

Interpreting and improving natural-language processing (in machines) with natural language-processing (in the brain)

Mariya Toneva

Machine Learning Department
Center for the Neural Basis of Cognition
Carnegie Mellon University
Pittsburgh, PA 15213
mariya@cmu.edu

Leila Wehbe

Machine Learning Department
Center for the Neural Basis of Cognition
Carnegie Mellon University
Pittsburgh, PA 15213
lwehbe@cmu.edu

Abstract

Neural network models for NLP are typically implemented without the explicit encoding of language rules and yet they are able to break one performance record after another. Despite much work, it is still unclear what the representations learned by these networks correspond to. We propose here a novel approach for interpreting neural networks that relies on the only processing system we have that does understand language: the human brain. We use brain imaging recordings of subjects reading complex natural text to interpret word and sequence embeddings from 4 recent NLP models - ELMo, USE, BERT and Transformer-XL. We study how their representations differ across layer depth, context length, and attention type. Our results reveal differences in the context-related representations across these models. Further, in the transformer models, we find an interaction between layer depth and context length, and between layer depth and attention type. We finally use the insights from the attention experiments to alter BERT: we remove the learned attention at shallow layers, and show that this manipulation improves performance on a wide range of syntactic tasks. Cognitive neuroscientists have already begun using NLP networks to study the brain, and this work closes the loop to allow the interaction between NLP and cognitive neuroscience to be a true cross-pollination.

1 Introduction

The large success of deep neural networks in NLP is perplexing when considering that unlike most other NLP approaches, neural networks are typically not informed by explicit language rules. Yet, neural networks are constantly breaking records in various NLP tasks from machine translation to sentiment analysis. Even more interestingly, it has been shown that word embeddings and language models trained on a large generic corpus and then optimized for downstream NLP tasks produce even better results than training the entire model only to solve this one task (Peters *et al.*, 2018; Howard and Ruder, 2018; Devlin *et al.*, 2018). These models seem to capture something generic about language. What representations do these models capture of their language input?

Different approaches have been proposed to probe the representations in the network layers through NLP tasks designed to detect specific linguistic information (Conneau *et al.*, 2018; Zhu *et al.*,

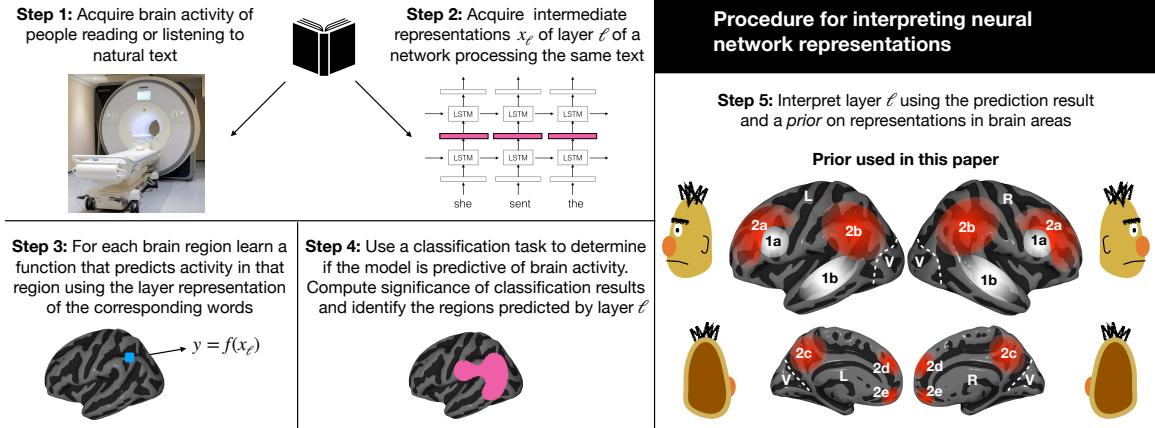


Figure 1: Diagram of approach and prior on brain function. The prior was constructed using the results of Lerner *et al.* (2011): regions in group 1 (white) process information related to isolated words and word sequences while group 2 (red) process only information related to word sequences (see 1.1). V indicates visual cortex. The drawing indicates the views of the brain with respect to the head. See Supplementary Table 2 for the names of the brain areas and Supplementary Section 7.3 for full description of the methods.

2018; Linzen *et al.*, 2016). Other approaches have attempted to offer a more theoretical assessment of how recurrent networks propagate information, or what word embeddings can represent (Peng *et al.*, 2018; Chen *et al.*, 2017; Weiss *et al.*, 2018). Most of this work has been centered around understanding the properties of sequential models such as LSTMs and RNNs, with considerably less work focused on non-sequential models such as transformers.

Using specific NLP tasks, word annotations or behavioral measures to detect if a type of information is present in a network-derived representation (such as a word embedding of an LSTM or a state vector of a transformer) can be informative. However, complex and arguably more interesting aspects of language, such as high level meaning, are difficult to capture in an NLP task or in behavioral measures. We therefore propose a novel approach for interpreting neural networks that relies on the only processing system we have that does understand language: the human brain. Indeed, the brain does represent complex linguistic information while processing language, and we can use brain activity recordings as a proxy for these representations. We can then relate the brain representations with neural network representations by learning a mapping from the latter to the former. We refer to this analysis as aligning the neural network representations with brain activity.

1.1 Proposed approach

We propose to look at brain activity of subjects reading naturalistic text as a source of additional information for interpreting neural networks. We use fMRI (functional Magnetic Resonance Imaging) and Magnetoencephalography (MEG) recordings of the brain activity of these subjects as they are presented text one word at a time. We present the same text to the NLP model we would like to investigate and extract representations from the intermediate layers of the network, given this text. We then learn an alignment between these extracted representations and the brain recordings corresponding to the same words to offer an evaluation of the information contained in the network representations.

Evaluating neural network representations with brain activity is a departure from existing studies that go the other way, using such an alignment to instead evaluate brain representations (Wehbe *et al.*, 2014a; Frank *et al.*, 2015; Hale *et al.*, 2018; Jain and Huth, 2018).

To align a layer ℓ representation with brain activity, we first learn a model that predicts the fMRI or MEG activity in every region of the brain (fig. 1). We determine the regions where this model is predictive of brain activity using a classification task followed by a significance test. If a layer representation can accurately predict the activity in a brain region r , then we conclude that the layer shares information with brain region r . We can thus make conclusions about the representation in layer ℓ based on our prior knowledge of region r .

Brain recordings have inherent, meaningful structure that is absent in network-derived representations. In the brain, different processes are assigned to specific locations as has been revealed by a large array of fMRI experiments. These processes have specific latencies and follow a certain order, which has been revealed by electrophysiology methods such as MEG. In contrast to the brain, a network-derived representation might encode information that is related to multiple of these processes without a specific organization. When we align that specific network representation with fMRI and MEG data, the result will be a decomposition of the representation into parts that correspond to different processes and should therefore be more interpretable. We can think of alignment with brain activity as a “demultiplexer” in which a single input (the network-derived representation) is decomposed into multiple outputs (relationship with different brain processes).

