

Mind the Gap: Aligning the Brain with Language Models Requires a Nonlinear and Multimodal Approach

Danny Dongyeop Han¹, Yunju Cho¹, Jiook Cha^{1,*}, Jay-Yoon Lee^{1,*}

¹Seoul National University

*Corresponding authors:

Jiook Cha (connectome@snu.ac.kr)
Jay-Yoon Lee (lee.jayyoon@snu.ac.kr)

February 19, 2025

Abstract

Self-supervised language and audio models effectively predict brain responses to speech. However, traditional prediction models rely on linear mappings from unimodal features, despite the complex integration of auditory signals with linguistic and semantic information across widespread brain networks during speech comprehension. Here, we introduce a nonlinear, multimodal prediction model that combines audio and linguistic features from pre-trained models (e.g., LLAMA, Whisper). Our approach achieves a 17.2% and 17.9% improvement in prediction performance (unnormalized and normalized correlation) over traditional unimodal linear models, as well as a 7.7% and 14.4% improvement, respectively, over prior state-of-the-art models. These improvements represent a major step towards future robust in-silico testing and improved decoding performance. They also reveal how auditory and semantic information are fused in motor, somatosensory, and higher-level semantic regions, aligning with existing neurolinguistic theories. Overall, our work highlights the often neglected potential of nonlinear and multimodal approaches to brain modeling, paving the way for future studies to embrace these strategies in naturalistic neurolinguistics research.

Schroeder, 2006). Furthermore, this process is inherently nonlinear, involving hierarchical and spatiotemporal transformations across distributed brain regions (Tuller et al., 2011). Understanding these complex mechanisms is crucial not only for advancing cognitive neuroscience but also for developing brain inspired artificial intelligence systems.

Language encoding models, which predict brain activity from speech stimuli, are a powerful tool for unraveling the neural processes of speech comprehension (Naselaris et al., 2011; Tang et al., 2023; Vaidya et al., 2022; LeBel et al., 2021; Jain & Huth, 2018; Goldstein et al., 2022). Unlike traditional contrast-based paradigms that rely on carefully controlled experiments, encoding models capitalize on naturalistic stimuli (e.g., spoken language) to capture real-world brain processing. This allows for a more comprehensive understanding of the brain’s activity in response to complex and ecologically valid tasks, offering greater generalizability compared to simplified, contrast-based tasks (Jain et al., 2024). Moreover, these models are increasingly used for in-silico experiments that help test brain function without collecting new data, (Jain et al., 2024; Bashivan et al., 2019; Wehbe et al., 2016) and to build decoding models for language comprehension (Tang et al., 2023).

Early encoding models primarily mapped simple acoustic features (e.g., spectrograms) to brain activity (de Heer et al., 2017). The introduction of word embeddings (Mikolov, 2013) enabled the incorporation of semantic information, revealing how meaning is represented across the brain (Huth et al., 2016). Recent advances in large language models (LLMs) and sophisticated audio models have further enriched these features, leading to substantial gains in prediction accuracy (Antonello et al., 2024; Vaidya et al., 2022). However, most studies still rely on linear mappings of unimodal features, which fail to capture two fundamental principles of neural language processing.

1. Introduction

The brain seamlessly deciphers spoken language, integrating auditory signals with linguistic and semantic meaning through a dynamic interplay of neural networks. This process relies on the brain’s capacity to combine information from multiple modalities (e.g., auditory, linguistic, and motor systems) (McGettigan et al., 2012; Ghazanfar &

First, the brain operates through nonlinear computations (Friston et al., 2000; Beniaguev et al., 2021; Tuller et al., 2011), enabling complex spatiotemporal relationships across neural systems. Second, speech comprehension is a multimodal process, needing integration of diverse information sources (voice, gesture, linguistic) across distributed neural networks (McGettigan et al., 2012). These principles, along with the Motor Theory of Speech Perception (Lieberman et al., 1967) and the Convergence Divergence Zone model (Damasio, 1989), highlight the importance of encoding models capturing nonlinear dynamics and multimodal integration to reveal the brain’s functional organization.

Although prior studies have explored multimodal models combining linguistic with visual features (Oota et al., 2022; Wang et al., 2022; Scotti et al., 2024), the integration of **auditory representations**—especially from advanced speech models like Whisper—remains underexplored. This gap is significant, as auditory information is central to natural speech comprehension. Recent work by Oota et al. (2023) elegantly shows that speech models, unlike text-based language models, capture brain activity patterns in auditory regions that cannot be explained by low-level stimulus features, underscoring the complementary nature of auditory and linguistic representations. Investigating how semantic and auditory features interact in the brain is therefore critical for advancing brain-aligned models of speech processing.

In this study, we address these gaps by introducing a nonlinear, multimodal encoding model that combines audio and semantic features extracted from advanced models like Whisper and LLAMA. Our contributions are as follows:

- **Our nonlinear multimodal approach yields marked prediction performance improvements, showing a 17.2% higher unnormalized correlation and 17.9% higher normalized correlation compared to standard unimodal linear models (Antonello et al., 2024), while surpassing previous state-of-the-art models—which rely on weighted averaging of linear unimodal predictions—by 7.7% and 14.4%, respectively (Appendix C).** This performance boost demonstrates that incorporating nonlinearity and multimodality is crucial for capturing the brain’s language processing mechanisms, promising more robust in-silico experiments and improved brain decoding capabilities.
- **We propose a novel spatiotemporal clustering analysis that tracks neural responses to semantic and auditory information over time, extending beyond traditional spatial-only approaches.** By analyzing relative error differences between semantic and audio encoding models, we demonstrate that nonlinear models achieve superior functional clustering compared to both linear encoders and standard connectivity analyses. This method reveals previously hidden patterns in

brain organization, showcasing how nonlinear encoding models better capture the spatiotemporal dynamics of language processing.

- **We provide novel evidence for distributed multimodal processing in speech comprehension through variance partitioning and systematic comparisons of prediction accuracy when adding audio or semantic features.** Most brain regions utilize overlapping information from semantic and audio features, with neither modality strictly dominating. While both make unique contributions, their relative influence varies hierarchically from early sensory to higher-order areas. This extends existing neurolinguistic theories (Lieberman et al., 1967; Damasio, 1989; Davis & Yee, 2021) by revealing how different brain regions engage with multiple aspects of speech input.

2. Method

2.1. MRI Data

We used a publicly available fMRI dataset (LeBel et al., 2023; Tang et al., 2023) of three subjects listening to approximately 20 hours of English podcast. The training data comprised 95 stories across 20 scanning sessions (approximately 33,000 time points). For testing, we used three held-out stories: one averaged across ten repetitions and two averaged across five repetitions each, with no session containing repeated test stimuli. Each voxel was normalized to zero mean and unit variance across time, ensuring consistent training and testing data with (Antonello et al., 2024).

