

---

# Contextual Feature Extraction Hierarchies Converge in Large Language Models and the Brain

---

**Gavin Mischler\***

Department of Electrical Engineering  
Columbia University  
New York, NY 10027  
gm2944@columbia.edu

**Yinghao Aaron Li\***

Department of Electrical Engineering  
Columbia University  
New York, NY 10027  
y14579@columbia.edu

**Stephan Bickel**

The Feinstein Institutes for Medical Research  
Northwell Health  
Manhasset, NY 11030  
sbickel@northwell.edu

**Ashesh D. Mehta**

The Feinstein Institutes for Medical Research  
Northwell Health  
Manhasset, NY 11030  
amehta@northwell.edu

**Nima Mesgarani†**

Department of Electrical Engineering  
Columbia University  
New York, NY 10027  
nima@ee.columbia.edu

## Abstract

Recent advancements in artificial intelligence have sparked interest in the parallels between large language models (LLMs) and human neural processing, particularly in language comprehension. While prior research has established similarities in the representation of LLMs and the brain, the underlying computational principles that cause this convergence, especially in the context of evolving LLMs, remain elusive. Here, we examined a diverse selection of high-performance LLMs with similar parameter sizes to investigate the factors contributing to their alignment with the brain’s language processing mechanisms. We find that as LLMs achieve higher performance on benchmark tasks, they not only become more brain-like as measured by higher performance when predicting neural responses from LLM embeddings, but also their hierarchical feature extraction pathways map more closely onto the brain’s while using fewer layers to do the same encoding. We also compare the feature extraction pathways of the LLMs to each other and identify new ways in which high-performing models have converged toward similar hierarchical processing mechanisms. Finally, we show the importance of contextual information in improving model performance and brain similarity. Our findings reveal the converging aspects of language processing in the brain and large language models and offer new directions for developing LLMs that align more closely with human cognitive processing.

---

\*These authors contributed equally to this work

†Corresponding author

# 1 Introduction

The intersection of artificial intelligence and neuroscience has emerged as a frontier of great interest, particularly in understanding how large language models (LLMs) and the human brain process language. Prior research has laid foundational work in this area, uncovering intriguing parallels in feature extraction and representational similarities between LLMs and neural responses during language processing. Studies [1, 2, 3, 4, 5, 6, 7, 8, 9] have demonstrated that the representations learned by LLMs can be linearly mapped to neural responses, suggesting that both LLMs and the brain utilize comparable features in language processing. However, these findings offer limited insight into the fundamental characteristics of LLMs that enable this brain-like processing.

Further investigations have delved into different aspects of LLMs to elucidate their resemblance to brain processes. Some studies [10, 11] support the predictive coding hypothesis in human language processing by finding stronger similarities with autoregressive LLMs. Others [8, 5, 12, 13] have explored various factors, such as the LLM language modeling performance, model size and capacity, and the generalizability of linguistic representations as indicators of brain-like processing. These studies imply that the quality of an LLM significantly contributes to its brain-like representations, yet the underlying reason for this similarity remain an open question. Is it merely a matter of scaling up the models [12], or do these models share fundamental computational principles that increasingly align well with the spoken language processing pathway in the human brain? This question is significant as it may suggest a potential shift in the paradigm of model optimization. Although both brains [14, 15, 16, 17] and LLMs [18, 19] process speech and language in hierarchical pathways, most studies have analyzed the similarity of their representations without detailed comparisons of the hierarchical processes through which they are created. Thus, it is still unclear whether brains and models arrive at similar representations through the same or different pathways.

We aim to answer these questions by examining the interplay between LLM performance, neural predictability, anatomical alignment, and contextual encoding, potentially paving the way toward models that perform with high accuracy and process language in a manner similar to the human brain. We examined 12 open-source, pre-trained LLMs, all uniform in size but diverse in their linguistic capabilities, particularly in language-understanding tasks such as reading comprehension. We recorded neural responses with intracranial electroencephalography (iEEG) in the auditory cortex and speech processing regions of neurosurgical patients as they listened to speech. We then predicted these neural responses from the embeddings extracted from each layer of the LLMs as they processed the same linguistic input. This approach allowed us to pinpoint which layers and aspects of the LLMs were most predictive of brain activity and explore how variations in model performance align with differences in neural prediction and their anatomical and functional correspondence. Our findings offer a fresh perspective on the evolving landscape of LLMs, providing insights that reveal more intricate parallels in language comprehension between artificial and biological systems, uncover new potential reasons for LLM performance differences, and point to a convergence in LLMs towards a more optimal, brain-like language processing system.

## 2 Results

### 2.1 Brain Similarity of Large Language Models

We studied 12 recent, popular, open-source LLMs, all with approximately 7 billion parameters. We evaluated each model on a suite of benchmark tasks to assess its language modeling performance, splitting these tasks into categories relevant to English language comprehension, specifically reading comprehension and commonsense reasoning as in [20] (see Methods for details). Overall LLM performance was estimated as the average score over these two categories.

Neural responses were recorded with invasive electrodes (intracranial EEG) from eight neurosurgical patients, with electrode placement determined by clinical need (Supplementary Fig. 1). The subjects listened to between 20 and 30 minutes of speech from various talkers, including stories voiced by voice actors and dialogues between characters. The text for each audio was fed into each LLM, and we extracted the causal embeddings of each word at every layer. We reduced these embeddings to 500 components with PCA, ensuring consistent dimensionality across models since the models were only approximately the same size to begin with. We used ridge regression to estimate the similarity of a model’s features to the brain (Fig. 1) [10, 3, 8]. We analyzed 707 electrodes which were responsive

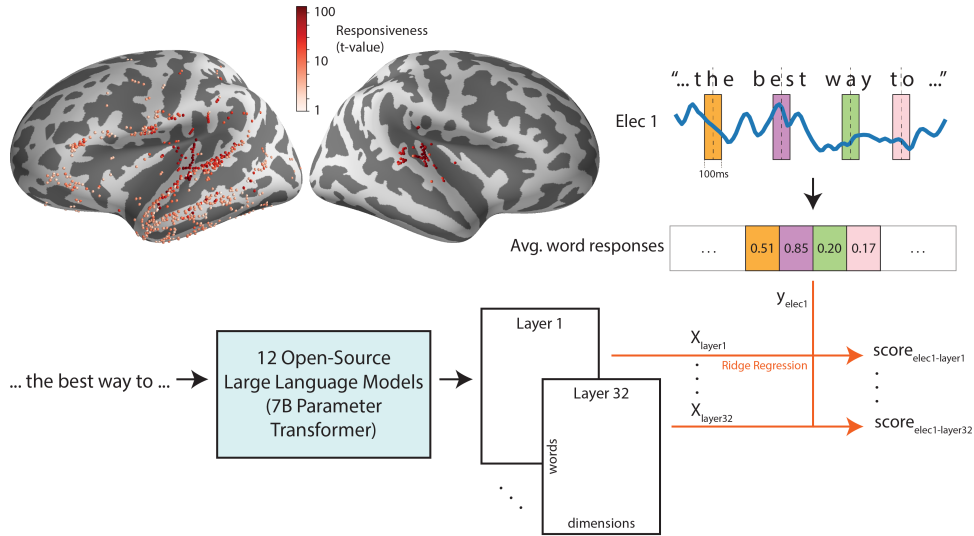


Figure 1: Mapping LLM embeddings to the brain. Speech responsive electrodes are shown on an inflated brain (shaded by their responsiveness t-value from a paired t-test between speech and silence). As subjects listened to speech, the average neural response in a 100ms window around a word center was used as a given electrode’s word response. The same text was fed to an LLM and the embeddings from all 32 layers were extracted. Ridge regression was used to predict the word responses from the LLM representations, producing a brain correlation score for each electrode-layer pair.

to speech, as determined by a t-test between responses to words and silence (FDR corrected,  $p < 0.05$  [21]). For each responsive electrode, we extracted the average high-gamma band envelope response in a 100ms window around the center of every word. Then, we fit cross-validated ridge regression models to predict these neural responses from the word embeddings and used the average prediction correlation on the withheld folds as the brain similarity with that electrode. Neither the number of principal components of the embeddings nor the window size used to compute the neural response to words significantly impacted the results (Supplementary Fig. 2).

Electrode-averaged brain similarity over each model’s layers is shown in Fig. 2A. With these latest LLMs, we confirm previous findings showing that neural responses can be predicted from model representations, and we find that brain similarity generally increases over layers and peaks in middle or later layers [3, 8]. Higher-performing LLMs also achieve higher peak brain scores (Pearson  $r = 0.92, p = 2.24 \times 10^{-5}$ ) (Fig. 2B), indicating that they extract more brain-like features from language.

