

## PERSPECTIVE

# Inferring the nature of linguistic computations in the brain

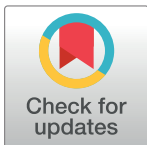
Sanne Ten Oever <sup>1,2,3</sup>, Karthikeya Kaushik<sup>1,2</sup>, Andrea E. Martin <sup>1,2\*</sup>

**1** Language and Computation in Neural Systems Group, Max Planck Institute for Psycholinguistics, Nijmegen, the Netherlands, **2** Donders Centre for Cognitive Neuroimaging, Radboud University, Nijmegen, the Netherlands, **3** Department of Cognitive Neuroscience, Faculty of Psychology and Neuroscience, Maastricht University, Maastricht, the Netherlands

\* [andrea.martin@mpi.nl](mailto:andrea.martin@mpi.nl)

## Abstract

Sentences contain structure that determines their meaning beyond that of individual words. An influential study by Ding and colleagues (2016) used frequency tagging of phrases and sentences to show that the human brain is sensitive to structure by finding peaks of neural power at the rate at which structures were presented. Since then, there has been a rich debate on how to best explain this pattern of results with profound impact on the language sciences. Models that use hierarchical structure building, as well as models based on associative sequence processing, can predict the neural response, creating an inferential impasse as to which class of models explains the nature of the linguistic computations reflected in the neural readout. In the current manuscript, we discuss pitfalls and common fallacies seen in the conclusions drawn in the literature illustrated by various simulations. We conclude that inferring the neural operations of sentence processing based on these neural data, and any like it, alone, is insufficient. We discuss how to best evaluate models and how to approach the modeling of neural readouts to sentence processing in a manner that remains faithful to cognitive, neural, and linguistic principles.



## OPEN ACCESS

**Citation:** Ten Oever S, Kaushik K, Martin AE (2022) Inferring the nature of linguistic computations in the brain. *PLoS Comput Biol* 18(7): e1010269. <https://doi.org/10.1371/journal.pcbi.1010269>

**Editor:** Daniel Bush, University College London, UNITED KINGDOM

**Published:** July 28, 2022

**Copyright:** © 2022 Ten Oever et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Funding:** AEM was supported by a Max Planck Research Group and a Lise Meitner Research Group “Language and Computation in Neural Systems” from the Max Planck Society, and by the Netherlands Organization for Scientific Research (NWO; grant 016.Vidi.188.029). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

## Introduction

Language is not contained in the physicality of speech, sign, or text; rather, the brain must construct meaning from the sensation of the physical input based on internal knowledge. To infer sentence meaning, we need to understand individual words as well as how these words structurally relate to each other. Ding and colleagues [1] showed that the human brain is sensitive to linguistic structure by presenting adjective-noun-verb-noun sentences containing a noun and a verb phrase (NP and VP, respectively). The stimuli are presented at a 4-Hz rate and contain only a 4-Hz acoustic signal. Nonetheless, they find that participants’ MEG responses show peaks not only at the acoustic (4-Hz) rate, but also at the sentence (1-Hz) and phrasal (2-Hz) rate. There has been a wide debate, based on this finding, about what the computational neural principles governing linguistic structure processing are.

One way to explain the Ding and colleagues data [1] is to assume that the brain has explicit models of linguistic structures that are imposed on the stimulus input ([2]; also see [3]). The 1-Hz patterns follow as the brain has an explicit response to sentences. Alternatively, a

response to sentence structure is extracted from the statistics in the stimulus input. Frank and Yang [4] showed that by concatenating distributed semantic representations of words, a 1-2-4-Hz pattern also emerges. As such, the authors state that the 1-Hz response does not require an explicit model of sentence structure and distributed associative representations in the brain can explain the data [4].

These two accounts create an interesting dichotomy that is unresolved. Which account provides a “better” mechanistic explanation of what the brain does? Here, we discuss some limitations on the interpretation of models of neural readouts related to linguistic structures. While the debate is open, it is important to shed light on how we should treat modeling of neural data in order to refrain from making unlicensed conclusions. This paper aims to serve as a didactic tool to improve future argumentation on what counts as evidence for hierarchical linguistic structure in neural readouts, and the resulting inferred computations. We have therefore structured our paper by posing (in our opinion) strong statements from each polar view that appear implicitly as well as explicitly in the literature.