2.2. Feature Extraction

We extracted the features from the stimuli by taking the hidden layer representations of various LLMs and audio models exposed to the same stimuli as the subject. For semantic feature extraction, we utilized LLAMA-1 (Touvron et al., 2023a) models with 7B, 13B, 33B, and 65B parameters, LLAMA-2 7B (Touvron et al., 2023b), and LLAMA-3 8B (Dubey et al., 2024). For audio feature extraction, we employed Whisper (Radford et al., 2023) models, including Tiny, Base, Small, Large, and Large v2 and v3. All models were obtained from Hugging Face (Wolf, 2019) and computed using half-precision (float16) for efficiency.

For LLAMA models, the stimuli were presented using a dynamically sized context window strategy (Antonello et al., 2024) to balance computational efficiency and contextual coherence (details are in Appendix B.1). For Whisper models, the audio stimuli (waveform) were processed using a sliding-window approach with a fixed window size of 16 seconds and a stride of 0.1 seconds. Features were extracted exclusively from the encoder portion of the model, as it processes only the raw audio waveform input. This ensured

that the extracted features accurately captured audio-specific representations relevant to the stimuli. Refer to (Antonello et al., 2024) for further details model contexts handling.

The following process adheres to (Antonello et al., 2024) for fair comparison. We temporally aligned the hidden states from the l^{th} layer of the language or audio models with fMRI acquisition times using Lanczos interpolation. To account for temporal delays between stimulus presentation and neural responses, we concatenated representations from the four preceding timepoints (2, 4, 6, and 8 seconds prior) for each fMRI acquisition timepoint, yielding a feature vector for each TR (see Appendix B.3). Unless stated otherwise, we extracted semantic features from the 12th layer of LLAMA-7B and audio features from the final encoder layer of Whisper Large V1. The layers were selected based on our findings that performance scaling with increasing LLM size, as reported in (Antonello et al., 2024), does not hold for LLAMA models of size $\geq 7\text{B}$ (see Appendix F).

2.3. Representations for fMRI Data

We predicted PCA-reduced fMRI representations, rather than the full voxel space, motivated by three benefits. First, PCA is a common dimensionality reduction method in fMRI analysis that helps prevent overfitting and has been widely applied in neuroimaging studies (Wang et al., 2010; Mourao-Miranda et al., 2005; López et al., 2011; Koutsouleris et al., 2009). This was crucial as speech comprehension engages the whole cortex and hence a vast number of voxels 80 – 90k (LeBel et al., 2023), far more than vision encoding $\approx 15\text{k}$ (Allen et al., 2022). In fact, linearly mapping the semantic stimulus (4×4096) to subject S1 require 1.3 billion parameters, whereas utilizing PCA (512 components) reduced this to 8.4 million, preventing overfitting. Second, PCA is effective at untangling the redundancy in brain data, as fMRI voxels are highly correlated. Studies have shown that masking up to 90% of voxels does not significantly impact fMRI decoding or reconstruction performance (Jabakhanji et al., 2022; Lin et al., 2022), suggesting that information is distributed and redundant. Lastly, PCA allows us to recover the original voxel space from predicted PCA embeddings, maintaining the interpretability of the model’s predictions.

In detail, we applied PCA to the aggregate fMRI response matrix $Y_{\text{org}} \in \mathbb{R}^{N_{\text{TR}} \times N_{\text{voxels}}}$, reducing its dimensionality to $Y_{\text{PCA}} \in \mathbb{R}^{N_{\text{TR}} \times N_{\text{PCA}}}$. N_{TR} , N_{voxels} each refers to the number of time points (TRs) and voxels respectively, and N_{PCA} was set to 512. The encoding model was trained to predict these PCA-reduced representations, and during evaluation, the predicted $\hat{Y}_{\text{PCA}}^{\text{test}}$ was reconstructed back to the original voxel space using inverse projection. This reconstructed output was then compared to the actual fMRI responses, $Y^{\text{test}} \in \mathbb{R}^{N_{\text{TR-test}} \times N_{\text{voxels}}}$. More details are provided in Appendix B.4.

2.4. Encoding Model

Previous research primarily employed linear regression to predict voxel responses from unimodal features (audio or semantic) (Tang et al., 2023; Huth et al., 2016; de Heer et al., 2017; LeBel et al., 2021; Jain & Huth, 2018; Schrimpf et al., 2021). Our study expands upon this approach by systematically investigating a range of encoding models varying in complexity and input modality to capture more nuanced relationships between stimulus representations and brain responses. We explored combinations of different stimulus representations, encoding model architectures, and response representations (as in Table 1). The following encoding architectures were used to assess the impact of complexity and nonlinearity (Details are in Appendix B.5):

- *Linear Regression (Linear)*: Following (Antonello et al., 2024), we used ridge regression.
- *Multi-Layer Perceptron (MLP)*: MLP with a single hidden layer of 256 units.
- *Multi-Layer Linear (MLLinear)*: MLP but without dropout, batch normalization, and with the identity activation function. This model serves as a reduced-rank linear regression, helping to isolate the effects of dimensionality reduction from nonlinearity.
- *Delayed Interaction MLP (DIMLP)*: Used for multimodal cases, this MLP variant processes each modality through separate 256-unit hidden layers before concatenation and final linear projection. This allows nonlinear processing within each modality while limiting cross-modal interaction to be linear, revealing the effects of nonlinear fusion of modalities.

2.5. Noise Ceiling and Normalized Correlation Coefficient

Due to the inherent noise in fMRI data, there is a theoretical upper limit to the amount of explainable variance an ideal encoding model can achieve, known as the noise ceiling. The noise ceiling for each voxel was estimated by applying an existing method (Schoppe et al., 2016) on ten responses of the same test story (see Appendix B.2). Afterwards, by dividing CC_{abs} , the correlation coefficient of the encoding model’s prediction with the ground truth fMRI signals (estimated as the average of test responses to the same stimuli), by CC_{max} , we obtain CC_{norm} , the normalized correlation coefficient. Due to the large number of voxels ($\approx 80,000$), random noise can cause certain voxels to have $CC_{\text{abs}} > CC_{\text{max}}$, resulting in $CC_{\text{norm}} > 1$. To prevent this, we regularized for noisy voxels by setting those with $CC_{\text{max}} < 0.25$ to 0.25 when computing CC_{norm} .

2.6. RED (Relative Error Difference)

We introduce a novel metric called the Relative Error Difference (RED). For each voxel v at time t , $RED(v, t) = |f_1(v, t) - y(v, t)| - |f_2(v, t) - y(v, t)|$, where $f_1(v, t)$ and $f_2(v, t)$ are the predictions from model 1 (LLAMA) and model 2 (Whisper), respectively, and $y(v, t)$ is the true fMRI response. A positive RED value indicates that model 2 outperforms model 1 at that voxel and timepoint, while a negative value indicates the opposite.

RED extends beyond traditional voxel-wise analyses ($f(v)$) that focus only on spatial patterns of brain activity. By preserving temporal information ($f(v, t)$), RED enables analysis of both spatial and temporal dynamics of neural processing. We leverage this spatio-temporal information in Section 3.1.2 to develop a novel approach for clustering ROIs based on their semantic and audio processing dynamics.