Similar to the layers of a model, the auditory and language processing pathway demonstrates hierarchical organization [14, 22, 15, 16]. The primary auditory cortex, the first point of auditory processing in the cortex, is centered around posteromedial Heschl’s gyrus (pmHG, or TE1.1) [23]. Since this is a common reference point in auditory cortical processing, we quantify the depth of each electrode in the brain’s spoken language processing pathway using its distance from this landmark [24, 25, 26]. Prior studies have found that deeper layers of LLMs correspond better to deeper language processing regions of the brain [8, 11, 27]. We confirm this result (Fig. 2C), but interestingly, we also find that better-performing LLMs peak in brain similarity at earlier layers compared to worse models (Pearson  $r = -0.81, p = 0.0013$ ) (Fig. 2D). This uncovers a new dimension in the evolution of LLMs: the progression of feature extraction over layers aligns differently with the brain for higher-performing versus lower-performing models.

## 2.2 Alignment of Language Processing Hierarchies Between Models and the Brain

Given that the layer-wise brain similarity appears different between good and bad models, we hypothesized that better models were not only learning more brain-like features, but that the progression of feature extraction within these models was different. Taking inspiration from an investigation

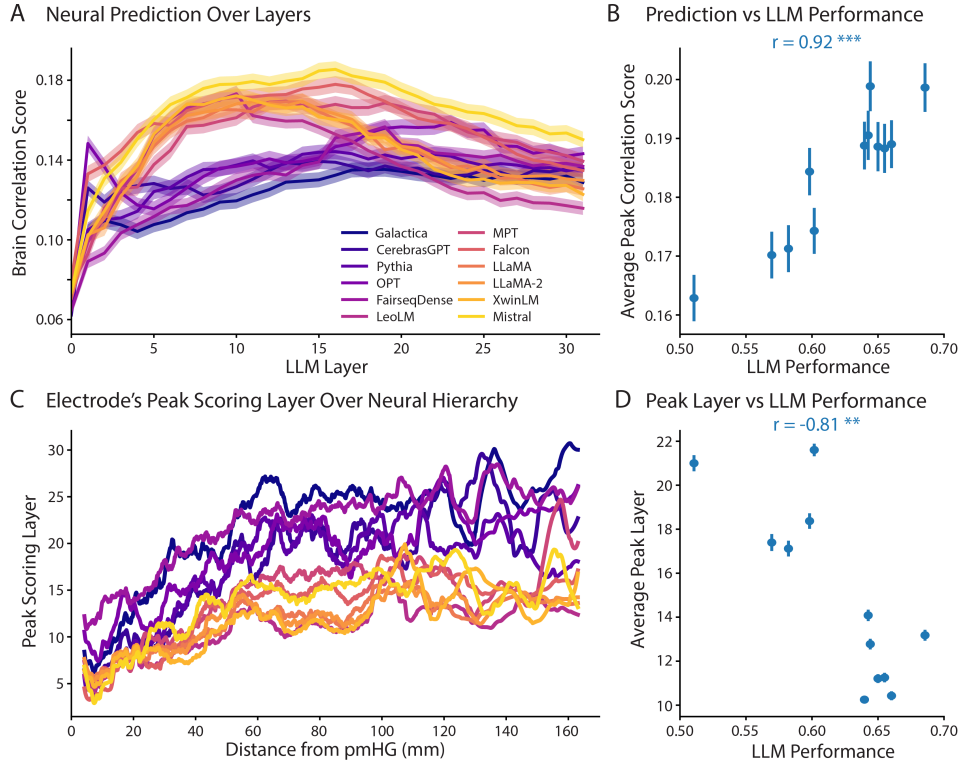


Figure 2: Peak brain correlations and layers relate to LLM performance. A) Average brain correlation over all electrodes for each LLM. LLMs are colored in order of their separately-measured benchmark performance, with blue/purple models performing the worst and yellow models performing the best. Shaded regions indicate standard error of the mean over electrodes. B) The peak correlation over all layers of a given model was computed for each electrode, then averaged over all electrodes. Bars indicate standard error of the mean over electrodes. Average peak correlation score is significantly related to LLM performance (Pearson  $r = 0.92$ ,  $p = 2.24 \times 10^{-5}$ ). Stars indicate statistical significance level thresholds of  $p < 0.05$ ,  $p < 0.01$ , and  $p < 0.001$  with \*, \*\*, and \*\*\*, respectively. C) The peak scoring layer of each model was computed for each electrode. Then electrodes were sorted by distance from pmHG and a sliding window average (centered,  $n = 50$ ) was taken across the electrodes of each model to compute the smoothed, local estimate of the most brain-like LLM layer. The peak scoring layer generally increases with distance from pmHG, and the better models (yellow) peak at lower layers compared to the worse models (blue/purple). D) The average peak layer for a given model over all electrodes is shown with bars indicating standard error of the mean. Average peak layer is significantly negatively related to LLM performance (Pearson  $r = -0.81$ ,  $p = 0.0013$ ).

of hierarchical correspondence between stages of visual cortex processing and image classification networks [28], we sought to compute the alignment between hierarchical feature extraction pathways in brains and models. Although the brain’s exact hierarchical processing stages, analogous to layers of a model, are not perfectly known, we again used the distance from pmHG to quantify the stages of hierarchical processing. We grouped electrodes into bins at 10mm intervals. Then, for each electrode, we normalized the brain similarity scores over layers. Finally, we averaged these layer-wise scores over the electrodes in a bin, producing a set of layer scores which are shown as a single row of the alignment matrix in Fig. 3A. We used the center of mass of this average brain similarity score over layers within each electrode bin to quantify the LLM layer most similar to a given stage of the brain’s hierarchy. Then, we compared the progression of these most-similar LLM layers to the bin distances along the hierarchy, visually finding that some models achieve a more linear increase in LLM layers over bins. We summarize the alignment between the language processing hierarchies of each LLM and the brain using the Pearson correlation between the layer center of mass in each bin and the hierarchical stage of each bin (i.e. the distance of each bin from pmHG) [28]. We

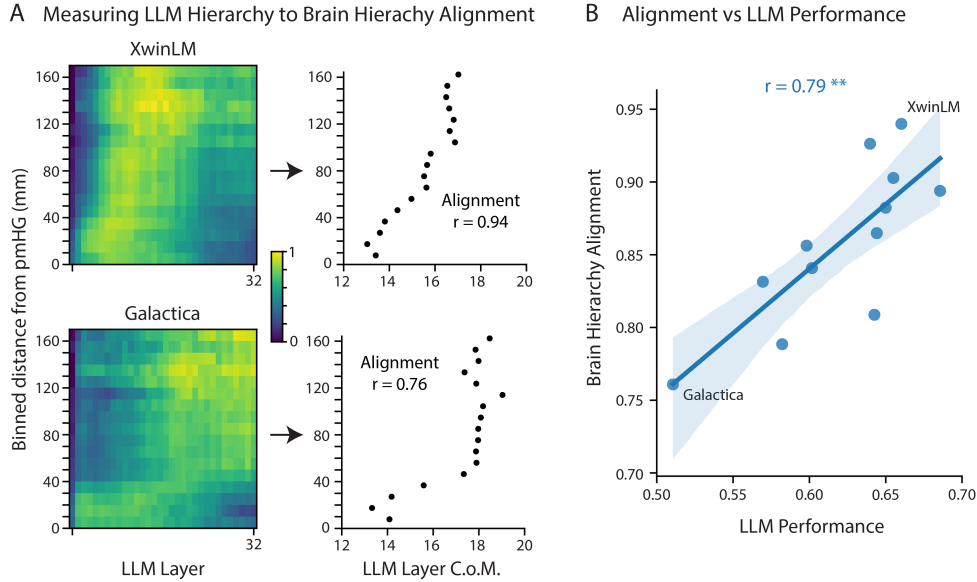


Figure 3: Better LLMs display more brain-like hierarchical processing. A) Examples of computing the brain hierarchy alignment are shown for two models: XwinLM (the model with the highest alignment score) and Galactica (the model with the lowest alignment score). Electrodes were first binned into a hierarchy by distance from pmHG. Within a bin, the correlations over all 32 layers were normalized between 0 and 1 and then averaged over electrodes in the bin, producing one row for each bin in the matrix on the left. The center of mass (C.o.M.) of the distribution of brain similarity scores over LLM layers for each bin was computed and plotted in the scatter plot to the right. The brain hierarchy alignment score was then computed as the Pearson correlation between LLM layer C.o.M. and distance from pmHG. B) A scatter plot of brain hierarchy alignment scores and LLM performance shows a significant positive correlation (Pearson  $r = 0.79$ ,  $p = 0.0021$ ,  $**$  indicates  $p < 0.01$ ). Line and shaded region shows linear regression fit and bootstrapped ( $n = 1000$ ) 95% confidence interval.

