

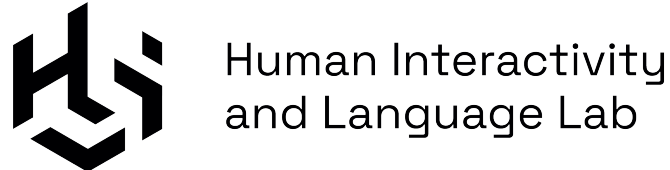
Semantic Space of Distributional Models as compared to Human’s Abstract Categorization and **Roots for Creative Association**

Urszula Kuczma, University of Warsaw

Joanna Rączaszek-Leonardi, University of Warsaw

Kristian Tylén, Aarhus University

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 952324



MAIN STUDY

Presented work shows the role of distributional semantics in divergent thinking / creativity study. It also indicates how such research can create potential for developing computational creativity models. The purpose of the main study was to discover:

- what strategies people use when associating concepts
- whether application of different strategies influences creative potential

In the main study participants were asked to perform a Verbal Fluency Task followed by an interview, to first obtain their associative path (manner in which they listed animals one by another) and then discuss with them what was the reason behind what they said in relevance to the previous response (if there was any conscious reason).

This task helped us to reveal certain strategies (Fig 3.) behind the process of associating. We confirmed a hypothesis that there is a correlation between the manner in which one uses these strategies and one's creative potential. This outcome indicates that the structure of semantic space and the manner of “moving” through it, also for digital agents, matters for the creativity of the outcome.