There doesn’t yet exist a unique theory of how the brain processes language that researchers agree upon (Hickok and Poeppel, 2007; Friederici, 2011; Hagoort, 2003). Because we don’t know which of the existing theories are correct, we abandon the theory based approach and adopt a fully data-driven approach. We focus on results from experiments that use naturalistic stimuli to derive our priors on the function of specific brain areas during language processing. These experiments have found that a set of regions in the temporo-parietal and frontal cortices are activated in language processing (Lerner *et al.*, 2011; Wehbe *et al.*, 2014b; Huth *et al.*, 2016; Blank and Fedorenko, 2017) and are collectively referred to as the language network (Fedorenko and Thompson-Schill, 2014). Using the results of Lerner *et al.* (2011) we subdivide this network into two groups of areas: group 1 is consistently activated across subjects when they listen to disconnected words or to complex fragments like sentences or paragraphs and group 2 is consistently activated only when they listen to complex fragments. We will use group 1 as our prior on brain areas that process information at the level of both short-range context (isolated words) and long-range context (multi-word composition), and group 2 as a prior on areas that process long-range context only. Fig. 1 shows a simple approximation of these areas on the Montreal Neurological Institute (MNI) template. Inspection of the results of Jain and Huth (2018) shows they corroborate the division of language areas into group 1 and group 2. Because our prior relies on experimental results and not theories of brain function, it is data-driven.

We use this setup to investigate a series of questions about the information represented in different layers of neural network models. We explore four recent models: ELMo, a language model by Peters *et al.* (2018), BERT, a transformer by Devlin *et al.* (2018), USE (Universal Sentence Encoder) a sentence encoder by Cer *et al.* (2018), and T-XL (Transformer-XL) a transformer that includes a recurrence mechanism by Dai *et al.* (2019). We investigate multiple questions about these networks. Is word-level specific information represented only at input layers? Does this differ across recurrent models, transformers and other sentence embedding methods? How many layers do we need to represent a specific length of context? Is attention affecting long range or short range context?

Intricacies As a disclaimer, we warn the reader that one should be careful while dealing with brain activity. Say a researcher runs a task T in fMRI (e.g. counting objects on the screen) and finds it activates region R , which is shown in another experiment to also be active during process P (e.g. internal speech). It is seductive to then infer that process P is involved during task T . This “reverse inference” can lead to erroneous conclusions, as region R can be involved in more than one task (Poldrack, 2006). To avoid this trap, we only interpret alignment between network-derived representations and brain regions if (1) the function of the region is well studied and we have some confidence on its function during a task similar to ours (e.g. the primary visual cortex processing letters on the screen or group 2 processing long range context) or (2) we show a brain region has overlap in the variance explained by the network-derived layer and by a specific process, in the same experiment. We further take sound measures for reporting results: we cross-validate our models and report results on unseen test sets, and present results for multiple participants. Another possible fallacy is to directly compare the performance of layers from different networks and conclude that one network performs better than the other: information is likely organized differently across networks and such comparisons are misleading. Instead we only perform controlled experiments where we look at one network and vary one parameter at a time, such as context length, layer depth or attention type.

1.2 Contributions

1. We present a new method to interpret network representations and a proof of concept for it.
2. We use our method to analyze and provide hypotheses about ELMo, BERT, USE and T-XL.
3. We find the middle layers of transformers are better at predicting brain activity than other layers. We find that T-XL’s performance doesn’t degrade as context is increased, unlike the other models’. We find that using uniform attention in early layers of BERT (removing the pretrained attention on the previous layer) leads to better prediction of brain activity.
4. We show that our brain-derived heuristic: using uniform attention in early layers of BERT (layers 1-6) leads to considerable improvement on benchmark NLP tasks (Marvin and Linzen, 2018) (up to 8% improvement in accuracy). This shows a transfer of knowledge from the brain to NLP tasks and validates our approach.

2 Related work on brains and language

Most work investigating language in the brain has been done in a controlled experiment setup where two conditions are contrasted (Friederici, 2011). These conditions typically vary in complexity (simple vs. complex sentences), vary in the presence or absence of a linguistic property (sentences vs. lists of words) or vary in the presence or absence of incongruities (e.g. semantic surprisal) (Friederici, 2011). A few researchers instead use naturalistic stimuli such as stories (Brennan *et al.*, 2010; Lerner *et al.*, 2011; Speer *et al.*, 2009; Wehbe *et al.*, 2014b; Huth *et al.*, 2016; Blank and Fedorenko, 2017). Some use predictive models of brain activity as a function of multi-dimensional features spaces describing the different properties of the stimulus (Wehbe *et al.*, 2014b; Huth *et al.*, 2016).

A few previous works have used neural network representations as a source of feature spaces to model brain activity. Wehbe *et al.* (2014b) aligned the MEG brain activity we use here with a

Recurrent Neural Network (RNN), trained on an online archive of Harry Potter Fan Fiction. The authors aligned brain activity with the context vector and the word embedding, allowing them to trace sentence comprehension at a word-by-word level. Jain and Huth (2018) aligned layers from a Long Short-Term Memory (LSTM) model to fMRI recordings of subjects listening to stories to differentiate between the amount of context maintained by each brain region. Other approaches rely on computing surprisal or cognitive load metrics using neural networks to identify processing effort in the brain, instead of aligning entire representations (Frank *et al.*, 2015; Hale *et al.*, 2018).

3 Approach

Network-derived Representations The approach we propose in this paper is general and can be applied to a wide variety of current NLP models. We present four case-studies of recent models that have very good performance on downstream tasks: ELMo, BERT, USE and T-XL.

- ELMo is a bidirectional language model that incorporates multiple layers of LSTMs. It can be used to derive contextualized embeddings by concatenating the LSTM output layers at that word with its non-contextualized embedding. We use a pretrained version of ELMo with 2 LSTM layers provided by Gardner *et al.* (2017).
- BERT is a bidirectional model of stacked transformers that is trained to predict whether a given sentence follows the current sentence, in addition to predicting a number of input words that have been masked (Devlin *et al.*, 2018). Upon release, this recent model achieved state of the art performance across a large array of NLP tasks, ranging from question answering to named entity recognition. We use a pretrained model provided by Hugging Face¹. We investigate the base BERT model, which has 12 layers, 12 attention heads, and 768 hidden units.
- USE is a method of encoding sentences into an embedding (Cer *et al.*, 2018) using a task similar to Skip-thought (Kiros *et al.*, 2015). USE is able to produce embeddings in the same space for single words and passages of text of different lengths. We use a version of USE from tensorflow hub trained with a deep averaging network² that has 512 dimensions.
- T-XL incorporates segment level recurrence into a transformer with the goal of capturing longer context than either recurrent networks or usual transformers (Dai *et al.*, 2019). We use a pretrained model provided by Hugging Face¹, with 19 layers and 1024 hidden units.