Table 1. Performance of encoding models across different modalities and architectures (sem=semantic, aud=audio). This table presents the average voxelwise r^2 and normalized correlation coefficient (CC_{norm}) values for various encoding models, comparing their ability to predict fMRI responses across different input modalities *semantic*, *audio* or *multimodal* and encoder architectures *Linear*, *MLinear*, *DIMLP* and *MLP*. Notably, MLP encoders consistently outperform linear models and their variants, highlighting the importance of incorporating nonlinearity for accurate fMRI prediction. r^2 is computed as $|r| * r$.

Input	Encoder	Output	r^2 (%) (Δr^2 %)	CC_{norm} (%) (ΔCC_{norm} %)
sem+aud	MLP	PCA	4.29 (+17.2)	34.32 (+17.9)
sem+aud	DIMLP	PCA	4.18 (+14.2)	32.59 (+11.9)
sem+aud	MLinear	PCA	4.10 (+12.0)	32.41 (+11.3)
sem+aud	Linear	voxels	4.10 (+12.0)	31.36 (+7.7)
sem+aud	Linear	PCA	3.87 (+5.7)	28.92 (-0.7)
sem+aud	MLP	voxels	3.83 (+4.6)	31.11 (+6.8)
sem	MLP	PCA	3.79 (+3.6)	30.89 (+6.1)
sem	MLinear	PCA	3.67 (+0.3)	29.95 (+2.8)
sem	Linear	voxels	3.66 (base)	29.12 (base)
sem	Linear	PCA	3.56 (-2.7)	26.88 (-7.7)
sem	MLP	voxels	3.36 (-8.2)	27.45 (-5.7)
aud	MLP	PCA	3.01 (-17.8)	29.01 (-0.4)
aud	MLP	voxels	2.89 (-21.0)	28.21 (-3.1)
aud	MLinear	PCA	2.89 (-21.0)	27.50 (-5.6)
aud	Linear	PCA	2.81 (-23.2)	26.71 (-8.3)
aud	Linear	voxels	2.77 (-24.3)	25.20 (-13.5)

3. Results

3.1. Nonlinear Encoders

3.1.1. NONLINEARITY, NOT REDUCED

DIMENSIONALITY ALONE, IMPROVES ENCODING PERFORMANCE

To examine the benefits of nonlinearity across feature hierarchies, we compared MLP and linear encoders across different layers of LLAMA and Whisper models (Figure 9 (Appendix)). Our findings reveal that MLP consistently outperformed linear models across all feature hierarchies and layer depths, supporting the notion that nonlinear transformations capture richer and more complex relationships in brain activity than linear mappings alone.

To disentangle the role of nonlinearity from dimensionality reduction, we compared the MLP encoder with two linear control models: “Linear,” which uses linear regression on PCA-reduced data, and “MLinear,” which mirrors the MLP architecture without nonlinear activation functions. As shown in Table 1, both Linear and MLinear models performed similarly to or worse than linear regression on the full voxel space (baseline model). These findings highlight the MLP’s ability to capture nonlinear relations drives its superior performance, not merely dimensionality reduction. Additionally, PCA proves essential for leveraging nonlinearity. MLP models predicting all voxels directly performed poorly, likely due to overfitting from the large voxel space (80–90k voxels compared to 512 PCA components).

3.1.2. NONLINEARITY IMPROVES BRAIN-WIDE PREDICTIONS AND SPATIO-TEMPORAL CLUSTERING

Nonlinear MLP models provide a crucial advantage over linear models by effectively capturing the complex nonlinear relationships inherent in brain activity during speech comprehension. Comparative brain maps in Appendix I.3 illustrate the superior performance of MLP encoders over linear encoders, with improvements in prediction accuracy distributed across the cortex. The MLP model shows significant gains in regions associated with semantic and auditory processing, such as the medial prefrontal cortex (mPFC), precuneus (PrCu), and lateral temporal cortex (LTC). These gains highlight the role of nonlinear interactions in accurately modeling brain activity, particularly in areas involved in higher-order language processing.

Furthermore, our hierarchical clustering analysis based on the Relative Error Difference (RED) between Whisper and LLAMA encoding models (Figure 1 and Appendix I.4) reveals two key insights: (1) RED enables functional clustering of brain regions where traditional functional connectivity fails, and (2) nonlinear encoders achieve superior functional grouping over linear ones, as shown by higher

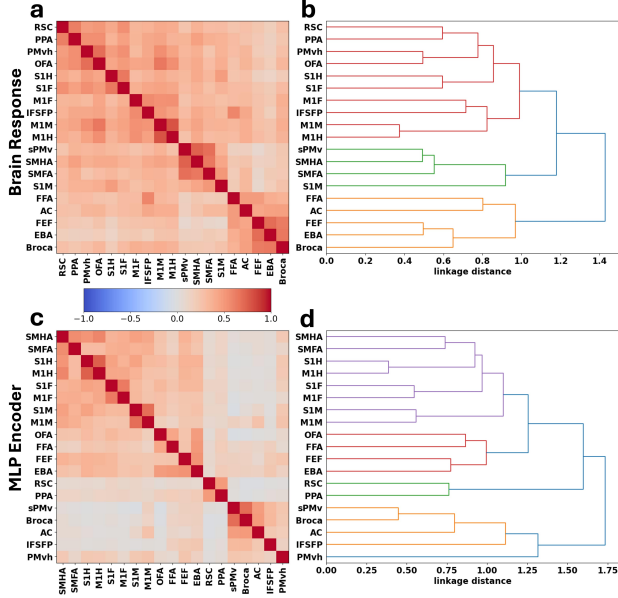


Figure 1. Spatio-temporal clustering analysis: (a,b) Functional connectivity matrix and hierarchical clustering dendrogram from raw fMRI correlations. (c,d) Correlation matrices and dendrograms from Relative Error Difference (RED) between semantic and audio encoding models using MLP encoders. Matrix values indicate regional similarity. Hierarchical clustering reveals brain region organization by response profiles. The nonlinear models (d) show clearer functional groupings than standard connectivity (b), quantified by higher modularity scores (see main text).

modularity Q values (nonlinear: 0.155, linear: 0.145, FC: 0.068). The MLP encoder’s clustering (Figure 1 (d)) demonstrates clear hierarchical functional organization: primary motor and somatosensory regions (M1/S1) first pair by body part (e.g., M1M/S1M) before forming broader motor and somatosensory clusters; higher visual regions cluster by function (face processing: OFA/FFA; scene processing and spatial navigation: PPA/RSC); and speech-language regions (sPMv/Broca/AC) form a cluster aligned with the dorsal stream pathway. This underscores RED’s potential and shows that nonlinear models capture structured spatiotemporal relationships in brain responses, aligning with known functional organization principles.