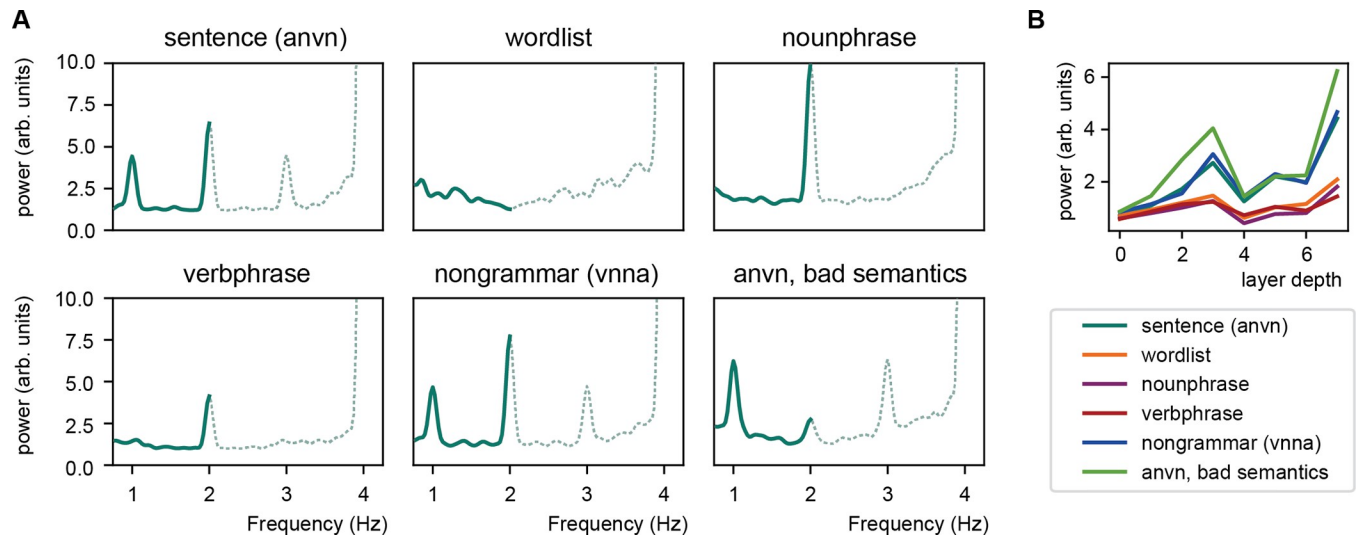
## Examples of unfounded inference

### **The output of artificial neural nets shows the 1-2-4 Hz without explicitly processing sentences, thus the brain does not explicitly process sentences**

Artificial neural networks (ANNs) of different architectures feature a collection of interconnected units that are trained according to a variety of procedures in order to extract statistical patterns in data. They can be trained to predict upcoming words, or can be used to create distributional semantic representations such as word2vec. Frank and Yang [4] showed that a 1-2-4-Hz pattern for sentences can be created by using word2vec representations and extracting the frequency content of the individual dimensions of this vector representation in the sentence. From this, the authors infer that it is not necessary for the brain to form an explicit representation of sentences when it creates the pattern found in the Ding and colleagues [1] data. Instead, they argue that because adjective-noun-verb-noun patterns follow each other, words can be processed sequentially, and the brain could simply be sensitive to the regularities on the word level, instead of creating an integrated representation of the sentence structure. Importantly, these results should not be interpreted as a proof that the brain does not form an abstract sentence representation (and Frank and Yang [4] also do not do this, but discourse in the field has interpreted this finding in this vein). Simply because one model can create the pattern by not using sentence representations does not imply that the brain does not reach the pattern in this way (i.e., the multiple realizability problem, see [5–9]). Indeed, Martin and Doumas [2] have shown that the 1-2-4-Hz pattern can also be created by a model that does encode sentence-like propositions by explicitly representing relationships between words and phrases in time [10]. Based on these two modeling approaches, and the Ding and colleagues [1] data alone, it is impossible via inductive inference to know whether the readout arise from the brain constructing explicit sentences representations or not.

