

Word class representations spontaneously emerge in a deep neural network trained on next word prediction

Kishore Surendra¹, Achim Schilling^{2,3}, Paul Stoewer^{3,4}, Andreas Maier⁴, and Patrick Krauss^{2,3,5}

¹University Hospital Hamburg, Germany

²Neuroscience Lab, University Hospital Erlangen, Germany

³Cognitive Computational Neuroscience Group, University Erlangen-Nuremberg, Germany

⁴Pattern Recognition Lab, University Erlangen-Nuremberg, Germany

⁵Linguistics Lab, University Erlangen-Nuremberg, Germany

February 16, 2023

Keywords: successor representations, cognitive maps, word classes, deep neural networks, text prediction, syntax, construction grammar, usage-based models of language acquisition, ChatGPT

Abstract

How do humans learn language, and can the first language be learned at all? These fundamental questions are still hotly debated. In contemporary linguistics, there are two major schools of thought that give completely opposite answers. According to Chomsky's theory of universal grammar, language cannot be learned because children are not exposed to sufficient data in their linguistic environment. In contrast, usage-based models of language assume a profound relationship between language structure and language use. In particular, contextual mental processing and mental representations are assumed to have the cognitive capacity to capture the complexity of actual language use at all levels. The prime example is syntax, i.e., the rules by which words are assembled into larger units such as sentences. Typically, syntactic rules are expressed as sequences of word classes. However, it remains unclear whether word classes are innate, as implied by universal grammar, or whether they emerge during language acquisition, as suggested by usage-based approaches. Here, we address this issue from a machine learning and natural language processing perspective. In particular, we trained an artificial deep neural network on predicting the next word, provided sequences of consecutive words as input. Subsequently, we analyzed the emerging activation patterns in the hidden layers of the neural network. Strikingly, we find that the internal representations of nine-word input sequences cluster according to the word class of the tenth word to be predicted as output, even though the neural network did not receive any explicit information about syntactic rules or word classes during training. This surprising result suggests, that also in the human brain, abstract representational categories such as word classes may naturally emerge as a consequence of predictive coding and processing during language acquisition.

Introduction

The question of how humans come to language is one of the oldest scientific problems [1]. According to the Greek historian Herodotus, already 2500 years ago the Egyptian pharaoh Psamtik sought to discover the origin of language. Therefore, he conducted an experiment with two children which he gave as newborn babies to a shepherd who should feed and care for them, but had the instruction not to speak to them. Psamtik hypothesized that the infants' first word would be uttered in the root language of all people. Consequently, as one of the children cried *'bekos'* which was the sound of the Phrygian word for "bread", Psamtik concluded that Phrygian was the root language of all humans because that [2]. Obviously, the assumption behind this cruel language deprivation experiment was that humans are born with innate words and their meanings, and that this root language is somehow 'over-ruled' during individual development and first language learning.

Nowadays, it is of course clear that words and meanings are not innate but rather learned during language acquisition [3], and that there is no causal relation between the signifier (sound pattern) and the signified (meaning) [4]. However, it is still highly debated to what extent language capacities are innate or must be learned.

According to Chomsky's theory of universal grammar, humans have an innate, genetically determined language faculty that e.g. distinguishes between different word classes such as nouns and verbs making it easier and faster for children to learn to speak [5–7]. In contrast, in cognitive linguistics and usage-based approaches, a profound relationship between language structure and language use is assumed [8–11]. In particular, contextual mental processing and mental representations are assumed to have the cognitive capacity to capture the complexity of actual language use at all levels [12–17]. According to Diessel, grammar is a "dynamic system of emergent structures" and it needs to be explained "how linguistic structures evolve" during language acquisition [18].

Predictive coding and processing are thought to be canonical computations of the human brain [19–22], in particular during speech and language processing which involves the prediction of which words come next [23]. In previous studies, we already demonstrated that efficient successor representations to form cognitive maps of space and language can be learned by artificial neural networks [24, 25]. In particular, we demonstrated how a neural network model can infer the underlying word classes of a simplified artificial language just by observing sequences of words, i.e. sentences, and without any prior knowledge about actual word classes or grammar. The emerging representations share important properties with network-like cognitive maps, enabling e.g. navigation in arbitrary abstract and conceptual spaces, and thereby broadly supporting domain-general cognition, as proposed by Bellmund et al. [26].

