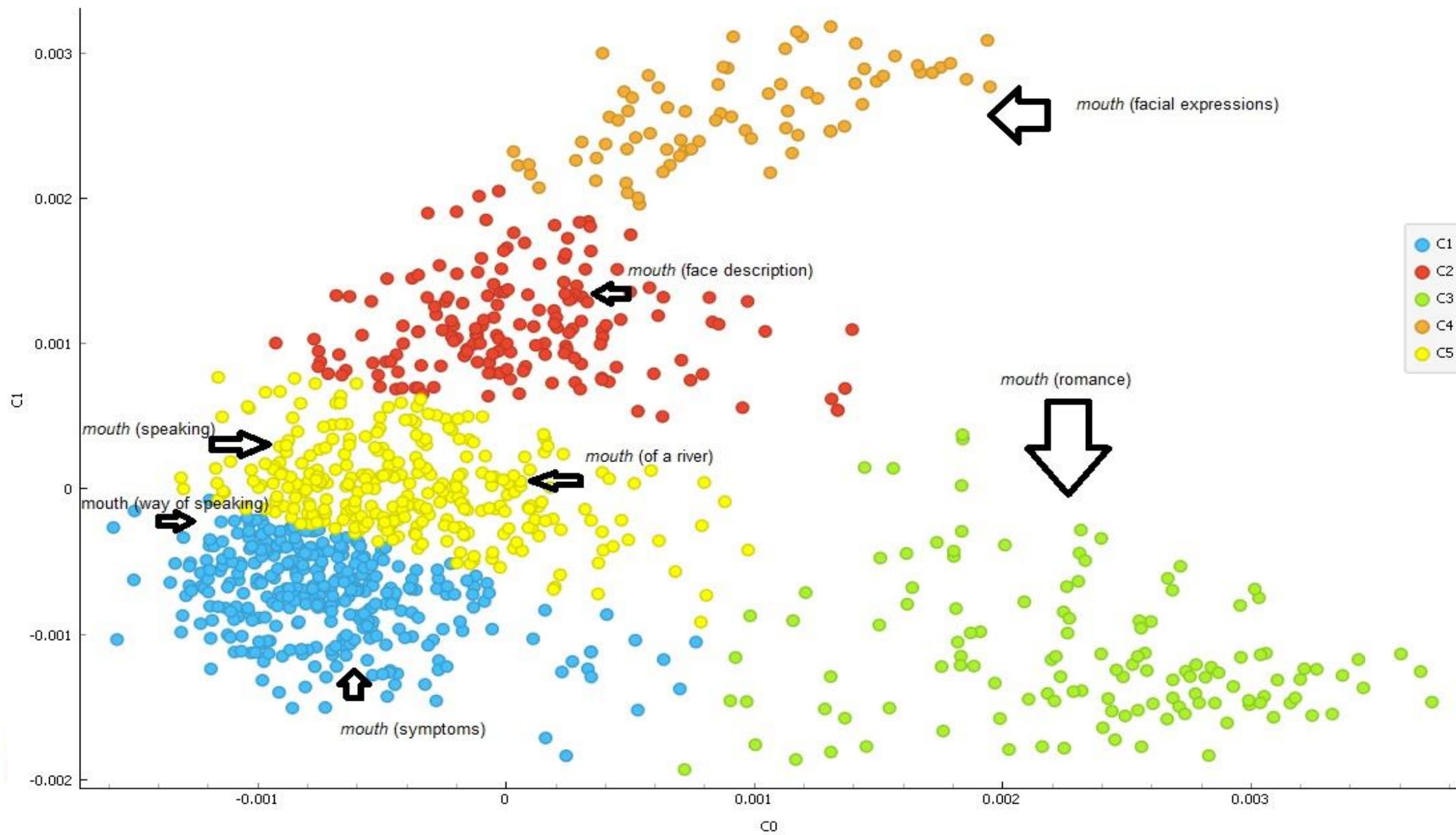
# Results | Silhouette scores for clustering *mouth*

