








Artificial Neural Network Language Models Predict Human Brain Responses to Language Even After a Developmentally Realistic Amount of Training

Eghbal A. Hosseini^{1,2} , Martin Schrimpf^{3,4} , Yian Zhang⁵, Samuel Bowman^{6,7,8} ,
Noga Zaslavsky^{1,2,9,10} , and Evelina Fedorenko^{1,2,3,11} 

¹Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA

²McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA, USA

³The MIT Quest for Intelligence Initiative, Cambridge, MA, USA

⁴Swiss Federal Institute of Technology, Lausanne, Switzerland

⁵Computer Science Department, Stanford University, Stanford, CA, USA

⁶Center for Data Science, New York University, New York, NY, USA

⁷Department of Linguistics, New York University, New York, NY, USA

⁸Department of Computer Science, New York University, New York, NY, USA

⁹K. Lisa Yang Integrative Computational Neuroscience (ICoN) Center, Massachusetts Institute of Technology, Cambridge, MA, USA

¹⁰Department of Language Science, University of California, Irvine, CA, USA

¹¹Speech and Hearing Bioscience and Technology Program, Harvard University, Boston, MA, USA

Keywords: language network, human behavior, artificial neural network, development, ANN-neural data alignment

ABSTRACT

Artificial neural networks have emerged as computationally plausible models of human language processing. A major criticism of these models is that the amount of training data they receive far exceeds that of humans during language learning. Here, we use two complementary approaches to ask how the models' ability to capture human fMRI responses to sentences is affected by the amount of training data. First, we evaluate GPT-2 models trained on 1 million, 10 million, 100 million, or 1 billion words against an fMRI benchmark. We consider the 100-million-word model to be developmentally plausible in terms of the amount of training data given that this amount is similar to what children are estimated to be exposed to during the first 10 years of life. Second, we test the performance of a GPT-2 model trained on a 9-billion-token dataset to reach state-of-the-art next-word prediction performance on the human benchmark at different stages during training. Across both approaches, we find that (i) the models trained on a developmentally plausible amount of data already achieve near-maximal performance in capturing fMRI responses to sentences. Further, (ii) lower perplexity—a measure of next-word prediction performance—is associated with stronger alignment with human data, suggesting that models that have received enough training to achieve sufficiently high next-word prediction performance also acquire representations of sentences that are predictive of human fMRI responses. In tandem, these findings establish that although *some* training is necessary for the models' predictive ability, a developmentally realistic amount of training (~100 million words) may suffice.

INTRODUCTION

A central objective in cognitive neuroscience is to develop models that can accurately predict human brain responses and behavior. In the neuroscience of language, some artificial neural

Citation: Hosseini, E. A., Schrimpf, M., Zhang, Y., Bowman, S., Zaslavsky, N., & Fedorenko, E. (2024). Artificial neural network language models predict human brain responses to language even after a developmentally realistic amount of training. *Neurobiology of Language*, 5(1), 43–63. https://doi.org/10.1162/nol_a_00137

DOI:
https://doi.org/10.1162/nol_a_00137

Supporting Information:
https://doi.org/10.1162/nol_a_00137

Received: 29 March 2023
Accepted: 9 January 2024

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

Corresponding Authors:
Eghbal A. Hosseini
ehosseini@mit.edu
Evelina Fedorenko
evelina9@mit.edu

Handling Editor:
Alessandro Lopopolo

Copyright: © 2024
Massachusetts Institute of Technology
Published under a Creative Commons
Attribution 4.0 International
(CC BY 4.0) license

network (ANN) language models were recently shown to be effective at predicting human brain activity and behavior during language processing (Caucheteux & King, 2022; Gauthier & Levy, 2019; Goldstein et al., 2022; Jain & Huth, 2018; Schrimpf et al., 2021; Toneva & Wehbe, 2019; Wilcox et al., 2020). For example, Schrimpf et al. (2021) examined the ability of over 40 language models to capture human responses to language and found that transformer architectures (Radford et al., 2019; Vaswani et al., 2017) fare best in aligning with human data. However, off-the-shelf models vary along many dimensions, making it difficult to unambiguously attribute any given model's success in aligning with human data to particular model properties (architecture, objective function, amount/kind of training data, etc.). Gaining insights into human linguistic mechanisms requires controlled experiments' on the models, where different properties are systematically manipulated (Hu et al., 2020; Kumar et al., 2022; Warstadt & Bowman, 2019). This is the approach we adopt here in order to investigate how the amount of training data affects model-to-human alignment.

One common criticism of ANN models as models of human language processing is that their training data size (often, billions of words) far surpasses the amount of language exposure in humans during their learning phase (Chang & Bergen, 2021; Dupoux, 2018; Frank, 2023; Linzen & Leonard, 2018; van Schijndel et al., 2019; see Warstadt & Bowman, 2022, for discussion). For example, Hart and Risley (1992) estimated that children are exposed to 3–11 million words each year, so by the time they turn 10 and possess adult-like linguistic competence, they are exposed to 30–110 million words. In contrast to a human child, who can learn a language from only ~100 million words (or less), many current models get orders of magnitude more training data (20,000 human years' worth for some models; Warstadt & Bowman, 2022). More recently, Gilkerson et al. (2017) estimated that by age of 10, the amount of language exposure is around 20 million words, and Frank (2023) put this estimate at between 9 and 110 million words (extrapolating from their estimate of between 200 and 400 million by age 20). Here, we ask whether this extensive amount of training is necessary for the models to acquire representations that are predictive of human brain responses during language sentence comprehension.

Prior studies on the effects of training data on the models' linguistic ability found that even with limited amounts of training data, models achieve considerable proficiency (Warstadt & Bowman, 2022). For example, Hu et al. (2020) and Zhang et al. (2020) report impressive syntactic generalizations in a BERT model (Devlin et al., 2018) trained on only millions of tokens (see also Huebner & Willits, 2021; Pannitto & Herbelot, 2020, for related evidence from a RoBERTa model trained on 5 million words of child-directed speech). Pérez-Mayos et al. (2021) find that a RoBERTa model (Liu et al., 2019) trained on 100 million words performs similarly to a model trained on 1 billion words on several syntactic benchmarks. These findings suggest that massive amounts of training may not be necessary for models to acquire certain aspects of linguistic competence. However, it is not known whether models trained on limited amounts of data can also predict human neural and behavioral responses to language.

Here, we evaluate how the amount of training data affects model-to-human alignment. In line with increasing emphasis in the field on robustness and replicability (Button et al., 2013; Ioannidis et al., 2014; Poldrack et al., 2017; Simmons et al., 2011), we adopt two complementary approaches (Figure 1). First, we investigate how well GPT-2 models (Radford et al., 2019) that are trained on different-sized datasets (1 million, 10 million, 100 million, or 1 billion words) to reach their best training task performance, predict human functional magnetic resonance imaging (fMRI) and behavioral responses to sentences. Second, we investigate how a GPT-2 model's ability to predict human fMRI and behavioral responses to sentences changes over the course of training on a large dataset to capture the "learning trajectory" of

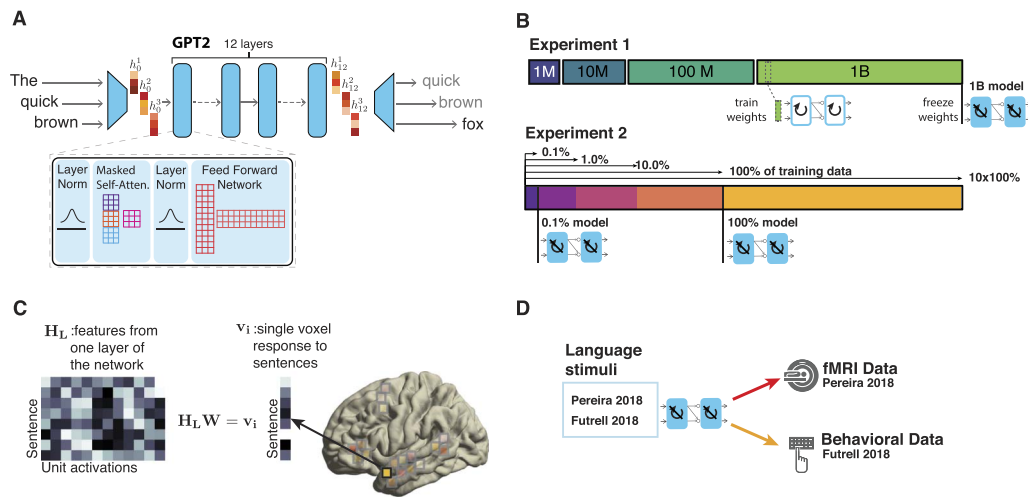


Figure 1. Methodological approach. (A) Unidirectional-attention transformer architecture. Text input is processed sequentially to predict the next likely token at each step. (B) The setup for Experiments 1 and 2. In Experiment 1, four models were trained using different-sized datasets, and for each model, the weights with the best validation perplexity were frozen and used in the model-to-brain comparison. In Experiment 2, the GPT-2 model was trained using a very large dataset, and the weights were frozen at different steps during training and used in the model-to-brain comparison. (C) Model representations were related to human representations by building a linear regression between unit activations for each layer of the model and voxel activity (in the language-selective network; Fedorenko et al., 2011) or reading times for the stimuli used in each of the benchmarks. This regression was then used to make predictions about human neural/behavioral responses for unseen language stimuli, and a Pearson correlation was computed between these predictions and the observed responses. (D) The general pipeline for predicting human brain and behavioral responses. For each benchmark, each model was exposed to the same language stimuli as humans, and the model-to-human match was evaluated as shown in C.

model-to-brain alignment. In addition, we also examine the role of model perplexity in the ability of a model to predict human responses. To foreshadow the key results, we find that models reach high performance in predicting human responses to sentences even with developmentally realistic amounts of training data.