In this follow-up study, we further address the question if abstract linguistic categories and structures can be learned from experienced language alone in a more complex and naturalistic linguistic task, i.e. word prediction in a natural language scenario. In particular, we trained an artificial deep neural network to predict the next word (successor) in a novel given the nine consecutive predecessor words as input. Subsequently, we analyzed the emerging activation patterns in the hidden layers of the neural network. Strikingly, we find that the internal representations of nine-word input sequences cluster according to the word class of the tenth word to be predicted as output, even though the neural network did not receive any explicit information about syntactic rules or word classes during training. This surprising result suggests, that also in the human brain, abstract representational categories such as word classes may naturally emerge as a consequence of predictive coding and processing of language input. Based on these findings we hypothesize that during language acquisition – which at least partly corresponds to learn to predict which word or utterance comes next –, word classes spontaneously emerge as clusters of successor representations of perceived utterances. We conclude that word classes need not to be innate to enable efficient language acquisition as suggested by universal grammar.

Methods

Data pre processing

The German novel *Gut gegen Nordwind* by Daniel Glattauer (© *Deuticke im Paul Zsolnay Verlag*, Wien 2006, published by *Deuticke Verlag* served as natural language text data for training and testing our model. The complete text data consists of a total number of 40460 tokens and 6117 types. Prior to further analysis, punctuation and special characters have been removed from the text corpus. Furthermore, repetitive words and extra white spaces have been removed, and all numbers have been replaced by a single word (cf. table 1). All words are converted to lower case to maintain uniformity, so that the same word occurring in a different case is considered as two tokens of the same type, instead of two different types. All words have been encoded as 384-dimensional word vectors using the word2vec embedding function from the python library *spaCy* [27]. Sequences of nine consecutive word vectors served as input, while one (the tenth) or two (tenth and eleventh) word vectors served as corresponding output. Finally, the all word vector sequences were split into a training (chapters 1 to 7 of the novel) and a test data set (chapters 8 and 9 of the novel).

Character/Word	Operation
Repetitive words: RE:, AW:, Eine, Zwei, ..., Stunden, Sekunden, Stunden..., später, Am nächsten, Kein Betreff, Betreff	Remove completely
Punctuation and other characters: ., ;, ?, %, &, ', ', !	Remove completely
Numbers: 18, 1, 500, ...	Replace with 'nummer'
Extra whitespaces	Replace with single space
E-mail	Replace with 'email'

Table 1: **Data cleaning.** Words, characters and their replacements during data cleaning.

Neural network architecture and training procedure

For the task at hand, i.e. to predict the tenth word (or the tenth and the eleventh word), given a prior sequence of nine words occurring in the corpus, recurrent neural networks (RNNs) are perfectly suited. Here, we implemented a neural network consisting of four bi-directional LSTM (long short-term memories) layers (with 128, 128, 64, and 64 neurons) followed by a flatten layer, and a dense output layer (384 neurons). The input consisted of sequences of nine 384-dimensional word embedding vectors generated as described above. The expected output is a single (or a sequence of two) 384-dimensional word embedding vectors. Weights were initialized using the Glorot uniform initialization, which is Keras's default initializer. As optimizer, we used Adam with a learning rate of 0.001 and as loss function we used mean-squared error. Training was performed for 100 epochs.

Word classes

Word classes were analysed by applying *part-of-speech (POS) tagging* [?, ?, ?] as implemented in the python library *spaCy* [27]. The used POS tags comprised the following 13 default word classes: 'NUM', 'VERB', 'ADJ', 'X', 'PART', 'NOUN', 'SCONJ', 'ADP', 'DET', 'PRON', 'CONJ', 'AUX' and 'ADV'. Their exact definitions can be found in [28]. Note that, during training the neural networks, we did not provide any information about word classes as input.