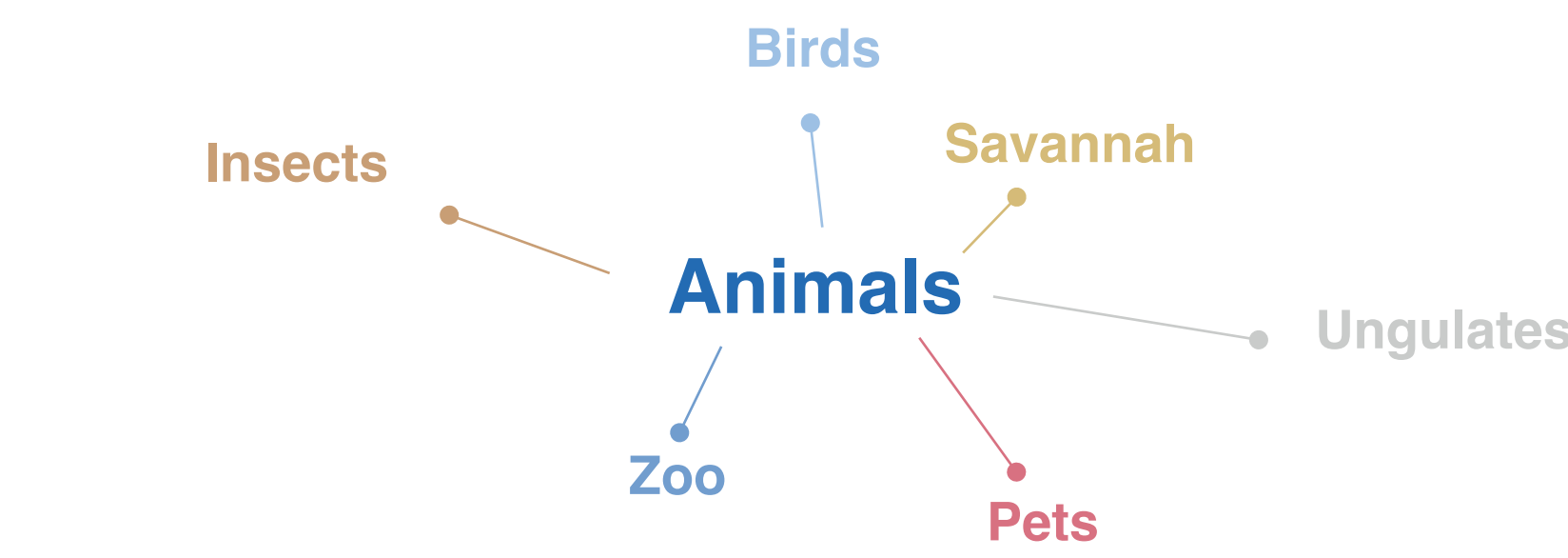


Fig 1. The domain chosen in the study

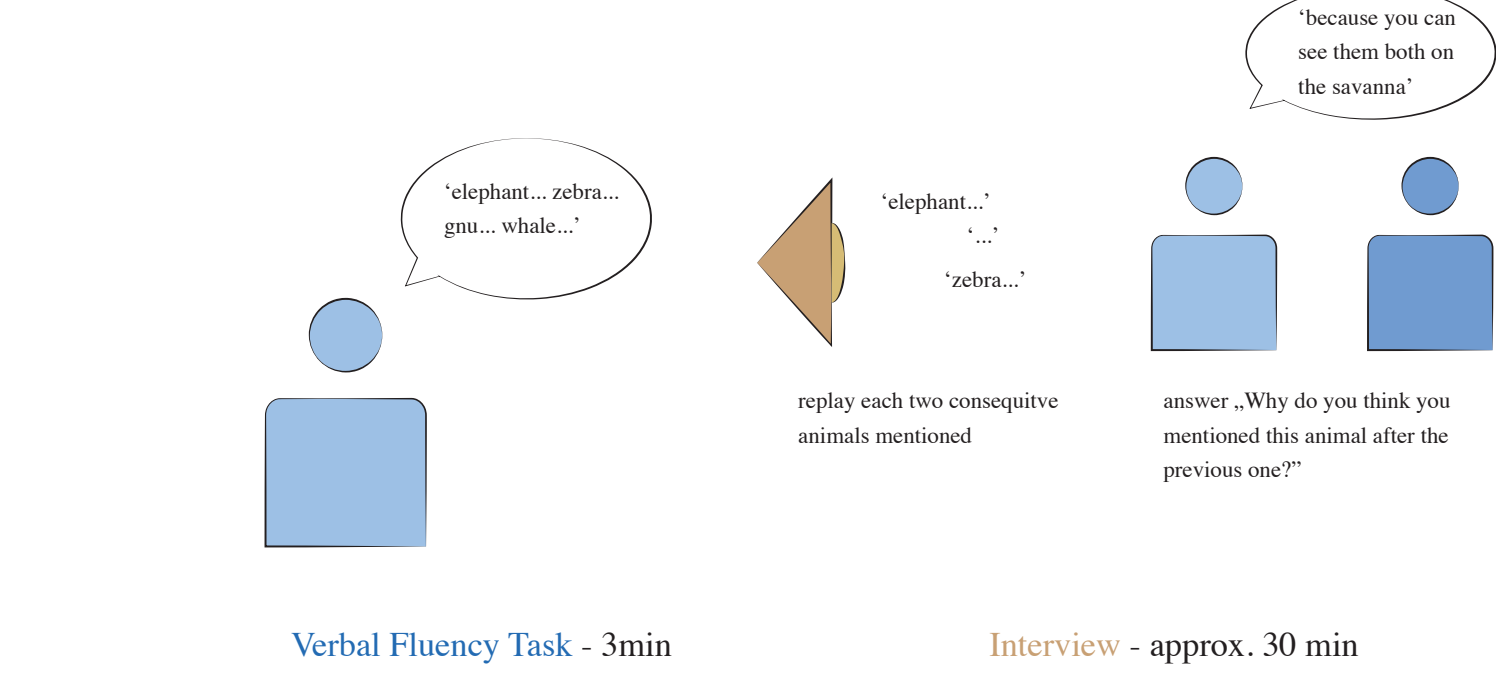