3.2. Multimodal Encoders

3.2.1. MULTIMODALITY PREDOMINANTLY CONTRIBUTES TO AND IMPROVES BRAIN-WIDE PREDICTIONS

Improvements from multimodal encoding are cortex-wide and extend beyond modality-specific processing regions. Voxelwise analyses in Figure 2 (a,b) show that adding audio features improves not only auditory areas but also primary

motor somatosensory regions. These improvements extend beyond expected auditory areas, with enhancements observed in the paracentral lobule, situated medially between the mPFC and Precuneus (PrCu), and in the occipital cortex (OC), reflecting the widespread impact of auditory information on cortical processing. Similarly, Figure 2 (c,d) shows that the addition of semantic features leads to improvements in most cortical areas, except certain parts of the auditory cortex (AC). These improvements are more pronounced in CC_{norm} visualizations (Appendix J.2), reinforcing the broad influence of semantic processing on neural activity beyond classical language areas.

This widespread improvement from multimodality is further amplified by nonlinearity. Comparing Figure 2 (b) with (a), and (d) with (c) reveals that MLP models not only enhance performance in regions that were already well predicted by multimodal linear encoders, but also in regions not initially benefiting from the added modality, such as the LTC when adding audio features, and LTC, mPFC, and OC when adding semantic features. This suggests multimodal MLP models can better exploit the additional modality through nonlinearity, which we discuss in Section 3.3.

To understand the source of these improvements, we conducted variance partitioning analysis, decomposing each voxel’s explained variance into modality-specific and joint components. Our analysis (Appendix L.2) reveals two key findings about brain-wide processing, with some regional exceptions such as the auditory cortex: first, the majority of explained variance comes from joint audio-semantic processing rather than from either modality alone; second, both semantic and audio features make unique contributions to explaining brain activity, though audio’s unique contribution is much smaller in magnitude across most regions. These results suggest that, for the majority of cortical areas, auditory and semantic models contain largely overlapping information, with semantic models providing additional unique predictive information beyond what audio models capture.

By assigning each voxel to its dominant predictive modality, we found that joint audio-semantic features dominate cortical representations (Figure 3 (a), shown for subject S1; plots for all subjects in Appendix L.3). This pattern is consistent across subjects, with semantic, audio, and joint features being the most attributable source for approximately 21.4%, 10.1%, and 68.5% of significantly predicted voxels across the cortex, respectively, as shown by ROI-wise analysis in Figure 3 (b) (subject-wise analyses in Appendix L.4).

Our findings align with and diverge from prior multimodal language studies. The cortex-wide improvement contrasts with (Antonello et al., 2024), which reported localized enhancements in AC and M1M with auditory features. Methodological differences may explain this: First, they used multiple Whisper layers, potentially adding redundancy, whereas

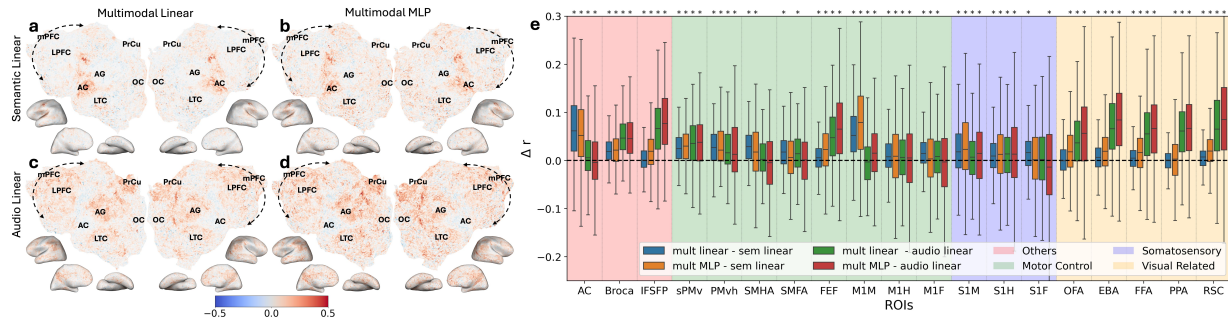


Figure 2. Multimodality improvement in encoding models. Panels (a)-(d) display voxelwise Δr values of a single subject, with warmer colors indicating regions where multimodal models outperform linear models. Each panel corresponds to the difference between voxel-wise predictions of the model in the corresponding column and the model in the corresponding row. E.g., panel (a) shows the difference between the Multimodal Linear and Semantic Linear models. (e) Box plot showing Δr across different regions of interest (ROIs), where the Δr values are aggregated over all subjects. *mult* and *sem* each refer to multimodal and semantic encoders. Asterisk* indicate ROI where $\Delta r > 0$ is statistically significant ($p < 0.05$). ROIs are grouped and color-coded by their functions. The boxes represent the range between the 25th and 75th percentiles, with the line inside showing the median. Whiskers extend to 1.5 times this range. (A complete list of ROI abbreviations are at Appendix A. Voxelwise and ROI-wise plots for each subject are in Figure 17, 18, and 21 in the Appendix).

our approach focuses on the final layer. Second, their linear stacked regression model averaged predictions from separate encoders, limiting modality interaction, while our concatenation allows direct interaction (Appendix C).

Our analysis also provides new insights into the nature of modality-specific representations. While our results align with (Oota et al., 2023) in that semantic models contain information beyond low-level features present in audio models, our results reveal a more nuanced picture - audio models, though showing smaller unique contributions, provide meaningful complementary information across multiple brain regions. This is evidenced by both improved prediction performance and non-zero unique variance contribution in our voxel-wise analysis. This apparent discrepancy may be because our voxel-wise approach might have captured finer-grained patterns of unique audio contributions that might be averaged out in their ROI-level analyses.

These patterns of distributed joint processing align with the Convergence-Divergence-Zone theory (Damasio, 1989), which posits that semantic information is integrated from multiple modalities across the cortex.

3.2.2. MULTIMODAL FUSION SUPPORTS AND EXTENDS NEUROLINGUISTIC THEORIES

Building on the brain-wide improvements observed, regions of interest (ROI) analyses reveal how multimodal integration supports and extends multiple neurolinguistic theories.

Speech Related Regions (AC, Broca, sPMv, M1M)

Our analysis reveals a systematic organization of speech processing that follows the auditory dorsal pathway, a key component of the dual-stream model of language processing

(Hickok & Poeppel, 2007). This pathway, extending from the auditory cortex (AC) through Broca’s area and sPMv to the primary motor cortex, exhibits distinct patterns of multimodal integration at each stage.

In early AC, voxel-wise variance partitioning shows that unique contributions from audio features dominate (Figure 3), reinforcing its role in processing low-level acoustic information. However, processing in broader AC regions shows a shift to joint audio-semantic representations, with 83.3% of significantly predicted voxels showing joint audio-semantic representation. The improved performance from adding auditory features (Figures 2 (a,b)) supports this hierarchical pattern, with earlier AC areas showing greater gains.

Moving along the dorsal pathway to Broca’s area and sPMv (superior ventral premotor speech area), we find predominant joint feature attribution (88.2% and 84.8% of voxels respectively) with improved predictions from the addition of either modality. This multimodal integration aligns with these regions’ role in speech planning and articulatory control—processes that require integrating acoustic targets with semantic content and motor programs (Gough et al., 2005; Nixon et al., 2004; de Heer et al., 2017; Glanz et al., 2018).