Silhouette scores and k-means clusters for *mouth*'s example sentences

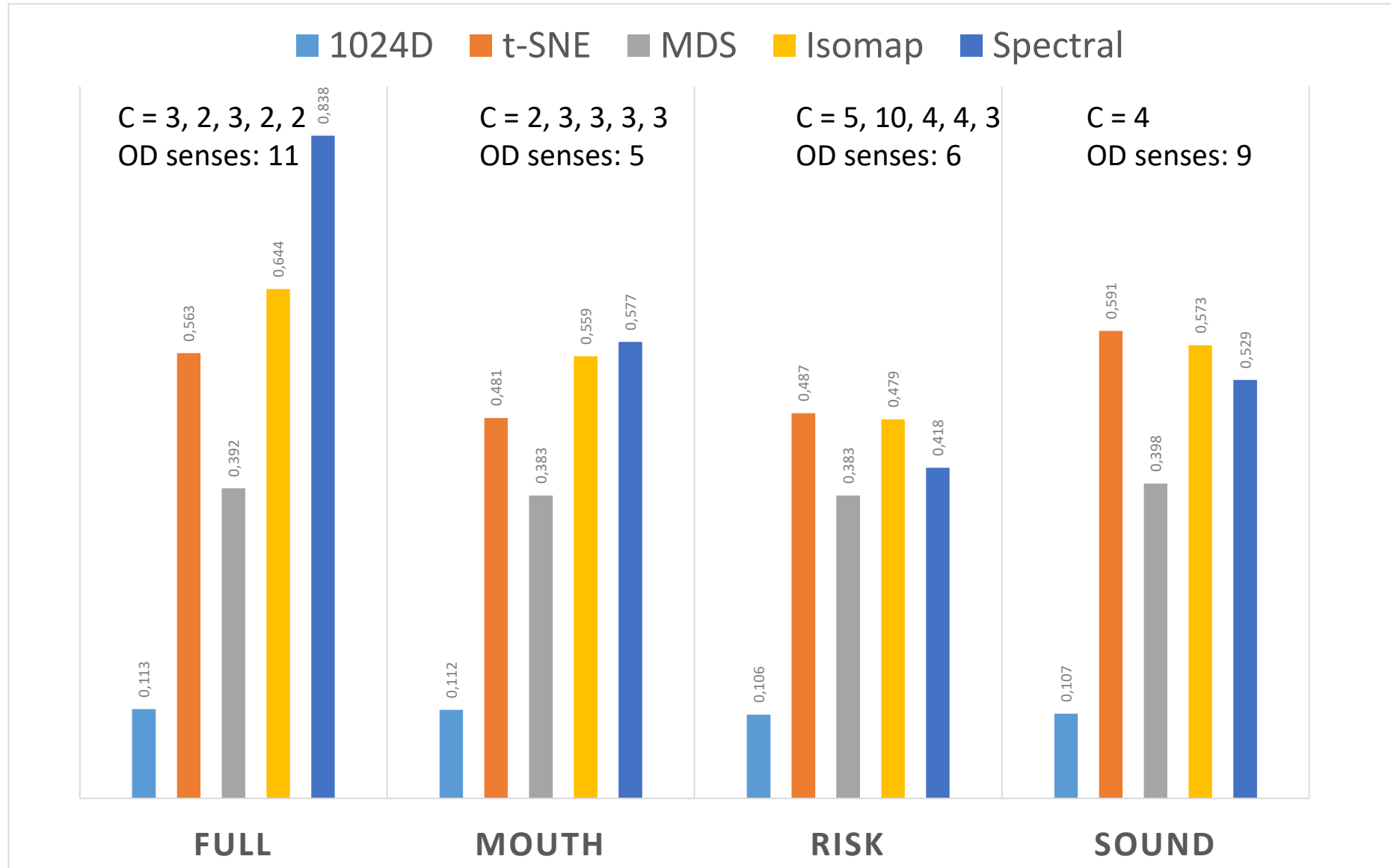