Fig 2. Scheme of the task

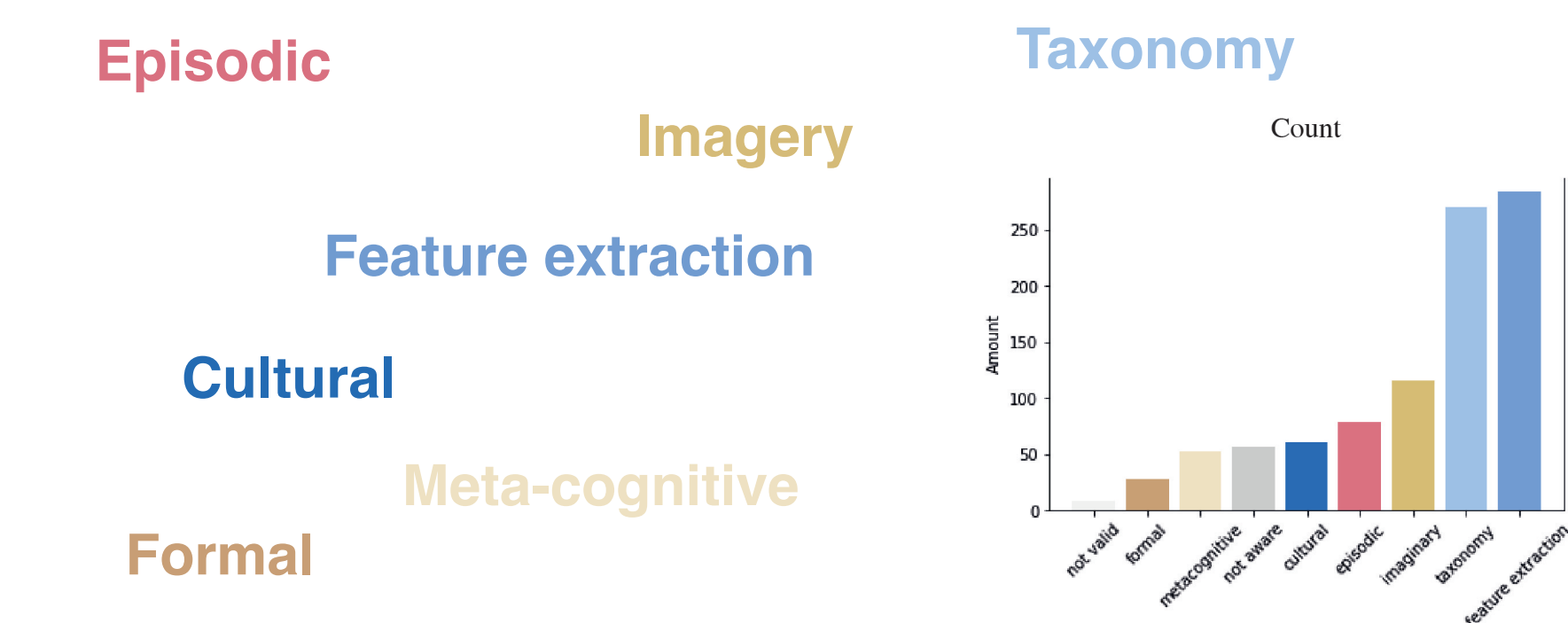
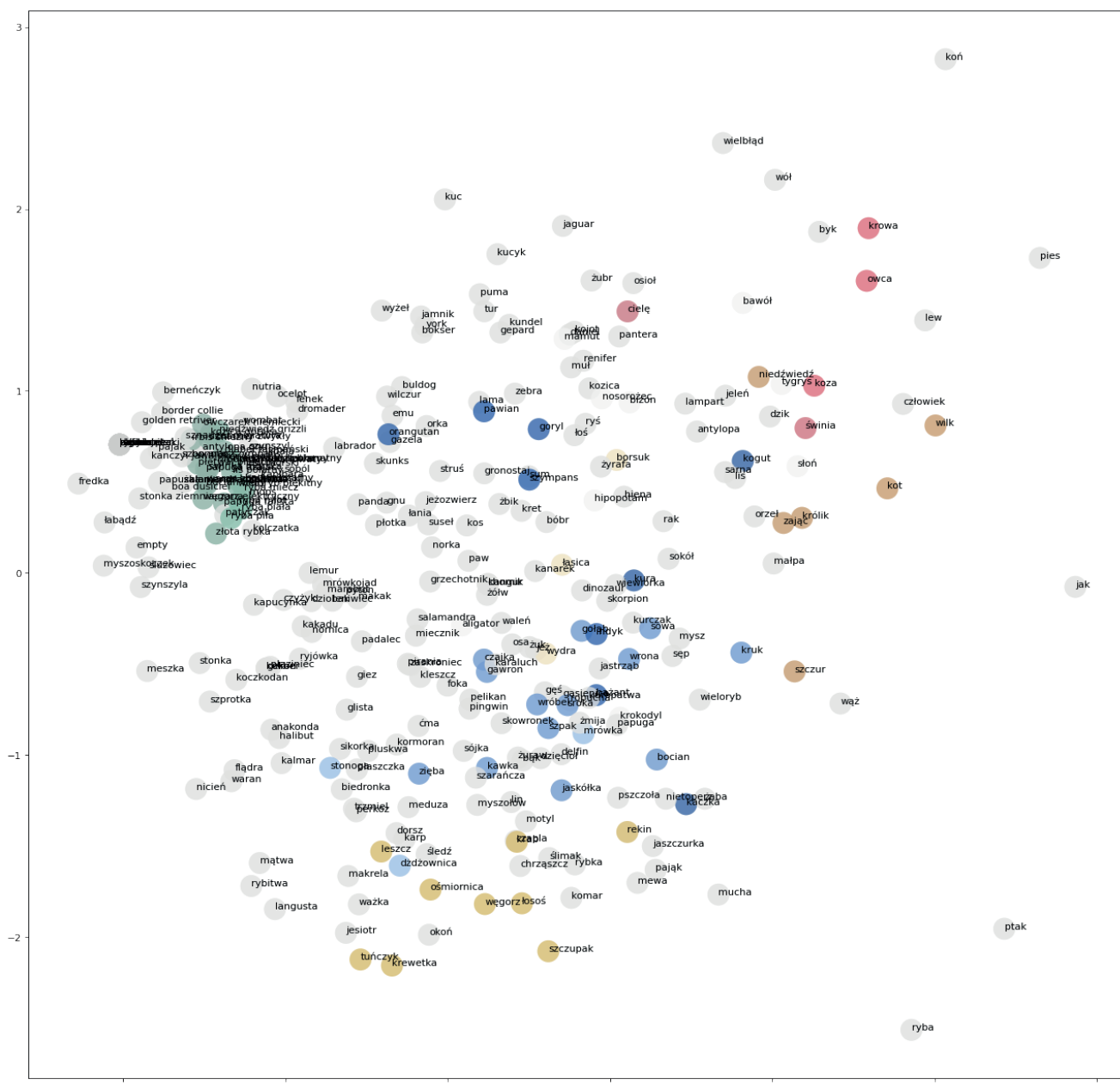