Multi-dimensional scaling

A frequently used method to generate low-dimensional embeddings of high-dimensional data is t-distributed stochastic neighbor embedding (t-SNE) [29]. However, in t-SNE the resulting low-

dimensional projections can be highly dependent on the detailed parameter settings [30], sensitive to noise, and may not preserve, but rather often scramble the global structure in data [31, 32]. In contrast to that, multi-Dimensional-Scaling (MDS) [33–36] is an efficient embedding technique to visualize high-dimensional point clouds by projecting them onto a 2-dimensional plane. Furthermore, MDS has the decisive advantage that it is parameter-free and all mutual distances of the points are preserved, thereby conserving both the global and local structure of the underlying data.

When interpreting patterns as points in high-dimensional space and dissimilarities between patterns as distances between corresponding points, MDS is an elegant method to visualize high-dimensional data. By color-coding each projected data point of a data set according to its label, the representation of the data can be visualized as a set of point clusters. For instance, MDS has already been applied to visualize for instance word class distributions of different linguistic corpora [37], hidden layer representations (embeddings) of artificial neural networks [38, 39], structure and dynamics of highly recurrent neural networks [40–43], or brain activity patterns assessed during e.g. pure tone or speech perception [37, 44], or even during sleep [45–48]. In all these cases the apparent compactness and mutual overlap of the point clusters permits a qualitative assessment of how well the different classes separate.

Generalized Discrimination Value (GDV)

We used the GDV to calculate cluster separability as published and explained in detail in [38]. Briefly, we consider N points $\mathbf{x}_{\mathbf{n}=1..N} = (x_{n,1}, \dots, x_{n,D})$, distributed within D -dimensional space. A label l_n assigns each point to one of L distinct classes $C_{l=1..L}$. In order to become invariant against scaling and translation, each dimension is separately z-scored and, for later convenience, multiplied with $\frac{1}{2}$:

$$s_{n,d} = \frac{1}{2} \cdot \frac{x_{n,d} - \mu_d}{\sigma_d}. \quad (1)$$

Here, $\mu_d = \frac{1}{N} \sum_{n=1}^N x_{n,d}$ denotes the mean, and $\sigma_d = \sqrt{\frac{1}{N} \sum_{n=1}^N (x_{n,d} - \mu_d)^2}$ the standard deviation of dimension d . Based on the re-scaled data points $\mathbf{s}_{\mathbf{n}} = (s_{n,1}, \dots, s_{n,D})$, we calculate the *mean intra-class distances* for each class C_l

$$\bar{d}(C_l) = \frac{2}{N_l(N_l-1)} \sum_{i=1}^{N_l-1} \sum_{j=i+1}^{N_l} d(\mathbf{s}_i^{(l)}, \mathbf{s}_j^{(l)}), \quad (2)$$

and the *mean inter-class distances* for each pair of classes C_l and C_m

$$\bar{d}(C_l, C_m) = \frac{1}{N_l N_m} \sum_{i=1}^{N_l} \sum_{j=1}^{N_m} d(\mathbf{s}_i^{(l)}, \mathbf{s}_j^{(m)}). \quad (3)$$

Here, N_k is the number of points in class k , and $\mathbf{s}_i^{(k)}$ is the i^{th} point of class k . The quantity $d(\mathbf{a}, \mathbf{b})$ is the euclidean distance between \mathbf{a} and \mathbf{b} . Finally, the Generalized Discrimination Value (GDV) is calculated from the mean intra-class and inter-class distances as follows:

$$\text{GDV} = \frac{1}{\sqrt{D}} \left[\frac{1}{L} \sum_{l=1}^L \bar{d}(C_l) - \frac{2}{L(L-1)} \sum_{l=1}^{L-1} \sum_{m=l+1}^L \bar{d}(C_l, C_m) \right] \quad (4)$$

whereas the factor $\frac{1}{\sqrt{D}}$ is introduced for dimensionality invariance of the GDV with D as the number of dimensions.