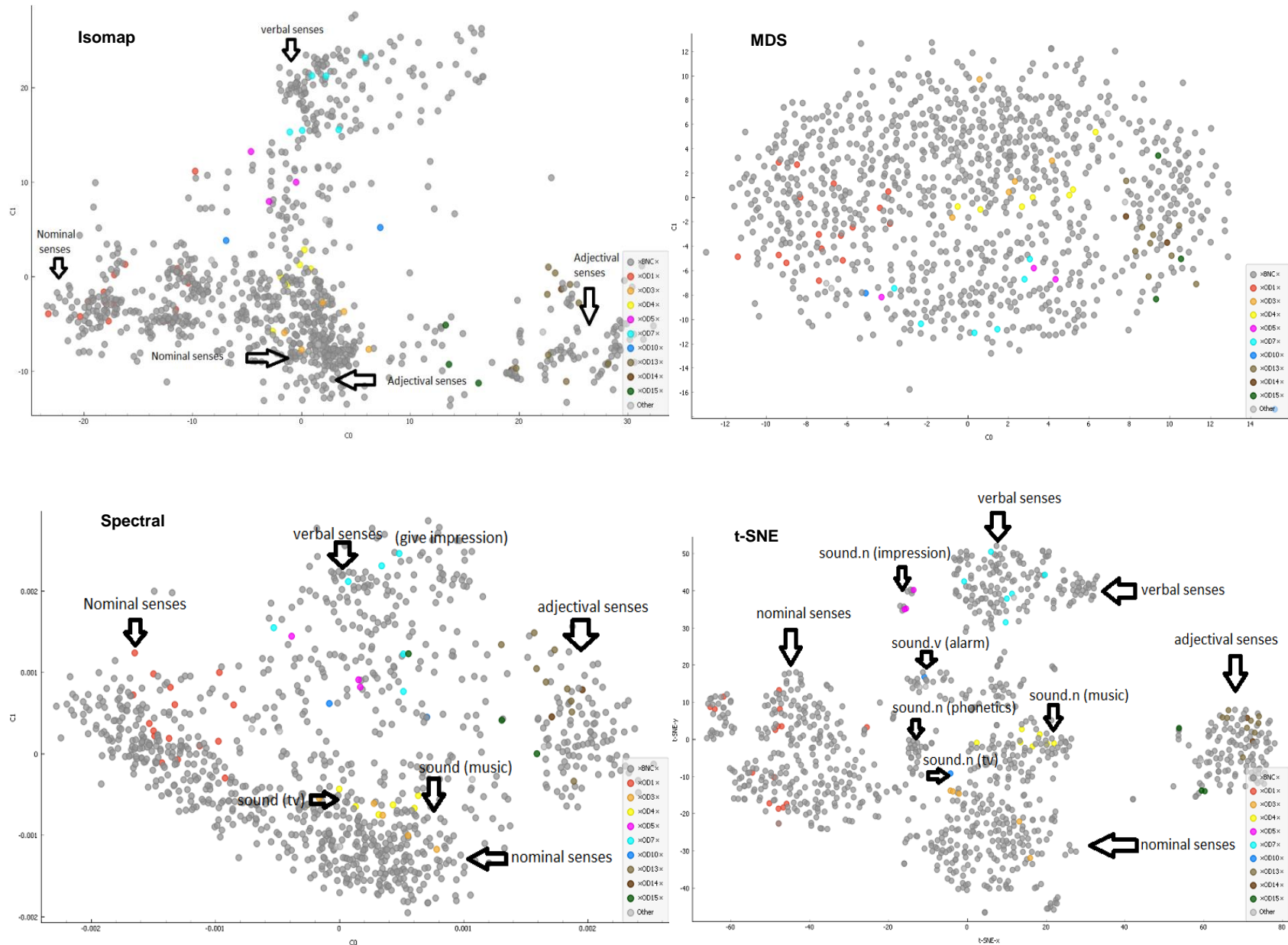
# Results | *mouth* in Spectral visualization & clusters for $k=5$