At the terminus of the dorsal pathway, M1M shows a strong contribution from auditory features, exceeding even AC, consistent with its role in executing speech articulation (32.4% of voxels) (Figure 3 (b)). This strong auditory presence in motor areas is further supported by substantial performance improvements when adding auditory features, reinforcing previous findings from (Wu et al., 2014) that highlight the coupling between auditory and motor processes in speech production.

These findings extend our understanding of speech model

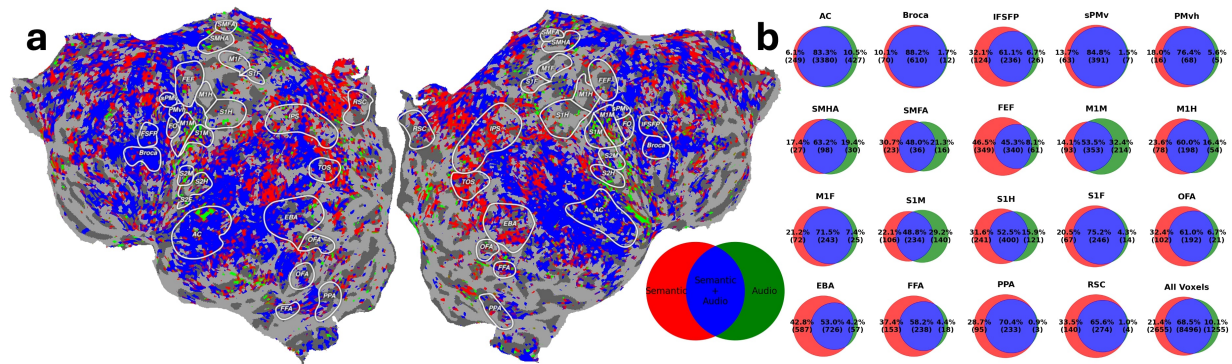


Figure 3. Visualization of most dominant feature type in brain activity predictions from variance partitioning analysis. (a) Voxel-wise plots from a single subject (S1) and (b) ROI-wise Venn diagrams showing which feature type (semantic: red, audio: green, joint: blue) explains the largest variance for each significantly predicted voxel ($q(\text{FDR}) < 0.01$) using MLP encoders. ROI results are aggregated across subjects with numbers indicating voxel percentages and counts.

representations. Our variance partitioning results align with previous findings that semantic models primarily predict AC activity by capturing low-level speech features (Oota et al., 2023). Our analysis also reveal some voxels show unique semantic contributions, and audio models capture distinct brain features beyond the typical scope of language models. The observed semantic contribution in AC, sPMv and Broca’s area aligns with prior findings (de Heer et al., 2017) and may be a general mechanism for language processing.

Motor and Somatosensory Areas: Embodied Speech Processing

The addition of audio or semantic features consistently improved predictions in cortical regions associated with motor control (green) and somatosensory processing (blue) (Figure 2 (e)). These improvements vary across ROIs: some benefit from the addition of semantic features (e.g., frontal eye field (FEF)), others from audio features (e.g., primary mouth motor cortex (M1M)), and some from both. Furthermore, variance partitioning analysis reveals that motor and somatosensory regions show unique contributions from both modalities - for instance, in M1M, audio features uniquely explain 32.4% of the variance while semantic features explain 14.1%, with 53.5% jointly explained. Similar patterns emerge across motor areas (SMHA, SMFA, FEF, M1H, M1F) and somatosensory regions (S1M, S1H, S1F), suggesting these regions process unique auditory and semantic information absent from their overlapping features.

These findings align with the long-standing Motor Theory of Speech Perception (Liberman et al., 1967; 1952; Poeppel & Assaneo, 2020), which posits that motor regions actively participate in simulating the articulatory movements necessary for speech production, thereby aiding comprehension. In particular, improvements from the addition of and the unique contribution from auditory features align with re-

search that discovered a tight coupling between auditory and motor-sensory processing (Skipper et al., 2005; Wu et al., 2014; Wilson et al., 2004).

These findings also suggests that semantic information plays a critical role in shaping activity within somatosensory regions. This suggests a broader involvement of these areas in speech comprehension than previously recognized. This aligns with the concept of embodiment of semantic memory, where the understanding of concepts is grounded in sensory and motor experiences and their memory in the neocortex (Binder & Desai, 2011). Our results align with (Nagata et al., 2022), who showed that the sensorimotor cortex is engaged in processing both concrete and abstract word semantics.

The enhancements in motor and sensory area predictions are more pronounced with MLP models, underscoring the importance of nonlinear interactions between auditory and semantic information. We explore this in more details in Section 3.3. See Appendix J for subject-wise plots.

Higher-Order Visual Areas: Multimodal Semantic Representations

Adding semantic features significantly enhances fMRI prediction accuracy in high-level visual areas like OFA (Pitcher et al., 2011), EBA (Downing et al., 2001), FFA (Kanwisher et al., 1997), PPA (Epstein & Kanwisher, 1998), and RSC (Vann et al., 2009) (Figure 2(e)). Variance partitioning (Figure 3 (b)) shows that these ROIs have largest contributions from semantic and joint features, suggesting text-derived semantics provide substantial predictive information for visual regions beyond audio features alone.

This finding aligns with prior studies demonstrating that visual and linguistic stimuli with similar semantic content elicit similar brain responses (Tang et al., 2024; Deniz et al., 2019; Devereux et al., 2013; Fairhall & Caramazza, 2013;

Popham et al., 2021; Huth et al., 2012; 2016). These studies, along with our results, support the convergence-divergence-zone theory (Popham et al., 2021; Damasio et al., 1996; 2004; Damasio, 1989), which posits semantic information from multiple modalities is integrated at points across the cortex, leading to a unified representation of semantic meaning. This model suggests that the brain constructs a modality-independent representation of semantics, drawing on information from vision, language, and other senses (Binder & Desai, 2011; Tang et al., 2023; 2024; Devereux et al., 2013; Fairhall & Caramazza, 2013; Martin, 2016).

Our study also provides novel evidence for the auditory modality’s contribution to this unified semantic representation. Variance partitioning (Figure 3 (b)) shows that auditory information accounts for approximately 5% of voxels in higher visual area ROIs. Adding audio features to our multimodal models resulted in a statistically significant performance increase in these ROIs (yellow) (Figure 2(e)), suggesting auditory information, such as tone of voice and environmental sounds, may provide unique semantic context not fully captured by visual or linguistic features alone.

The consistent observation that multimodal fusion, particularly with nonlinear models, enhances prediction accuracy emphasizes the brain’s use of complex, nonlinear computations to combine information from different modalities for a holistic understanding of language. Subject-wise ROI prediction differences are visualized in Figure 21 (Appendix).