Fig 3. Possible strategies of associations (non-exhaustive)

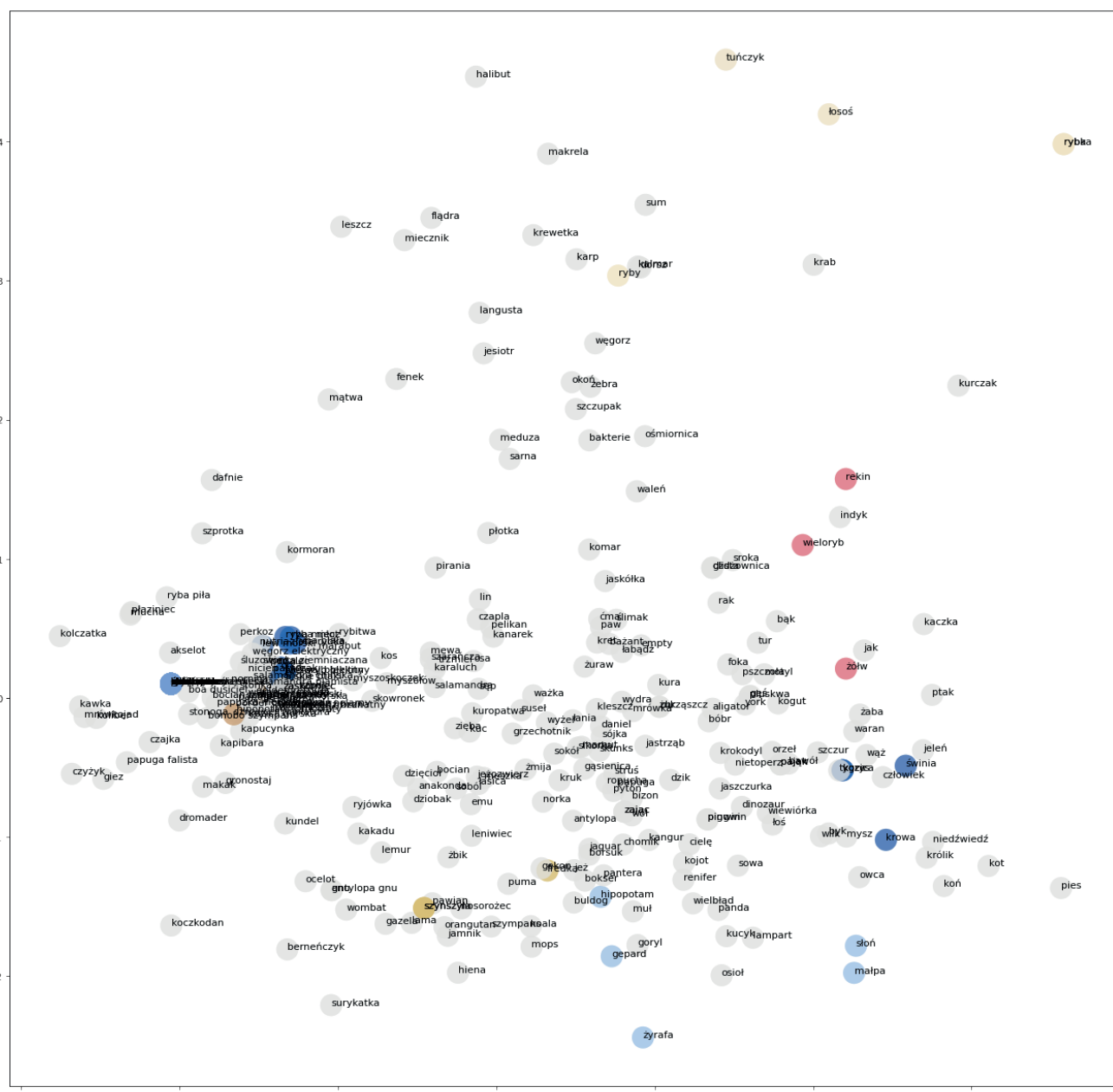
MODEL COMPARISON

Distributional Semantic Model was essential in the study to compare the distances between the concepts (flexibility of responses). In the past this measure was specified by independent judges, now it is mostly obtained through cosine similarity of vectors. We compared several currently available models and the resulting semantic spaces can be seen below (visualization through PCA dimension reduction [6]). To check validity of the models in terms of approximating human concepts categorization we used clustering method - HDBSCAN [5].