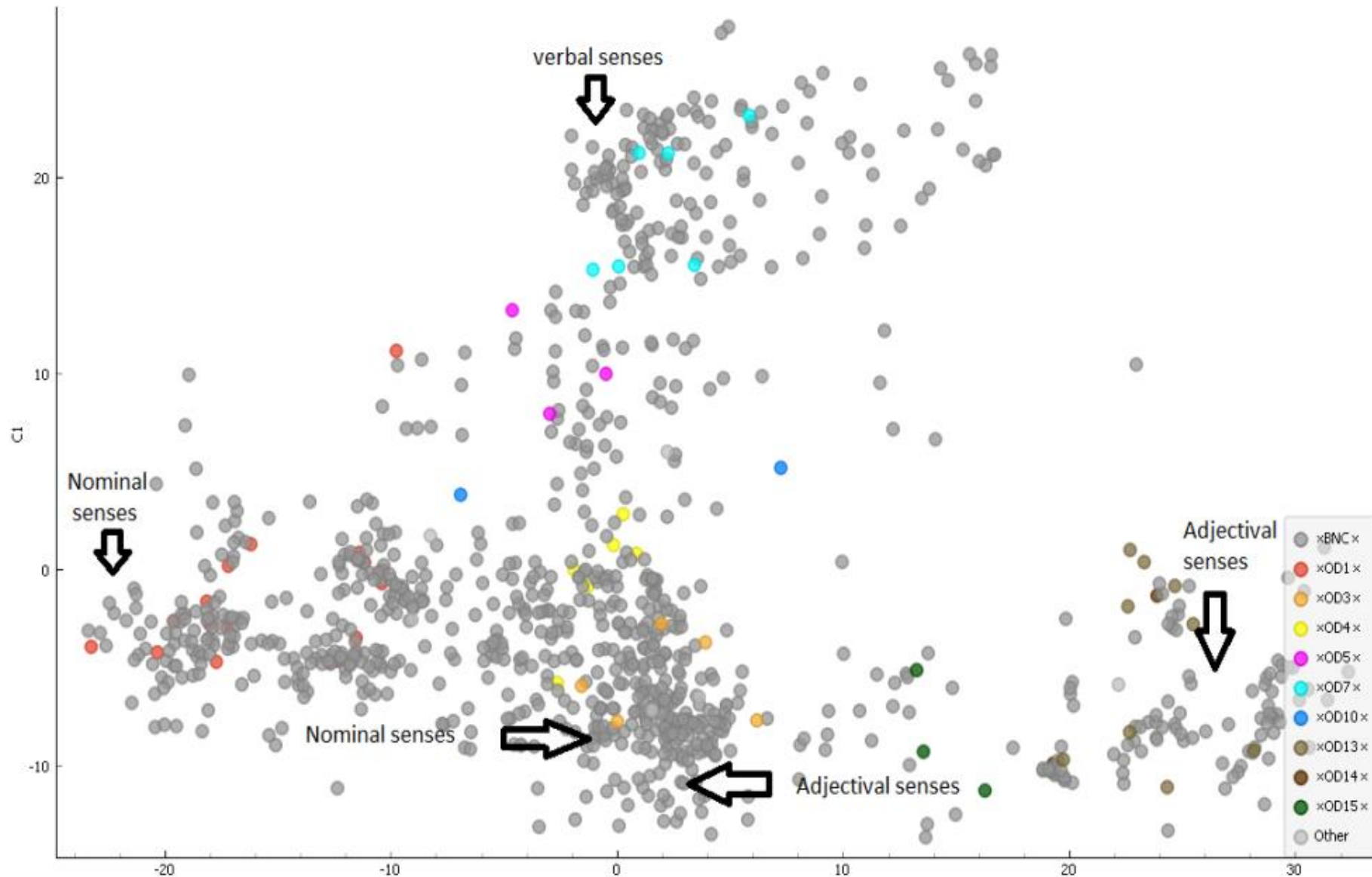
# Results | The highest Silhouette scores for the four words before and after dimension reduction