### **Artificial neural networks by definition do not explicitly process sentence structure**

Just as the multiple realizability problem exists across different classes of computational models, it also exists within a single class of models such as ANNs. Indeed, ANNs processing language can achieve a 1-2-4-Hz pattern under vastly different architectures and underlying model goals. Some of these models process sentence structure explicitly, but others do not. As show above, word2vec representations can create a 1-2-4-Hz pattern without directly



**Fig 1. Output of the Berkeley parser.** (A) Spectra of the output of the parser using different versions of the Ding and colleagues stimuli (a = adjective, n = noun, v = verb). Dashed lines indicate output from the FFT when interpolating the data by inserting zeros after every word (up-sampling the data from 4 to 8 Hz). (B) 1-Hz response across the different layers of the parser model.

<https://doi.org/10.1371/journal.pcbi.1010269.g001>

modeling abstract sentence representations [4]. However, one can also create an ANN whose purpose is to annotate sentences with their linguistic constituents. As these models output how different syntactic structures organize a text, they definitionally contain explicit representations of sentences. One such model is the Berkeley Neural Network Parser, which is trained to parse sentences in syntactic units [11,12]. We investigated how the Berkeley parser behaves using the Ding and colleagues [1] stimuli (Fig 1; for details see [https://github.com/sannetenover/2022\\_SimLinguisticInferences](https://github.com/sannetenover/2022_SimLinguisticInferences)). What is evident is that just like in [4] or [2], the Berkeley parser also peaks at the sentence rate (Fig 1A). It does so for sentences, but not for word lists, or noun and verb phrases. Interestingly, the model also shows a peak when presented with ungrammatical sentences with uninterpretable semantics, but in which syntactic structure is preserved (viz., shuffling the words within the original sentences; also see [13]). What can also be appreciated is that sentence representations are strongest at the deepest, most abstract layer (Fig 1B). This is expected, as the ultimate goal of the parser is to explicitly represent the sentences at the highest hierarchical level. Since sentence input in the Ding and colleagues data occurs at 1-Hz, the output should contain the 1-Hz representation. While many commonly known models are not explicitly provided with syntactic instruction, viz., word2vec, it cannot be assumed that all ANN architectures definitionally are, as that depends on the input and instructions to the model: Parsers explicitly leverage a representation of syntax. What is more, the above simulation shows that both models—either explicitly instructed to parse the sentence or instructed with a non-grammatical next-word prediction goal—can show the 1-2-4-Hz pattern.

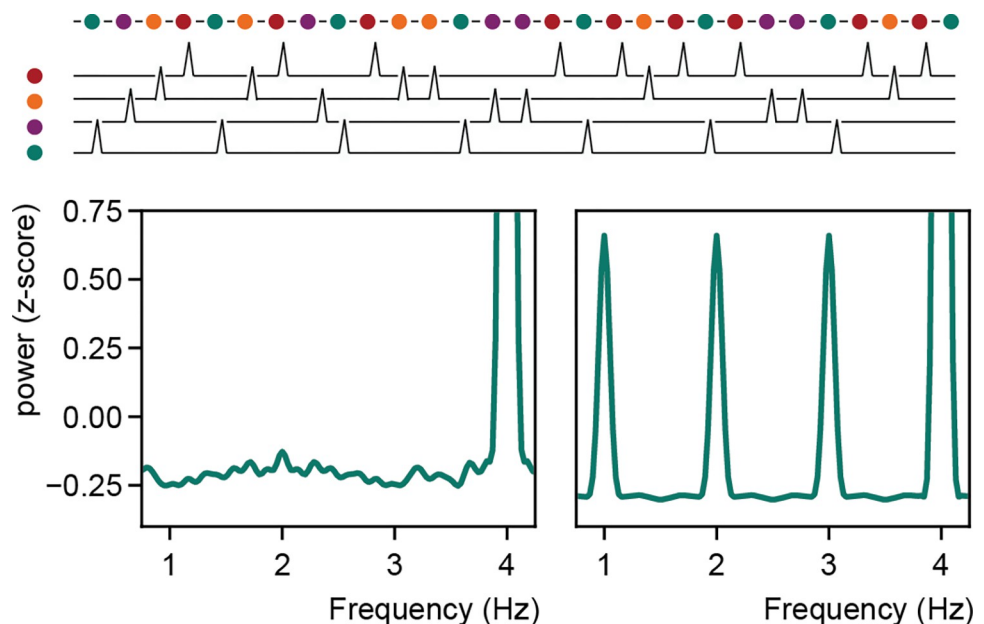
### The output shows a 1-2-4-Hz pattern, so the computation also has to