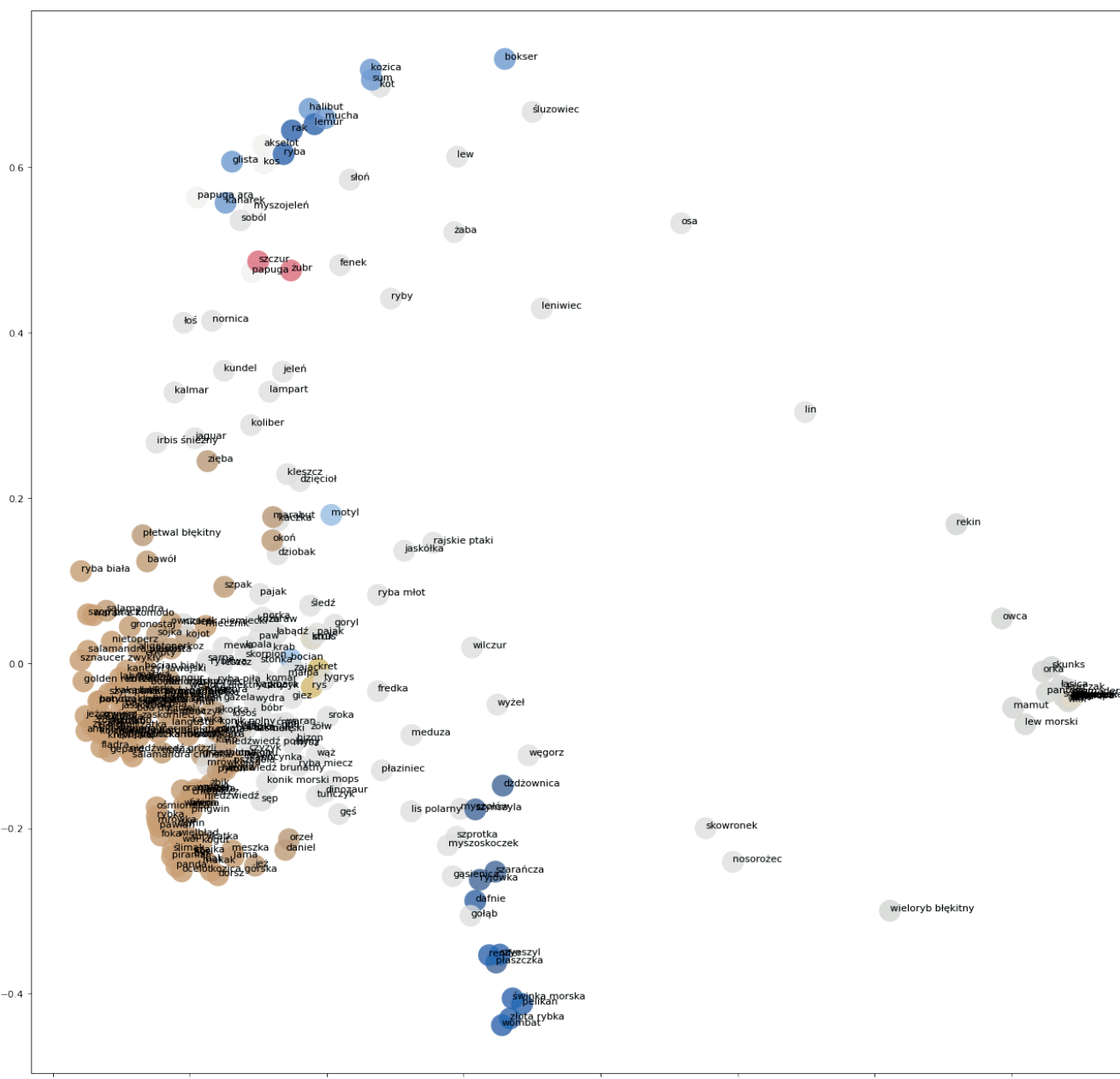
The results show that static models like GloVe seem more appropriate for this type of task, producing more easily identifiable conceptual clusters. Nevertheless, it seems that computationally obtained semantic spaces, based purely on linguistic input, might not be sufficient to accurately represent complexity of human abstraction (generalization and feature extraction). The character of discovered strategies of association indicate possibility of creating multi-modal concept space as solution.