Note that the GDV is invariant with respect to a global scaling or shifting of the data (due to the z-scoring), and also invariant with respect to a permutation of the components in the N -dimensional data vectors (because the euclidean distance measure has this symmetry). The GDV is zero for completely overlapping, non-separated clusters, and it becomes more negative as the separation increases. A GDV of -1 signifies already a very strong separation.

Code Implementation

The models were coded in Python. The neural networks were designed using the Keras [49] and Keras-RL [50] libraries. Mathematical operations were performed with numpy [51] and scikit-learn [52] libraries. Visualizations were realised with matplotlib [53] and networkX [54]. For natural language processing we used SpaCy [27].

Results

Next word prediction

We trained a neural network on next word prediction using sequences of nine consecutive word vectors as input. The trained network was tested with sequences of nine words not used for training. The resulting neural activation of each layer was read out and the corresponding activation vectors were projected onto a 2-dimensional plane using MDS. All projected points were then color coded according to the word class of the subsequent word of the corresponding input sequence. Word classes were assessed after training using POS tagging and did not serve as input during training. While layer 1 shows a random distribution of the data points ??, we find a remarkably strong clustering according to world classes in the last layer of the neural network ?. This is also confirmed by the corresponding GDV curve across the layers ?. This means that the neural network organizes its internal representations of input word sequences according to the word class of the next word to be predicted.

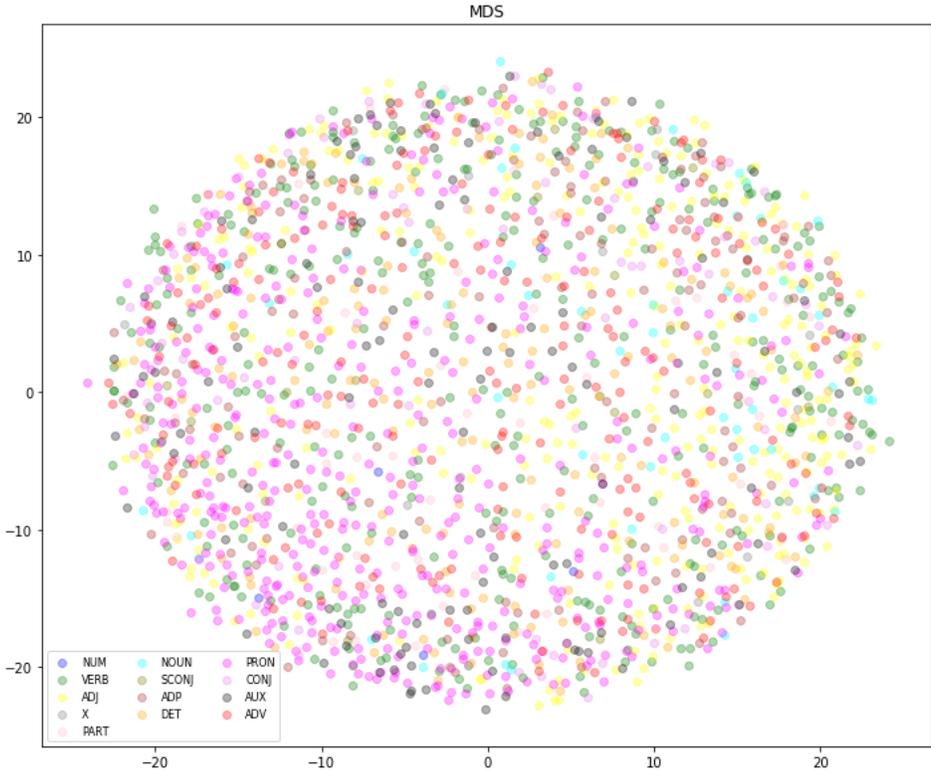


Figure 1: Layer 1 results of neural network testing and projection onto a 2-dimensional plane using MDS, with color coding according to subsequent word class. The points are randomly distributed.