MATERIALS AND METHODS

Human Datasets (Benchmarks)

Primary benchmark: fMRI dataset

We used fMRI data from two experiments (Experiments 2 and 3 in Pereira et al., 2018). This benchmark, hereafter called Pereira2018, is identical to the one reported in Schrimpf et al. (2021). Experiment 2 ($n = 9$ native English speakers) consisted of 384 sentences across 96 Wikipedia-style passages, spanning 24 broad topics (professions, clothing, musical instruments, etc.). There were four passages per topic (e.g., passages about a clarinet, an accordion, a piano, and a violin for the musical instruments topic). Each passage consisted of four sentences, and the sentences varied in length between seven and 18 words. Experiment 3 ($n = 9$ native English speakers) consisted of 243 sentences across 72 passages, which were a mix of Wikipedia-style passages and short narratives. Each passage consisted of three or four sentences, and the sentences varied in length between five and 20 words. The stimuli for both experiments were constructed so as to span a broad range of topic areas. In both experiments, each sentence was presented on the screen for 4 s, followed by 4 s of fixation, and each participant read the materials three times across three fMRI scanning sessions. The responses were averaged across the three repetitions to derive a single response per sentence.

Furthermore, in each participant, the analysis was restricted to a set of voxels that were identified as language-responsive in an independent extensively validated language *localizer* task (Fedorenko et al., 2010). In the localizer task, participants read sentences and list of nonwords (e.g., “blook” or “cre”) in a standard blocked design. Each item consisted of 12 words/nonwords that were presented one at a time at the rate of 450 ms per word/nonword. Each sentence/nonword list was followed by a simple button-press task (“Press a button when you see a picture of a finger on a button”). Each trial lasted 6 s. Each block consisted of three trials and lasted 18 s. Each scanning run consisted of 16 experimental blocks (8 per condition) and six fixation blocks and lasted 358 s. Each participant completed two runs, and the order of conditions was counterbalanced across runs. The localizer is available for download at (<https://evlab.mit.edu/funcloc/>). The contrast between sentences and nonword lists has been shown to robustly identify the frontotemporal language-selective network of brain areas (Fedorenko et al., 2011; Lipkin et al., 2022). These areas support language comprehension across modalities (listening, reading, etc.) and have been established to be sensitive to both word meanings and syntactic structure processing (e.g., Fedorenko et al., 2010; Fedorenko et al., 2020). For each participant, we selected the top 10% of most localizer-responsive voxels within a set of 12 broad masks (6 in each hemisphere) that cover inferior frontal and lateral temporal cortex (the masks were derived from an independent set of 220 participants who performed the language localizer task and are available at <https://evlab.mit.edu/funcloc/>). Thus, the fMRI benchmark consists of—for each of the two experiments—a set of language responsive voxels in each participant, and for each voxel, we have an estimate of the blood oxygen level dependent (BOLD) response to each of 384 sentences (Experiment 2 in Pereira et al., 2018) or 243 sentences (Experiment 3 in Pereira et al., 2018).

Secondary benchmark: Behavioral (reading-times) dataset

We used self-paced reading data from (Futrell et al., 2018). Similar to the fMRI benchmark, this benchmark, hereafter called Futrell2018, is identical to the one reported in Schrimpf et al. (2021). 179 native adult English-speaking participants (recruited through Mechanical Turk) read stories that were presented one word at a time; with each button press, the current word would disappear in place of the new word (e.g., Just et al., 1982). The time it took a participant to move to the next word, $n + 1$, was used as a measure of comprehension difficulty at word n . The stories were based on existing stories but were edited in a way so as to increase the frequency of rare words and constructions, including constructions that are known to cause comprehension difficulty (see Futrell et al., 2018, for details). The stories consisted of 33–64 sentences, and contained between 938 and 1,089 words. Each of 179 participants read between five and 10 stories and answered comprehension questions at the end of each story; each story was read by 82–98 participants. Following Futrell et al. (2018), we excluded reading times outside of the [100 ms, 3,000 ms] range. Note that we report the results for this benchmark in the Supporting Information, available at https://doi.org/10.1162/nol_a_00137 (see Supplementary Figure 3) because, as will be discussed below, we found that diverse variants of control (untrained) language models predict these responses well above chance, which suggests that model predictivity is unlikely to be related to the representation of the linguistic stimuli.

Artificial Neural Network Models

We used two different implementations of a GPT-2-style model. For Experiment 1, where a model was trained on a dataset with a controlled number of words, we used the GPT-NEOX library, which is a distributed training framework that uses the DeepSpeed library (Aminabadi

et al., 2022; Black et al., 2022). We used a unidirectional-attention transformer model (GPT-2; Radford et al., 2019) with 12 layers and an embedding layer which was learned during training. Each layer had a size of 768 units and consisted of four main blocks (Figure 1A): (i) first layer normalization, (ii) self-attention, (iii) second layer normalization, and (iv) the feedforward layer. The final layer consisted of a linear projection with a sigmoid nonlinearity that mapped hidden states into probabilities over the dictionary. The context size was 1,024 tokens. To test whether our results would generalize to bidirectional-attention transformer architectures, we additionally used publicly available miniBERTa models that were trained on the same datasets as the GPT-2 models (Zhang et al., 2020; <https://huggingface.co/nyu-ml>). Note, we did not include the model trained on the smallest (1 million words) dataset, for which Zhang et al. (2020) used a smaller-sized model, which therefore would not be directly comparable to the other models. The miniBERTas use the same design as the RoBERTa base model (Liu et al., 2019)—a bidirectional-attention model with 12 layers, each 768 units in size, and a context size of 512 tokens. Importantly, RoBERTa has the same number of parameters as GPT-2 (125 million), allowing for a relatively controlled comparison of uni- and bidirectional architectures.

For Experiment 2, to investigate model training dynamics with a very large dataset, where during the early stages of the training the model continues to see new input (cf. doing multiple passes through a smaller-size training corpus as in Experiment 1), we used GPT-2 model weights from a publicly available model from the Hugging Face Transformers library (<https://huggingface.co/stanford-crfm>). The model has a similar architecture to the GPT-2 model used in Experiment 1.

Model Training

Training datasets

For Experiment 1, we combined the BookCorpus (Zhu et al., 2015) and English Wikipedia (Liu et al., 2019; Zhu et al., 2015) with a 1:3 ratio. We then created four different datasets with 1 million, 10 million, 100 million, and 1 billion words. These were used for training both the GPT-2 models and the miniBERTa models. For Experiment 2, we used a model that was trained on the OpenWebText corpus (Gokaslan & Cohen, 2019) with more than 9 billion tokens.

To characterize the training corpora, we counted the number of unique tokens, token bi-grams, token tri-grams and token four-grams for different dataset sizes in Experiment 1 and for different checkpoints in Experiment 2. For this analysis, the data were tokenized as in the Penn Treebank corpus (Marcus et al., 1993). In particular, contractions were split (e.g., they're → they + re) and punctuation marks were treated as separate tokens. Afterwards, we counted unique occurrences of tokens, tokens bi-grams, and so on. As shown in Supplementary Figure 2B and D, the number of all n -grams increases with corpus size (Experiment 1) and for later checkpoints (Experiment 2), and the percentage of unique tokens relative to total tokens is always higher in Experiment 2 compared to Experiment 1. (See also Supplementary Figure 10 for an illustration of the training dynamics in Experiment 1 vs. Experiment 2; note, quantifying the variability in syntactic structures is more challenging given the size of the corpora and the fact that they are not parsed/POS-tagged.)

In addition, following a reviewer's request, we tested whether the experimental materials from the human benchmarks were present in the training corpora. We found that none of the sentences from the fMRI benchmark were present in any of the training corpora, and only a very small number of sentences from the behavioral benchmark were present in the training corpora. In particular, in Experiment 1, three of the sentences from the behavioral benchmark

were found in the 1 million word training dataset, two sentences in the 10 million word training dataset, four sentences in the 100 million word training dataset, and 17 sentences in the 1 billion word training dataset; in Experiment 2, four sentences were found in the training dataset.

Training procedure

For Experiment 1, to train the GPT-2 models, we used standard initialization from the GPT-NEOX library and standard training parameters (Radford et al., 2019; see Supplementary Figure 1 for details). After training, the model weights with the smallest validation perplexity were selected for evaluation on the human benchmarks. The smallest validation perplexity was reached after 1,000 steps for the 1 million word dataset (1,024 tokens * 128 batches * 1,000 steps = 131,072,000 tokens total), after 2,000 steps for the 10 million word dataset (1,024 tokens * 128 batches * 2,000 steps = 262,144,000 tokens), after 14,250 steps for the 100 million word dataset (1,024 tokens * 128 batches * 14,250 steps = 1,867,776,000 tokens), and after 310,000 steps for the 1 billion word dataset (1,024 tokens * 128 batches * 310,000 steps = 60,948,480,000 tokens). To train the miniBERTa models, we used standard initialization and training parameters from the Hugging Face Transformers library (Liu et al., 2019).

Predictivity of the trained models was compared to that of untrained models. Here, in addition to the untrained GPT-2 model available from the Hugging Face library, we created an alternative untrained GPT-2 model in order to investigate the effects of different weight initializations on the alignment between model representations and human neural responses and reading behavior, and thus to isolate the effects of model architecture alone (i.e., the units and the patterns of connections among them) on predictivity. This version implemented the same unidirectional mask as the trained models and the other untrained model, but all the weights were set to a Gaussian distribution with a fixed mean and standard deviation (mean: 0, standard deviation: 0.02 for the layer normalization, self-attention, and feedforward layer weights; see Supplementary Figure 6 for a detailed comparison with the Hugging Face initialization parameters).

For Experiment 2, the GPT-2 model was trained with standard initialization and training parameters until it reached state-of-the-art perplexity values. We selected several checkpoints at which we extracted model representations from each layer for evaluation on the human benchmarks. The model was trained on 16 GPUs, with eight batches per GPU, and updates were performed after four gradient accumulations. As a result, a *training step* constituted 1,024 (tokens) * 8 (batches) * 16 (GPUs) * 4 (gradient accumulations) = 524,288 tokens. Given that the tokenized OpenWebText corpus contains 9,036,044,288 tokens, it takes close to 20,000 training steps (specifically, 17.2K steps) to do one complete pass over the corpus. The checkpoints were selected in a logarithmic manner: 0, 0.1% (20 training steps, which corresponds to 524,288 tokens * 20 steps = 10,485,760 tokens), 1.0% (200 steps, which corresponds to 524,288 tokens * 200 steps = 104,857,600 tokens), 10% (2K steps, which corresponds to 524,288 tokens * 2,000 steps = 1,048,576,000 tokens), 100% (20K steps, which corresponds to 524,288 tokens * 20,000 steps = 10,485,760,000 tokens), and 10 × 100% (200K steps, which corresponds to 524,288 tokens * 200,000 steps = 104,857,600,000 tokens).