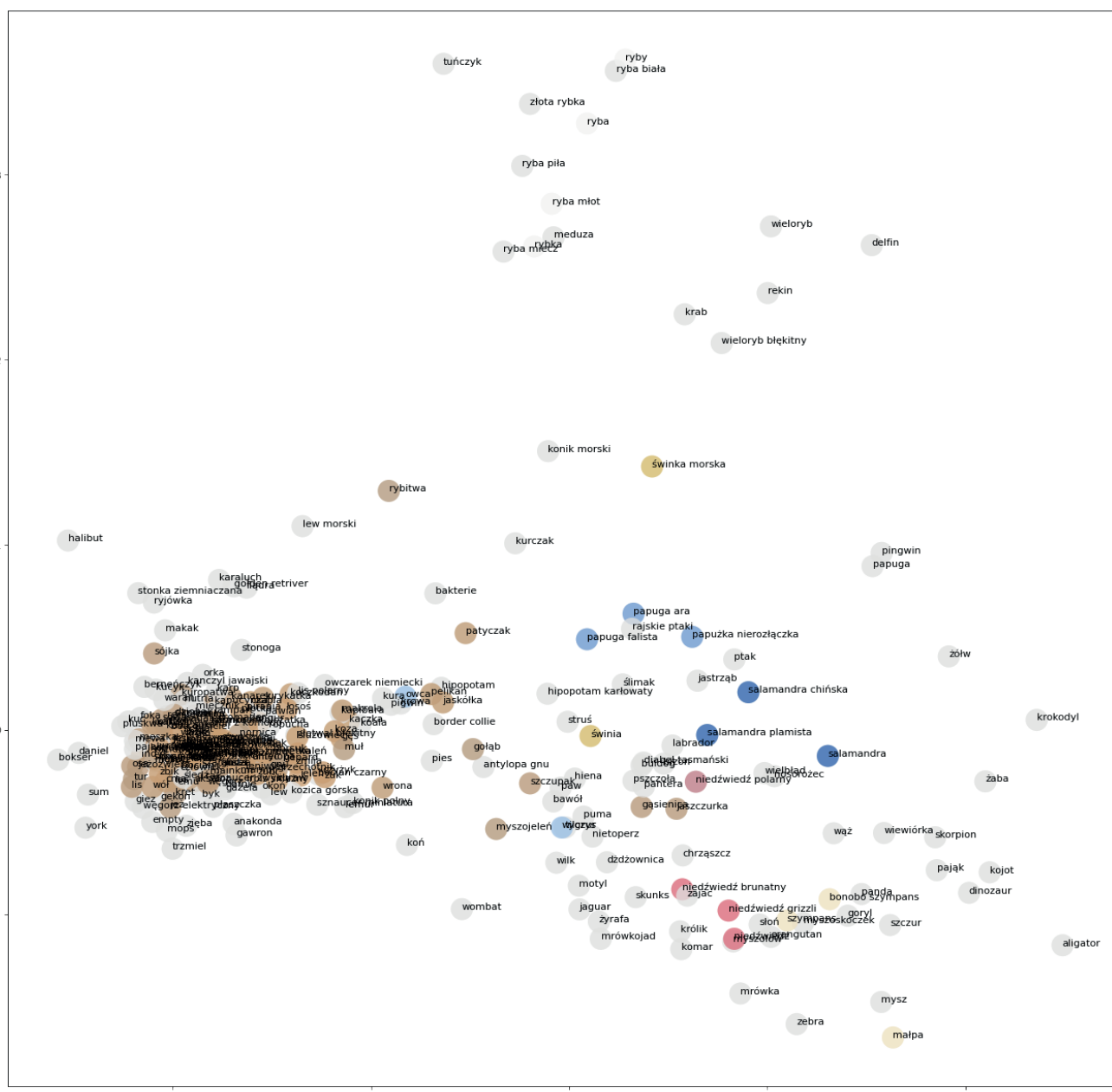
GloVe polish embeddings
glove_100_3_polish.txt [1]
11 Clusters recognized e.g.
Some farm mammals
Some apes
9 unrecognised words
1 Unclear cluster
1 Cluster including only compounds



GloVe from translation
glove.42B.100d.txt [2]
10 Clusters recognized e.g.:
Some animals from
Some fish
21 unrecognised words
2 Unclear clusters
1 Cluster including only compounds



BERT polish - averaged last layers
allegro / herbert-large-cased [3]
9 Clusters recognized e.g.:
Some forest animals
Some flying animals (only)
0 unrecognised words
7 Unclear clusters



Sentence Transformers
paraphrase-multilingual-MiniLM-L12-v2 [4]
8 Clusters recognized e.g.:
Some types of bears (3)
Some types of parrots (3)
0 unrecognised words
2 Unclear clusters

DIGITAL VS. HUMAN

Clustering between domains and within domains

As we could see above, the models are not fully capable of recognizing sub-domains, especially since many of the words from a given category can be fairly interchangeable grammatically and contextually. One of the things that distinguishes a human semantic space from a distributional model is the level of details. Human knowledge structure can possess vast hierarchy of elements, e.g. taxonomy system (fig.4).