illustrate this alignment computation for XwinLM and Galactica, two models which achieve the highest and lowest hierarchy alignment scores, respectively (Fig. 3A). These models also display a stark difference in benchmark performance, with Galactica being the lowest performing LLM. The alignment scores reveal that the better model (XwinLM) exhibits a feature extraction progression more consistent with the brain from early to late-stage processing compared to the bad model. This brain alignment is also highly correlated with LLM performance on the benchmark evaluation tasks (Pearson  $r = 0.79$ ,  $p = 0.0021$ ) (Fig. 3B). We find the same result when using electrode latency to measure the stages of the brain’s hierarchical processing, rather than distance from pmHG (Pearson  $r = 0.89$ ,  $p = 0.0001$ ) (Supplementary Fig. 3), which demonstrates that this finding holds for other estimates of the stages of the cortical hierarchy. Additionally, to ensure this effect was not the result of a single subject overpowering the distribution, we separated the even- and odd-numbered subjects and performed the analysis again, finding that brain hierarchy alignment was significantly correlated with LLM performance for each group (Pearson correlation, even subjects:  $r = 0.79$ ,  $p = 0.0022$ , odd subjects:  $r = 0.81$ ,  $p = 0.0013$ ) (Supplementary Fig. 4). Overall, these findings demonstrate that better-performing LLMs extract features using a hierarchy that more linearly aligns with the brain’s hierarchical language processing pathway.

To perform model-to-model comparisons, we used centered kernel alignment (CKA) [29], a method analogous to canonical correlation analysis (CCA) with a nonlinear kernel, which is able to capture similarity between high-dimensional representations like neural network embeddings. We computed the CKA similarity between all pairs of layers for all pairs of models. Thus, each pair of models creates a layer-by-layer similarity matrix describing their embedding similarity. High similarity along the diagonal indicates that the two models extract similar features at the same layers. Higher similarity offset from the diagonal indicates that one model exhibits a delay in extracting similar

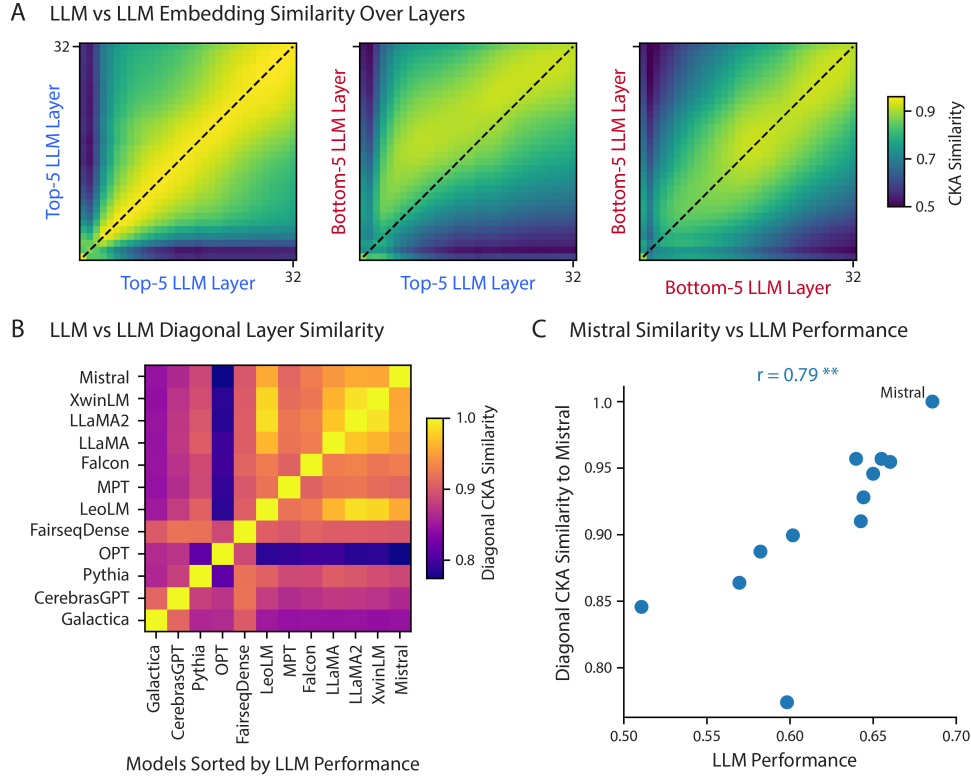


Figure 4: Comparing feature extraction hierarchies between LLMs. A) Layer-by-layer similarity matrices were computed using CKA for every pair of LLMs. LLMs were labeled as either top-5, bottom-5, or excluded, depending on their sorted performance on our LLM benchmark evaluation. Then, similarity matrices between all pairs of top-5 LLMs were averaged and displayed as the “top-5 versus top-5” average similarity matrix in the top left. The same was done to create similarity matrices between the average “top-5 versus bottom-5” as well as the average “bottom-5 versus bottom-5” LLMs. Visually, the “top-5 versus top-5” similarity matrix is highly diagonal, while “bottom-5 versus bottom-5” is less similar in early layers. The “top-5 versus bottom-5” similarity matrix shows an offset diagonal, indicating a delay in feature extraction for the bottom-5 models compared to the top-5. B) Model-by-model diagonal similarity was computed as the average along the diagonal of their similarity matrix. Models are arranged in sorted order of LLM benchmark performance from worst to best. This visually confirms that the best models are fairly similar to each other in layer-wise feature extraction, while worse models are less similar to each other and less similar to the best models. C) The diagonal similarity of each model with Mistral, the best performing LLM, is plotted against the LLM performance, showing a strong positive relationship (Pearson  $r = 0.79$ ,  $p = 0.0022$ , \*\* indicates  $p < 0.01$ ).

features. When grouping these similarity matrices by the top-5 and bottom-5 models based on LLM benchmark performance and averaging within a group, an interesting pattern emerges (Fig. 4A). We find that the top models exhibit a high degree of similarity to each other along the diagonal. On the other hand, the worst models are much less similar to each other in their early layers, and even in their later layers they are less consistent than the top-5-to-top-5 model pairs. Finally, comparing top-5 models to bottom-5 models reveals a striking offset in maximum similarity from the diagonal. This suggests that bad models require more layers to reach a similar level of feature extraction as good models. We summarize the layer-wise feature extraction similarity between each pair of models using the average CKA similarity along the diagonal in their CKA similarity matrix (Fig. 4B). The plot demonstrates that the top-5 models are indeed more similar to each other, with a sub-block of high similarity emerging among the top few models. Since Mistral is the best performing LLM, we look at the diagonal similarity to Mistral of each model and find that a more Mistral-like feature extraction progression correlates strongly with LLM performance (Pearson  $r = 0.79$ ,  $p = 0.0022$ )

(Fig. 4C). These results reveal new distinctions between the embeddings of LLMs and suggest that inefficient feature extraction or poor early-layer learning in bad models may contribute to their worse performance and lower brain similarity.

### 2.3 Contextual Content Supports Brain Hierarchy Alignment

Since the contextual nature of LLM features is critical for their brain similarity compared with non-contextual representations [10, 3, 30], we hypothesized that the amount of contextual information used by a model may also play a key role in determining the alignment between hierarchical feature extraction pathways of LLMs and the brain. We extracted limited-context embeddings from the LLMs by restricting their causal attention mechanism to a certain window of the previous text. Transformer architecture LLMs use tokenizers to separate text into discrete units, so we supplied the models with only the most recent  $N$  tokens, sweeping  $N$  over a range of values from 1 to 100. A single token input gives the model no context at all. For reference, Mistral’s tokenizer averages 1.15 tokens per word in our stimulus corpus. We then repeated our analysis of model-brain hierarchical alignment (as previously shown in Fig 3B) by computing the correlation between brain hierarchy alignment and LLM performance at each limited context window length. While this correlation is positive for all but the 1-token case, it is only significant for long contextual window lengths of 50 tokens and above (Fig. 5A). This suggests that the brain alignment of LLMs critically depends on the amount of contextual information the model is able to see, which then influences its hierarchical feature extraction mechanism.