We investigate how the representations of all four networks change as we provide varying lengths of context. We compute the representations $x_{\ell,k}$ in each available intermediate layer ($\ell \in \{1, 2\}$ for ELMo; $\ell \in \{1,..12\}$ for BERT; ℓ is the output embedding for USE; $\ell \in \{1,..19\}$ for T-XL). We compute $x_{l,k}$ for word w_n by passing the most recent k words ($w_{n-k+1},..,w_n$) through the network.

fMRI and MEG data In this paper we use fMRI and MEG data which have complementary strengths. fMRI is sensitive to the change in oxygen level in the blood that is a consequence to neural activity, it has high spatial resolution (2-3mm) and low temporal resolution (multiple seconds). MEG measures the change in the magnetic field outside the skull due to neural activity, it has low spatial

¹<https://github.com/huggingface/pytorch-pretrained-BERT/>

²<https://tfhub.dev/google/universal-sentence-encoder/2>

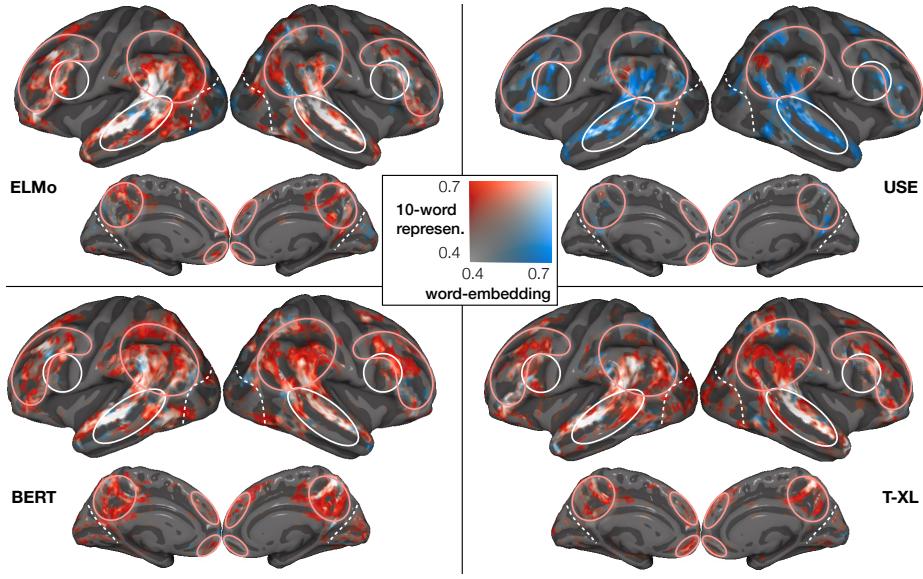


Figure 2: Comparison between the prediction performance of two network representations from each model: a 10-word representation corresponding to the 10 most recent words shown to the participant (Red) and a word-embedding corresponding to the last word (Blue). Areas in white are well predicted from both representations. These results align to a fair extent with our prior: group 2 areas (red outlines) are mostly predicted by the longer context representations while areas 1b (lower white outlines) are predicted by both word-embeddings and longer context representations.

resolution (multiple cm) and high temporal resolution (up to 1KHz). We use fMRI data published by Wehbe *et al.* (2014b). 8 subjects read chapter 9 of *Harry Potter and the Sorcerer’s stone* Rowling (2012) which was presented one word at a time for a fixed duration of 0.5 seconds each, and 45 minutes of data were recorded. The fMRI sampling rate (TR) was 2 seconds. The same chapter was shown by Wehbe *et al.* (2014a) to 3 subjects in MEG with the same rate of 0.5 seconds per word. Details about the data and preprocessing can be found in Supplementary Section 7.2.

Encoding models and comparisons For each type of network-derived representation $x_{\ell,k}$, we estimate an encoding model that takes $x_{\ell,k}$ as input and predicts the brain recording associated with reading the same words that were used to derive $x_{\ell,k}$. We estimate a function f , such that $f(x_{\ell,k}) = y$, where y is the brain activity recorded with either MEG or fMRI. We follow previous work (Sudre *et al.*, 2012; Wehbe *et al.*, 2014b,a; Nishimoto *et al.*, 2011; Huth *et al.*, 2016) and model f as a linear function, regularized by the ridge penalty. We test each encoding model by using it in a classification task on held-out data, in a four-fold cross-validation setting. The classification task is to predict which of two sets of words was being read based on the respective feature representations of these words (Mitchell *et al.*, 2008; Wehbe *et al.*, 2014b,a). The chance accuracy in this classification task is 0.5. The final fMRI results are reported on the MNI template, and we use pycortex to visualize them Gao *et al.* (2015). See Supplementary Section 7.3 for detailed methods descriptions.

Proof of concept Since MEG signals are faster than the rate of word presentation, they are more appropriate to study the components of word embeddings than the slow fMRI signals that cannot

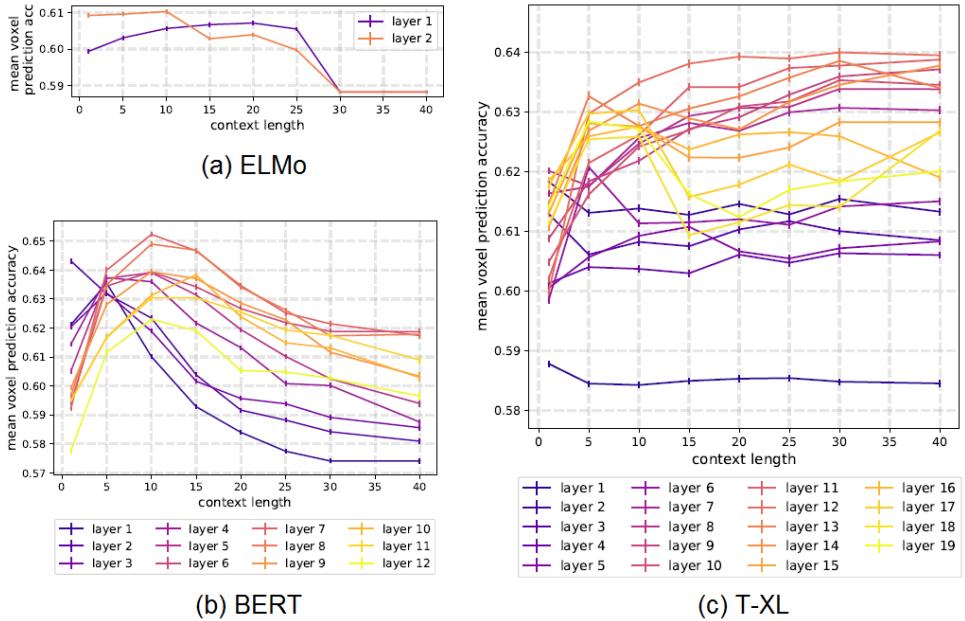


Figure 3: Performance of encoding models for all hidden layers in ELMo, BERT, and T-XL as the amount of context provided to the network is increased. Transformer-XL is the only model that continues to increase performance as the context length is increased. In all networks, the middle layers perform the best for contexts longer than 15 words. The deepest layers across all networks show a sharp increase in performance at short-range context (fewer than 10 words), followed by a decrease in performance.

be attributed to individual words. We know that a word embedding learned from a text corpus is likely to contain information related to the number of letters and part of speech of a word. We show in Supplementary Section 7.4 that the number of letters of a word and its ELMo embedding predict a shared portion of brain activity early on (starting 100ms after word onset) in the back of the MEG helmet, over the visual cortex. Indeed, this region and latency are when we expect the visual information related to a word to be processed (Sudre *et al.*, 2012). Further, a word’s part of speech and its ELMo embedding predict a shared portion of brain activity around 200ms after word onset in the left front of the MEG sensor. Indeed, we know from electrophysiology studies that part of speech violations incur a response around 200ms after word onset in the frontal lobe (Frank *et al.*, 2015). We conclude from these experiments that the ELMo embedding contains information about the number of letters and the part of speech of a word. Since we knew this from the onset, this experiment serves as a proof of concept for using our approach to interpret information in network representations.