Analyses

Model comparison to the fMRI benchmark

We followed the approach in Schrimpf et al. (2021). In particular, we first extracted the representation for all the sentences that were used in the human fMRI experiments from each

layer of each model. For each experiment, we split the stimuli into five ~equal-size batches (Experiment 1: three batches of size 76 sentences and two batches of size 78; Experiment 2: two batches of size 48 sentences and three batches of size 49 sentences) and used ~80% of the data to build an ordinary least squares regression model between model unit activations and voxel-level responses in the language network (defined by an extensively validated language localizer task, as described in *Primary Benchmark: fMRI Dataset*, above; Fedorenko et al., 2010; Lipkin et al., 2022). For both experiments, we selected the representation of the last word in each sentence (for multi-token words, we averaged the representations across the composite tokens). We then applied the regression to the left-out ~20% of sentences to generate predictions for BOLD responses in each voxel and compared these predictions against the observed BOLD responses using Pearson correlation. This procedure was iterated across the data folds, leaving out a different 20% each time, and the Pearson values were averaged across these iterations in each voxel of each participant. For each participant averaging across the two experiments (Experiments 2 and 3 in Pereira et al., 2018), we obtained a single score by taking the median Pearson value across the language-responsive voxels. A reliable positive Pearson correlation value indicates that the model is able to predict fMRI responses with a linear transformation. The resulting Pearson correlation values are divided by the ceiling value computed by estimating how well a best possible model of an average human would predict fMRI responses (Schrimpf et al., 2021). This value was estimated to be 0.32 for the Pereira et al. (2018) dataset. (We acknowledge that the issue of how to compute a noise ceiling remains an open issue in the field.)

Statistical testing was performed on the scores from the best model layer (as determined in Schrimpf et al., 2021) and took the form of independent two-sample *t* tests. In particular, participant scores for each model in Experiment 1 ($n = 5$ models: untrained, 1M, 10M, 100M, and 1B) or each checkpoint of a model in Experiment 2 ($n = 6$ checkpoints: untrained, 0.1% of training steps, 1% of training steps, 10% of training steps, 100% of training steps, and $10 \times 100\%$ of training steps) were compared to the scores for the fully trained model (from Schrimpf et al., 2021). The resulting *p* values were Bonferroni-corrected for the number of comparisons (5 and 6 comparisons in Experiments 1 and 2, respectively).

Model comparison to the behavioral benchmark

Similar to the fMRI benchmark, we followed the approach in Schrimpf et al. (2021). In particular, we first extracted the representation for all the stimuli (individual words) that were used in the behavioral experiment from the last layer of each model. For each participant, we split the total words (which varied across participants depending on the number of stories that a participant read) into five ~equal-sized batches, and used ~80% of the data to build an ordinary least squares regression model between model unit activations and reading times. In dividing the data into batches, we ensured that (a) the same word did not appear in the training versus the test set, and (b) for any given sentence, the words were divided as evenly as possible between the training and the test set. We then applied the regression to the left-out ~20% of words to generate predictions for reading times and compared these predictions against the observed reading times using Pearson correlation. This procedure was iterated across the data folds, leaving out a different 20% each time, and the Pearson values were averaged across these iterations to obtain a single score per participant. A reliable positive Pearson correlation value indicates that the model is able to predict reading times with a linear transformation. The resulting Pearson correlation values are divided by the ceiling value computed by estimating how well a best possible model of an average human would predict reading times (Schrimpf et al., 2021). This value was estimated to be 0.76 for the Futrell et al. (2018) dataset.

Statistical testing was performed on the scores from the last model layer and took the form of independent two-sample *t* tests, with a Bonferroni correction for the number of tests, similar to the fMRI benchmark.

Model perplexity

Following standard practice (e.g., Jelinek et al., 1977), we used perplexity as a measure of model performance on the language prediction tasks (next-word prediction for the GPT-2 models and missing-word prediction for the miniBERTa models). Perplexity (PPL) is defined as:

$$\text{PPL} = 2^{H(x)}$$

$$H(x) = -\frac{1}{N} \sum_{i=1}^N \log_2 P(x_i)$$

where H is entropy, and x denotes the tokens.

For both experiments, we used the test set from the wikitext-103-raw-v1 dataset (Merity et al., 2016) to compute perplexity. Perplexity was computed using a context size of 1,024 tokens and a stride of 512 tokens.

RESULTS

Models Trained on Small Corpora Predict Human Responses

We started by examining the performance of a unidirectional transformer model trained on the standard language modeling task in predicting human fMRI responses during the processing of sentences. Specifically, we tested a GPT-2 architecture (Figure 1A), which has previously been shown to best capture human neural and behavioral responses (Schrimpf et al., 2021). We trained four independent models on 1 million, 10 million, 100 million, and 1 billion words, respectively (Supplementary Figure 1). A training dataset of size 100 million words is comparable to the amount of language input that children have been estimated to get during the first decade of life (Frank, 2023; Hart & Risley, 1992). Of course, the *nature* of the language input is still quite different between models and children both with respect to the content and the modality—text only for models vs. multimodal input for children. We return to this point in the Discussion. After training, we selected the checkpoint with the best perplexity on the validation set and tested how well the model representations capture human neural (fMRI) responses to sentences in the language-selective network (Fedorenko et al., 2011) and human behavioral responses in a self-paced reading task (Figure 1D).

For the Pereira2018 benchmark (Pereira et al., 2018; Schrimpf et al., 2021), we observed a consistent increase in performance with an increase in the size of the training set (Figure 2A; see Supplementary Figure 5 for evidence—for this and the behavioral benchmark—of robustness of this pattern to seed choice during model initialization; cf. Mehrer et al., 2021; see Frank et al., 2015, and Aurnhammer & Frank, 2019, for similar findings from earlier, pre-transformer models, including *n*-gram, RNN, and phrase-structure grammar models). Critically, however, the model trained on just 100 million words already exhibits fMRI response predictivity that is similar to that of the fully trained GPT-2 model as reported in Schrimpf et al. (2021), with no significant difference in predictivity values ($p = 0.99$; the data frames are available at OSF: see Data and Code Availability Statement). The model trained on 1 billion words also does not differ from the fully trained model in predictivity ($p = 0.82$). In contrast, the predictivity of the untrained model (the version with the Hugging Face initialization parameters)