The 1-Hz peak in the Ding and colleagues [1] paper could imply that there is some integrative process happening that is slower than the individually presented words [14–17]. This interpretation is appealing as we have the inherent feeling that we integrate sequential words into a hierarchical representation that is the sentence, and we know that we arrive at integrated sentence meanings dictated by syntactic structure. However, nothing about the neural data shows

that this integration necessarily has to occur. Indeed, in similar visual steady-state evoked response studies random pictures are presented and every  $n$ th picture represents some specific class like faces [18–20]. If one is sensitive to the class, one finds a peak at the facial rate (Fig 2). One finds the peak as every  $n$ th stimulus something different happens. The computation does not necessarily happen at this rate, but the peak occurs at this rate as some event happens at that rate. Such a pattern has been shown for faces [18], emotions [21], single words [22], or arbitrarily created changes transitional probabilities [23]. That it is not necessary that the computation happens at the extracted rate is what Frank and Yang [4] argue in the case of sentences, and is also what happens for the output of the Berkeley parser when using ungrammatical sentences (Fig 1). In sum, one has to be cautious in interpreting effects at slow rates as an integrative process when a process simply repeating itself at a slower rate could generate the effect.

### Predictive value is always the ultimate goal in cognitive and neural science

Modern ANN methods have a strong focus on predicting data as best as possible. Intuitively, this seems a good starting point for scientific practice; however, there are serious issues with putting predictive value on a pedestal. While the goal for many ANN methods is to predict new data, in the cognitive and neural sciences, the goal is to explain the capacities of a system and to understand the mechanisms that lead to observed data. The only way to do this is to pose hypotheses about brain mechanisms, constrained through abductive inference based on what we know to an acceptable degree of verisimilitude about language and the brain, and test explicitly for these mechanisms [7,8,15,24]. Some forms of data science applications simply replace the thing we aim to understand (e.g., the brain, behavior, the capacity of a system for human language) with another black box (e.g., an ANN model; [25,26]). A simple classical cognitive box model can add to the level of understanding of brain operations, even if it does not do a great job at quantitatively predicting the data. To use an ANN to explain brain mechanisms,



**Fig 2. SSVEP studies do not assume any integration of responses.** When multiple category stimuli (here different colors) are presented with all a stereotypical response, a low frequency response will occur if any of them is presented at a specific rate (here green). Left: all items randomly presented. Right: single item is presented at a 1-Hz rate.

<https://doi.org/10.1371/journal.pcbi.1010269.g002>

at a minimum one must show that the operations and parameters of a complex model are related to brain operations [27–29]. This last step is often difficult. Thus, it often is still very valuable to work with models that provide explicit claims about the features that influence the model (e.g., encoding models such as regression).

### **When a model is trained on more data it is more realistic**

It is easy to be impressed by large language models that are trained on an incredible amount of data and have high predictive power. As cognitive scientists and neuroscientists, our aim is to explain brain computation, and we therefore must stay true to the data that the brain processes. It has often been argued that the brain cannot be like an ANN because it does not have access to so much training data [26,30]. Sometimes, learning a simple rule can account for some transferable skills the brain possesses, while contemporary language models and deep nets more broadly, with more data, training, energy use, and compute, have a harder time with task transfer [31–33]. However, this state of affairs does not exclude that for other tasks, the brain might employ a strong associative approach. When interpreting computational models, one must determine whether the brain realistically had access to, and has the capacity for processing the data put into a computational model to obtain a desired result. If the required data size and computational power are super- or inhuman, rule learning may be a form of computational compression that is more explanatory in the context of cognitive modeling.