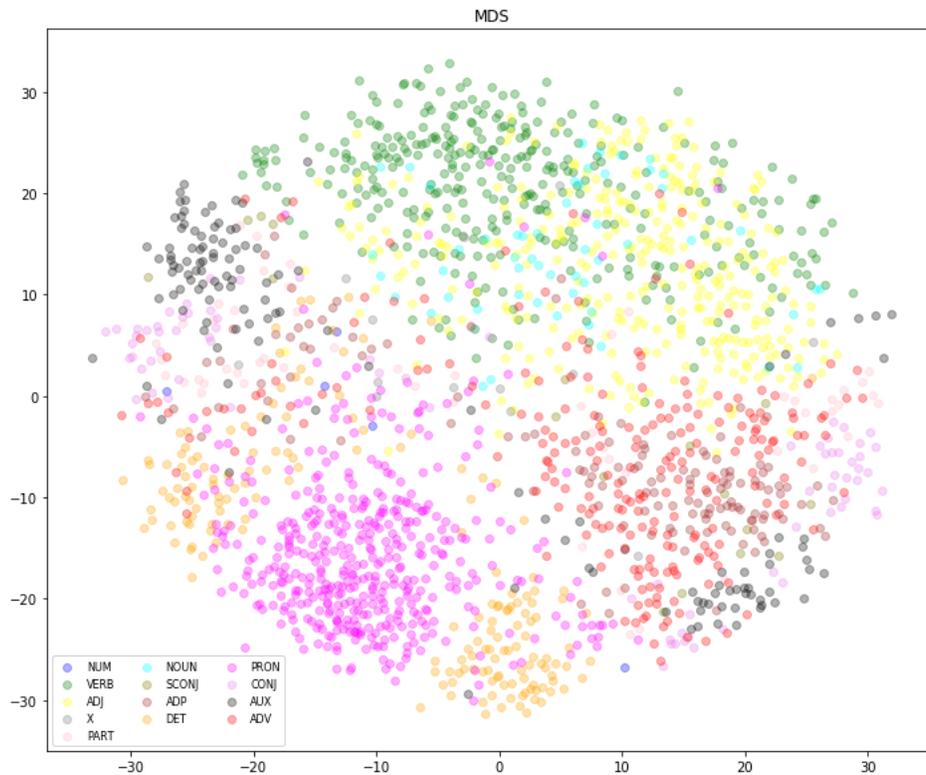


Figure 2: Last layer results of neural network testing and projection onto a 2-dimensional plane using MDS, with color coding according to subsequent word class. The final layer shows strong clustering by word class, indicating that the neural network organizes internal representations based on the predicted subsequent word's class.

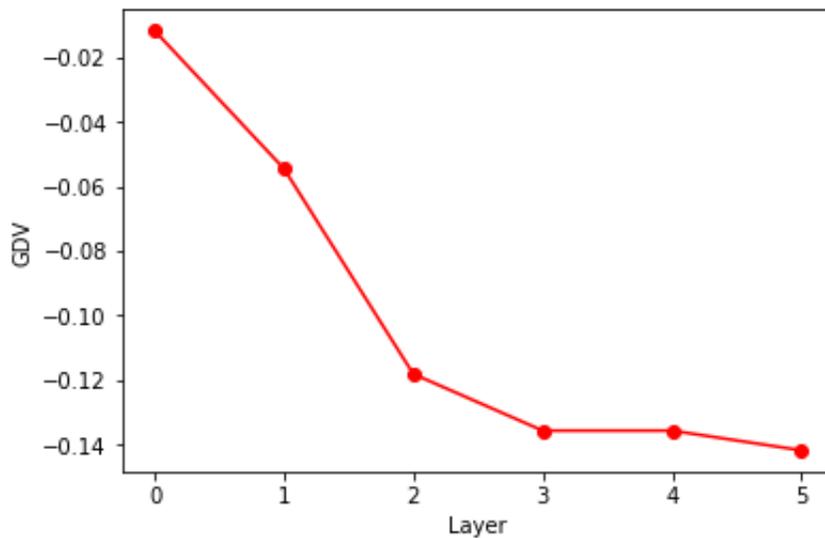


Figure 3: GDV curve across layers of the neural network. The decline of the GDV indicates that the neural network has learned to cluster internal representations according to the subsequent word's class with increasingly strong clustering from input to output layer.