Since the correlation between LLM performance and hierarchical alignment is strongly positive for long context lengths, we expected that better-performing models would be better at incorporating contextual information into their language representations. To test this, we quantified the amount of contextual information present in a model’s embeddings by measuring how much its embeddings changed when contextual information was added to the input. We measured the CKA difference ( $1 - \text{similarity}_{CKA}(\text{full-context}, \text{1-token})$ ) of the embeddings of each layer when given the full context compared to the first-layer embeddings when given only a 1-token limited context window. We refer to the average of this CKA difference over all layers as the contextual content of the model’s representations. We find that this contextual content is positively correlated with LLM performance (Spearman  $r = 0.66, p = 0.020$ ) (Fig. 5B). Additionally, it is very strongly correlated with brain similarity (Spearman  $r = 0.84, p = 0.0006$ ) (Fig. 5C). These findings indicate that contextual information plays a crucial role in natural language processing in both natural and artificial language models, and contextual feature extraction enables brain hierarchy alignment in LLMs.

We further investigated the impact of contextual information on neural similarity by computing how much each LLM’s peak similarity score with a given electrode changed when the models were given the full context versus no context (1 token). We then averaged this difference over all LLMs for each electrode and plotted the electrodes on the brain, finding that being given the extra context more greatly improved similarity scores with electrodes in higher-level language processing areas (Fig. 5D). Averaging electrodes within major anatomical regions further quantifies this result, as we find higher average context effects on brain correlation score within the higher-level linguistic-processing area of inferior frontal gyrus (IFG) [31] compared to sensory regions like Heschl’s gyrus (HG) and superior temporal gyrus (STG) (Wilcoxon rank-sum test,  $p < 0.05$ ) (Fig. 5E). Interestingly, the articulatory region of subcentral gyrus, which has also been implicated in high-level linguistic processing [32], displays the highest average score improvement, but due to its high variance it does not meet statistical significance. These results show that contextual information becomes more critical in determining brain similarity further along the spoken language processing hierarchy, which supports previous investigations of high-level linguistic feature encoding in more downstream regions [33, 34]. This finding strengthens the notion that both the brain and LLMs are extracting context along their hierarchies, and that LLMs need contextual information to achieve brain similarity in downstream processing regions. Taken together, our analyses reveal that high-performing LLMs not only extract representations of language that are similar to the brain, but they also use hierarchical feature extraction pathways which more strongly align with that of the brain due to contextual information processing abilities, a finding that uncovers new ways in which the best LLMs are continuously converging toward the brain.



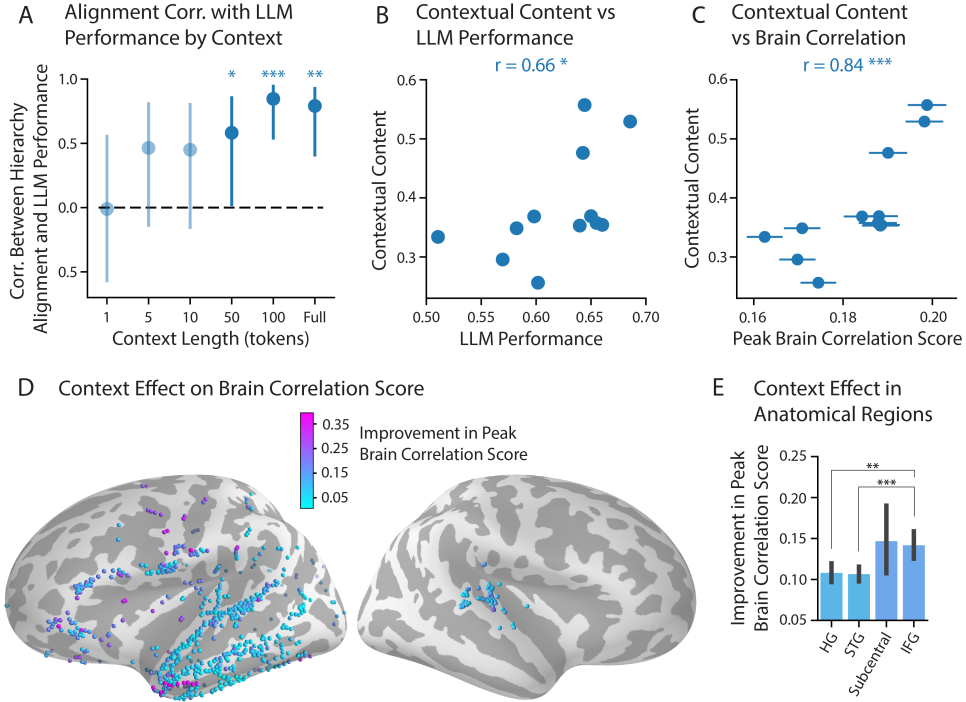


Figure 5: Effect of contextual information. A) Using embeddings from the LLMs when given a certain limited number of the previous tokens as a context window, we performed the analysis of brain hierarchy alignment again. The correlation between LLM performance and brain hierarchy alignment is illustrated by each dot with 95% confidence interval bars, showing only a significant correlation for long contextual windows. Stars illustrate the significance level of the correlation, with \*, \*\*, and \*\*\* indicating chance levels below 0.05, 0.01, and 0.001, respectively. B) The contextual content of a model’s representations is plotted against its benchmark performance, showing a positive correlation between the two (Spearman  $r = 0.66, p = 0.020$ ). C) Contextual content of each model is plotted against its average peak brain similarity over electrodes, showing a strong correlation (Spearman  $r = 0.84, p = 0.0006$ ). Horizontal lines show standard error of the mean over electrodes for brain similarity. D) Electrodes plotted on the FreeSurfer average inflated brain, colored by the effect of contextual information on peak brain similarity score. E) Bar plot of the average context effect on peak brain similarity score for electrodes within four main anatomical regions along the linguistic hierarchy. Each bar is colored by its value according to the same colormap as used on the brain plot, and error bars show standard error of the mean. Stars indicate significant differences between a pair of regions (Wilcoxon rank-sum test).

### 3 Discussion

We explored LLMs and their alignment with neural responses during language processing, uncovering several key findings. Firstly, we observed a clear correlation between the language task performance of LLMs and their accuracy in predicting neural responses in the auditory cortex, with higher-performing models exhibiting greater functional alignment with the speech cortex. Secondly, we showed that the models with higher performance on benchmark tasks achieved peak predictive accuracy in earlier layers. In contrast, lower-performing models exhibited a delayed representation, necessitating deeper layers to approach similar levels of brain prediction accuracy. Finally, our study highlights the crucial role of contextual information in both LLMs and brain processing, where the contextual window’s size significantly influenced the difference between better and worse models, with the availability of long-range contextual information driving the high-performing LLMs closer to the brain’s hierarchical pathway. These findings uncover fundamental principles in language processing, highlighting the critical role of hierarchical structure and contextual dependencies in language which give rise to convergent processing strategies in both artificial and biological systems.



### 3.1 Hierarchical Processing and Inter-Model Comparisons

We found that better-performing LLMs exhibit a more brain-like hierarchy of layers, offering new insights into their language processing. While previous studies have revealed similarities in the hierarchical stages found in the brain and deep neural networks for linguistic [11, 8, 27], acoustic [35, 36], visual [37, 38, 39], and imagined stimuli [40], a distinct approach in our study is the inter-model comparison within a consistent architectural framework. In related work analyzing deep neural networks for vision tasks, recent evidence [28] has shown that better performance can create a less brain-like progression of feature extraction in models when compared to the visual cortex, suggesting that the complex architectures of high-performing image processing networks have steered them away from neural alignment. By examining LLMs based on a single architecture, the stacked transformer decoder [41], we uncover differences in their alignment with the brain’s hierarchical stages during language comprehension. Transformer language models use contextual features to encode linguistic, syntactic, and positional structures [42, 43], and increasingly high-level and context-specific features arise throughout a model’s layers [18, 19]. This may be partly because later layers bind linguistic structures over longer contexts [44]. The crucial observation that such models display brain-like hierarchies resonates with neurobiological findings of hierarchical organization in the auditory and language-related cortex [14, 22, 33, 17, 15, 16, 45, 46]. The convergence of the two systems highlights language’s inherent hierarchical structure as we increasingly form larger units of representation, from articulatory features to phonemes, syllables, words, sentences, and phrases [34, 47, 48]. Our results demonstrate that as LLMs have achieved higher performance, they have done so using feature extraction pathways that more closely resemble the human brain.