3.3. Nonlinear and Multimodal Encoders

3.3.1. NONLINEAR INTERACTIONS BETWEEN MODALITIES ENHANCE FMRI PREDICTIONS

To assess the role of nonlinear cross-modal interactions, we developed a Delayed Interaction MLP (DIMLP), which processes audio and semantic features separately before a final linear fusion stage. This contrasts with MLP, which allows full nonlinear interactions across modalities. This design enables a direct comparison of within-modality nonlinearity (DIMLP) vs. cross-modal nonlinear interactions (MLP).

Both DIMLP and MLP outperform linear models (Table 1). DIMLP, incorporating only within-modality nonlinearity, yields a 2.0% gain over the linear model (from 4.10% average r^2 to 4.18%). But the standard MLP, allowing full nonlinear interactions, achieves a further 2.6% gain (from 4.18% to 4.29%). These results suggest that both forms of nonlinearity enhance encoding performance, but cross-modal nonlinear interactions contribute most significantly,

This conclusion is further supported by voxelwise analysis (Appendix K and Figure 23 in Appendix). While DIMLP improves prediction accuracy across brain regions compared to the linear model, the transition to a standard MLP leads to further, cortex-wide enhancements. This suggests that non-

linear interactions between audio and semantic features are essential for modeling the complex, distributed neural representations underlying speech comprehension (see Appendix K).

ROI-wise analysis (Appendix Figure 24) shows regional variation in nonlinearity’s benefits. Multimodal MLP consistently matches or outperforms DIMLP and often surpasses linear models. Notably, motor (e.g., M1M) and somatosensory regions (e.g., S1M) benefit significantly from nonlinear cross-modal interactions, showing their role in complex multimodal processing during speech comprehension.

4. Discussion and Conclusion

This study underscores the transformative potential of nonlinear, multimodal encoding models for advancing our understanding of speech comprehension in the brain. By introducing a nonlinear Multi-layer Perceptron (MLP) and integrating audio and linguistic features, we achieved a 14.4% increase in mean normalized correlation across the cortex compared to previous state-of-the-art (Antonello et al., 2024), predicting 34.3% of the brain’s explainable variance.

A key finding is that nonlinearity is fundamental to neural speech processing - nonlinear models outperformed linear approaches across all network layers, with improvements stemming from nonlinearity rather than dimensionality reduction alone as shown by linear control models. This cortex-wide enhancement reveals the brain’s reliance on nonlinear computation, further supported by our novel RED analysis showing improved hierarchical clustering of brain regions, with higher modularity (0.155) than linear models (0.145) and traditional connectivity measures (0.068).

Our second key finding is that speech comprehension involves inherent multimodal fusion across the cortex. Adding either audio or semantic features improved predictions cortex-wide, while variance partitioning showed 68.5% of significantly predicted voxels are best explained by joint audio-language processing rather than either modality alone. Through ROI-wise analyses of both variance partitioning and performance improvements, we provide support for key neurolinguistic theories including the Motor Theory of Speech Perception (Lieberman et al., 1967), Convergence-Divergence Zone model (Damasio, 1989), and embodied semantics (Davis & Yee, 2021), highlighting the brain’s reliance on distributed multimodal fusion.

Our nonlinear encoding approach has two main limitations. First, insufficient dataset size currently constrains model complexity, leading to overfitting when adding hidden layers or using RNNs and Transformers (Appendix D). Given data scaling benefits in linear encoders (Antonello et al., 2024) and how a large dataset such as the Natural Scenes Dataset (Allen et al., 2022) enabled deep learning breakthroughs

in visual encoding and decoding (Adeli et al., 2023; Scotti et al., 2024), larger language fMRI datasets are needed to fully harness the potential of deep learning and drive further advancements. Second, while nonlinear encoders offer strong performance gains, they create new interpretability challenges. While variance partitioning and RED-based clustering offer preliminary insights, further innovations such as RSA (Kriegeskorte et al., 2008) and novel feature attribution (Oota et al., 2023) are necessary. Moreover, nonlinear models offer unique interpretative possibilities, as shown by (Yang et al., 2023) in memory vision encoding.

In conclusion, our study demonstrates that nonlinear, multi-modal encoding models are crucial for understanding brain speech comprehension. Addressing dataset size and model interpretability limitations will be key to advancing brain aligned AI, enabling models that better reflect the hierarchical and distributed nature of neural processing. These insights have implications for neural representation learning, deep learning interpretability, and brain computer interfaces.

5. Impact Statement

This paper presents work whose goal is to advance the field of Neuroscience and Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Adeli, H., Minni, S., and Kriegeskorte, N. Predicting brain activity using transformers. *bioRxiv*, pp. 2023–08, 2023.
- Akiba, T., Sano, S., Yanase, T., Ohta, T., and Koyama, M. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 2623–2631, 2019.
- Allen, E. J., St-Yves, G., Wu, Y., Breedlove, J. L., Prince, J. S., Dowdle, L. T., Nau, M., Caron, B., Pestilli, F., Charest, I., et al. A massive 7t fmri dataset to bridge cognitive neuroscience and artificial intelligence. *Nature neuroscience*, 25(1):116–126, 2022.
- Antonello, R., Vaidya, A., and Huth, A. Scaling laws for language encoding models in fmri. *Advances in Neural Information Processing Systems*, 36, 2024.
- Bashivan, P., Kar, K., and DiCarlo, J. J. Neural population control via deep image synthesis. *Science*, 364(6439): eaav9436, 2019.
- Beniaguev, D., Segev, I., and London, M. Single cortical neurons as deep artificial neural networks. *Neuron*, 109(17):2727–2739, 2021.
- Binder, J. R. and Desai, R. H. The neurobiology of semantic memory. *Trends in cognitive sciences*, 15(11):527–536, 2011.
- Bonnasse-Gahot, L. and Pallier, C. fmri predictors based on language models of increasing complexity recover brain left lateralization. *arXiv preprint arXiv:2405.17992*, 2024.
- Damasio, A. R. The brain binds entities and events by multiregional activation from convergence zones. *Neural computation*, 1(1):123–132, 1989.
- Damasio, H., Grabowski, T. J., Tranel, D., Hichwa, R. D., and Damasio, A. R. A neural basis for lexical retrieval. *Nature*, 380(6574):499–505, 1996.
- Damasio, H., Tranel, D., Grabowski, T., Adolphs, R., and Damasio, A. Neural systems behind word and concept retrieval. *Cognition*, 92(1-2):179–229, 2004.
- Davis, C. P. and Yee, E. Building semantic memory from embodied and distributional language experience. *Wiley Interdisciplinary Reviews: Cognitive Science*, 12(5): e1555, 2021.
- de Heer, W. A., Huth, A. G., Griffiths, T. L., Gallant, J. L., and Theunissen, F. E. The hierarchical cortical organization of human speech processing. *Journal of Neuroscience*, 37(27):6539–6557, 2017.
- Deniz, F., Nunez-Elizalde, A. O., Huth, A. G., and Gallant, J. L. The representation of semantic information across human cerebral cortex during listening versus reading is invariant to stimulus modality. *Journal of Neuroscience*, 39(39):7722–7736, 2019.
- Devereux, B. J., Clarke, A., Marouchos, A., and Tyler, L. K. Representational similarity analysis reveals commonalities and differences in the semantic processing of words and objects. *Journal of Neuroscience*, 33(48):18906–18916, 2013.
- Downing, P. E., Jiang, Y., Shuman, M., and Kanwisher, N. A cortical area selective for visual processing of the human body. *Science*, 293(5539):2470–2473, 2001.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Epstein, R. and Kanwisher, N. A cortical representation of the local visual environment. *Nature*, 392(6676):598–601, 1998.
- Fairhall, S. L. and Caramazza, A. Brain regions that represent amodal conceptual knowledge. *Journal of Neuroscience*, 33(25):10552–10558, 2013.