Discussion

The results of our study provide evidence that abstract linguistic categories, such as word classes, can emerge spontaneously in neural representations of linguistic input. This finding challenges the

notion that the ability to recognize and categorize words by their grammatical function is innate and hardwired in the human brain, as proposed by Chomsky's theory of universal grammar. Our results suggest that language acquisition involves, at least in part, the learning of predictive structures and categories based on statistical regularities in the input, rather than relying solely on innate linguistic knowledge. This is consistent with the view that language is a complex adaptive system shaped by both biological and environmental factors. It is interesting to note that the clustering of input sequences by word class is evident only in the last layer of the neural network, suggesting that the network may gradually learn and refine more abstract and complex features of language as information flows through its layers. This finding is consistent with the hierarchical nature of language processing, in which higher-level representations build on lower-level representations. One potential application of our findings is in natural language processing, where understanding the organization of neural representations of language input can help improve language modeling, machine translation, and other related tasks. In addition, our study provides a starting point for further investigations into the neural mechanisms underlying language acquisition and processing. In conclusion, our study provides compelling evidence that neural networks can spontaneously learn to organize their internal representations of language input according to abstract linguistic categories such as word classes. Our results support the view that language acquisition is a complex and dynamic process that relies on both innate mechanisms and statistical learning from environmental input.

Acknowledgments

We are grateful to the publisher *Deuticke Verlag* for the permission to use the novel *Gut gegen Nordwind* by Daniel Glattauer for the present and future studies.

This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation): grants KR 5148/2-1 (project number 436456810), KR 5148/3-1 (project number 510395418) and GRK 2839 (project number 468527017) to PK, and grant SCHI 1482/3-1 (project number 451810794) to AS, and by the Emerging Talents Initiative (ETI) of the University Erlangen-Nuremberg (grant 2019/2-Phil-01 to PK).

Author contributions

KS, AS, and PK performed computer simulations. PK designed the study. PK, AS and AM supervised the study. KS prepared the figures. All authors discussed the results and wrote the manuscript.

Competing interests

The authors declare no competing financial interests.

References

- [1] Dieter E Zimmer. *So kommt der Mensch zur Sprache: über Spracherwerb, Sprachentstehung und Sprache & Denken*, volume 16. Heyne TB, 1986.
- [2] Henry Creswicke Rawlinson, John Gardner Wilkinson, et al. *The history of Herodotus*, volume 1. 1861.
- [3] Helen Goodluck. *Language acquisition: A linguistic introduction*. Basil Blackwell, 1991.
- [4] Ferdinand De Saussure. *Course in general linguistics*. Columbia University Press, 2011.