4 Interpreting long-range contextual representations

Integrated contextual information in ELMo, BERT, and T-XL One question of interest in NLP is how successfully a model is able to integrate context into its representations. We investigate whether the four NLP models we consider are able to create an integrated representation of a text sequence by comparing the performance of encoding models trained with two kinds of representations: a

token-level word-embedding corresponding to the most recent word token a participant was shown and a 10-word representation corresponding to the 10 most recent words. For each of the models with multiple layers (all but USE), this 10-word representation was derived from a middle layer in the network (layer 1 in ELMo, layer 7 in BERT, and layer 11 in T-XL). We present the comparisons across the four models in figure 2, where only significantly predicted voxels for each of the 8 subjects were included with the false discovery rate controlled at level 0.05 (see Supplementary Section 7.3.1 for more details). We observe similarities in the word-embedding performances across all models, which all predict the brain activity in the left and right group 1b regions and to some extent in group 1a regions. We also observe differences in the longer context representations between USE and the rest of the models:

- ELMo, BERT, and T-XL long context representations predict subsets of both group 1 regions and group 2 regions. Most parts that are predicted by the word-embedding are also predicted by the long context representations (almost no blue voxels). We conclude that the long context representations most probably include information about both the long range context and the very recent word embeddings. These results may be due to the fact that all these models are at least partially trained to predict a word at a given position. They must encode long range information and also local information that can predict the appropriate word.
- USE long context representations predict the activity in a much smaller subset of group 2 regions. The low performance of the USE vectors might be due to the deep averaging which might be composing words in a crude manner. The low performance in predicting group 1 regions is most probably because USE computes representations at a sentence level and does not have the option of retaining recent information like the other models. USE long context representations therefore only have long range information.

Relationship between layer depth and context length We investigate how the performances of ELMo, BERT, and T-XL change at different layers as they are provided varying size of contexts. The results are shown in figure 3. We observe that in all networks, the middle layers perform the best for contexts longer than 15 words. In addition, the deepest layers across all networks show a sharp increase in performance at short-range context (fewer than 10 words), followed by a decrease in performance. We further observe that T-XL is the only model that continues to increase performance as the context length is increased. T-XL was designed to represent long range information better than a usual transformer and our results suggest that it does. Finally, we observe that layer 1 in BERT behaves differently from the first layers in the other two networks. In figure 4, we show that when we instead examine the increase in performance of all subsequent layers from the performance of the first layer, the resulting context-layer relationships resemble the ones in T-XL. This suggests that BERT layer 1 combines the information from the token-level embeddings in a way that limits the retention of longer context information in the layer 1 representations.

Effect of attention on layer representation We further investigate the effect of attention across different layers by measuring the negative impact that removing its learned attention has on its brain prediction performance. Specifically we replaced the learned attention with uniform attention over the representations from the previous layer. More concretely, to alter the attention pattern at a single layer in BERT, for each attention head $h_i = \text{Attn}_i(QW_i^Q, KW_i^K, VW_i^V)$, we replace the pretrained

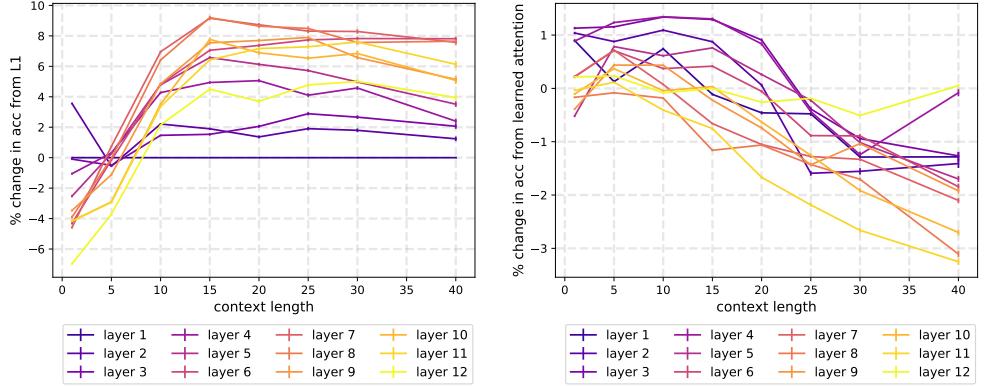


Figure 4: Change in encoding model performance of BERT layers from the performance of the first layer. When we adjust for the performance of the first layer, the performance of the remaining layers resemble that of T-XL more closely, as shown in Figure 3.

Figure 5: Change in encoding model performance of BERT layer l when the attention in layer l is made uniform. The performance of deep layers, other than the output layer, is harmed by the change in attention. Shallow layers benefit from the uniform attention for context lengths up to 25 words.

parameter matrices W_i^Q , W_i^K , and W_i^V for this layer, such that the attention $Attn(Q, K, V)$, defined as $\text{softmax}(QK/\sqrt{d_k})^T V$ (Vaswani *et al.*, 2017), yields equal probability over the values in value matrix V (here d_k denotes the dimensionality of the keys and queries). To this end, for a single layer, we replace W_i^Q and W_i^K with zero-filled matrices and W_i^V with the identity matrix. We only alter a single layer at a time, while keeping all other parameters of the pretrained BERT fixed. In figure 5, we present the change in performance of each layer with uniform attention when compared to pretrained attention. The performance of deep layers, other than the output layer, is harmed by the change in attention. However, surprisingly and against our expectations, shallow layers benefit from the uniform attention for context lengths up to 25 words.

5 Applying insight from brain interpretations to NLP tasks

After observing that the layers in the first half of the pretrained BERT model benefit from uniform attention for predicting brain activity, we test how the same alterations affect the performance of BERT on natural language processing tasks. We evaluate on tasks that do not require fine-tuning beyond the available pretrained models to ensure that there is an opportunity to transfer the insight from the brain interpretations of the pretrained BERT model. To this end, we evaluate on a range of syntactic tasks proposed by Marvin and Linzen (2018), that have been previously used to quantify BERT’s syntactic capabilities (Goldberg, 2019). We adopt the evaluation protocol of Goldberg (2019), in which BERT is first fed a complete sentence where the single focus verb is masked (e.g. [CLS] the game that the guard hates [MASK] bad .), then the prediction for the masked position is obtained using the pretrained language-modeling head, and lastly the accuracy is obtained by comparing the scores for the original correct verb (e.g. is) to the score for the incorrect verb (e.g. are). We make the attention in layers 1 through 6 in BERT uniform, a single layer at a

condition	uni L1	uni L2	uni L6	uni L11	base	count
simple	1.00	1.00	1.00	0.98	1.00	120
in a sentential complement	0.83	0.83	0.83	0.83	0.83	1440
short VP coordination	0.88	0.90	0.91	0.88	0.89	720
long VP coordination	0.96	0.97	1.00	0.96	0.98	400
across a prepositional phrase	0.86	0.93	0.88	0.82	0.85	19440
across a subject relative clause	0.83	0.83	0.85	0.83	0.84	9600
across an object relative clause	0.87	0.91	0.92	0.86	0.89	19680
across an obj. relative clause (no that)	0.87	0.80	0.87	0.84	0.86	19680
in an obj. relative clause	0.97	0.95	0.91	0.93	0.95	15960
in an object relative clause (no that)	0.83	0.72	0.74	0.72	0.79	15960
reflexive anaphora: simple	0.91	0.94	0.99	0.95	0.94	280
reflexive anaphora: in a sent. complem.	0.88	0.85	0.86	0.85	0.89	3360
reflexive anaphora: across rel. clause	0.79	0.84	0.79	0.76	0.80	22400

Table 1: Comparison of performance of pretrained BERT (‘base’) against models with uniform attention in layers 1, 2, 6, and 11 on syntactic tasks by Marvin and Linzen (2018).

time while keeping the remaining parameters fixed as described in Section 4, and evaluate on the 13 tasks. We present the results of altering layers 1,2, and 6 in Table 1. We observe that the altered models outperform the pretrained model (‘base’) in 10 of the 13 tasks and achieve parity in two of the remaining 3 tasks. Performance of altering layers 3-5 is similar and is presented in Supplementary Table 3. We contrast the performance of these layers with that of a model with uniform attention at layer 11, which is the model that suffers the most from this change for predicting the brain activity as shown in Figure 5. We observe that this model also performs poorly on the NLP tasks as it performs on par or worse than the base model in 12 of the 13 tasks.