- Friston, K. J., Mechelli, A., Turner, R., and Price, C. J. Nonlinear responses in fmri: the balloon model, volterra kernels, and other hemodynamics. *NeuroImage*, 12(4): 466–477, 2000.
- Ghazanfar, A. A. and Schroeder, C. E. Is neocortex essentially multisensory? *Trends in cognitive sciences*, 10(6): 278–285, 2006.
- Glanz, O., Derix, J., Kaur, R., Schulze-Bonhage, A., Auer, P., Aertsen, A., and Ball, T. Real-life speech production and perception have a shared premotor-cortical substrate. *Scientific reports*, 8(1):8898, 2018.
- Goldstein, A., Zada, Z., Buchnik, E., Schain, M., Price, A., Aubrey, B., Nastase, S. A., Feder, A., Emanuel, D., Cohen, A., et al. Shared computational principles for language processing in humans and deep language models. *Nature neuroscience*, 25(3):369–380, 2022.
- Gough, P. M., Nobre, A. C., and Devlin, J. T. Dissociating linguistic processes in the left inferior frontal cortex with transcranial magnetic stimulation. *Journal of Neuroscience*, 25(35):8010–8016, 2005.
- Hickok, G. and Poeppel, D. The cortical organization of speech processing. *Nature reviews neuroscience*, 8(5): 393–402, 2007.
- Huth, A. G., Nishimoto, S., Vu, A. T., and Gallant, J. L. A continuous semantic space describes the representation of thousands of object and action categories across the human brain. *Neuron*, 76(6):1210–1224, 2012.
- Huth, A. G., De Heer, W. A., Griffiths, T. L., Theunissen, F. E., and Gallant, J. L. Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature*, 532(7600):453–458, 2016.
- Jabakhanji, R., Vigotsky, A. D., Bielefeld, J., Huang, L., Baliki, M. N., Iannetti, G., and Apkarian, A. V. Limits of decoding mental states with fmri. *Cortex*, 149:101–122, 2022.
- Jain, S. and Huth, A. Incorporating context into language encoding models for fmri. *Advances in neural information processing systems*, 31, 2018.
- Jain, S., Vo, V. A., Wehbe, L., and Huth, A. G. Computational language modeling and the promise of in silico experimentation. *Neurobiology of Language*, 5(1):80–106, 2024.
- Kanwisher, N., McDermott, J., and Chun, M. M. The fusiform face area: A module in human extrastriate cortex specialized for face perception. *The Journal of Neuroscience*, 17(11):4302–4311, 1997.
- Koutsouleris, N., Meisenzahl, E. M., Davatzikos, C., Bottlender, R., Frodl, T., Scheuerecker, J., Schmitt, G., Zetzsche, T., Decker, P., Reiser, M., et al. Use of neuroanatomical pattern classification to identify subjects in at-risk mental states of psychosis and predict disease transition. *Archives of general psychiatry*, 66(7):700–712, 2009.
- Kriegeskorte, N., Mur, M., and Bandettini, P. A. Representational similarity analysis-connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, 2:249, 2008.
- LeBel, A., Jain, S., and Huth, A. G. Voxelwise encoding models show that cerebellar language representations are highly conceptual. *Journal of Neuroscience*, 41(50): 10341–10355, 2021.
- LeBel, A., Wagner, L., Jain, S., Adhikari-Desai, A., Gupta, B., Morgenthal, A., Tang, J., Xu, L., and Huth, A. G. A natural language fmri dataset for voxelwise encoding models. *Scientific Data*, 10(1):555, 2023.
- Liberman, A. M., Delattre, P., and Cooper, F. S. The role of selected stimulus-variables in the perception of the unvoiced stop consonants. *The American journal of psychology*, pp. 497–516, 1952.
- Liberman, A. M., Cooper, F. S., Shankweiler, D. P., and Studdert-Kennedy, M. Perception of the speech code. *Psychological review*, 74(6):431, 1967.
- Lin, S., Sprague, T., and Singh, A. Redundancy and dependency in brain activities. *Shared Visual Representations in Human & Machine Intelligence*, 2022.
- López, M., Ramírez, J., Górriz, J. M., Álvarez, I., Salas-Gonzalez, D., Segovia, F., Chaves, R., Padilla, P., Gómez-Río, M., Initiative, A. D. N., et al. Principal component analysis-based techniques and supervised classification schemes for the early detection of alzheimer’s disease. *Neurocomputing*, 74(8):1260–1271, 2011.
- Loshchilov, I. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Martin, A. Grapes—grounding representations in action, perception, and emotion systems: How object properties and categories are represented in the human brain. *Psychonomic bulletin & review*, 23:979–990, 2016.
- McGettigan, C., Faulkner, A., Altarelli, I., Obleser, J., Baverstock, H., and Scott, S. K. Speech comprehension aided by multiple modalities: behavioural and neural interactions. *Neuropsychologia*, 50(5):762–776, 2012.
- Mikolov, T. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