- [5] Noam Chomsky. On the nature, use and acquisition of language. In *Language and Meaning in Cognitive Science*, pages 13–32. Routledge, 2012.
- [6] Noam Chomsky. *Aspects of the Theory of Syntax*, volume 11. MIT press, 2014.
- [7] Charles Yang, Stephen Crain, Robert C Berwick, Noam Chomsky, and Johan J Bolhuis. The growth of language: Universal grammar, experience, and principles of computation. *Neuroscience & Biobehavioral Reviews*, 81:103–119, 2017.
- [8] Adele E Goldberg. *Constructions: A construction grammar approach to argument structure*. University of Chicago Press, 1995.
- [9] Adele E Goldberg. Constructions: A new theoretical approach to language. *Trends in cognitive sciences*, 7(5):219–224, 2003.
- [10] Michael Tomasello. *Constructing a language: A usage-based theory of language acquisition*. Harvard university press, 2005.
- [11] Ronald W Langacker. Cognitive grammar. *Basic Readings*, page 29, 2008.
- [12] Joan Bybee, Revere Perkins, William Pagliuca, et al. *The Evolution of Grammar: Tense, Aspect, and Modality in the Languages of the World*. University of Chicago Press, 1994.
- [13] Paul J Hopper and Joan L Bybee. Frequency and the emergence of linguistic structure. *Frequency and the Emergence of Linguistic Structure*, pages 1–502, 2001.
- [14] Joan L Bybee. Usage-based theory and exemplar representations of constructions. 2013.
- [15] Holger Diessel, Ewa Dabrowska, and Dagmar Divjak. Usage-based construction grammar. *Cognitive linguistics*, 2:50–80, 2019.
- [16] Adele Goldberg and Adele E Goldberg. *Explain me this*. Princeton University Press, 2019.
- [17] Hans-Jörg Schmid. *The dynamics of the linguistic system: Usage, conventionalization, and entrenchment*. Oxford University Press, 2020.
- [18] Holger Diessel. 14. usage-based construction grammar. In *Handbook of cognitive linguistics*, pages 296–322. De Gruyter Mouton, 2015.
- [19] James M Kilner, Karl J Friston, and Chris D Frith. Predictive coding: an account of the mirror neuron system. *Cognitive processing*, 8(3):159–166, 2007.
- [20] Andre M Bastos, W Martin Usrey, Rick A Adams, George R Mangun, Pascal Fries, and Karl J Friston. Canonical microcircuits for predictive coding. *Neuron*, 76(4):695–711, 2012.
- [21] Georg B Keller and Thomas D Mrsic-Flogel. Predictive processing: a canonical cortical computation. *Neuron*, 100(2):424–435, 2018.
- [22] Achim Schilling, William Sedley, Richard Gerum, Claus Metzner, Konstantin Tziridis, Andreas Maier, Holger Schulze, Fan-Gang Zeng, Karl J Friston, and Patrick Krauss. Predictive coding and stochastic resonance: Towards a unified theory of auditory (phantom) perception. *arXiv preprint arXiv:2204.03354*, 2022.
- [23] Armine Garibyan, Achim Schilling, Claudia Boehm, Alexandra Zankl, and Patrick Krauss. Neural correlates of linguistic collocations during continuous speech perception. *bioRxiv*, 2022.

- [24] Paul Stoewer, Christian Schlieker, Achim Schilling, Claus Metzner, Andreas Maier, and Patrick Krauss. Neural network based successor representations to form cognitive maps of space and language. *Scientific Reports*, 12(1):1–13, 2022.
- [25] Paul Stoewer, Achim Schilling, Andreas Maier, and Patrick Krauss. Neural network based formation of cognitive maps of semantic spaces and the emergence of abstract concepts. *arXiv preprint arXiv:2210.16062*, 2022.
- [26] Jacob LS Bellmund, Peter Gärdenfors, Edvard I Moser, and Christian F Doeller. Navigating cognition: Spatial codes for human thinking. *Science*, 362(6415), 2018.
- [27] AI Explosion. spacy-industrial-strength natural language processing in python. *URL: https://spacy.io*, 2017.
- [28] Slav Petrov, Dipanjan Das, and Ryan McDonald. A universal part-of-speech tagset. *arXiv preprint arXiv:1104.2086*, 2011.
- [29] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [30] Martin Wattenberg, Fernanda Viégas, and Ian Johnson. How to use t-sne effectively. *Distill*, 1(10):e2, 2016.
- [31] Catalina A Vallejos. Exploring a world of a thousand dimensions. *Nature biotechnology*, 37(12):1423–1424, 2019.
- [32] Kevin R Moon, David van Dijk, Zheng Wang, Scott Gigante, Daniel B Burkhardt, William S Chen, Kristina Yim, Antonia van den Elzen, Matthew J Hirn, Ronald R Coifman, et al. Visualizing structure and transitions in high-dimensional biological data. *Nature biotechnology*, 37(12):1482–1492, 2019.
- [33] Warren S Torgerson. Multidimensional scaling: I. theory and method. *Psychometrika*, 17(4):401–419, 1952.
- [34] Joseph B Kruskal. Nonmetric multidimensional scaling: a numerical method. *Psychometrika*, 29(2):115–129, 1964.
- [35] Joseph B Kruskal and Myron Wish. *Multidimensional scaling*, volume 11. Sage, 1978.
- [36] Michael AA Cox and Trevor F Cox. Multidimensional scaling. In *Handbook of data visualization*, pages 315–347. Springer, 2008.
- [37] Achim Schilling, Rosario Tomasello, Malte R Henningsen-Schomers, Alexandra Zankl, Kishore Surendra, Martin Haller, Valerie Karl, Peter Uhrig, Andreas Maier, and Patrick Krauss. Analysis of continuous neuronal activity evoked by natural speech with computational corpus linguistics methods. *Language, Cognition and Neuroscience*, 36(2):167–186, 2021.
- [38] Achim Schilling, Andreas Maier, Richard Gerum, Claus Metzner, and Patrick Krauss. Quantifying the separability of data classes in neural networks. *Neural Networks*, 139:278–293, 2021.
- [39] Patrick Krauss, Claus Metzner, Nidhi Joshi, Holger Schulze, Maximilian Traxdorf, Andreas Maier, and Achim Schilling. Analysis and visualization of sleep stages based on deep neural networks. *Neurobiology of sleep and circadian rhythms*, 10:100064, 2021.
- [40] Patrick Krauss, Alexandra Zankl, Achim Schilling, Holger Schulze, and Claus Metzner. Analysis of structure and dynamics in three-neuron motifs. *Frontiers in Computational Neuroscience*, 13:5, 2019.