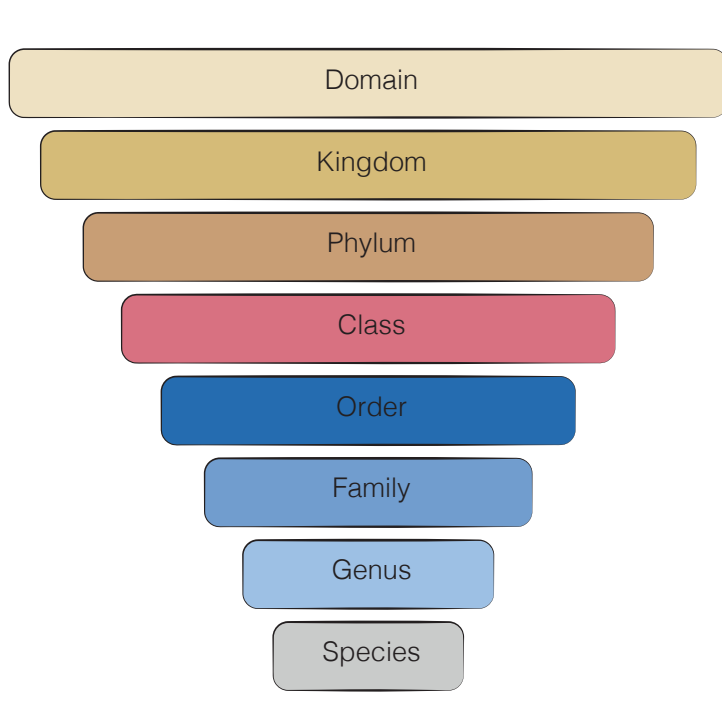


Fig 4. Example of hierarchical categorization of concepts

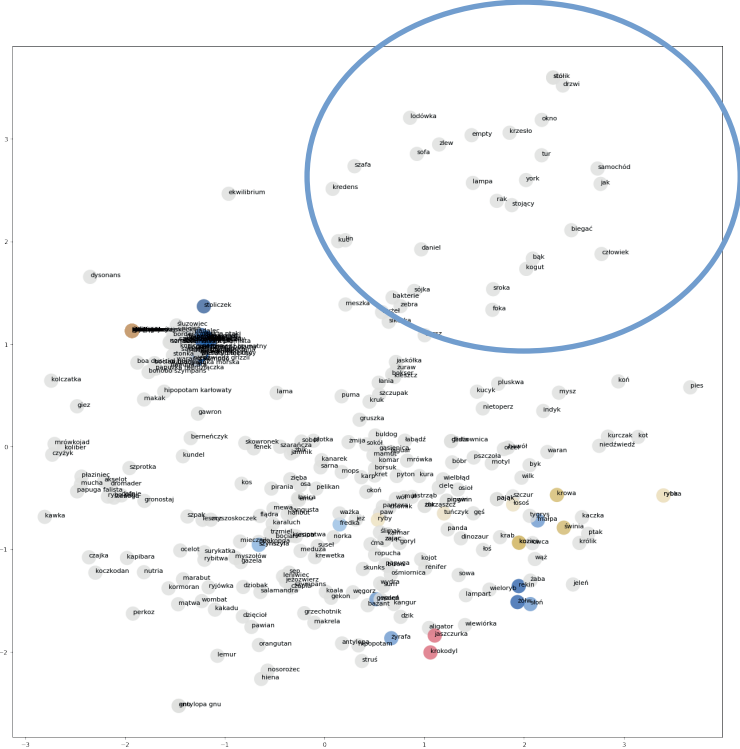


Fig 5. Position of non-clustered noise words in the GloVe Polish semantic space

Individual differences

The other important aspect of human semantic space is its uniqueness. Each person has a different semantic space (with more emphasis e.g. on experience, vision or language). What is more, semantic space of an individual can also change, based on time or external factors. The semantic space used commonly to compare between participants is more of a social product of such individual spaces.