6 Discussion

We introduced an approach to use brain activity recordings of subjects reading naturalistic text to interpret different representations derived from neural networks. We used MEG to show that the (non-contextualized) word embedding of ELMo contains information about word length and part of speech as a proof of concept. We used fMRI to show that different network representation (for ELMo, USE, BERT, T-XL) encode information relevant to language processing at different context lengths. USE long-range context representations perform differently from the other models and do not also include short-range information. The transformer models (BERT and T-XL) both capture the most brain-relevant context information in their middle layers. T-XL, by combining both recurrent properties and transformer properties, has representations that don’t degrade in performance when very long context is used, unlike pure recurrent models or transformers like ELMo or BERT.

We found that uniform attention on the previous layer actually improved the brain prediction performance of the shallow layers (layers 1-6) over using learned attention. After this observation, we modified attention at shallow layers of BERT, ran a set of NLP tasks and saw that the modified networks do perform better on a vast majority of tasks. We were thus able to extract from brain data a piece of knowledge about a neural network that can be materialized as a heuristic and that leads to better performance on NLP tasks: “uniform attention is better for shallow layers”.

Future work We hope that as naturalistic brain experiments become more popular and data more widely shared, aligning brain activity with neural networks will become a research area. Our next steps are to expand the analysis using MEG to uncover new aspects of word-embeddings and to derive more informative fMRI brain priors that contain specific conceptual information that is linked to brain areas, and use them to study the high level semantic information in network representations.

References

- Arias-Castro, E., Chen, S., *et al.* (2017). Distribution-free multiple testing. *Electronic Journal of Statistics*, **11**(1), 1983–2001.
- Barber, R. F., Candès, E. J., *et al.* (2015). Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, **43**(5), 2055–2085.
- Blank, I. and Fedorenko, E. (2017). Domain-general brain regions do not track linguistic input as closely as language-selective regions. *Journal of Neuroscience*, pages 3642–16.
- Brennan, J., Nir, Y., Hasson, U., Malach, R., Heeger, D., and Pylkkänen, L. (2010). Syntactic structure building in the anterior temporal lobe during natural story listening. *Brain and language*.
- Cer, D., Yang, Y., Kong, S.-y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-Cespedes, M., Yuan, S., Tar, C., *et al.* (2018). Universal sentence encoder. *arXiv preprint arXiv:1803.11175*.
- Chen, Y., Gilroy, S., Knight, K., and May, J. (2017). Recurrent neural networks as weighted language recognizers. *arXiv preprint arXiv:1711.05408*.
- Conneau, A., Kruszewski, G., Lample, G., Barrault, L., and Baroni, M. (2018). What you can cram into a single vector: Probing sentence embeddings for linguistic properties. *arXiv preprint arXiv:1805.01070*.
- Dai, Z., Yang, Z., Yang, Y., Cohen, W. W., Carbonell, J., Le, Q. V., and Salakhutdinov, R. (2019). Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Fedorenko, E. and Thompson-Schill, S. L. (2014). Reworking the language network. *Trends in cognitive sciences*, **18**(3), 120–126.
- Fischl, B. (2012). Freesurfer. *Neuroimage*, **62**(2), 774–781.
- Frank, S. L., Otten, L. J., Galli, G., and Vigliocco, G. (2015). The erp response to the amount of information conveyed by words in sentences. *Brain and language*, **140**, 1–11.
- Friederici, A. D. (2011). The brain basis of language processing: from structure to function. *Physiological reviews*, **91**(4), 1357–1392.
- Gao, J. S., Huth, A. G., Lescroart, M. D., and Gallant, J. L. (2015). Pycortex: an interactive surface visualizer for fmri. *Frontiers in neuroinformatics*, **9**, 23.
- Gardner, M., Grus, J., Neumann, M., Tafjord, O., Dasigi, P., Liu, N. F., Peters, M., Schmitz, M., and Zettlemoyer, L. S. (2017). AllenNLP: A deep semantic natural language processing platform.
- Goldberg, Y. (2019). Assessing bert’s syntactic abilities. *arXiv preprint arXiv:1901.05287*.

- Hagoort, P. (2003). How the brain solves the binding problem for language: a neurocomputational model of syntactic processing. *Neuroimage*, **20**, S18–S29.
- Hale, J., Dyer, C., Kuncoro, A., and Brennan, J. R. (2018). Finding syntax in human encephalography with beam search. *arXiv preprint arXiv:1806.04127*.
- Hickok, G. and Poeppel, D. (2007). The cortical organization of speech processing. *Nature Reviews Neuroscience*, **8**(5), 393–402.
- Howard, J. and Ruder, S. (2018). Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 328–339.
- Huth, A. G., de Heer, W. A., Griffiths, T. L., Theunissen, F. E., and Gallant, J. L. (2016). Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature*, **532**(7600), 453–458.
- Jain, S. and Huth, A. (2018). Incorporating context into language encoding models for fmri. *bioRxiv*, page 327601.
- Kay, K. N., Naselaris, T., Prenger, R. J., and Gallant, J. L. (2008). Identifying natural images from human brain activity. *Nature*, **452**(7185), 352.
- Kiros, R., Zhu, Y., Salakhutdinov, R. R., Zemel, R., Urtasun, R., Torralba, A., and Fidler, S. (2015). Skip-thought vectors. In *Advances in neural information processing systems*, pages 3294–3302.
- Kriegeskorte, N., Goebel, R., and Bandettini, P. (2006). Information-based functional brain mapping. *Proceedings of the National Academy of Sciences*, **103**(10), 3863–3868.
- Lerner, Y., Honey, C. J., Silbert, L. J., and Hasson, U. (2011). Topographic mapping of a hierarchy of temporal receptive windows using a narrated story. *The Journal of Neuroscience*, **31**(8), 2906–2915.
- Linzen, T., Dupoux, E., and Goldberg, Y. (2016). Assessing the ability of lstms to learn syntax-sensitive dependencies. *arXiv preprint arXiv:1611.01368*.
- Marvin, R. and Linzen, T. (2018). Targeted syntactic evaluation of language models. *arXiv preprint arXiv:1808.09031*.
- Mitchell, T. M., Shinkareva, S. V., Carlson, A., Chang, K.-M., Malave, V. L., Mason, R. A., and Just, M. A. (2008). Predicting human brain activity associated with the meanings of nouns. *science*, **320**(5880), 1191–1195.
- Nishimoto, S., Vu, A., Naselaris, T., Benjamini, Y., Yu, B., and Gallant, J. (2011). Reconstructing visual experiences from brain activity evoked by natural movies. *Current Biology*.
- Peng, H., Schwartz, R., Thomson, S., and Smith, N. A. (2018). Rational recurrences. *arXiv preprint arXiv:1808.09357*.
- Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*.
- Poldrack, R. A. (2006). Can cognitive processes be inferred from neuroimaging data? *Trends in cognitive sciences*, **10**(2), 59–63.
- Rabinovich, M., Ramdas, A., Jordan, M. I., and Wainwright, M. J. (2017). Optimal rates and tradeoffs in multiple testing. *arXiv preprint arXiv:1705.05391*.