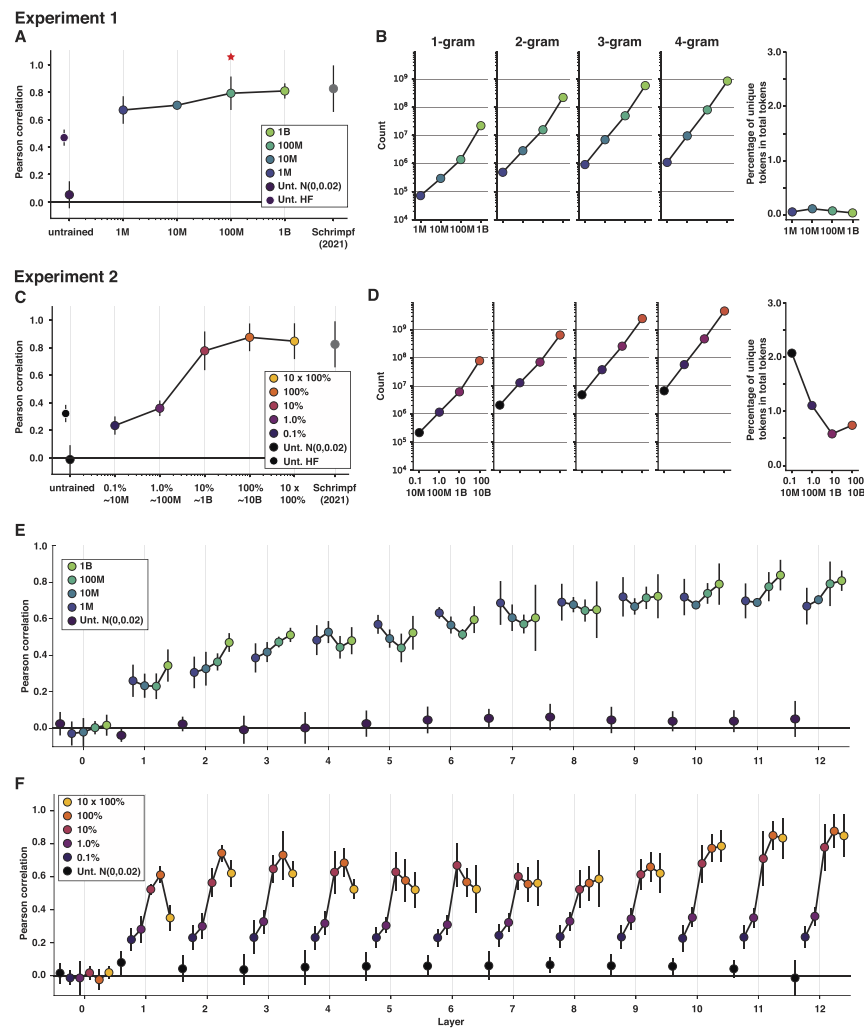


Figure 2. Model performance on the fMRI (Pereira2018) benchmark as a function of training. (A) Experiment 1 results: performance (normalized predictivity) of the best-performing GPT-2 layer, as reported in Schrimpf et al. (2021), in predicting language-responsive voxels’ activation in the Pereira2018 benchmark. The results are shown for (i) two versions of an untrained (Unt.) model (initialized in two different ways: Unt. N(0,0.02) corresponds to the untrained model initialized with a mean of 0 and a standard deviation of 0.02, and Unt. HF corresponds to the untrained model initialized with the Hugging Face parameters; black dots); (ii) four models trained on datasets of different sizes (1 million, 10 million, 100 million, and 1 billion words; blue-to-green dots connected by a line; the model trained on a developmentally plausible amount of data—100 million—is marked with a red asterisk; see also Supplementary Figure 4 for the results for a 50 million word model); and (iii) a fully trained model, as reported in Schrimpf et al. (2021; gray dots). Here, in Figure 3A, and in Figure 3, we computed a median score across participants and divided it by an estimated ceiling value to get a normalized score, and we computed a median absolute deviation over participants for use as error bars. (B) The number of unique tokens (1 gram), token bi-grams, token tri-grams, and token four-grams in each training dataset. There is at least a logarithmic increase in the counts with the increase in the dataset size. The rightmost panel shows the percentage of unique tokens relative to all tokens for each dataset. For Experiment 1, the total number of tokens is computed based on the number of training steps that was needed to reach best validation loss (see Materials and Methods; see also Supplementary Figure 1 for the illustration of the training dynamics). (C) Experiment 2 results: performance of the best-performing GPT-2 layer, as reported in Schrimpf et al. (2021), in predicting language-responsive voxels’ activation in the Pereira2018 benchmark. The results are shown for (i) two versions of an untrained model (initialized in two different ways, as in Figure 2A; see Materials and Methods; black dots); (ii) a model trained on a large dataset examined at different points during the training (0.1%, 1.0%, 10%, 100%, and 10 × 100% of training steps; purple-to-yellow dots connected by a line); and (iii) a fully trained model, as reported in Schrimpf et al. (2021; gray dots). (D) Same as in Figure 2B but for the OpenWebText training dataset. (E–F) Exploratory analyses of individual model layers: performance of the 12 GPT-2 model layers in predicting human neural responses in the Pereira2018 benchmark in (E) Experiment 1 and (F) Experiment 2. Layer 0 is the token embedding layer, and layer 12 is the last layer. The results are shown for (i) an untrained (Unt.) model (with the Gaussian initialization; black dots); and (ii) four models trained on datasets of different size (blue-to-green dots connected by a line in A) or a model trained on a large dataset examined at different points during the training (purple-to-yellow dots connected by a line in B).

and the models trained on 1 million and 10 million words is significantly below the predictivity of the fully trained model (p values < 0.0001 , 0.007 , and 0.016 , respectively; here and elsewhere, the values are Bonferroni-corrected, as described in Analyses).

The untrained model performance differs between the two versions (see Training Procedure, above). The version initialized with the standard Hugging Face parameters performs well above chance ($p < 0.0001$), as reported in Schrimpf et al. (2021; see also Caucheteux & King, 2022), but the version initialized with the alternative parameters (all weights set to a normal distribution with a mean of 0 and a standard deviation of 0.02) performs around 0 (not significantly different from 0; $p = 0.14$; Figure 2A).

The results also generalize, to some degree, to a bidirectional transformer model (mini-BERTa; Liu et al. 2019; Supplementary Figure 2). In particular, similar to the GPT-2 models, we observed a consistent increase in model performance with an increase in the training dataset size, which suggests that this pattern is robust to architecture. However, the 100 million word model still performs below the fully trained model. This difference between the GPT-2 and miniBERTa models in the amount of training they require to align with human data is likely due to the difference in the directionality of the attention mechanisms, with unidirectional-attention mechanisms being more sample efficient. Generalizing these results to other minimally different variants of uni- vs. bidirectional-attention transformer models will help strengthen this conclusion.

In exploratory analyses, in addition to examining the language network as an integrated system, we examined the effects of the amount of training data on the models' ability to predict fMRI responses in individual frontal and temporal language functional regions of interest (ROIs) for a total of six language fROIs in the left hemisphere (LH), and six homotopic regions in the right hemisphere (RH; e.g., Lipkin et al., 2022). The results are shown in Supplementary Figure 9. The overall pattern was similar across all language fROIs, including between the LH inferior frontal gyrus fROI and the LH post-temporal fROI (which have been argued by some to differ functionally; e.g., Friederici, 2018; Hagoort, 2019; cf. Fedorenko et al., 2020). The overall predictivity was lower in the RH than the LH language fROIs ($p \ll 0.0001$ for all models in Experiment 1 and all checkpoints in Experiment 2), in line with past findings (e.g., Schrimpf et al., 2021; Tuckute et al., 2024).

We also investigated the patterns of model-to-brain alignment across model layers. Prior work in vision (Geiger et al., 2020; Storrs et al., 2021) has suggested that training affects model performance differently across layers, with early layers already reaching close to maximal performance with a limited amount of training, but later layers continuing to benefit from increasingly more training. In line with these prior observations, for the Pereira2018 benchmark, we observed that for layers 4–9, performance peaks for the 1 million word model, and for the last three layers (layers 10–12), a consistent improvement in performance is observed with larger datasets (Figure 2E). This observation echoes prior work showing that later layers build more contextualized representation of linguistic stimuli and better capture syntactic and compositional semantic aspects of the linguistic signal (Belinkov et al., 2017; Hewitt & Manning, 2019; Tenney et al., 2019), to which the language brain regions are also deeply sensitive (e.g., Fedorenko et al., 2010; Fedorenko et al., 2020; Pallier et al., 2011; Shain et al., 2023).

The general pattern of results was also similar for the secondary, behavioral benchmark (Futrell2018; Supplementary Figure 3A): The predictivity of the untrained model and the model trained on 1 million words is significantly below the predictivity of the fully trained model (p values < 0.0001 and $p = 0.00016$, respectively); and the predictivity of the models trained on 10 million words, 100 million words, and 1 billion words does not

significantly differ from that of the fully trained model (p values > 0.05). However, because both of the untrained models achieve reliably above-zero predictivity on the Futrell2018 benchmarks ($ps < 0.0001$), model performance is unlikely to be related to the representation of linguistic stimuli. As a result, we present these findings in Supporting Information, for completeness.

Models Trained on a Small Portion of a Massive Corpus Predict Human Responses

In the previous section, we investigated how models that are trained on small corpora (until they reach their best performance on the target language modeling task) perform in predicting human data. However, humans, including children learning a language, are continuously exposed to new words and constructions (see Supplementary Figure 10). To better simulate such scenarios, as well as to evaluate the robustness of the results to approach, we examined how the ability of a model to predict human fMRI responses to sentences changes over time as the model is being trained on a very large corpus, similar to (Caucheteux & King, 2022). To do so, we used a GPT-2 model that was trained on a corpus consisting of over 9 billion tokens and selected several checkpoints during the training process (0.1%, 1.0%, 10.0%, 100%, and $10 \times 100\%$ of training steps, where 100% of training steps approximately equal one complete pass over the full dataset; see Materials and Methods). At each of these checkpoints, we tested how well the model representations capture human responses to sentences.

For the Pereira2018 benchmark, the performance of the fully trained model (i.e., $10 \times 100\%$ of training steps) closely matches the results reported in Schrimpf et al. (2021; where the Hugging Face version of the model was used; cf. the GPT-NEOX library version here), with no significant difference in predictivity ($p = 0.67$). This result shows that model-to-human alignment is robust to the details of model implementation, as one would hope. Critically, mirroring the results from Experiment 1, we observed a consistent increase in how well the model predicts fMRI responses to sentences until the model reaches the 10% checkpoint, at which point the performance plateaus. Critically, the predictivity of the models trained on 10% or 100% of the training steps does not significantly differ from the predictivity of the fully trained model (p values > 0.05). In contrast, the predictivity of the untrained model (the version with the Hugging Face initialization parameters) and models trained on 0.1% or 1.0% of the training steps is significantly below that of a fully trained model ($ps < 0.001$; Figure 2C; see Supplementary Figure 4 for evidence of robustness to seed choice during model initialization).

The slight decrease in performance with more training (from 100% to $10 \times 100\%$) suggests that more training does not necessarily lead to better alignment with human brain data, although it is possible that this result is due to the relatively spatially and temporally coarse nature of our neural measurements. In particular, a response in a given fMRI voxel reflects an average activity of a large population (a few hundred thousand) of neurons, and the activity is averaged over multiple seconds, which necessarily obscures the fast dynamics of language processing. It is possible that for finer-grained neural data, such as intracranial recordings (electrocorticography or stereo electroencephalography, or EEG), we might continue to see improvements with more training.

In exploratory analyses of the individual model layers, we observed that performance shows a consistent increase across layers up to the 1.0% checkpoint. After that, the early and middle layers show a drop in performance from the 10% checkpoint to the $10 \times 100\%$ checkpoint, whereas in the later layers, performance increases from the 10% to the 100%

checkpoint and then reaches a plateau (Figure 2F). Additionally, as in Experiment 1, and in line with prior work in vision (e.g., Geiger et al., 2020; Storrs et al., 2021), earlier layers reach close to maximal performance earlier in the training (at the 1% checkpoint), whereas later layers reach their peak close to the 10% checkpoint (Figure 2F).

The pattern of results for the secondary, behavioral benchmark (Futrell2018) closely follows the pattern that we observed with limited-size training datasets in Experiment 1, with predictivity reaching a plateau after the 1% checkpoint (Supplementary Figure 3B). The predictivity of the untrained model and the models trained on 0.1% or 1% of the training steps is significantly below the predictivity of the fully trained model (p values < 0.001 , < 0.001 , and 0.0015 , respectively); and the predictivity of the models trained on 10%, 100%, and $10 \times 100\%$ of the training steps does not significantly differ from that of the fully trained model (p values > 0.05).

Model Perplexity Predicts Model Performance

For ANN language models, perplexity (a measure of performance on the next-word prediction task; see Analyses) is a reliable predictor of model performance on diverse NLP benchmarks (e.g., Brown et al., 2020; Radford et al., 2019). Schrimpf et al. (2021) further found that off-the-shelf models that perform better on the next-word prediction task are also better able to capture human neural and behavioral responses (see Antonello & Huth, 2024, for evidence that a similar relationship obtains for other tasks), in line with prior work showing a similar relationship in pre-transformer models between the amount of training and the ability of a model to predict ERP responses (Aurnhammer & Frank, 2019; Frank et al., 2015; cf. Pasquiou et al., 2022). Here, we examined the relationship between model perplexity and its ability to predict human fMRI responses for models that only differ in the size of the training corpus and for a model at different stages of training, in order to test whether better performance on the next-word prediction task is associated with representations that are more strongly predictive of human neural responses to language.