# Results | *sound* in 4 visualizations (overview)

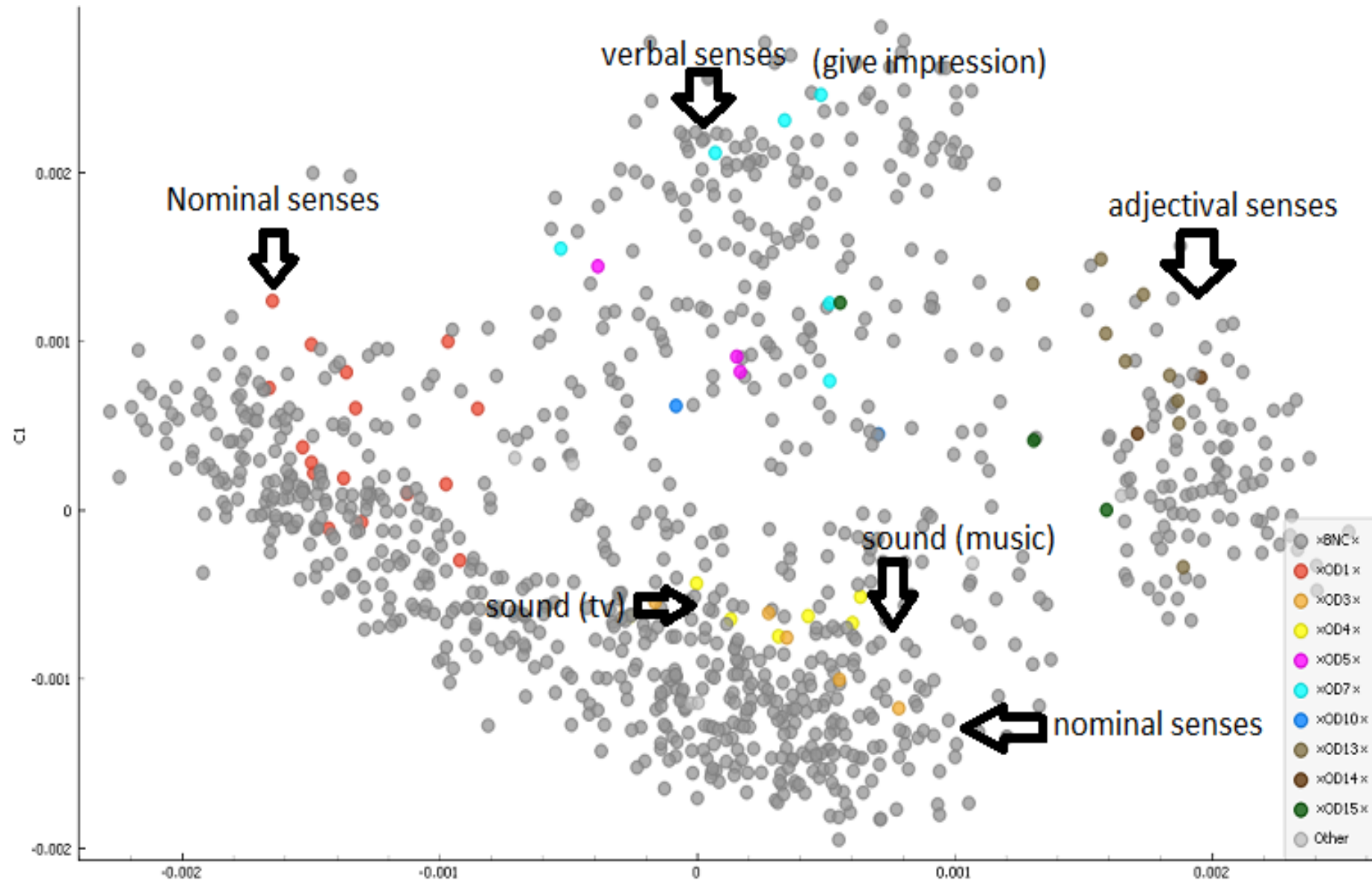


# Results | *sound* in Isomap



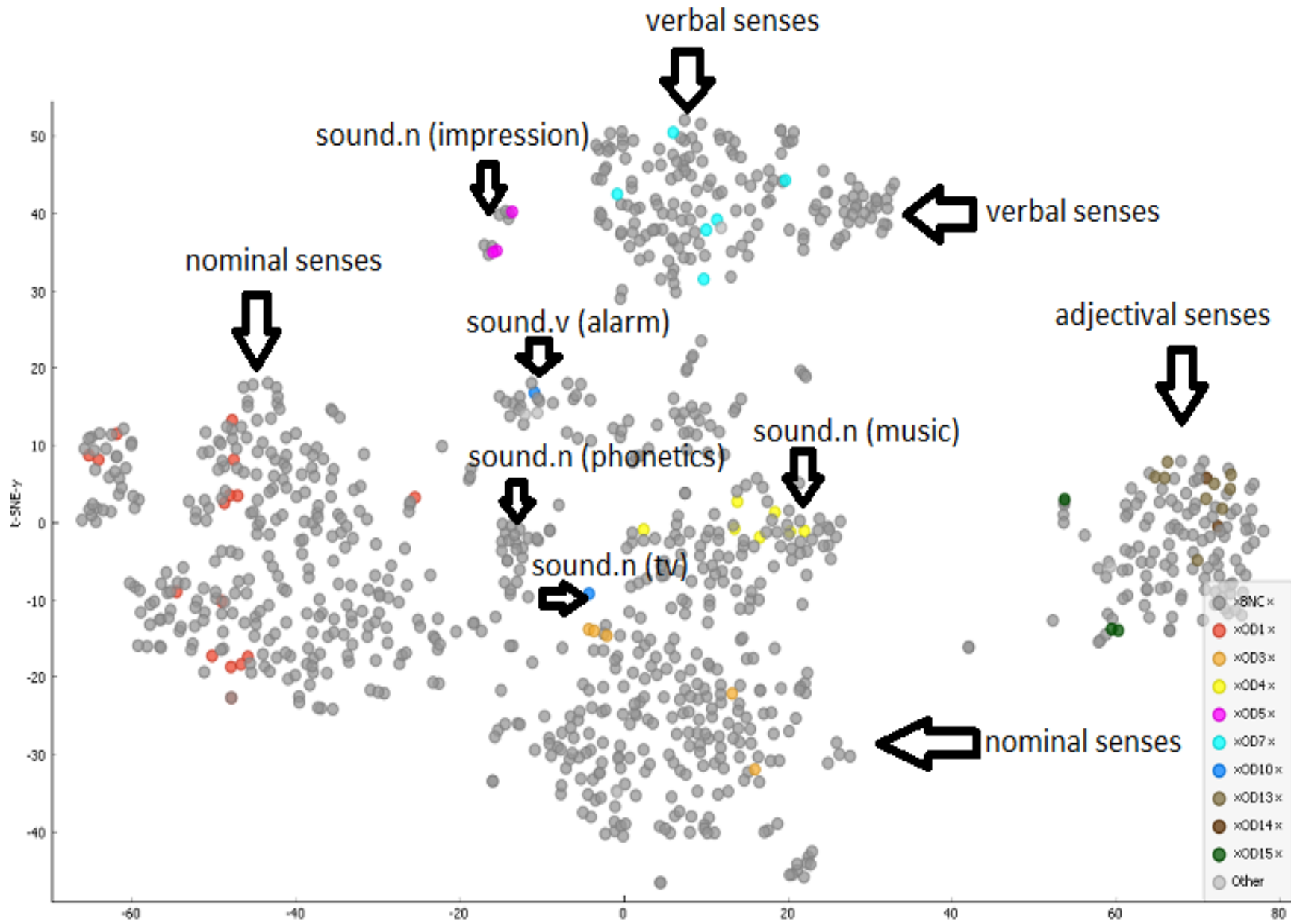


# Results | *sound* in Spectral





# Results | *sound* in t-SNE



# Parameter choices for the dimension reduction methods

Dimension reduction	Settings
<b>t-SNE</b>	perplexity = 20 (also tested: 10, 30) distance = Euclidian (also tested: Manhattan, Chebychev) initialization = PCA max. iterations = 3000 learning rate = 200
<b>MDS</b>	initialization = PCA max. iterations = 5000
<b>Isomap</b>	neighbours = 20
<b>Spectral</b>	affinity = RBF kernel (also tested: Nearest neighbour)



# Discussion

- In our experiments, unsupervised separation between the metaphoric, metonymic and literal senses of words such as *mouth* and *sound*, based on the distributional features of the word uses, is reasonably good.
- The uses of words with relevance to specific semantic fields (e.g., *risk* in financial domains, *mouth* to make facial expressions, *full* with relevance to emotions) stood out in the automatically generated clusters.