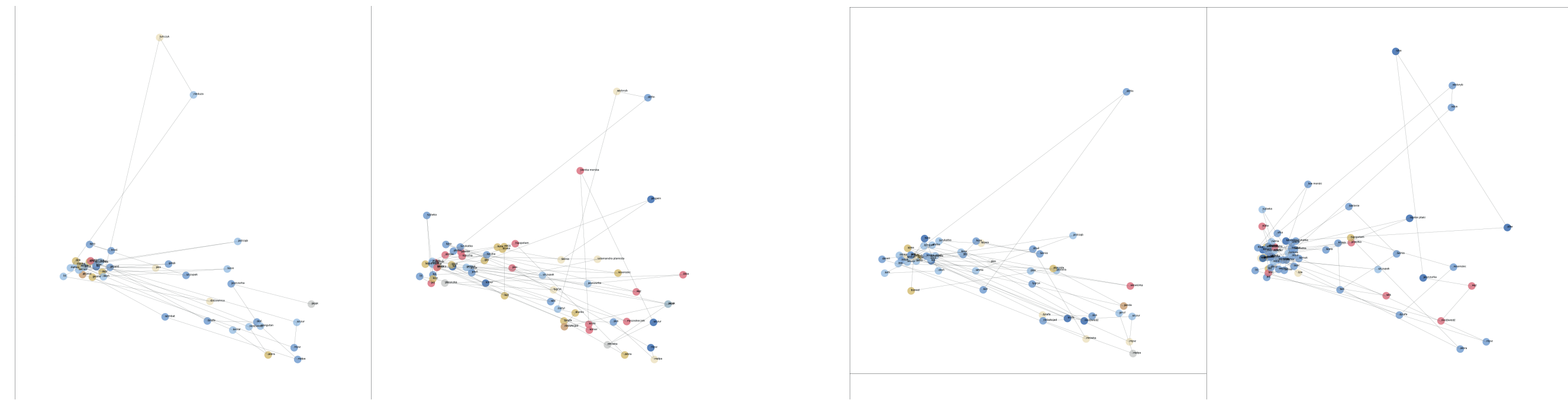


Fig 6. Example of differences within semantic space between participants

CONCLUSIONS

In this work we assessed the applicability of distributional semantics models to reflect human semantic knowledge. We discovered certain shortcomings, but also possibilities for improvement. We propose to advance from semantic space representation to 'conceptual' space representation. It is especially vital in the research of computational creativity, where both extensive exploration of the semantic space and validation of relevance of the responses is needed.

Such 'conceptual' space could possibly be obtained via combination of modalities - creating similarity distances not only through linguistic representation, but also visual or auditory. This is one of the topics we would like to expand our research with. Another aspect is creating a "surprise" factor differentiating individual human semantic spaces, namely episodic memories. This would allow for less obvious connections and possibly spark more of the creative exploration, present in the initial phase of ideation.

References

1. Dadas S., (2019) A repository of Polish NLP resources
2. Pennington J., Socher R., Manning C. D. (2014). GloVe: Global Vectors for Word Representation
3. Mroczkowski, R., Rybak, P., Wróblewska, A., and Gawlik, I. (2021). HerBERT: Efficiently Pretrained Transformer-based Language Model for Polish
4. Reimers N., Gurevych I. (2019) Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks
5. McInnes L., Healy J., Astels S. (2017), hdbscan: Hierarchical density based clustering, *Journal of Open Source Software*, 2(11), 205,
6. Pedregosa et al. (2011) Scikit-learn: Machine Learning

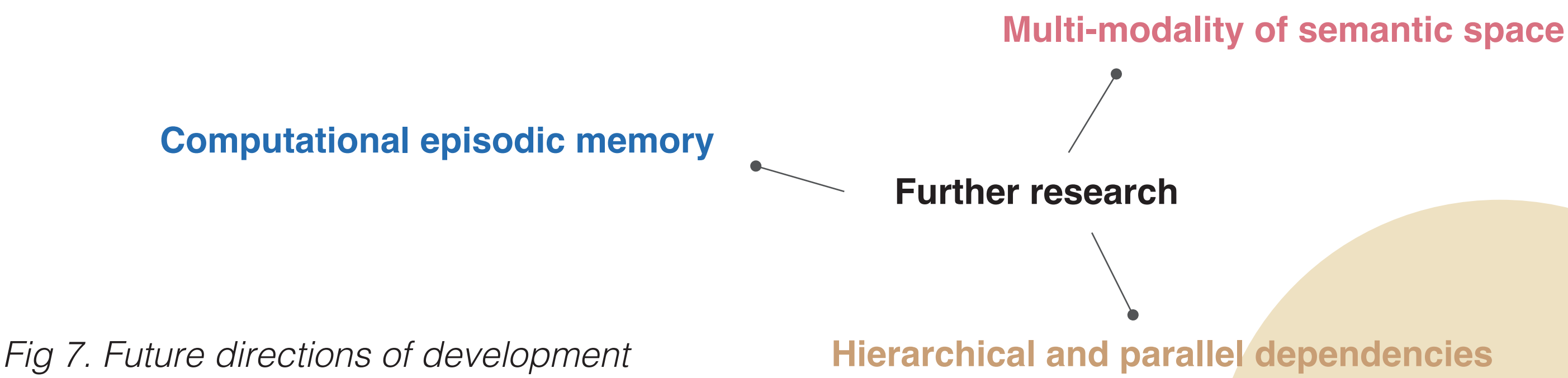


Fig 7. Future directions of development

For further questions please contact Urszula Kuczma at urszulakuczma@gmail.com
This research was possible thanks to the funding from Traincrease Program Horizon 2020.