As expected, perplexity is lower (i.e., the ability to predict upcoming words is better) for models that are trained on larger datasets (Figure 3A) and for a given model at the later stages of training (Figure 3B; see Supplementary Figure 3C and D for the results on the behavioral benchmark). Critically, across both Experiments 1 and 2, we observed a consistent relationship between perplexity and neural predictivity, such that lower perplexity is associated with higher predictivity. However, once a model reaches a certain level of perplexity, further improvements in the model's ability to predict the next word are no longer associated with increases in predictivity, in line with recent findings (Oh & Schuler, 2022, 2023).

DISCUSSION

In this work, we investigated the relationship between the amount of training data and predictivity of fMRI responses to sentences for transformer-based ANN language models. Our study makes several contributions: (1) Even when trained on a developmentally realistic amount of data, transformer language models align with human data; (2) alignment between untrained ANN language models and human fMRI responses is strongly affected by the initial unit weight configuration; and (3) model perplexity predicts brain scores.

Performance on Developmentally Realistic Amount of Training Data

Using an fMRI benchmark (Pereira et al., 2018), we established that even with a developmentally realistic amount of training data (~100 million words, comparable to what humans

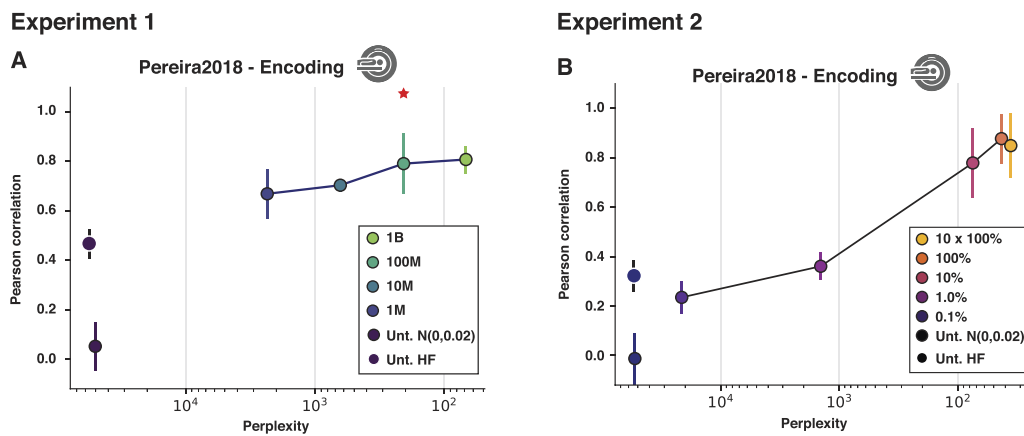


Figure 3. Relationship between model perplexity and model ability to predict human brain responses to sentences. (A) Experiment 1 results: the relationship between perplexity, i.e., the model's ability to predict the next token in an independent dataset (wikitext-103-raw-v1), shown on the x-axis, with lower values corresponding to better performance, and model performance in predicting language-responsive voxels' activation in the Pereira2018 fMRI benchmark. The results are shown for (i) two versions of an untrained model (black dots; see Figure 2 caption for details); and (ii) four models trained on datasets of different sizes (1M, 10M, 100M, and 1B words; blue-to-green dots connected by a line; the model trained on a developmentally plausible amount of data—100 million words—is marked with a red asterisk). (B) Experiment 2 results: the relationship between perplexity and model performance in predicting language-responsive voxels' activation in the Pereira2018 fMRI benchmark. The results are shown for (i) two versions of an untrained model (black dots; see Figure 2 caption for details); and (ii) a model trained on a large dataset examined at different points during the training (0.1%, 1.0%, 10%, 100%, and 10 × 100% of training steps; purple-to-yellow dots connected by a line).

get during the first 10 years of life; Frank, 2023; Hart & Risley, 1992), a GPT-2 model achieves near-maximal predictivity of fMRI responses to sentences. This effect generalizes to a different model architecture (a bidirectional-attention transformer: RoBERTa), although, compared to GPT-2, such models appear to be less sample efficient, requiring more training data to achieve peak predictivity. (The result also generalizes to a behavioral reading-times benchmark (Futrell et al., 2018); high performance of untrained models on this benchmark prompted us to move these findings to the Supporting Information.) In a complementary approach, we showed that when trained on a large dataset, a GPT-2 model already achieves near-maximal predictivity with only 10% of the training steps, well before a full pass over the dataset.

These results align with prior work in vision. For example, Geiger et al. (2020) found that even a small amount of training can result in model representations that are predictive of neural responses in macaques. Moreover, the logarithmic nature of the increase in predictivity between a model trained on 1 million tokens and a model trained on 1 billion tokens aligns with prior NLP results (e.g., see Kaplan et al., 2020, for evidence of a logarithmic relationship between training data size and the loss in training, and between model size and loss), as well as with vision research (e.g., see Geiger et al., 2020, for evidence of a logarithmic relationship between training data size and predictivity of neural firing rates).

The key implication of these findings is that although large language models are trained on vast amounts of data (and performance on some NLP benchmarks continues to improve with more training), this large amount of training is not necessary for these models to acquire representations that are predictive of human brain responses and behavior. The fact that ANN models trained on a developmentally plausible amount of data can accurately capture fMRI responses to sentences helps address one of the most common criticisms of these models as models of human language processing.

Performance of Untrained Models

By relating different versions of untrained models to human fMRI responses, this work clarifies the contributions of architecture to the models' ability to predict neural responses to linguistic input. Schrimpf et al. (2021; see also Caucheteux & King, 2022; Pasquiou et al., 2022) have found that untrained models predict fMRI responses quite well, albeit worse than trained models. They speculated that good performance of untrained models might be due to the smoothing of word embeddings across layers in a way that enables the embeddings to capture some aspects of statistical regularities of language, perhaps something as general as nearby words being likely to be related to one another. However, what counts as untrained is important to clarify.

Untrained models come with a particular setting of their unit weights. A particular weight configuration may get "baked into" a model during the process of model development, aimed at maximizing learning efficiency for the target task. Such potential biases in initial, pre-trained weights may be akin to innate, evolution-shaped, aspects of brain structure, which may filter information in specific ways as it travels within or across brain areas, even before any learning of the input regularities has occurred (e.g., Zador, 2019). We showed that initializing a model with a normal distribution for all weights leads to the model being unable to predict fMRI response to sentences (predictivity is at ~0; of course, such a model is also unable to perform the next-word prediction task). This inability to predict fMRI responses for models initialized with a normal distribution is not due to the lack of activity propagation across layers, as shown in Supplementary Figure 8B. We also showed that the standard deviation of weight initialization only has minimal effect on the predictivity for untrained model (Supplementary Figure 8D).

In summary, the ability of untrained models to predict fMRI responses to sentences reported in previous studies should not be taken as evidence that model architecture alone (i.e., the units and the patterns of connections among them) can capture human neural responses to linguistic input, or at least, it should be acknowledged that these effects are due to the particular pre-trained weight configurations. Furthermore, if a model can (at least partially) match human data with a few bits of information in the form of the initialization parameters (see Supplementary Figure 8C for evidence that above-baseline predictivity for some initializations may result from the representations for different sentences being more similar), then any results at that alignment level or below for trained models are not meaningful and we should focus on progress beyond that alignment level. Another implication is that future attempts to align trained ANN models with human data should generalize their findings across different weight initializations (Mehrer et al., 2020).

Model Perplexity

In line with Schrimpf et al.'s (2021) claim that models that perform better on next-word prediction are better at predicting brain data (see also Caucheteux & King, 2022; cf. Antonello & Huth, 2024), we found that model perplexity for different amounts of the training data is a good proxy for model performance in predicting responses to sentences. We observed this relationship both in Experiment 1, where we varied the size of the training dataset, and in Experiment 2, where we tested model representations at different points during the training on a large dataset. These findings provide further evidence that optimizing for predictive representations—through training the models on the next-word prediction task—may be critical for ANN models to acquire representations that are predictive of human responses to linguistic input.

This finding aligns well with earlier work, which showed that surprisal (how predictable a word is from the preceding context), which is closely related to perplexity, is generally predictive of human behavioral responses (e.g., Smith & Levy, 2013) and neural responses, as estimated with EEG (e.g., Aurnhammer & Frank, 2019; Frank et al., 2015; Rabovsky et al., 2018), magnetoencephalography (Brodbeck et al., 2022; Heilbron et al., 2022), fMRI (Brennan et al., 2016; Heilbron et al., 2022; Henderson et al., 2016; Lopopolo et al., 2017; Shain et al., 2020; Willems et al., 2016), or intracranially (Goldstein et al., 2022) during language processing. However, as recently shown in Tuckute et al. (2024), representations from language models achieve substantially higher predictivity for fMRI response to sentences than more traditional surprisal metrics based on n -gram counts or probabilistic context-free grammar parser probabilities.

One recent study (Pasquiou et al., 2022) did not observe a relationship between perplexity and model ability to predict human fMRI responses to linguistic input. We speculate that the lack of this relationship in Pasquiou et al.'s data may relate to the use of an extended-narrative stimulus (i.e., the entire *The Little Prince* book) rather than single sentences or short passages. The overall low encoding performance for such stimuli imposes a ceiling on the correlations between model-to-brain alignment and model perplexity (or other variables), making it difficult to differentiate among models. Alternatively, humans and models may use different information for predicting upcoming words, especially in extended linguistic stimuli (Oh & Schuler, 2022).

Why models struggle with predicting neural responses to long narratives is a separate and important question. We offer a speculation. In the human brain, division of labor exists between (i) the language-selective network, which integrates information within clauses/sentences but does not track longer-range contexts (e.g., Blank & Fedorenko, 2020), and (ii) the default network(s) (Buckner & DiNicola, 2019), which integrates information over extended temporal contexts (Lerner et al., 2011). Importantly, the default network does not operate over word sequences; instead, the information that this system represents is likely abstract, as evidenced by the fact that it processes long contexts in both linguistic and nonlinguistic stimuli (e.g., Baldassano et al., 2017; Simony et al., 2016). As a result, the ANN language models (like those used in current work and in Pasquiou et al. (2022)) may simply lack representations that are sufficiently abstract (not directly tied to the stimulus, i.e., to the word sequences) to match those in the default network, perhaps because language models eventually have to “go back” to specific words in order to perform the next-word prediction task. Some of the newer models, like GPT-3, seem to be able to handle a greater degree of abstraction (Brown et al., 2020) and thus may be promising for future attempts to capture human neural responses to long and complex linguistic stimuli.