### 3.2 Feature Extraction Efficiency and Contextual Processing

A significant finding of our study is the delayed feature extraction observed in less effective LLMs compared to their higher-performing counterparts. This delay, particularly evident in the early processing stages within transformer models, suggests a slower buildup of relevant linguistic and contextual information [19]. The implications of this observation are multifaceted. Firstly, it challenges the conventional emphasis on the final layers of LLMs [10], instead drawing attention to the critical role of initial layers in efficient language processing [13]. This shift in focus aligns with emerging neuroscience research that underscores the significance of early-stage processing in the human brain for complex cognitive tasks like language processing [46, 34, 48]. Secondly, this delayed representation in less effective models offers insights into potential inefficiencies in their training or design. Given the architectural similarity of models in our study, the variance in feature extraction efficiency among models may reflect differences in training strategies [49] and data quality [50, 51, 20], providing insights for future LLM model development. As LLMs have evolved in recent years, improvements in dataset size and cleanliness as well as architectural changes to increase context length have come along with their performance improvements, and our results show that these improvements have also given rise to greater brain similarity. Furthermore, the observation that higher-performing models utilize early layers more effectively and peak in their brain similarity in middle layers rather than later layers raises intriguing questions about the role of subsequent layers. It is possible that these later layers are engaged in next-level contextual integration and feature extraction, potentially analogous to higher-order stimulus integration to support cognitive functions in the human brain [52, 53]. Alternatively, this finding could point to a limitation in our current methodologies, such as limited iEEG coverage, the simplicity of the speech comprehension task, or the fact that LLMs are not explicitly trained to perform comprehension, but rather next-word prediction, which is slightly different from the speech listening comprehension task the subjects performed. Our iEEG recordings include broad coverage of speech processing regions, especially acoustic sensory regions like HG and STG, which, although critical for spoken language processing, represent a slightly different aspect of linguistic feature extraction than the token-level processing that transformer architecture LLMs begin with. Answering these questions is crucial for enriching our understanding of artificial language processing.

The influence of contextual information on brain similarity and LLM benchmark scores also points to specific avenues that may improve model performance on language tasks. Ensuring that models are able to extract long context windows, such as by using architectures that allow for long context windows [54] and utilizing training data that is rich in long context information, could enhance LLM performance further beyond simply scaling up a model’s parameter size. Transformer-based LLMs have been shown to suffer from unequal contextual information extraction when the prior

context occurs at different distances from the target [55], supporting the notion that improving the robustness of modern LLMs to varying context lengths may lead to performance improvements. Our investigation offers a unique lens through which to view the parallels and divergences between machine learning and human cognitive development.

### 3.3 Convergence to Brain-Like Models for Human-Level Artificial General Intelligence

The convergence of LLMs and human speech processing may suggest that certain fundamental principles underlying efficient language processing might be common to both artificial and biological systems. The human brain’s language capabilities have developed as an adaptive response to complex communication needs, optimizing for efficiency and versatility [56]. Our findings suggest that LLM architectures and processing strategies are gravitating towards these same principles, mimicking the brain’s evolutionary adaptations for language. LLMs are trained without consideration for brain similarity, yet they have become increasingly brain-like in their feature extraction and hierarchical processing. Brain-like processing may represent an optimal solution to language modeling found by evolution [57], although subject to biological constraints, and our results suggest that modern LLM training focused on performance optimization may have placed these models on a similar path. In our study, Mistral, the top-performing model, stands as a prime example of this convergence, where the degree of similarity of a model’s embeddings to those of Mistral is highly correlated with performance and brain similarity. This evolution towards an optimal brain-like model offers an intriguing suggestion regarding artificial general intelligence (AGI). While not clearly defined, AGI can be quantified as human-level performance on a broad set of benchmarks [58]. Our findings suggest that developing models mimicking human neural processing strategies [59], rather than solely focusing on augmenting computational power or diversifying learning algorithms [60], could accelerate the development of models that behave on par with human performance. Hence, brain similarity could be a useful evaluation and optimization metric for future model development.

Our research marks a significant stride in understanding the parallels between large language models and human brain processes in language comprehension, by revealing the intricate relationship between internal model representation, model performance, and neural predictive accuracy. Our findings enhance the understanding of LLMs and offer new insights into the cognitive mechanisms underlying human language processing.

## 4 Methods

### 4.1 Human Intracranial Recordings

Eight subjects undergoing clinical evaluation for drug-resistant epilepsy participated in the study. Electrodes were implanted intracranially (iEEG) with the clinical goal of identifying epileptogenic foci for surgical removal. Any electrodes showing signs of epileptiform discharges, as identified by an epileptologist, were not analyzed in this study. Prior to electrode implantation, all subjects provided written informed consent for research participation. The research protocol was approved by the institutional review board at North Shore University Hospital.

Subjects listened to naturalistic recordings of voice actors reading passages from stories and conversations. To ensure the subjects were paying attention to the stimuli, one of the voices in the recording occasionally directed a question at the listener directly, or the stories were paused and the subject was asked a question, to check their understanding. The subjects were able to effectively answer each question. These pauses separated the stimulus into separate passages.

The envelope of the high-gamma band (70-150 Hz) of the raw neural recordings was computed using the Hilbert transform [61] and downsampled to 100 Hz. This signal was used as the neural response due to its correlation with neuronal firing rates [62, 63] and its common use in auditory neuroscience research [64, 65]. We restricted our analysis to speech-responsive electrodes, which we estimated using a t-test between each electrode’s response to the first second of the stimulus compared to last second of silence preceding it (FDR corrected,  $p < 0.05$  [21]), which left 707 electrodes for analysis. We extracted average word responses from each electrode by taking the average high-gamma signal value in a 100ms window around the midpoint of each word.

## 4.2 Large Language Models

We analyzed 12 LLMs of approximately 7 billion parameters downloaded from Hugging Face and implemented with its Transformers library [66], including the most recent and most popular open-source LLMs. We selected these models by searching the Hugging Face Hub for 7 billion parameter models, then using as many of the trending or most-downloaded models that we were able to run without issue.

We computed two similar evaluation metrics to those used by LLaMA 2 [20]: Reading Comprehension and Commonsense Reasoning. As measures of English language understanding, these are both highly related to the listening comprehension task which was performed by the human subjects in the study. As in [20], these metrics were created by averaging the model’s performance on a certain set of related tasks. All individual benchmarks were computed for each model using the Language Model Evaluation Harness [67] on Github.

- Reading Comprehension - This metric was the average 0-shot performance of a model on SQuAD 2.0 [68] and BoolQ [69].
- Commonsense Reasoning - This metric consists of the average 0-shot performance on OpenBookQA [70], PIQA [71], HellaSwag [72], and WinoGrande [73].

Overall LLM Performance was computed as the average Reading Comprehension and Commonsense Reasoning scores.

The models used, and their benchmark performance and overall LLM performance scores, are shown in Table 1.

Models Used	Reading Comprehension	Commonsense Reasoning	LLM Performance
Galactica-6.7B [74]	0.486	0.535	0.511
CerebrasGPT-6.7B [75]	0.565	0.575	0.570
Pythia-6.9B [76]	0.568	0.597	0.582
OPT-6.7B [77]	0.581	0.616	0.598
FairseqDense-6.7B [78]	0.575	0.628	0.602
LeoLM-7B [79]	0.634	0.646	0.640
MPT-7B [80]	0.620	0.665	0.643
Falcon-7B [81]	0.619	0.669	0.644
LLaMA-7B [82]	0.626	0.674	0.650
LLaMA2-7B [20]	0.639	0.671	0.655
XwinLM-7B [83]	0.648	0.673	0.660
Mistral-7B [84]	0.669	0.703	0.686

Table 1: All models used in the study, along with their computed benchmark performances.

In order to extract LLM embeddings for each stimulus passage (approximately 30-60 seconds when spoken), we fed the text to the model and extracted the embeddings of each layer when given a causal attention mask. When limiting the contextual window of the model, the attention mask was truncated to only include the most recent  $N$  tokens. For multi-token words, we used the embedding of the last token in the word. Thus, for each passage, we extracted a tensor of embeddings of shape  $(L_{layers}, N_{words}, D_{dimensions})$  from each model.