### **Statistical models are “simpler” or “more parsimonious” than hierarchical models**

Frank and Yang [4] argue that representing the Ding and colleagues [1] data with only distributed lexical representations as in word2vec, compared to hierarchical representations [2]—which need both lexical and syntactic representations—is more parsimonious. However, word2vec is created by using a huge corpus of data with thousands of connections between many layers to squeeze the distributional information into a vector. Thus, even if word2vec by itself could be parsimonious, creating it is clearly not. Moreover, a word2vec representation is far from purely lexical as distributional patterns latently contain indirect syntactic information. Word2vec is created by forming associations with all possible neighboring words. Therefore, during the creation of a word2vec representation, the model uses sentence or at least neighboring word information (although not necessarily in a hierarchical manner). While indeed, creating the Ding and colleagues [1] pattern in neural readout cannot be said to require building a syntactic structure during sentence processing, word2vec can create the 1-Hz pattern precisely because it puts syntactic classes in neighboring positions in its vector representation. In any case, if one assumes the brain has a word2vec representation, it needed to construct it by integrating information across the sentence like the ANN does when creating the word2vec.

### **How can we conclude anything**

It is difficult to conclude whether statistical or hierarchical models can explain better the Ding and colleagues [1] data as both model classes can be made to recreate and predict the neural readout. Even though it is difficult to know whether a computational implementation is the one the brain uses [34], computational modeling still provides a means to be explicit about cognitive and neural processes at hand and thereby improve theory building [7]. The field will progress, but only if the outcomes of models and empirical data are interpreted, compared, and synthesized in a careful manner.

The starting point of every model should always be the explanandum: the process or capacity of a system that we aim to explain [24,35,36]. For cognitive neuroscientists, this is the brain operation that is, or that leads to or comprises, a cognitive operation. In trying to model sentence processing, deriving a processing mechanism that yields compositionality (viz., compositionality entails that the meaning of sentences is determined by that of individual words and the rules used to combine them [6,37]) seems absolutely crucial (or, alternatively providing a reason why compositionality is not necessary for sentence processing) for a theory of language representation and processing. The structured meaning that we experience from sentences is also the explanandum of cognitive models of language; if a model does not provide functionally equivalent “linguistic” output, then it does not account for the phenomenon, nor the capacity set out to be modeled. If one ignores these core principles of sentence processing, but goes directly toward trying to explain the neural data without context for abduction, then the model does not stay true to its epistemic goal.

If our goal is to explain neural computation, then similarly as above, only models that stay true to possible neurophysiological implementations of neural computation can ultimately be valid. One needs to have reasonable assumptions about how neuronal populations or ensembles could perform the computation; this should be explicitly stated, argued, and considered. If one models sentence processing as a hierarchical linguistic process, it is a reasonable starting point to investigate hierarchical, temporal [2,38,39], and anatomical organization in brain computation [40,41]. If one models sentence processing as a sequential process, one has to be explicit about the sequential steps taken in the brain. But one cannot simply assume that because a model fits the data, the brain computation is equivalent to what is instantiated in the model [8].

In an ideal world, a model of sentence processing should be able to explain not only the frequency-tagging data as generated by Ding and colleagues [1], but also stay faithful to the vast psycholinguistic and event-related brain potential (ERP) literature on how sentences are processed differently from word lists or unstructured acoustic input [42–44]. From this literature, we know that the brain is sensitive to the semantic and syntactic properties of sentences and discourses, indicating that the brain integrates words into phrases and sentences (and thus, does not treat words as independent units). This literature makes it unlikely that the only thing the brain does is create a word-by-word word2vec-like representation [4], but instead words need to be integrated (either sequentially or hierarchically). Putting too much emphasis on a single experimental finding can result in epistemic myopia; a simple way to safeguard against this is to require that theoretical and neurocomputational models explain a wide range of findings. More importantly, the ultimate goal of a cognitive model of language representation and processing is to explain the human capacity for natural language and language behavior. While experiments provide tailored designs that can discriminate between competing models, they are often intentionally far from natural circumstances in order to better orthogonalize factors. Exclusive focus on stimuli that are never encountered in daily life (such as in many experimental studies) and only become relevant in a lab setting can often obscure the original goal of the model, which is to explain a cognitive process or capacity, not only a specific dataset. It is therefore necessary that experimental approaches are complemented with studies using naturalistic paradigms [45,46] that more directly investigate the operations of interest, namely naturalistic spoken language comprehension.