Limitations and Future Directions

In general, it is challenging to compare the amount of training data that a model gets to the amount of linguistic input that a child gets. The consequences of a single token of input for a computational model depend on many aspects of the model's architecture, training setup, and so on; and of course, the fact that, for smaller datasets, the same dataset is repeated multiple times during the training varies drastically from what humans experience. All the results should therefore be interpreted in light of these limitations.

Furthermore, we have here focused on the effects of the *amount* of training data on the ANN language models' ability to capture human responses to language. However, the *nature* of the training data is, no doubt, also important. For example, training models on data that are

similar to what children are exposed to could lead to improved neural predictivity (Chang & Bergen, 2021; Warstadt & Bowman, 2019). Indeed, this approach has been shown to improve vision models' ability to capture primate neural responses (Mehrer et al., 2021). It will also be important to investigate the role of the learning algorithms that the models use and their training objective, as both likely affect the representations that the models learn (e.g., see Zhuang et al., 2022, for evidence from vision). Specifically, Zhuang et al. (2022) showed that in an object categorization task, the negative sampling objective function, which maximizes the similarity between objects in the same category while minimizing the similarity between objects in different categories in the internal representation of the model, can alleviate model failure in capturing human visual behavior, which occurs under the standard objective function. This failure is due to the presence of categories that are infrequent in the training data, and this finding can be relevant for language, which also contains infrequent elements (words and constructions) amid more common ones.

Another issue that we did not investigate here is the nature of various training parameters (e.g., learning rate, batch size, randomization of context, length of context, etc.). Such parameters can affect model performance on NLP tasks and, possibly, their ability to predict human neural or behavioral responses to language. However, we suspect that the influence of these parameters would be relatively minimal given the evidence from prior work that model size, dataset size, and amount of computing power are the main contributors to model performance after training, as measured in loss for predicting the next token (Kaplan et al., 2020), and that model size is the main contributor to model performance in predicting fMRI responses to language (Antonello et al., 2023).

Another aspect of the ANN models that is important for building accurate models of human language processing is the model architecture. We here generalized our training effects across uni- and bidirectional-attention transformers, but a systematic investigation of the effects of diverse architectural parameters (e.g., the number and size of layers, number of attention heads) on the models' ability to predict human responses to language would be valuable. Tightly controlled comparisons between different classes of model architectures are more challenging but creating numerous model variants all trained on the same dataset (e.g., Storrs et al., 2021) could enable identification of architectural motifs that are essential for a good match with human neural and behavioral data.

Perhaps the biggest limitation of this and related work is the obscurity of both the model representations and human neural representations. It is not known what aspects of the representations change as the models are trained on increasingly more data (aside from knowing that these changes lead to improved performance on the next-word prediction task), and how exactly these changes in the models' representations of linguistic input impact their ability to predict brain activation or behavioral processing difficulty. Some recent work has begun to attempt isolating the aspects of model representations that affect model-to-brain alignment. For example, Kauf et al. (2024) performed a series of experiments where model representations were obtained for different perturbations of a linguistic stimulus (e.g., scrambling the word order or dropping/replacing some of the words) and then related to neural representations of an intact stimulus in order to see which perturbations affect model representations negatively. They found that word-level and compositional semantic information appears to be more important than information related to the syntactic structure in the model-to-brain alignment. Still however, much about the details of how models vs. humans represent and process linguistic stimuli remains to be discovered. Another exciting recent approach (Sexton & Love, 2022) is to replace a part of a model (the one best aligned with human neural

responses) with fMRI signals to test whether parts of the model representation that align with neural data affect model behavior.

In future work, we aim to address these gaps in order to build increasingly more accurate and interpretable models of language processing in the brain.

ACKNOWLEDGMENTS

We are grateful to Josh McDermott and members of the Fedorenko Lab (especially Carina Kauf, Cory Shain, Greta Tuckute, and Chengxu Zhuang) for helpful discussions and comments on the drafts of the manuscript; to Stella Biderman (EleutherAI) for help with the setup of Experiment 1; and to Jason Bolton, Laurel Orr, and Siddharth Karamcheti for help with the setup of Experiment 2.

FUNDING INFORMATION

Evelina Fedorenko, National Institute of Neurological Disorders and Stroke (<https://dx.doi.org/10.13039/1000000065>), Award ID: U01-NS121471. Eghbal A. Hosseini, McGovern Institute for Brain Research, Massachusetts Institute of Technology (<https://dx.doi.org/10.13039/100019335>), Award ID: Friends of McGovern graduate fellowship. Martin Schrimpf, McGovern Institute for Brain Research, Massachusetts Institute of Technology (<https://dx.doi.org/10.13039/100019335>), Award ID: Friends of McGovern graduate fellowship. Noga Zaslavsky, McGovern Institute for Brain Research, Massachusetts Institute of Technology (<https://dx.doi.org/10.13039/100019335>), Award ID: K. Lisa Young ICoN Center post-doctoral fellowship. Evelina Fedorenko, National Institute on Deafness and Other Communication Disorders (<https://dx.doi.org/10.13039/100000055>), Award ID: DC016607. Evelina Fedorenko, National Institute on Deafness and Other Communication Disorders (<https://dx.doi.org/10.13039/100000055>), Award ID: R01-DC016950. Evelina Fedorenko, McGovern Institute for Brain Research, Massachusetts Institute of Technology (<https://dx.doi.org/10.13039/100019335>). Evelina Fedorenko, Simons Center for the Social Brain, Massachusetts Institute of Technology (<https://dx.doi.org/10.13039/100018792>). Evelina Fedorenko, Massachusetts Institute of Technology, Award ID: Middleton Professorship.

AUTHOR CONTRIBUTIONS

Eghbal A. Hosseini: Conceptualization: Lead; Data curation: Lead; Formal analysis: Lead; Investigation: Lead; Methodology: Lead; Software: Lead; Validation: Lead; Visualization: Lead; Writing – original draft: Lead; Writing – review & editing: Lead. **Martin Schrimpf:** Methodology: Supporting; Software: Supporting; Writing – review & editing: Supporting. **Yian Zhang:** Methodology: Supporting; Resources: Supporting; Software: Supporting. **Samuel Bowman:** Resources: Supporting; Writing – review & editing: Supporting. **Noga Zaslavsky:** Conceptualization: Equal; Supervision: Supporting; Writing – review & editing: Supporting. **Evelina Fedorenko:** Conceptualization: Equal; Investigation: Supporting; Writing – original draft: Equal; Writing – review & editing: Equal.

DATA AND CODE AVAILABILITY STATEMENT

The human benchmarks and the code for relating model representations to the benchmarks are publicly available at <https://github.com/mschrimpf/neural-nlp>. The code for reproducing the main figures in this study is available at https://github.com/eghbalhosseini/ann_brain_alignment and accompanying data are available at <https://osf.io/6bd85/>. For Experiment 1,

the GPT-2 models and the representations extracted for all the benchmarks are available at <https://osf.io/6bd85/>; the miniBERTa models are available upon request; the checkpoints are available at <https://huggingface.co/nyu-ml>. (The training corpora used for Experiment 1 have copyright restrictions so cannot be made publicly available.) For Experiment 2, the model checkpoints are available at <https://huggingface.co/stanford-crfm>; the training corpus is available at <https://huggingface.co/datasets/Skylion007/openwebtext>.