## 4.3 Ridge Regression Mapping from Embeddings to Neural Responses

We performed PCA to reduce the dimensionality of each model’s embeddings to 500 components. For a given model, PCA was performed for each layer separately. Then, we fit 10-fold cross-validated ridge regression models to predict the average word responses from each layer’s embeddings, sweeping over a range of regularization parameters for each training fold, using scikit-learn’s RidgeCV model [85].

#### 4.4 Electrode Localization and Brain Plotting

Each subject’s electrode positions were mapped to the subject’s brain using `iELVis` [86] to perform co-registration between pre- and post-implant MRI scans. Then, the subject-specific electrode locations were mapped to the FreeSurfer average brain [87]. Euclidean distance from posteromedial HG (TE1.1) [23] was computed in this average brain, since TE1.1 is a landmark of primary auditory cortex [24, 45, 25, 26]. When visualizing electrodes on the average brain, all subdural electrodes were snapped to the nearest surface point.

#### 4.5 Comparing LLMs with Centered Kernel Alignment

To estimate the similarity between high-dimensional embeddings of different models, we used CKA [29], a similarity metric which is related to CCA but has been shown to perform well in high-dimensional scenarios between neural network features. We used the RBF kernel to allow for nonlinear similarity measurement. For a given pair of models, we computed the CKA similarity between the embeddings of one layer of the first model with another layer of the second model. Iterating over all pairs of layers for those two models produced a single similarity matrix. These similarity matrices were then grouped by whether they described a comparison between two models in the top-5 of all LLMs for benchmark performance, one model in the top-5 and the other in the bottom-5, or two models in the bottom-5, and then averaged.

#### 4.6 Data and Code Availability

Although the iEEG recordings used in this study cannot be made publicly available, they can be requested from the author [N.M.]. Code for preprocessing neural recordings, including extracting the high-gamma envelope and identifying responsive electrodes is available in the `naplib-python` package [88].

### Acknowledgement

This work was funded by the National Institutes of Health, the National Institute on Deafness and Other Communication Disorders, and the National Science Foundation Graduate Research Fellowship Program. The funders had no role in study design, data collection and analysis, decision to publish or preparation of the manuscript.

### References

- [1] Mariya Toneva and Leila Wehbe. Interpreting and improving natural-language processing (in machines) with natural language-processing (in the brain). *Advances in neural information processing systems*, 32, 2019.
- [2] Samira Abnar, Lisa Beinborn, Rochelle Choenni, and Willem Zuidema. Blackbox meets blackbox: Representational similarity and stability analysis of neural language models and brains. *arXiv preprint arXiv:1906.01539*, 2019.
- [3] Martin Schrimpf, Idan Asher Blank, Greta Tuckute, Carina Kauf, Eghbal A Hosseini, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. The neural architecture of language: Integrative modeling converges on predictive processing. *Proceedings of the National Academy of Sciences*, 118(45):e2105646118, 2021.
- [4] Charlotte Caucheteux and Jean-Rémi King. Language processing in brains and deep neural networks: computational convergence and its limits. *BioRxiv*, pages 2020–07, 2020.
- [5] Eghbal A Hosseini, Martin Schrimpf, Yian Zhang, Samuel Bowman, Noga Zaslavsky, and Evelina Fedorenko. Artificial neural network language models align neurally and behaviorally with humans even after a developmentally realistic amount of training. *BioRxiv*, pages 2022–10, 2022.
- [6] Andrew James Anderson, Douwe Kiela, Jeffrey R Binder, Leonardo Fernandino, Colin J Humphries, Lisa L Conant, Rajeesh DS Raizada, Scott Grimm, and Edmund C Lalor. Deep artificial neural networks reveal a distributed cortical network encoding propositional sentence-level meaning. *Journal of Neuroscience*, 41(18):4100–4119, 2021.

- [7] Charlotte Caucheteux, Alexandre Gramfort, and Jean-Remi King. Disentangling syntax and semantics in the brain with deep networks. In *International conference on machine learning*, pages 1336–1348. PMLR, 2021.
- [8] Charlotte Caucheteux and Jean-Rémi King. Brains and algorithms partially converge in natural language processing. *Communications biology*, 5(1):134, 2022.
- [9] Jingyuan Sun, Shaonan Wang, Jiajun Zhang, and Chengqing Zong. Neural encoding and decoding with distributed sentence representations. *IEEE Transactions on Neural Networks and Learning Systems*, 32(2):589–603, 2020.
- [10] Ariel Goldstein, Zaid Zada, Eliav Buchnik, Mariano Schain, Amy Price, Bobbi Aubrey, Samuel A Nastase, Amir Feder, Dotan Emanuel, Alon Cohen, et al. Shared computational principles for language processing in humans and deep language models. *Nature neuroscience*, 25(3):369–380, 2022.
- [11] Charlotte Caucheteux, Alexandre Gramfort, and Jean-Rémi King. Evidence of a predictive coding hierarchy in the human brain listening to speech. *Nature human behaviour*, 7(3):430–441, 2023.
- [12] Richard Antonello, Aditya Vaidya, and Alexander G Huth. Scaling laws for language encoding models in fmri. *arXiv preprint arXiv:2305.11863*, 2023.
- [13] Richard Antonello and Alexander Huth. Predictive coding or just feature discovery? an alternative account of why language models fit brain data. *Neurobiology of Language*, pages 1–16, 2023.
- [14] Gregory Hickok and David Poeppel. The cortical organization of speech processing. *Nature reviews neuroscience*, 8(5):393–402, 2007.
- [15] Uri Hasson, Eunice Yang, Ignacio Vallines, David J Heeger, and Nava Rubin. A hierarchy of temporal receptive windows in human cortex. *Journal of Neuroscience*, 28(10):2539–2550, 2008.
- [16] Yulia Lerner, Christopher J Honey, Lauren J Silbert, and Uri Hasson. Topographic mapping of a hierarchy of temporal receptive windows using a narrated story. *Journal of Neuroscience*, 31(8):2906–2915, 2011.
- [17] Nai Ding, Lucia Melloni, Aotian Yang, Yu Wang, Wen Zhang, and David Poeppel. Characterizing neural entrainment to hierarchical linguistic units using electroencephalography (eeg). *Frontiers in human neuroscience*, 11:481, 2017.
- [18] Kawin Ethayarajh. How contextual are contextualized word representations? comparing the geometry of bert, elmo, and gpt-2 embeddings. *arXiv preprint arXiv:1909.00512*, 2019.
- [19] Ian Tenney, Dipanjan Das, and Ellie Pavlick. Bert rediscovers the classical nlp pipeline. *arXiv preprint arXiv:1905.05950*, 2019.
- [20] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [21] Sture Holm. A simple sequentially rejective multiple test procedure. *Scandinavian journal of statistics*, pages 65–70, 1979.
- [22] Tatyana O Sharpee, Craig A Atencio, and Christoph E Schreiner. Hierarchical representations in the auditory cortex. *Current opinion in neurobiology*, 21(5):761–767, 2011.
- [23] Patricia Morosan, Jorg Rademacher, Axel Schleicher, Katrin Amunts, Thorsten Schormann, and Karl Zilles. Human primary auditory cortex: cytoarchitectonic subdivisions and mapping into a spatial reference system. *Neuroimage*, 13(4):684–701, 2001.
- [24] Simon Baumann, Christopher I Petkov, and Timothy D Griffiths. A unified framework for the organization of the primate auditory cortex. *Frontiers in systems neuroscience*, 7:11, 2013.
- [25] Sam V Norman-Haignere and Josh H McDermott. Neural responses to natural and model-matched stimuli reveal distinct computations in primary and nonprimary auditory cortex. *PLoS biology*, 16(12):e2005127, 2018.
- [26] Gavin Mischler, Menoua Keshishian, Stephan Bickel, Ashesh D Mehta, and Nima Mesgarani. Deep neural networks effectively model neural adaptation to changing background noise and suggest nonlinear noise filtering methods in auditory cortex. *NeuroImage*, 266:119819, 2023.