Ultimately, every new discovery such as the Ding and colleagues [1] study leads to new questions. Models should be updated when new experimental data that can differentiate between proposed accounts (see, e.g., [13,47,48]) becomes available. Similarly, models must be explicit about the operations that hierarchical or sequential syntactic structure building draws upon and how they might operate in the brain in order to create predictions that can be

experimentally tested [15,39]. Only through systematic (theoretical) modeling one can generate model predictions which then can be experimentally tested, though prediction alone is not sufficient to explain how and why a behavior, capacity, or phenomenon is the way it is. Careful modeling can be followed by new experimental paradigms that put both principles and predictions to the test, which may lead to different types of evidence than the Ding and colleagues [1] study can provide by itself (see e.g., [45,46]). This type of data-to-model-to-data cycle catalyzes theoretical models, as well as theories of brain computation, and can disentangle different accounts of sentence processing [7,49]. For example, some studies have fit predictions from different computational models onto electroencephalographic (EEG) and functional magnetic resonance imaging (fMRI) data and have shown a better fit for hierarchical rather than sequential models for natural sentence processing [50,51]. These studies provide a first means to push the debate forward and highlight the value of comparing different computational models with each other to advance our understanding of the brain.

In sum, we believe that computational modeling is vital for understanding sentence processing through the lens of neural readouts (e.g., frequency-tagging, mutual information, phase synchronization, and connectivity signals). In order to be explanatory, models must compute linguistically-sufficient representations via explicit cognitive operations that are specified within a system architecture that stays faithful to neurophysiological principles. Furthermore, we should not be bewitched by the predictive power of our models, but instead should judge whether the proposed model is a likely model of brain function, and whether it explains more than just a single dataset. Only by generating explicit models we can differentiate, and ultimately integrate, sequential and hierarchical accounts, as well as arrive at modeling the neurocomputational principles on which these processes operate. There is no doubt that we need more modeling, as well as experimental, theoretical, and formal work, in order to reach these goals.

## References

1. Ding N, Melloni L, Zhang H, Tian X, Poeppel D. Cortical tracking of hierarchical linguistic structures in connected speech. *Nat Neurosci*. 2016; 19(1):158–64. <https://doi.org/10.1038/nn.4186> PMID: 26642090
2. Martin AE, Doumas LA. A mechanism for the cortical computation of hierarchical linguistic structure. *PLoS Biol*. 2017; 15(3):e2000663. <https://doi.org/10.1371/journal.pbio.2000663> PMID: 28253256
3. Kazanina N, Tavano A. What Neural Oscillations Can (not) Do for Syntactic Structure. *Building*. 2021.
4. Frank SL, Yang J. Lexical representation explains cortical entrainment during speech comprehension. *PLoS ONE*. 2018; 13(5):e0197304. <https://doi.org/10.1371/journal.pone.0197304> PMID: 29771964
5. Block NJ, Fodor JA. What Psychological States are Not. *Philos Rev*. 1972; 81(2):159–81.
6. Fodor JA, Pylyshyn ZW. Connectionism and cognitive architecture: A critical analysis. *Cognition*. 1988; 28(1–2):3–71. [https://doi.org/10.1016/0010-0277\(88\)90031-5](https://doi.org/10.1016/0010-0277(88)90031-5) PMID: 2450716
7. Guest O, Martin AE. How computational modeling can force theory building in psychological science. *Perspect Psychol Sci*. 2021. <https://doi.org/10.1177/1745691620970585> PMID: 33482070
8. Guest O, Martin AE. On logical inference over brains, behaviour, and artificial neural networks. *PsyArXiv [Preprint]*. 2021. <https://doi.org/10.31234/osf.io/tbmog>
9. Pylyshyn ZW. *Computation and cognition: Toward a foundation for cognitive science*. The MIT Press; 1986.
10. Doumas LA, Hummel JE, Sandhofer CM. A theory of the discovery and predication of relational concepts. *Psychol Rev*. 2008; 115(1):1. <https://doi.org/10.1037/0033-295X.115.1.1> PMID: 18211183
11. Kitaev N, Cao S, Klein D. Multilingual constituency parsing with self-attention and pre-training. *arXiv [Preprint]*. arXiv:181211760. 2018.
12. Kitaev N, Klein D. Constituency parsing with a self-attentive encoder. *arXiv [Preprint]*. arXiv:180501052. 2018.
13. Burroughs A, Kazanina N, Houghton C. Grammatical category and the neural processing of phrases. *Sci Rep*. 2021; 11(1):1–10.