REFERENCES

- Aminabadi, R. Y., Rajbhandari, S., Zhang, M., Awan, A. A., Li, C., Li, D., Zheng, E., Rasley, J., Smith, S., Ruwase, O., & He, Y. (2022). DeepSpeed Inference: Enabling efficient inference of transformer models at unprecedented scale. *ArXiv*. <https://doi.org/10.48550/arXiv.2207.00032>
- Antonello, R., & Huth, A. (2024). Predictive coding or just feature discovery? An alternative account of why language models fit brain data. *Neurobiology of Language*, *5*(1), 64–79. https://doi.org/10.1162/nol_a_00087
- Antonello, R., Vaidya, A., & Huth, A. G. (2023). Scaling laws for language encoding models in fMRI. *ArXiv*. <https://doi.org/10.48550/arXiv.2305.11863>
- Aurnhammer, C., & Frank, S. L. (2019). Evaluating information-theoretic measures of word prediction in naturalistic sentence reading. *Neuropsychologia*, *134*, Article 107198. <https://doi.org/10.1016/j.neuropsychologia.2019.107198>, PubMed: 31553896
- Baldassano, C., Chen, J., Zadbood, A., Pillow, J. W., Hasson, U., & Norman, K. A. (2017). Discovering event structure in continuous narrative perception and memory. *Neuron*, *95*(3), 709–721. <https://doi.org/10.1016/j.neuron.2017.06.041>, PubMed: 28772125
- Belinkov, Y., Durrani, N., Dalvi, F., Sajjad, H., & Glass, J. (2017). What do neural machine translation models learn about morphology? *ArXiv*. <https://doi.org/10.48550/arXiv.1704.03471>
- Black, S., Biderman, S., Hallahan, E., Anthony, Q., Gao, L., Golding, L., He, H., Leahy, C., McDonell, K., Phang, J., Pieler, M., Prashanth, U. S., Purohit, S., Reynolds, L., Tow, J., Wang, B., & Weinbach, S. (2022). GPT-NeoX-20B: An open-source autoregressive language model. *ArXiv*. <https://doi.org/10.48550/arXiv.2204.06745>
- Blank, I. A., & Fedorenko, E. (2020). No evidence for differences among language regions in their temporal receptive windows. *NeuroImage*, *219*, Article 116925. <https://doi.org/10.1016/j.neuroimage.2020.116925>, PubMed: 32407994
- Brennan, J. R., Stabler, E. P., Van Wagenen, S. E., Luh, W.-M., & Hale, J. T. (2016). Abstract linguistic structure correlates with temporal activity during naturalistic comprehension. *Brain and Language*, *157–158*, 81–94. <https://doi.org/10.1016/j.bandl.2016.04.008>, PubMed: 27208858
- Brodbeck, C., Bhattasali, S., Cruz Heredia, A. A. L., Resnik, P., Simon, J. Z., & Lau, E. (2022). Parallel processing in speech perception with local and global representations of linguistic context. *eLife*, *11*, Article e72056. <https://doi.org/10.7554/eLife.72056>, PubMed: 35060904
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., ... Amodei, D. (2020). Language models are few-shot learners. *ArXiv*. <https://doi.org/10.48550/arXiv.2005.14165>
- Buckner, R. L., & DiNicola, L. M. (2019). The brain's default network: Updated anatomy, physiology and evolving insights. *Nature Reviews Neuroscience*, *20*(10), 593–608. <https://doi.org/10.1038/s41583-019-0212-7>, PubMed: 31492945
- Button, K. S., Ioannidis, J. P. A., Mokrysz, C., Nosek, B. A., Flint, J., Robinson, E. S. J., & Munafò, M. R. (2013). Power failure: Why small sample size undermines the reliability of neuroscience. *Nature Reviews Neuroscience*, *14*(5), 365–376. <https://doi.org/10.1038/nrn3475>, PubMed: 23571845
- Caucheteux, C., & King, J.-R. (2022). Brains and algorithms partially converge in natural language processing. *Communications Biology*, *5*(1), Article 134. <https://doi.org/10.1038/s42003-022-03036-1>, PubMed: 35173264
- Chang, T. A., & Bergen, B. K. (2021). Word acquisition in neural language models. *ArXiv*. <https://doi.org/10.48550/arXiv.2110.02406>
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. *ArXiv*. <https://doi.org/10.48550/arXiv.1810.04805>
- Dupoux, E. (2018). Cognitive science in the era of artificial intelligence: A roadmap for reverse-engineering the infant language-learner. *Cognition*, *173*, 43–59. <https://doi.org/10.1016/j.cognition.2017.11.008>, PubMed: 29324240
- Fedorenko, E., Behr, M. K., & Kanwisher, N. (2011). Functional specificity for high-level linguistic processing in the human brain. *Proceedings of the National Academy of Sciences of the United States of America*, *108*(39), 16428–16433. <https://doi.org/10.1073/pnas.1112937108>, PubMed: 21885736
- Fedorenko, E., Blank, I. A., Siegelman, M., & Mineroff, Z. (2020). Lack of selectivity for syntax relative to word meanings throughout the language network. *Cognition*, *203*, Article 104348. <https://doi.org/10.1016/j.cognition.2020.104348>, PubMed: 32569894
- Fedorenko, E., Hsieh, P.-J., Nieto-Castañón, A., Whitfield-Gabrieli, S., & Kanwisher, N. (2010). New method for fMRI investigations of language: Defining ROIs functionally in individual subjects. *Journal of Neurophysiology*, *104*(2), 1177–1194. <https://doi.org/10.1152/jn.00032.2010>, PubMed: 20410363
- Frank, M. C. (2023). Bridging the data gap between children and large language models. *PsyArXiv*. <https://doi.org/10.31234/osf.io/qzbgx>
- Frank, S. L., Otten, L. J., Galli, G., & Vigliocco, G. (2015). The ERP response to the amount of information conveyed by words in sentences. *Brain and Language*, *140*, 1–11. <https://doi.org/10.1016/j.bandl.2014.10.006>, PubMed: 25461915
- Friederici, A. D. (2018). The neural basis for human syntax: Broca's area and beyond. *Current Opinion in Behavioral Sciences*, *21*, 88–92. <https://doi.org/10.1016/j.cobeha.2018.03.004>
- Futrell, R., Gibson, E., Tily, H. J., Blank, I., Vishnevetsky, A., Piantadosi, S., & Fedorenko, E. (2018). The natural stories corpus. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)* (pp. 76–82). European Language Resources Association.

- Gauthier, J., & Levy, R. (2019). Linking artificial and human neural representations of language. *ArXiv*. <https://doi.org/10.48550/arXiv.1910.01244>
- Geiger, F., Schrimpf, M., Marques, T., & DiCarlo, J. J. (2020). Wiring up vision: Minimizing supervised synaptic updates needed to produce a primate ventral stream. *BioRxiv*. <https://doi.org/10.1101/2020.06.08.140111>
- Gilkerson, J., Richards, J. A., Warren, S. F., Montgomery, J. K., Greenwood, C. R., Kimbrough Oller, D., Hansen, J. H. L., & Paul, T. D. (2017). Mapping the early language environment using all-day recordings and automated analysis. *American Journal of Speech-Language Pathology*, 26(2), 248–265. https://doi.org/10.1044/2016_AJSLP-15-0169, PubMed: 28418456
- Gokaslan, A., & Cohen, V. (2019). *OpenWebText corpus* [Dataset].
- Goldstein, A., Zada, Z., Buchnik, E., Schain, M., Price, A., Aubrey, B., Nastase, S. A., Feder, A., Emanuel, D., Cohen, A., Jansen, A., Gazula, H., Choe, G., Rao, A., Kim, C., Casto, C., Fanda, L., Doyle, W., Friedman, D., ... Hasson, U. (2022). Shared computational principles for language processing in humans and deep language models. *Nature Neuroscience*, 25(3), 369–380. <https://doi.org/10.1038/s41593-022-01026-4>, PubMed: 35260860
- Hagoort, P. (2019). The neurobiology of language beyond single-word processing. *Science*, 366(6461), 55–58. <https://doi.org/10.1126/science.aax0289>, PubMed: 31604301
- Hart, B., & Risley, T. R. (1992). American parenting of language-learning children: Persisting differences in family-child interactions observed in natural home environments. *Developmental Psychology*, 28(6), 1096–1105. <https://doi.org/10.1037/0012-1649.28.6.1096>
- Heilbron, M., Armeni, K., Schoffelen, J.-M., Hagoort, P., & de Lange, F. P. (2022). A hierarchy of linguistic predictions during natural language comprehension. *Proceedings of the National Academy of Sciences of the United States of America*, 119(32), Article e2201968119. <https://doi.org/10.1073/pnas.2201968119>, PubMed: 35921434
- Henderson, J. M., Choi, W., Lowder, M. W., & Ferreira, F. (2016). Language structure in the brain: A fixation-related fMRI study of syntactic surprisal in reading. *NeuroImage*, 132, 293–300. <https://doi.org/10.1016/j.neuroimage.2016.02.050>, PubMed: 26908322
- Hewitt, J., & Manning, C. D. (2019). A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)* (pp. 4129–4138). Association for Computational Linguistics. <https://doi.org/10.18653/v1/N19-1419>
- Hu, J., Gauthier, J., Qian, P., Wilcox, E., & Levy, R. P. (2020). A systematic assessment of syntactic generalization in neural language models. *ArXiv*. <https://doi.org/10.48550/arXiv.2005.03692>
- Huebner, P. A., & Willits, J. A. (2021). Scaffolding input promotes atomic organization in the recurrent neural network language model. In *Proceedings of the 25th Conference on Computational Natural Language Learning* (pp. 408–422). Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.conll-1.32>
- Ioannidis, J. P. A., Munafò, M. R., Fusar-Poli, P., Nosek, B. A., & David, S. P. (2014). Publication and other reporting biases in cognitive sciences: Detection, prevalence, and prevention. *Trends in Cognitive Sciences*, 18(5), 235–241. <https://doi.org/10.1016/j.tics.2014.02.010>, PubMed: 24656991
- Jain, S., & Huth, A. (2018). Incorporating context into language encoding models for fMRI. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), *Proceedings of the 32nd International Conference on Neural Information Processing Systems* (pp. 6628–6637). Curran Associates.
- Jelinek, F., Mercer, R. L., Bahl, L. R., & Baker, J. K. (1977). Perplexity—A measure of the difficulty of speech recognition tasks. *Journal of the Acoustical Society of America*, 62(S1), S63. <https://doi.org/10.1121/1.2016299>
- Just, M. A., Carpenter, P. A., & Woolley, J. D. (1982). Paradigms and processes in reading comprehension. *Journal of Experimental Psychology: General*, 111(2), 228–238. <https://doi.org/10.1037/0096-3445.111.2.228>, PubMed: 6213735
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., & Amodei, D. (2020). Scaling laws for neural language models. *ArXiv*. <https://doi.org/10.48550/arXiv.2001.08361>
- Kauf, C., Tuckute, G., Levy, R., Andreas, J., & Fedorenko, E. (2024). Lexical-semantic content, not syntactic structure, is the main contributor to ANN-brain similarity of fMRI responses in the language network. *Neurobiology of Language*, 5(1), 7–42. https://doi.org/10.1162/nol_a_00116
- Kumar, S., Sumers, T. R., Yamakoshi, T., Goldstein, A., Hasson, U., Norman, K. A., Griffiths, T. L., Hawkins, R. D., & Nastase, S. A. (2022). Reconstructing the cascade of language processing in the brain using the internal computations of a transformer-based language model. *BioRxiv*. <https://doi.org/10.1101/2022.06.08.495348>
- Lerner, Y., Honey, C. J., Silbert, L. J., & Hasson, U. (2011). Topographic mapping of a hierarchy of temporal receptive windows using a narrated story. *Journal of Neuroscience*, 31(8), 2906–2915. <https://doi.org/10.1523/JNEUROSCI.3684-10.2011>, PubMed: 21414912
- Linzen, T., & Leonard, B. (2018). Distinct patterns of syntactic agreement errors in recurrent networks and humans. *ArXiv*. <https://doi.org/10.48550/arXiv.1807.06882>
- Lipkin, B., Tuckute, G., Affourtit, J., Small, H., Mineroff, Z., Kean, H., Jouravlev, O., Rakocevic, L., Pritchett, B., Siegelman, M., Hoeflin, C., Pongos, A., Blank, I. A., Struhl, M. K., Ivanova, A., Shannon, S., Sathe, A., Hoffmann, M., Nieto-Castañón, A., & Fedorenko, E. (2022). Probabilistic atlas for the language network based on precision fMRI data from >800 individuals. *Scientific Data*, 9(1), Article 529. <https://doi.org/10.1038/s41597-022-01645-3>, PubMed: 36038572
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *ArXiv*. <https://doi.org/10.48550/arXiv.1907.11692>
- Lopopolo, A., Frank, S. L., van den Bosch, A., & Willems, R. M. (2017). Using stochastic language models (SLM) to map lexical, syntactic, and phonological information processing in the brain. *PLOS ONE*, 12(5), Article e0177794. <https://doi.org/10.1371/journal.pone.0177794>, PubMed: 28542396
- Marcus, M. P., Santorini, B., & Marcinkiewicz, M. A. (1993). Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2), 313–330.
- Mehrer, J., Spoerer, C. J., Jones, E. C., Kriegeskorte, N., & Kietzmann, T. C. (2021). An ecologically motivated image dataset for deep learning yields better models of human vision. *Proceedings of the National Academy of Sciences of the United States of America*, 118(8), Article e2011417118. <https://doi.org/10.1073/pnas.2011417118>, PubMed: 33593900
- Mehrer, J., Spoerer, C. J., Kriegeskorte, N., & Kietzmann, T. C. (2020). Individual differences among deep neural network models. *Nature Communications*, 11(1), Article 5725. <https://doi.org/10.1038/s41467-020-19632-w>, PubMed: 33184286
- Merity, S., Xiong, C., Bradbury, J., & Socher, R. (2016). Pointer sentinel mixture models. *ArXiv*. <https://doi.org/10.48550/arXiv.1609.07843>