- [27] Sreejan Kumar, Theodore R Sumers, Takateru Yamakoshi, Ariel Goldstein, Uri Hasson, Kenneth A Norman, Thomas L Griffiths, Robert D Hawkins, and Samuel A Nastase. Reconstructing the cascade of language processing in the brain using the internal computations of a transformer-based language model. *BioRxiv*, pages 2022–06, 2022.
- [28] Soma Nonaka, Kei Majima, Shuntaro C Aoki, and Yukiyasu Kamitani. Brain hierarchy score: Which deep neural networks are hierarchically brain-like? *IScience*, 24(9), 2021.
- [29] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pages 3519–3529. PMLR, 2019.
- [30] Charlotte Caucheteux, Alexandre Gramfort, and Jean-Rémi King. Deep language algorithms predict semantic comprehension from brain activity. *Scientific reports*, 12(1):16327, 2022.
- [31] Sergi G Costafreda, Cynthia HY Fu, Lucy Lee, Brian Everitt, Michael J Brammer, and Anthony S David. A systematic review and quantitative appraisal of fmri studies of verbal fluency: role of the left inferior frontal gyrus. *Human brain mapping*, 27(10):799–810, 2006.
- [32] Sophie Arana, André Marquand, Annika Hultén, Peter Hagoort, and Jan-Mathijs Schoffelen. Sensory modality-independent activation of the brain network for language. *Journal of neuroscience*, 40(14):2914–2924, 2020.
- [33] Jingwei Sheng, Li Zheng, Bingjiang Lyu, Zhehang Cen, Lang Qin, Li Hai Tan, Ming-Xiong Huang, Nai Ding, and Jia-Hong Gao. The cortical maps of hierarchical linguistic structures during speech perception. *Cerebral cortex*, 29(8):3232–3240, 2019.
- [34] Menoua Keshishian, Serdar Akkol, Jose Herrero, Stephan Bickel, Ashesh D Mehta, and Nima Mesgarani. Joint, distributed and hierarchically organized encoding of linguistic features in the human auditory cortex. *Nature Human Behaviour*, 7(5):740–753, 2023.
- [35] Bruno L Giordano, Michele Esposito, Giancarlo Valente, and Elia Formisano. Intermediate acoustic-to-semantic representations link behavioral and neural responses to natural sounds. *Nature Neuroscience*, 26(4):664–672, 2023.
- [36] Greta Tuckute, Jenelle Feather, Dana Boebinger, and Josh H McDermott. Many but not all deep neural network audio models capture brain responses and exhibit correspondence between model stages and brain regions. *Plos Biology*, 21(12):e3002366, 2023.
- [37] Nikolaus Kriegeskorte. Deep neural networks: a new framework for modeling biological vision and brain information processing. *Annual review of vision science*, 1:417–446, 2015.
- [38] Radoslaw Martin Cichy, Aditya Khosla, Dimitrios Pantazis, Antonio Torralba, and Aude Oliva. Comparison of deep neural networks to spatio-temporal cortical dynamics of human visual object recognition reveals hierarchical correspondence. *Scientific reports*, 6(1):27755, 2016.
- [39] Nicholas J Sexton and Bradley C Love. Reassessing hierarchical correspondences between brain and deep networks through direct interface. *Science Advances*, 8(28):eabm2219, 2022.
- [40] Tomoyasu Horikawa and Yukiyasu Kamitani. Hierarchical neural representation of dreamed objects revealed by brain decoding with deep neural network features. *Frontiers in computational neuroscience*, 11:4, 2017.
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [42] Joe O’Connor and Jacob Andreas. What context features can transformer language models use? *arXiv preprint arXiv:2106.08367*, 2021.
- [43] Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. What does bert look at? an analysis of bert’s attention. *arXiv preprint arXiv:1906.04341*, 2019.
- [44] David Skrill and Samuel Victor Norman-Haignere. Large language models transition from integrating across position-yoked, exponential windows to structure-yoked, power-law windows. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [45] Sam V Norman-Haignere, Laura K Long, Orrin Devinsky, Werner Doyle, Ifeoma Irobunda, Edward M Merricks, Neil A Feldstein, Guy M McKhann, Catherine A Schevon, Adeen Flinker, et al. Multiscale temporal integration organizes hierarchical computation in human auditory cortex. *Nature human behaviour*, 6(3):455–469, 2022.

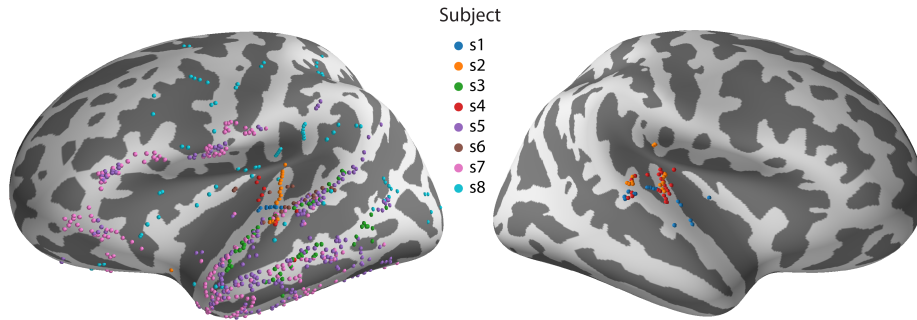
- [46] Wendy A de Heer, Alexander G Huth, Thomas L Griffiths, Jack L Gallant, and Frédéric E Theunissen. The hierarchical cortical organization of human speech processing. *Journal of Neuroscience*, 37(27):6539–6557, 2017.
- [47] Giovanni M Di Liberto, Jingping Nie, Jeremy Yeaton, Bahar Khalighinejad, Shihab A Shamma, and Nima Mesgarani. Neural representation of linguistic feature hierarchy reflects second-language proficiency. *Neuroimage*, 227:117586, 2021.
- [48] Xue L Gong, Alexander G Huth, Fatma Deniz, Keith Johnson, Jack L Gallant, and Frédéric E Theunissen. Phonemic segmentation of narrative speech in human cerebral cortex. *Nature communications*, 14(1):4309, 2023.
- [49] Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*, 2023.
- [50] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [51] Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. *arXiv preprint arXiv:2107.06499*, 2021.
- [52] Alexander G Huth, Wendy A De Heer, Thomas L Griffiths, Frédéric E Theunissen, and Jack L Gallant. Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature*, 532(7600):453–458, 2016.
- [53] Elliot Murphy, Kiefer J Forseth, Cristian Donos, Kathryn M Snyder, Patrick S Rollo, and Nitin Tandon. The spatiotemporal dynamics of semantic integration in the human brain. *Nature Communications*, 14(1):6336, 2023.
- [54] Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang, Prajjwal Bhargava, Rui Hou, Louis Martin, Rashi Rungta, Karthik Abinav Sankararaman, Barlas Oguz, et al. Effective long-context scaling of foundation models. *arXiv preprint arXiv:2309.16039*, 2023.
- [55] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*, 2023.
- [56] Steven Pinker and Paul Bloom. Natural language and natural selection. *Behavioral and brain sciences*, 13(4):707–727, 1990.
- [57] Terrence William Deacon. *The symbolic species: The co-evolution of language and the brain*. Number 202. WW Norton & Company, 1997.
- [58] Ben Goertzel. Artificial general intelligence: concept, state of the art, and future prospects. *Journal of Artificial General Intelligence*, 5(1):1, 2014.
- [59] Lin Zhao, Lu Zhang, Zihao Wu, Yuzhong Chen, Haixing Dai, Xiaowei Yu, Zhengliang Liu, Tuo Zhang, Xintao Hu, Xi Jiang, et al. When brain-inspired ai meets agi. *Meta-Radiology*, page 100005, 2023.
- [60] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [61] Erik Edwards, Maryam Soltani, Won Kim, Sarang S Dalal, Srikantan S Nagarajan, Mitchel S Berger, and Robert T Knight. Comparison of time–frequency responses and the event-related potential to auditory speech stimuli in human cortex. *Journal of neurophysiology*, 102(1):377–386, 2009.
- [62] Supratim Ray and John HR Maunsell. Different origins of gamma rhythm and high-gamma activity in macaque visual cortex. *PLoS biology*, 9(4):e1000610, 2011.
- [63] Mitchell Steinschneider, Yonatan I Fishman, and Joseph C Arezzo. Spectrotemporal analysis of evoked and induced electroencephalographic responses in primary auditory cortex (a1) of the awake monkey. *Cerebral Cortex*, 18(3):610–625, 2008.
- [64] Nima Mesgarani, Connie Cheung, Keith Johnson, and Edward F Chang. Phonetic feature encoding in human superior temporal gyrus. *Science*, 343(6174):1006–1010, 2014.