- Rowling, J. (2012). *Harry Potter and the Sorcerer's Stone*. Harry Potter US. Pottermore Limited.
- Speer, N., Reynolds, J., Swallow, K., and Zacks, J. (2009). Reading stories activates neural representations of visual and motor experiences. *Psychological Science*, **20**(8), 989–999.
- Sudre, G., Pomerleau, D., Palatucci, M., Wehbe, L., Fyshe, A., Salmelin, R., and Mitchell, T. (2012). Tracking neural coding of perceptual and semantic features of concrete nouns. *NeuroImage*, **62**(1), 451–463.
- Taulu, S. and Simola, J. (2006). Spatiotemporal signal space separation method for rejecting nearby interference in meg measurements. *Physics in Medicine & Biology*, **51**(7), 1759.
- Taulu, S., Kajola, M., and Simola, J. (2004). Suppression of interference and artifacts by the signal space separation method. *Brain topography*, **16**(4), 269–275.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Wehbe, L., Vaswani, A., Knight, K., and Mitchell, T. M. (2014a). Aligning context-based statistical models of language with brain activity during reading. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 233–243, Doha, Qatar. Association for Computational Linguistics.
- Wehbe, L., Murphy, B., Talukdar, P., Fyshe, A., Ramdas, A., and Mitchell, T. M. (2014b). Simultaneously uncovering the patterns of brain regions involved in different story reading subprocesses. *PLOS ONE*, **9**(11): e112575.
- Wehbe, L., Ramdas, A., Steorts, R. C., and Shalizi, C. R. (2015). Regularized brain reading with shrinkage and smoothing. *Annals of Applied Statistics*, **9**(4), 1997–2022.
- Weiss, G., Goldberg, Y., and Yahav, E. (2018). On the practical computational power of finite precision rnns for language recognition. *arXiv preprint arXiv:1805.04908*.
- Zhu, X., Li, T., and de Melo, G. (2018). Exploring semantic properties of sentence embeddings. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 632–637.

7 Supplementary Materials

7.1 Brain areas included in prior

1a	Inferior Frontal Gyrus
1b	Middle/Superior Temporal
2a	Lateral Middle/Superior Frontal
2b	Supramarginal Gyrus / Posterior Superior Temporal / Angular Gyrus
2c	Precuneus
2d	Medial Superior Frontal
2e	Medial Orbito-Frontal

Table 2: Name of regions of interest in fig. 1 of main manuscript. Regions were approximated from the results of (Lerner *et al.*, 2011).

7.2 Data Preprocessing

We use fMRI data of 8 subjects reading chapter 9 of *Harry Potter and the Sorcerer’s Stone* (Rowling, 2012), collected and made available online by Wehbe *et al.* (2014b)³. Words were presented one at a time at a rate of 0.5s each. fMRI data was acquired at a rate of 2s per image, i.e. the repetition time (TR) is 2s. The images were comprised of $3 \times 3 \times 3\text{mm}$ voxels. The data for each subject was slice-time and motion corrected using SPM8 (Kay *et al.*, 2008), then detrended and smoothed with a 3mm full-width-half-max kernel. The brain surface of each subject was reconstructed using Freesurfer (Fischl, 2012), and a grey matter mask was obtained. The Pycortex software (Gao *et al.*, 2015) was used to handle and plot the data. For each subject, 25000-31000 cortical voxels were kept.

The same paradigm was recorded for 3 subjects using MEG by the authors of Wehbe *et al.* (2014a) and shared upon our request. This data was recorded at 306 sensors organized in 102 locations around the head. MEG records the change in magnetic field due to neuronal activity and the data we used was sampled at 1kHz, then preprocessed using the Signal Space Separation method (SSS) (Taulu *et al.*, 2004) and its temporal extension (tSSS) (Taulu and Simola, 2006). The signal in every sensor was downsampled into 25ms non-overlapping time bins. For each of the 5176 word in the chapter, we therefore obtained a recording for 306 sensors at 20 time points after word onset (since each word was presented for 500ms).

7.3 Encoding Models

7.3.1 fMRI

Ridge regularization is used to estimate the parameters of a linear model that predicts the brain activity y^i in every fMRI voxel i as a linear combination of a particular layer representation x^ℓ . For each output dimension (voxel), the Ridge regularization parameter is chosen independently by nested cross-validation. We use Ridge regression because of its computational efficiency and because of the results of Wehbe *et al.* (2015) showing that for fMRI data, as long as proper regularization is used and the regularization parameter is chosen by cross-validation for each voxel independently, different regularization techniques lead to similar results. Indeed, Ridge regression is indeed a common regularization technique used for building predictive fMRI (Mitchell *et al.*, 2008; Nishimoto *et al.*, 2011; Wehbe *et al.*, 2014b; Huth *et al.*, 2016).

³<http://www.cs.cmu.edu/afs/cs/project/theo-73/www/plosone/>

For every voxel i , a model is fit to predict the signals $y^i = [y_1^i, y_2^i, \dots, y_n^i]$, where n is the number of time points, as a function of the representation derived from layer ℓ of a network. The words presented to the participants are first grouped by the TR interval in which they were presented. Then, the features of layer ℓ of the words in every group are averaged to form a sequence of features $x^\ell = [x_1^\ell, x_2^\ell, \dots, x_n^\ell]$ which are aligned with the brain signals. The models are trained to predict the signal at time t , y_t , using the concatenated vector z_t^ℓ formed of $[x_{t-1}^\ell, x_{t-2}^\ell, x_{t-3}^\ell, x_{t-4}^\ell]$. The features of the words presented in the previous volumes are included in order to account for the lag in the hemodynamic response that fMRI records. Indeed, the response measured by fMRI is an indirect consequence of brain activity that peaks about 6 seconds after stimulus onset, and the solution of expressing brain activity as a function of the features of the preceding time points is a common solution for building predictive models (Nishimoto *et al.*, 2011; Wehbe *et al.*, 2014b; Huth *et al.*, 2016).

For each given subject and each layer ℓ , we perform a cross-validation procedure to estimate how predictive that layer is of brain activity in each voxel i . For each fold:

- The fMRI data Y and feature matrix $Z^\ell = z_1^\ell, z_2^\ell, \dots, z_n^\ell$ are split into corresponding train and validation matrices and these matrices are individually normalized (to get a mean of 0 and standard deviation of 1 for each voxel across time), ending with train matrices Y^R and $Z^{R,\ell}$ and validation matrices Y^V and $Z^{V,\ell}$.
- Using the train fold, a model $w^{i,\ell}$ is estimated as:

$$\arg \min_{w^{i,\ell}} ||y^{R,i} - Z^{R,\ell} w^{i,\ell}||_2^2 + \lambda^i ||w^{i,\ell}||_2^2$$

A ten-fold nested cross-validation procedure is first used to identify the best λ^i for every voxel i that minimizes nested cross-validation error. $w^{i,\ell}$ is then estimated using λ^i on the entire training fold.