- Oh, B.-D., & Schuler, W. (2022). Why does surprisal from larger transformer-based language models provide a poorer fit to human reading times? *ArXiv*. <https://doi.org/10.48550/arXiv.2212.12131>
- Oh, B.-D., & Schuler, W. (2023). Transformer-based language model surprisal predicts human reading times best with about two billion training tokens. *ArXiv*. <https://doi.org/10.48550/arXiv.2304.11389>
- Pallier, C., Devauchelle, A.-D., & Dehaene, S. (2011). Cortical representation of the constituent structure of sentences. *Proceedings of the National Academy of Sciences of the United States of America*, 108(6), 2522–2527. <https://doi.org/10.1073/pnas.1018711108>, PubMed: 21224415
- Pannitto, L., & Herbelot, A. (2020). Recurrent babbling: Evaluating the acquisition of grammar from limited input data. In *Proceedings of the 24th Conference on Computational Natural Language Learning* (pp. 165–176). Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.conll-1.13>
- Pasquiou, A., Lakretz, Y., Hale, J., Thirion, B., & Pallier, C. (2022). Neural language models are not born equal to fit brain data, but training helps. *ArXiv*. <https://doi.org/10.48550/arXiv.2207.03380>
- Pereira, F., Lou, B., Pritchett, B., Ritter, S., Gershman, S. J., Kanwisher, N., Botvinick, M., & Fedorenko, E. (2018). Toward a universal decoder of linguistic meaning from brain activation. *Nature Communications*, 9(1), Article 963. <https://doi.org/10.1038/s41467-018-03068-4>, PubMed: 29511192
- Pérez-Mayos, L., Ballesteros, M., & Wanner, L. (2021). How much pretraining data do language models need to learn syntax? *ArXiv*. <https://doi.org/10.48550/arXiv.2109.03160>
- Poldrack, R. A., Baker, C. I., Durnez, J., Gorgolewski, K. J., Matthews, P. M., Munafò, M. R., Nichols, T. E., Poline, J.-B., Vul, E., & Yarkoni, T. (2017). Scanning the horizon: Towards transparent and reproducible neuroimaging research. *Nature Reviews Neuroscience*, 18(2), 115–126. <https://doi.org/10.1038/nrn.2016.167>, PubMed: 28053326
- Rabovsky, M., Hansen, S. S., & McClelland, J. L. (2018). Modelling the N400 brain potential as change in a probabilistic representation of meaning. *Nature Human Behaviour*, 2(9), 693–705. <https://doi.org/10.1038/s41562-018-0406-4>, PubMed: 31346278
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. In *Better language models and their implications*. OpenAI Blog. https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf
- Rohde, D. L., & Plaut, D. C. (1999). Language acquisition in the absence of explicit negative evidence: How important is starting small? *Cognition*, 72(1), 67–109. [https://doi.org/10.1016/S0010-0277\(99\)00031-1](https://doi.org/10.1016/S0010-0277(99)00031-1), PubMed: 10520565
- Schrimpf, M., Blank, I. A., Tuckute, G., Kauf, C., Hosseini, E. A., Kanwisher, N., Tenenbaum, J. B., & Fedorenko, E. (2021). The neural architecture of language: Integrative modeling converges on predictive processing. *Proceedings of the National Academy of Sciences of the United States of America*, 118(45), Article e2105646118. <https://doi.org/10.1073/pnas.2105646118>, PubMed: 34737231
- Sexton, N. J., & Love, B. C. (2022). Reassessing hierarchical correspondences between brain and deep networks through direct interface. *Science Advances*, 8(28), Article eabm2219. <https://doi.org/10.1126/sciadv.abm2219>, PubMed: 35857493
- Shain, C., Blank, I. A., van Schijndel, M., Schuler, W., & Fedorenko, E. (2020). fMRI reveals language-specific predictive coding during naturalistic sentence comprehension. *Neuropsychologia*, 138, Article 107307. <https://doi.org/10.1016/j.neuropsychologia.2019.107307>, PubMed: 31874149
- Shain, C., Kean, H., Casto, C., Lipkin, B., Affourtit, J., Siegelman, M., Mollica, F., & Fedorenko, E. (2023). Graded sensitivity to structure and meaning throughout the human language network. *BioRxiv*. <https://doi.org/10.1101/2021.11.12.467812>
- Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science*, 22(11), 1359–1366. <https://doi.org/10.1177/0956797611417632>, PubMed: 22006061
- Simony, E., Honey, C. J., Chen, J., Lositsky, O., Yeshurun, Y., Wiesel, A., & Hasson, U. (2016). Dynamic reconfiguration of the default mode network during narrative comprehension. *Nature Communications*, 7, Article 12141. <https://doi.org/10.1038/ncomms12141>, PubMed: 27424918
- Smith, N. J., & Levy, R. (2013). The effect of word predictability on reading time is logarithmic. *Cognition*, 128(3), 302–319. <https://doi.org/10.1016/j.cognition.2013.02.013>, PubMed: 23747651
- Storrs, K. R., Kietzmann, T. C., Walther, A., Mehrer, J., & Kriegeskorte, N. (2021). Diverse deep neural networks all predict human inferior temporal cortex well, after training and fitting. *Journal of Cognitive Neuroscience*, 33(10), 2044–2064. https://doi.org/10.1162/jocn_a_01755, PubMed: 34272948
- Tenney, I., Das, D., & Pavlick, E. (2019). BERT rediscovers the classical NLP pipeline. *ArXiv*. <https://doi.org/10.48550/arXiv.1905.05950>
- Toneva, M., & Wehbe, L. (2019). Interpreting and improving natural-language processing (in machines) with natural language-processing (in the brain). In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, & R. Garnett (Eds.), *Proceedings of the 33rd International Conference on Neural Information Processing Systems* (pp. 14954–14964). Curran Associates.
- Tuckute, G., Sathe, A., Srikant, S., Taliaferro, M., Wang, M., Schrimpf, M., Kay, K., & Fedorenko, E. (2024). Driving and suppressing the human language network using large language models. *Nature Human Behaviour*, 8(3), 544–561. <https://doi.org/10.1038/s41562-023-01783-7>, PubMed: 38172630
- van Schijndel, M., Mueller, A., & Linzen, T. (2019). Quantity doesn't buy quality syntax with neural language models. *ArXiv*. <https://doi.org/10.48550/arXiv.1909.00111>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. *ArXiv*. <https://doi.org/10.48550/arXiv.1706.03762>
- Warstadt, A., & Bowman, S. R. (2019). Linguistic analysis of pre-trained sentence encoders with acceptability judgments. *ArXiv*. <https://doi.org/10.48550/arXiv.1901.03438>
- Warstadt, A., & Bowman, S. R. (2022). What artificial neural networks can tell us about human language acquisition. *ArXiv*. <https://doi.org/10.48550/arXiv.2208.07998>
- Wilcox, E. G., Gauthier, J., Hu, J., Qian, P., & Levy, R. (2020). On the predictive power of neural language models for human real-time comprehension behavior. *ArXiv*. <https://doi.org/10.48550/arXiv.2006.01912>
- Willems, R. M., Frank, S. L., Nijhof, A. D., Hagoort, P., & van den Bosch, A. (2016). Prediction during natural language comprehension. *Cerebral Cortex*, 26(6), 2506–2516. <https://doi.org/10.1093/cercor/bhv075>, PubMed: 25903464
- Zador, A. M. (2019). A critique of pure learning and what artificial neural networks can learn from animal brains. *Nature*

- Communications*, 10(1), Article 3770. <https://doi.org/10.1038/s41467-019-11786-6>, PubMed: 31434893
- Zhang, Y., Liu, H., Li, H.-S., Warstadt, A., & Bowman, S. R. (2020). *The MiniBERTs: Testing what RoBERTa learns with varying amounts of pretraining*. *Cilvr at NYU*. <https://wp.nyu.edu/cilvr/2020/07/02/the-minibertas-testing-what-roberta-learns-with-varying-amounts-of-pretraining/>
- Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., & Fidler, S. (2015). Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. *ArXiv*. <https://doi.org/10.48550/arXiv.1506.06724>
- Zhuang, C., Xiang, V., Bai, Y., Jia, X., Turk-Browne, N., Norman, K., DiCarlo, J. J., & Yamins, D. L. K. (2022). How well do unsupervised learning algorithms model human real-time and life-long learning? In *Advances in Neural Information Processing Systems 35: 36th Conference on Neural Information Processing Systems (NeurIPS 2022)* (pp. 22628–22642). Curran Associates.