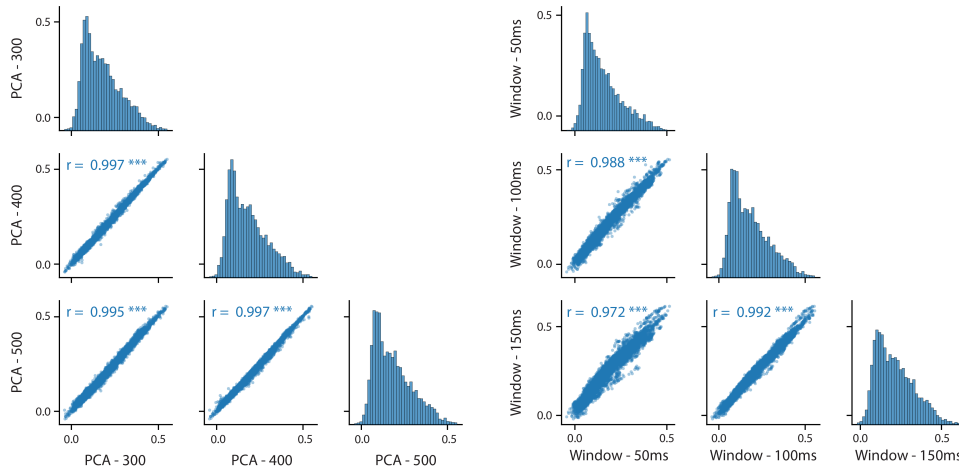
- [65] Kristofer E Bouchard, Nima Mesgarani, Keith Johnson, and Edward F Chang. Functional organization of human sensorimotor cortex for speech articulation. *Nature*, 495(7441):327–332, 2013.
- [66] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- [67] Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, September 2021.
- [68] Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions for squad. *arXiv preprint arXiv:1806.03822*, 2018.
- [69] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*, 2019.
- [70] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018.
- [71] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020.
- [72] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- [73] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- [74] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022.
- [75] Nolan Dey, Gurpreet Gosal, Hemant Khachane, William Marshall, Ribhu Pathria, Marvin Tom, Joel Hestness, et al. Cerebras-gpt: Open compute-optimal language models trained on the cerebras wafer-scale cluster. *arXiv preprint arXiv:2304.03208*, 2023.
- [76] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR, 2023.
- [77] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [78] Mikel Artetxe, Shruti Bhosale, Naman Goyal, Todor Mihaylov, Myle Ott, Sam Shleifer, Xi Victoria Lin, Jingfei Du, Srinivasan Iyer, Ramakanth Pasunuru, et al. Efficient large scale language modeling with mixtures of experts. *arXiv preprint arXiv:2112.10684*, 2021.
- [79] Laion leolm: Linguistically enhanced open language model. <https://huggingface.co/LeoLM/leo-hessianai-13b>. Accessed: 2023-10-01.
- [80] MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable llms. <https://www.mosaicml.com/blog/mpt-7b>, 2023.
- [81] Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, et al. The falcon series of open language models. *arXiv preprint arXiv:2311.16867*, 2023.
- [82] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

- [83] Xwin-lm: Powerful, stable, and reproducible llm alignment. <https://huggingface.co/Xwin-LM/Xwin-LM-7B-V0.2>. Accessed: 2023-10-01.
- [84] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [85] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.
- [86] David M Groppe, Stephan Bickel, Andrew R Dykstra, Xiuyuan Wang, Pierre Mégevand, Manuel R Mercier, Fred A Lado, Ashesh D Mehta, and Christopher J Honey. ielvis: An open source matlab toolbox for localizing and visualizing human intracranial electrode data. *Journal of neuroscience methods*, 281:40–48, 2017.
- [87] Bruce Fischl, André Van Der Kouwe, Christophe Destrieux, Eric Halgren, Florent Ségonne, David H Salat, Evelina Busa, Larry J Seidman, Jill Goldstein, David Kennedy, et al. Automatically parcellating the human cerebral cortex. *Cerebral cortex*, 14(1):11–22, 2004.
- [88] Gavin Mischler, Vinay Raghavan, Menoua Keshishian, and Nima Mesgarani. naplib-python: Neural acoustic data processing and analysis tools in python. *Software Impacts*, 17:100541, 2023.

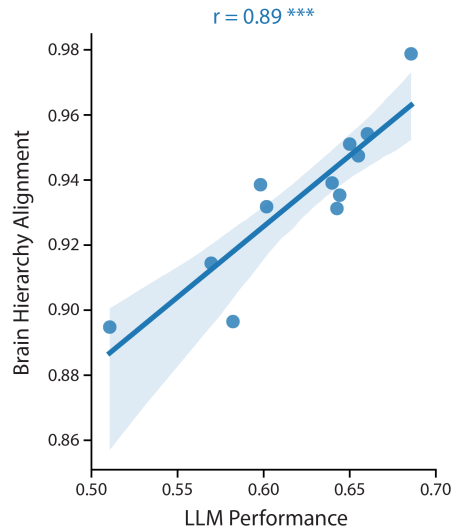
## 5 Supplementary Figures



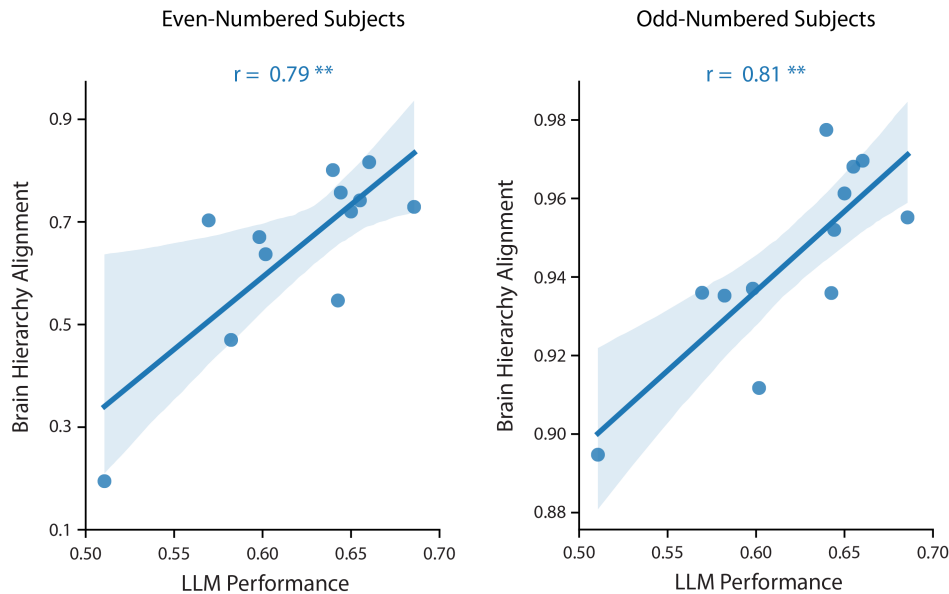
Supplementary Figure 1: Subject-wise electrode localization. Electrodes are plotted on the inflated Freesurfer average brain and are colored by their corresponding subject identity.



Supplementary Figure 2: Effect of regression hyperparameters on scores. The left plot shows the pairwise effects on the peak brain similarity scores when altering the number of principal components of the LLM embeddings used for computing scores with ridge regression, keeping a 100ms window size constant. The right plot shows the pairwise effects of altering the width of the averaging window around word centers for estimating neural responses to words, keeping the PCA dimensionality of 500 constant. Along each plot's diagonal is the marginal distribution for that hyperparameter setting. The off-diagonal plots display scatter plots of all the peak-scores for all models together for one hyperparameter setting against another. Each dot represents the peak brain correlation score for one model-electrode pair. All pairs of settings produce scores which are highly correlated, as written in each subplot (Pearson correlation, \*\*\* indicates  $p < 0.001$ ).



Supplementary Figure 3: Hierarchy alignment by model when using electrode lag instead of distance to estimate neural hierarchy. We used the electrode lag, instead of distance from primary auditory cortex, to bin electrodes into a hierarchy with a bin-width of 40ms. We estimated electrode lag using the peak of a 1D temporal receptive field fitted for each electrode to predict its response from the acoustic envelope of the stimulus sound. We then performed the same analysis as shown in Fig. 3, reproducing Fig. 3B with new brain hierarchy alignment for each model. These alignment values are similarly significantly correlated with LLM performance (Pearson  $r = 0.89$ ,  $p = 0.0001$ ).



Supplementary Figure 4: Hierarchy alignment patterns hold for partial subject groupings. Splitting the electrodes based on whether they came from even- or odd-numbered subjects, we performed the same analyses as in Fig. 3B. Both subject groups show that brain hierarchy alignment is significantly correlated with LLM performance (Pearson correlations in figure, even  $p = 0.0022$ , odd  $p = 0.0013$ ) demonstrating that this effect is not the result of a single outlier subject.