- In almost all cases, Silhouette scoring for 2D representations recommended fewer categories than the number of Oxford Dictionary sense categories. Some dictionary distinctions were preserved within the sub-clusters (e.g., *sound* of music vs. *sound* of TV and radio), but others were lost (e.g. the four verbal senses of *risk*).

# Conclusion

- The BERT-based, distributionally-motivated clusters did not correspond to the number of dictionary senses, but they did show BERT's sensitivity to semantic and syntactic similarities between word uses.
- Before dimension reduction, Silhouette scores of the  $k$ -means clusters were low, and so was the qualitative cohesion between the sentences in the cluster.
- Visualizing BERT representations in 2-dimensional spaces using Spectral, t-SNE and Isomap showed quantitative and qualitative improvements that can be beneficial to lexicographers. Not only the Silhouette scores of the clusters increased, but also semantic and syntactic similarities appeared in the clusters.

# Conclusion

- MDS was inferior to the 3 remaining manifold learning algorithms in our case study.
- These visualizations can be helpful in enriching dictionary entries with additional, corpus-based examples; the closest BNC sentences to the dictionary examples mostly reflected very similar semantic and syntactic patterns.
- In our charts, we also saw thematically-motivated clusters of BNC sentences that were ignored during exemplification of the OD headword (e.g., the uses of the word *mouth* in romantic literature).



# Acknowledgement

This publication was supported  
by the *University of Debrecen*  
*Faculty of Humanities Scholarly*  
*Fund*.



# References

- Demsar J., Curk, T., Erjavec, A., Gorup, C., Hocevar, T., Milutinovic, M., Mozina, M., Polajnar, M., Toplak, M., Staric, A., Stajdohar, M., Umek, L., Zagar, L., Zbontar, J., Zitnik, M., Zupan, B. (2013.) Orange: Data Mining Toolbox in Python. *Journal of Machine Learning Research*, 14, pp. 2349–2353.
- Harris, Z. S. (1954) Distributional Structure. *Word*, 10:2-3, pp. 146–162, DOI: 10.1080/00437956.1954.11659520
- Rousseeuw, P. J. (1987). Silhouettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis. *Journal of Computational and Applied Mathematics*, 20, pp. 53–65.
- Tenenbaum, J. B., de Silva, V. & Langford, J. C. (2000). A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500), pp. 2319–2323.
- van der Maaten, L. & Hinton, G. (2008). Visualizing Data Using t-SNE. *Journal of Machine Learning Research*, 1, pp. 1–48.