- [41] Patrick Krauss, Karin Prebeck, Achim Schilling, and Claus Metzner. Recurrence resonance” in three-neuron motifs. *Frontiers in computational neuroscience*, page 64, 2019.
- [42] Patrick Krauss, Marc Schuster, Verena Dietrich, Achim Schilling, Holger Schulze, and Claus Metzner. Weight statistics controls dynamics in recurrent neural networks. *PloS one*, 14(4):e0214541, 2019.
- [43] Claus Metzner, Marius E Yamakou, Dennis Voelkl, Achim Schilling, and Patrick Krauss. Quantifying and maximizing the information flux in recurrent neural networks. *arXiv preprint arXiv:2301.12892*, 2023.
- [44] Patrick Krauss, Claus Metzner, Achim Schilling, Konstantin Tziridis, Maximilian Traxdorf, Andreas Wollbrink, Stefan Rampp, Christo Pantev, and Holger Schulze. A statistical method for analyzing and comparing spatiotemporal cortical activation patterns. *Scientific reports*, 8(1):1–9, 2018.
- [45] Patrick Krauss, Achim Schilling, Judith Bauer, Konstantin Tziridis, Claus Metzner, Holger Schulze, and Maximilian Traxdorf. Analysis of multichannel eeg patterns during human sleep: a novel approach. *Frontiers in human neuroscience*, 12:121, 2018.
- [46] Maximilian Traxdorf, Patrick Krauss, Achim Schilling, Holger Schulze, and Konstantin Tziridis. Microstructure of cortical activity during sleep reflects respiratory events and state of daytime vigilance. *Somnologie*, 23(2):72–79, 2019.
- [47] Claus Metzner, Achim Schilling, Maximilian Traxdorf, Konstantin Tziridis, Andreas Maier, Holger Schulze, and Patrick Krauss. Classification at the accuracy limit: facing the problem of data ambiguity. *Scientific Reports*, 12(1):22121, 2022.
- [48] Claus Metzner, Achim Schilling, Maximilian Traxdorf, Holger Schulze, Konstantin Tziridis, and Patrick Krauss. Extracting continuous sleep depth from eeg data without machine learning. *arXiv preprint arXiv:2301.06755*, 2023.
- [49] François Chollet et al. Keras. <https://keras.io>, 2015. Last visited: February 16, 2023.
- [50] Matthias Plappert. keras-rl. <https://github.com/keras-rl/keras-rl>, 2016. Last visited: February 16, 2023.
- [51] Charles R. Harris, K. Jarrod Millman, St’efan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fern’andez del R’io, Mark Wiebe, Pearu Peterson, Pierre G’erard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020.
- [52] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [53] J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007.
- [54] Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics, and function using networkx. In Gaël Varoquaux, Travis Vaught, and Jarrod Millman, editors, *Proceedings of the 7th Python in Science Conference*, pages 11 – 15, Pasadena, CA USA, 2008.