- The predictions for each voxel on the validation fold are obtained as $p^\ell = Z^{V,\ell} w^{i,\ell}$.
- A classification task is then performed to assess the prediction performance of the learned model. This classification task is based on searchlight classification (Kriegeskorte *et al.*, 2006), in which a sliding window groups each voxel with its immediate neighbors in the 3D grid of voxels. We perform a more accurate searchlight analysis we refer to as cortical-searchlight. We are interested only in the grey matter voxels (which contain neurons) and these comprise the most external part of the brain: the cortical sheet. The cortical sheet of each subject is highly folded, and voxels that lie in a neighborhood on the sheet are not necessarily neighbors in the 3D grid of voxels. Using the reconstructed cortical sheet of each subject, we estimate for each cortical voxel a surrounding neighborhood by including the voxels adjacent to it on the cortical sheet, and the voxels adjacent to those voxels. See figure 6. We use for each voxel i this neighborhood of voxels N^i with $|N^i| = k^i$ in a classification task.
- For each voxel i , we use the signals predicted for layer ℓ to classify a contiguous chunk of real data of length 20TRs. Since fMRI data is noisy, performance using a single TR will be close to chance accuracy and will therefore have low power and will not be informative for our purpose. Indeed, for this reason most experiments using predictive fMRI models test them on a part of the experiment that is repeated multiple times (Kay *et al.*, 2008; Nishimoto *et al.*, 2011; Huth *et al.*, 2016). These repetitions are then averaged into one test set which is predicted, and this less noisy average leads to better prediction accuracy. The experiment we are using however doesn't have any repetitions and not specific test set, and therefore by raise the number of TRs and classify 20TRs at a time, we are able to improve the classification accuracy. Wehbe *et al.* (2014b) have shown that classification accuracy reaches a plateau after around 15 TRs and we pick 20TRs for good measure. The classification task takes an unlabeled chunk of real data of size $20 \times k$ and two possible predicted data chunks of the same size, one being the predicted data corresponding to the same time, and another randomly chosen chunk. Euclidean distance is computed between the real chunk and the two predicted chunks, and the closest chunk is chosen. This is repeated a large number of times and average accuracy is computed at each voxel.

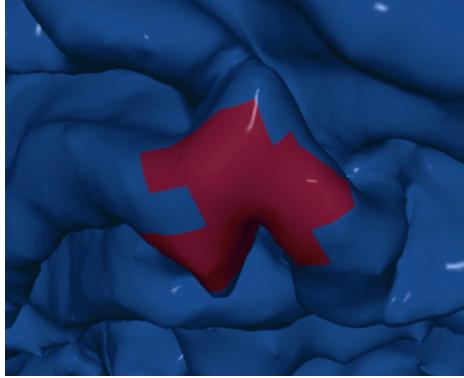


Figure 6: Example neighborhood estimated using the cortical sheet and not the 3D grid of voxels.

The above steps are repeated for each of the four cross-validation folds and average accuracy is obtained for each voxel i for layer ℓ , for each subject.

We use a new empirical based method to compute statistical significance that relies on the distribution of average accuracies over a subject’s brain to estimate the False Discovery Proportion (FDP). The voxel accuracies belong to two distributions: either the voxel has chance accuracy or the voxel is truly predicted by the corresponding layer ℓ . Average chance accuracy for our binary balanced task is 0.5, however the accuracies due to chance performance might have a varying distribution around 0.5. The accuracies above 0.5 are a mixture of predicted voxels and voxels with chance performance. We assume that chance performance is symmetrically distributed around 0.5, and we use the set of accuracies that are less than 0.5—which we consider to be in the chance distribution—to estimate the distribution of chance accuracies above 0.5. We want to find a set of voxels where to reject the null hypothesis such that the FDP is ≤ 0.05 . For that purpose we find the smallest margin δ , $0 < \delta < 0.5$ such that:

$$\widehat{\text{FDP}} = \frac{1 + \#\{\text{voxel } s.t. \text{ accuracy} \leq 0.5 - \delta\}}{1 \vee \#\{\text{voxel } s.t. \text{ accuracy} \geq 0.5 + \delta\}} \leq q$$

where $q = 0.05$, by starting at $\delta = 0.001$ and increasing it in increments of 0.001, stopping when $\widehat{\text{FDP}} \leq 0.05$ or the limit is reached. This approach is adapted from the Barber-Candès approach which has been proposed and analyzed by Barber *et al.* (2015); Arias-Castro *et al.* (2017); Rabinovich *et al.* (2017), and shown to control the False Discovery Rate (FDR) at level q when δ_{final} is chosen as a threshold. We reject the null hypothesis for all voxels where the accuracy is $\geq 0.5 + \delta_{\text{final}}$.

To combine results across different subjects, we use pycortex (Gao *et al.*, 2015) to transform each subject to the Montreal Neurological Institute (MNI) space, the most commonly used template space in fMRI. We can then average the results of different participants.

7.3.2 MEG

MEG data is sampled faster than the rate of word presentation, so for each word, we have 20 times points recorded at 306 sensors. Ridge regularization is similarly used to estimate the parameters of a linear model that predicts the brain activity $y^{i,\tau}$ in every MEG sensor i at time τ after word onset. For each output dimension (sensor/time tuple i, τ), the Ridge regularization parameter is chosen independently by nested cross-validation.

For every sensor/time tuple i, τ , a model is fit to predict the signals $y^{i,\tau} = [y_1^{i,\tau}, y_2^{i,\tau}, \dots, y_n^{i,\tau}]$, where n is the number of words in the story, as a function of the representation derived from layer ℓ of a network. We use as input the word vector x^ℓ without the delays we used in fMRI because the MEG recordings capture instantaneous consequences of brain activity (change in the magnetic field). The models are trained to predict the signal at word t , $y_t^{i,\tau}$, using the vector x_t^ℓ .

For each each given subject and each layer ℓ , we perform a cross-validation procedure to estimate how predictive that layer is of brain activity in each voxel i . For each fold:

- The MEG data Y and feature matrix $X^\ell = x_1^\ell, x_2^\ell, \dots, x_n^\ell$ are split into corresponding train and validation matrices and these matrices are individually normalized (to get a mean of 0 and standard deviation of 1 for each voxel across time), ending with train matrices Y^R and $X^{R,\ell}$ and validation matrices Y^V and $Z^{V,\ell}$.
- Using the train fold, a model $w^{(i,\tau)\ell}$ is estimated as:

$$\arg \min_{w^{(i,\tau)\ell}} \|y^{(i,\tau),R} - X^{R,\ell} w^{(i,\tau)\ell}\|_2^2 + \lambda^{(i,\tau)} \|w^{(i,\tau)\ell}\|_2^2$$

A ten-fold nested cross-validation procedure is first used to identify the best $\lambda^{(i,\tau)}$ for every sensor, time-point tuple (i, τ) that minimizes nested cross-validation error. $w^{(i,\tau)\ell}$ is then estimated using $\lambda^{(i,\tau)}$ on the entire training fold.