14. Ghitza O. "Acoustic-driven oscillators as cortical pacemaker": a commentary on Meyer, Sun & Martin (2019). *Lang Cogn Neurosci*. 2020; 35(9):1100–5.
15. Martin AE. A compositional neural architecture for language. *J Cogn Neurosci*. 2020:1–20. [https://doi.org/10.1162/jocn\\_a\\_01552](https://doi.org/10.1162/jocn_a_01552) PMID: 32108553
16. Meyer L, Sun Y, Martin AE. Synchronous, but not entrained: Exogenous and endogenous cortical rhythms of speech and language processing. *Lang Cogn Neurosci*. 2019:1–11.
17. Rimmele JM, Poeppel D, Ghitza O. Acoustically Driven Cortical  $\delta$  Oscillations Underpin Prosodic Chunking. *eNeuro*. 2021; 8(4). <https://doi.org/10.1523/ENEURO.0562-20.2021> PMID: 34083380
18. Ales JM, Farzin F, Rossion B, Norcia AM. An objective method for measuring face detection thresholds using the sweep steady-state visual evoked response. *J Vis*. 2012; 12(10):18. <https://doi.org/10.1167/12.10.18> PMID: 23024355
19. Norcia AM, Appelbaum LG, Ales JM, Cottareau BR, Rossion B. The steady-state visual evoked potential in vision research: A review. *J Vis*. 2015; 15(6):4. <https://doi.org/10.1167/15.6.4> PMID: 26024451
20. Zoefel B, Ten Oever S, Sack AT. The involvement of endogenous neural oscillations in the processing of rhythmic input: More than a regular repetition of evoked neural responses. *Front Neurosci*. 2018; 12:95. <https://doi.org/10.3389/fnins.2018.00095> PMID: 29563860
21. Schettino A, Porcu E, Gundlach C, Keitel C, Müller MM. Rapid processing of neutral and angry expressions within ongoing facial stimulus streams: Is it all about isolated facial features? *PLoS ONE*. 2020; 15(4):e0231982. <https://doi.org/10.1371/journal.pone.0231982> PMID: 32330160
22. De Rosa M, Ktori M, Vidal Y, Bottini R, Crepaldi D. Frequency-based neural discrimination in fast periodic visual stimulation. *Cortex*. 2022; 148:193–203. <https://doi.org/10.1016/j.cortex.2022.01.005> PMID: 35180482
23. Henin S, Turk-Browne NB, Friedman D, Liu A, Dugan P, Flinker A, et al. Learning hierarchical sequence representations across human cortex and hippocampus. *Sci Adv*. 2021; 7(8):eabc4530. <https://doi.org/10.1126/sciadv.abc4530> PMID: 33608265
24. van Rooij I, Blokpoel M. Formalizing verbal theories. *Soc Psychol*. 2020.
25. Kay KN. Principles for models of neural information processing. *Neuroimage*. 2018; 180:101–9. <https://doi.org/10.1016/j.neuroimage.2017.08.016> PMID: 28793238
26. Marcus G. Deep learning: A critical appraisal. *arXiv [Preprint]*. arXiv:180100631. 2018.
27. Khaligh-Razavi S-M, Kriegeskorte N. Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS Comp Biol*. 2014; 10(11):e1003915.
28. Kubilius J, Bracci S, Op de Beeck HP. Deep neural networks as a computational model for human shape sensitivity. *PLoS Comp Biol*. 2016; 12(4):e1004896. <https://doi.org/10.1371/journal.pcbi.1004896> PMID: 27124699
29. Cichy RM, Kaiser D. Deep neural networks as scientific models. *Trends Cogn Sci*. 2019; 23(4):305–17. <https://doi.org/10.1016/j.tics.2019.01.009> PMID: 30795896
30. Dupoux E. Cognitive science in the era of artificial intelligence: A roadmap for reverse-engineering the infant language-learner. *Cognition*. 2018; 173:43–59. <https://doi.org/10.1016/j.cognition.2017.11.008> PMID: 29324240
31. George D, Lehrach W, Kansky K, Lázaro-Gredilla M, Laan C, Marthi B, et al. A generative vision model that trains with high data efficiency and breaks text-based CAPTCHAs. *Science*. 2017; 358(6368). <https://doi.org/10.1126/science.aag2612> PMID: 29074582
32. Lake B, Baroni M, editors. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. *International conference on machine learning*; 2018: PMLR.
33. Lappin S. *Deep learning and linguistic representation*. CRC Press; 2021.
34. Putnam H. Psychological predicates. *Art, mind, and religion*. 1967; 1:37–48.
35. van Rooij I, Baggio G. Theory before the test: How to build high-verisimilitude explanatory theories in psychological science. *Perspect Psychol Sci*. 2021; 16(4):682–97. <https://doi.org/10.1177/1745691620970604> PMID: 33404356
36. Cummins R. "How does it work?" versus "what are the laws?": Two conceptions of psychological explanation. *Explanation and cognition*. 2000:117–44.
37. Partee B. Montague grammar and transformational grammar. *Linguist Inq*. 1975:203–300.
38. Lakatos P, Shah AS, Knuth KH, Ulbert I, Karmos G, Schroeder CE. An oscillatory hierarchy controlling neuronal excitability and stimulus processing in the auditory cortex. *J Neurophysiol*. 2005; 94(3):1904–11. <https://doi.org/10.1152/jn.00263.2005> PMID: 15901760