- The predictions for each sensor, time-point tuple (i, τ) on the validation fold are obtained as $p^\ell = X^{V,\ell} w^{(i,\tau)\ell}$.
- A classification task is then performed to assess the prediction performance of the learned model. This classification task also pools spatially: we use the 3 sensors at each location, pooling across all the subjects, ending up with 102 classifications at 20 time-points. By pooling the data in each sensor location across subjects, we increase the signal-to-noise ratio.
- For each sensor location s and time-point τ , we use the signals predicted from layer ℓ for the three sensors at time-point τ after word onset to classify a set of 20 words. Since MEG data is noisy, performance using a single word will be close to chance accuracy and will therefore have low power and will not be informative for our purpose. Indeed, for this reason most experiments using predictive MEG models test them on a part of the experiment that is repeated multiple times (Sudre *et al.*, 2012). These repetitions are then averaged into one test set which is predicted, and this less noisy average leads to better prediction accuracy. The experiment we are using however doesn't have any repetitions and not specific test set, and therefore by raising the number of words and classify 20 words at a time, we are able to improve the classification accuracy. We use the value of 20 words from Wehbe *et al.* (2014a).

The above steps are repeated for each of the four cross-validation folds and average accuracy is obtained for each sensor location, time-point tuple (s, τ) for layer ℓ , for each subject.

In our proof of concept experiment, we run an analysis in which we try to find, using the classification task outlined here, classification accuracy that is common both to a word embedding ℓ and to other features of a word such as a one-hot vector encoding its part of speech. This analysis is a proxy for finding the shared explained variance between the vectors, which we can call A and B. We concatenate A and B into a vector (representing $A \cup B$). We run the classification analysis using A , B and $(A \cup B)$. We then estimate the shared accuracy as: $A + B - A \cup B$.

7.4 MEG results as proof of concept

We use MEG to provide a proof of concept of our approach. We know that single word non-contextualized embeddings likely have information about the part-of-speech and the length of a word. We will show here how our approach can recover this information from brain activity as a proof-of-concept. We use MEG to study word embeddings because unlike fMRI we can access the brain activity to reading a single word. We know from the Neuroscience literature that MEG activity can be related to the length of the current word Sudre *et al.* (2012) and its part of speech Frank *et al.* (2015) at different times. We investigate whether word length and part-of-speech (POS) information is also present in the non-contextualized embedding by computing the

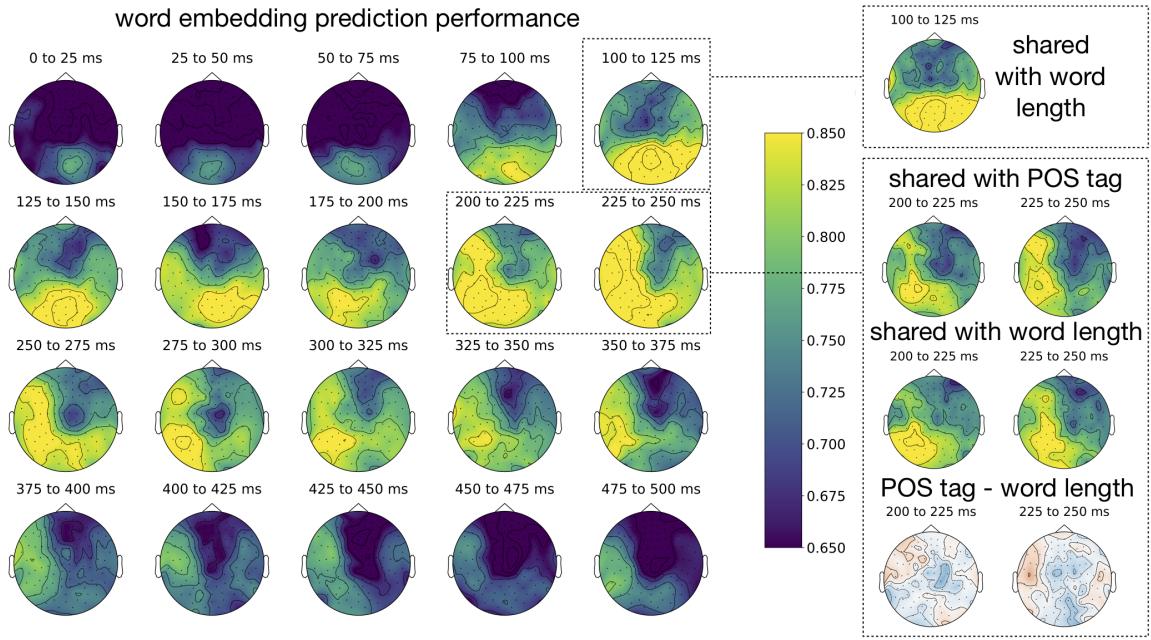


Figure 7: Performance of ELMo current word embedding at predicting MEG activity at each sensor location and time point, compared with the performance shared with word length and Part-Of-Speech (POS) tags. Around 200-250ms, the word embedding predicts a part of the activity at the top of the helmet, and this is shared mostly with the POS tags and not with word length (see bottom-right comparison). Indeed, we know from electrophysiology studies studies that POS violations incur a response around 200ms after word onset in the front of the brain Frank *et al.* (2015), which aligns with our analysis. We hypothesize from these results that the word-embedding contains both word length and POS information.

shared performance ($A \cap B$) between the pairs of features (A and B) as $A + B - A \cup B$ as explained in the previous section.

We present the results in Figure 7. The current word embedding is able to predict activity as the current word is being perceived starting at the back of the sensor helmet (generally on top of the visual cortex) around 100ms. This is when we expect the visual signal to start reaching the visual cortex. Indeed, we see that the word-embedding and the word length have overlap in the activity they predict in the visual cortex at that time. Gradually, the areas predicted by the word embedding move forward in the brain towards areas known to be involved in more high level aspects of reading. Around 200-250ms, we see the word embedding predicts a part of the activity at the top of the helmet, and this is shared mostly with the POS tags and not with word length (see bottom-right comparison). Indeed, we know from electrophysiology studies studies that POS violations incur a response around 200ms after word onset in the front of the brain Frank *et al.* (2015), which aligns with our analysis. From these results we can hypothesize that the word-embedding contains both word length and POS information, as was expected.

7.5 Complete Attention Results

condition	uni L1	uni L2	uni L3	uni L4	uni L5	uni L6	base	count
simple	1.00	1.00	0.96	1.00	0.99	1.00	1.00	120
in a sentential complement	0.83	0.83	0.83	0.83	0.84	0.83	0.83	1440
short VP coordination	0.88	0.90	0.91	0.88	0.88	0.91	0.89	720
long VP coordination	0.96	0.97	0.95	0.95	0.96	1.00	0.98	400
across a prepositional phrase	0.86	0.93	0.88	0.86	0.80	0.88	0.85	19440
across a subject relative clause	0.83	0.83	0.84	0.84	0.83	0.85	0.84	9600
across an object relative clause	0.87	0.91	0.90	0.86	0.83	0.92	0.89	19680
across an object relative clause (no that)	0.87	0.80	0.75	0.72	0.75	0.87	0.86	19680
in an object relative clause	0.97	0.95	0.96	0.92	0.91	0.91	0.95	15960
in an object relative clause (no that)	0.83	0.72	0.70	0.69	0.74	0.74	0.79	15960
reflexive anaphora: simple	0.91	0.94	0.99	0.98	1.00	0.99	0.94	280
reflexive anaphora: in a sent. complem.	0.88	0.85	0.88	0.87	0.86	0.86	0.89	3360
reflexive anaphora: across a rel. clause	0.79	0.84	0.82	0.68	0.66	0.79	0.80	22400

Table 3: Performance of models with uniformly-altered attention in layers 1-6 in BERT on a range of syntactic tasks by Marvin and Linzen (2018). ‘Base’ refers to pretrained BERT.