39. Martin AE, Dumas LA. Predicate learning in neural systems: using oscillations to discover latent structure. *Curr Opin Behav Sci.* 2019; 29:77–83.
40. Felleman DJ, Van Essen DC. Distributed hierarchical processing in the primate cerebral cortex. *Cereb Cortex.* 1991; 1(1):1–47. <https://doi.org/10.1093/cercor/1.1.1-a> PMID: 1822724
41. Hickok G, Poeppel D. The cortical organization of speech processing. *Nat Rev Neurosci.* 2007; 8(5):393–402. <https://doi.org/10.1038/nrn2113> PMID: 17431404
42. Coulson S, King JW, Kutas M. Expect the unexpected: Event-related brain response to morphosyntactic violations. *Lang Cognit Process.* 1998; 13(1):21–58.
43. Hagoort P, Hald L, Bastiaansen M, Petersson KM. Integration of word meaning and world knowledge in language comprehension. *Science.* 2004; 304(5669):438–41. <https://doi.org/10.1126/science.1095455> PMID: 15031438
44. Kutas M, Federmeier KD. Electrophysiology reveals semantic memory use in language comprehension. *Trends Cogn Sci.* 2000; 4(12):463–70. [https://doi.org/10.1016/s1364-6613\(00\)01560-6](https://doi.org/10.1016/s1364-6613(00)01560-6) PMID: 11115760
45. Kaufeld G, Bosker HR, Ten Oever S, Alday PM, Meyer AS, Martin AE. Linguistic structure and meaning organize neural oscillations into a content-specific hierarchy. *J Neurosci.* 2020; 40(49):9467–75. <https://doi.org/10.1523/JNEUROSCI.0302-20.2020> PMID: 33097640
46. Keitel A, Gross J, Kayser C. Perceptually relevant speech tracking in auditory and motor cortex reflects distinct linguistic features. *PLoS Biol.* 2018; 16(3):e2004473. <https://doi.org/10.1371/journal.pbio.2004473> PMID: 29529019
47. Glushko A, Poeppel D, Steinhauer K. Overt and covert prosody are reflected in neurophysiological responses previously attributed to grammatical processing. *bioRxiv.* 2020.
48. Tavano A, Blohm S, Knoop CA, Muralikrishnan R, Scharinger M, Wagner V, et al. Neural harmonics of syntactic structure. *bioRxiv.* 2021:2020. 04. 08.031575.
49. Hale JT, Campanelli L, Li J, Bhattasali S, Pallier C, Brennan JR. Neurocomputational Models of Language Processing. *Annu Rev Linguist.* 2022; 8:427–46.
50. Brennan JR, Stabler EP, Van Wagenen SE, Luh W-M, Hale JT. Abstract linguistic structure correlates with temporal activity during naturalistic comprehension. *Brain Lang.* 2016; 157:81–94. <https://doi.org/10.1016/j.bandl.2016.04.008> PMID: 27208858
51. Brennan JR, Hale JT. Hierarchical structure guides rapid linguistic predictions during naturalistic listening. *PLoS ONE.* 2019; 14(1):e0207741. <https://doi.org/10.1371/journal.pone.0207741> PMID: 30650078