- Mourao-Miranda, J., Bokde, A. L., Born, C., Hampel, H., and Stetter, M. Classifying brain states and determining the discriminating activation patterns: support vector machine on functional mri data. *NeuroImage*, 28(4):980–995, 2005.
- Nagata, K., Kunii, N., Shimada, S., Fujitani, S., Takasago, M., and Saito, N. Spatiotemporal target selection for intracranial neural decoding of abstract and concrete semantics. *Cerebral Cortex*, 32(24):5544–5554, 2022.
- Naselaris, T., Kay, K. N., Nishimoto, S., and Gallant, J. L. Encoding and decoding in fmri. *Neuroimage*, 56(2):400–410, 2011.
- Nixon, P., Lazarova, J., Hodinott-Hill, I., Gough, P., and Passingham, R. The inferior frontal gyrus and phonological processing: an investigation using rtms. *Journal of cognitive neuroscience*, 16(2):289–300, 2004.
- Oota, S. R., Arora, J., Rowtula, V., Gupta, M., and Bapi, R. S. Visio-linguistic brain encoding. *arXiv preprint arXiv:2204.08261*, 2022.
- Oota, S. R., Çelik, E., Deniz, F., and Toneva, M. Speech language models lack important brain-relevant semantics. *arXiv preprint arXiv:2311.04664*, 2023.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Pitcher, D., Walsh, V., and Duchaine, B. The role of the occipital face area in the cortical face perception network. *Experimental brain research*, 209:481–493, 2011.
- Poeppl, D. and Assaneo, M. F. Speech rhythms and their neural foundations. *Nature reviews neuroscience*, 21(6):322–334, 2020.
- Popham, S. F., Huth, A. G., Bilenko, N. Y., Deniz, F., Gao, J. S., Nunez-Elizalde, A. O., and Gallant, J. L. Visual and linguistic semantic representations are aligned at the border of human visual cortex. *Nature neuroscience*, 24(11):1628–1636, 2021.
- Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., and Sutskever, I. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pp. 28492–28518. PMLR, 2023.
- Schoppe, O., Harper, N. S., Willmore, B. D., King, A. J., and Schnupp, J. W. Measuring the performance of neural models. *Frontiers in computational neuroscience*, 10:10, 2016.
- Schrimpf, M., Blank, I. A., Tuckute, G., Kauf, C., Hosseini, E. A., Kanwisher, N., Tenenbaum, J. B., and Fedorenko, E. The neural architecture of language: Integrative modeling converges on predictive processing. *Proceedings of the National Academy of Sciences*, 118(45):e2105646118, 2021.
- Scotti, P., Banerjee, A., Goode, J., Shabalin, S., Nguyen, A., Dempster, A., Verlinde, N., Yundler, E., Weisberg, D., Norman, K., et al. Reconstructing the mind’s eye: fmri-to-image with contrastive learning and diffusion priors. *Advances in Neural Information Processing Systems*, 36, 2024.
- Skipper, J. I., Nusbaum, H. C., and Small, S. L. Listening to talking faces: motor cortical activation during speech perception. *Neuroimage*, 25(1):76–89, 2005.
- Tang, J., LeBel, A., Jain, S., and Huth, A. G. Semantic reconstruction of continuous language from non-invasive brain recordings. *Nature Neuroscience*, 26(5):858–866, 2023.
- Tang, J., Du, M., Vo, V., Lal, V., and Huth, A. Brain encoding models based on multimodal transformers can transfer across language and vision. *Advances in Neural Information Processing Systems*, 36, 2024.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Tuller, B., Nguyen, N., Lancia, L., and Vallabha, G. K. Non-linear dynamics in speech perception. *Nonlinear Dynamics in Human Behavior*, pp. 135–150, 2011.
- Vaidya, A. R., Jain, S., and Huth, A. Self-supervised models of audio effectively explain human cortical responses to speech. In *International Conference on Machine Learning*, pp. 21927–21944. PMLR, 2022.
- Vann, S. D., Aggleton, J. P., and Maguire, E. A. What does the retrosplenial cortex do? *Nature reviews neuroscience*, 10(11):792–802, 2009.
- Wang, A. Y., Kay, K., Naselaris, T., Tarr, M. J., and Wehbe, L. Incorporating natural language into vision models improves prediction and understanding of higher visual cortex. *BioRxiv*, pp. 2022–09, 2022.

-
- Wang, Y., Fan, Y., Bhatt, P., and Davatzikos, C. High-dimensional pattern regression using machine learning: from medical images to continuous clinical variables. *Neuroimage*, 50(4):1519–1535, 2010.
- Wehbe, L., Huth, A. G., Deniz, F., Gao, J., Kieseler, M.-L., and Gallant, J. L. Bold predictions: Automated simulation of fmri experiments. *NeurIPS Demonstr. Track*, 2016.
- Wilson, S. M., Saygin, A. P., Sereno, M. I., and Iacoboni, M. Listening to speech activates motor areas involved in speech production. *Nature neuroscience*, 7(7):701–702, 2004.
- Wolf, T. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- Wu, Z.-M., Chen, M.-L., Wu, X.-H., and Li, L. Interaction between auditory and motor systems in speech perception. *Neuroscience bulletin*, 30:490–496, 2014.
- Yang, H., Gee, J., and Shi, J. Memory encoding model. *arXiv preprint arXiv:2308.01175*, 2023.

A. Abbreviations of Brain Areas and Regions of Interest (ROIs)

Brain Areas are abbreviated as follows :

- **AC**: Auditory Cortex
- **AG**: Angular Gyrus
- **LPFC**: Lateral Prefrontal Cortex
- **LTC**: Lateral Temporal Cortex
- **mPFC**: Medial Prefrontal Cortex
- **OC**: Occipital Cortex
- **PrCu**: Precuneus

The ROIs are abbreviated as follows :

- **AC**: Auditory Cortex
- **AG**: Angular Gyrus
- **Broca**: Broca’s Area
- **EBA**: Extrastriate Body Area
- **FFA**: Fusiform Face Area
- **FEF**: Frontal Eye Field
- **IFSFP**: Inferior Frontal Sulcus Face Patch
- **LPFC**: Lateral Prefrontal Cortex
- **LTC**: Lateral Temporal Cortex
- **M1F**: Primary Motor Cortex - Foot
- **M1H**: Primary Motor Cortex - Hand
- **M1M**: Primary Motor Cortex - Mouth
- **mPFC**: Medial Prefrontal Cortex
- **OC**: Occipital Cortex
- **OFA**: Occipital Face Area
- **PMvh**: Ventral Premotor Hand Area
- **PPA**: Parahippocampal Place Area
- **PrCu**: Precuneus
- **RSC**: Retrosplenial Cortex
- **S1F**: Primary Somatosensory Cortex - Foot
- **S1H**: Primary Somatosensory Cortex - Hand
- **S1M**: Primary Somatosensory Cortex - Mouth
- **sPMv**: Superior Ventral Premotor Speech Area
- **SMFA**: Supplementary Motor Foot Area
- **SMHA**: Supplementary Motor Hand Area

B. Details of Implementation

B.1. LLAMA Feature Extraction Strategy

LLAMA feature extraction was done in a dynamical window size manner for efficiency. Initially, the context window grew incrementally as tokens were added, up to a maximum of 512 tokens, after which the window was reset to a new context of 256 tokens. This approach avoided memory overheads associated with processing the entire tokenized text while maintaining sufficient contextual information for accurate semantic representation.

B.2. Noise Ceiling (CC_{max}) and Normalized Correlation (CC_{norm}) Calculation

For each voxel, the maximum correlation coefficient is estimated as $CC_{max} = (\sqrt{1 + \frac{NP}{SP \times N}})^{-1}$, where N is the number of repeats (10 in our case), NP is the noise power or unexplainable variance, and SP is the amount of variance that could be explained by an ideal predictive model.

B.3. Resampling the hidden state of LLMs to fMRI time points

After giving the language/audio model the same input as the subject, we temporally aligned the hidden states of its l^{th} layer corresponding to a given i^{th} token (last token of the i^{th} word for language models), $H_l^i(S_{\{k|k \leq i\}}) \in \mathbb{R}^{d_{\text{model}}^l}$ (aggregate shape of $\mathbb{R}^{N_{\text{token}} \times d_{\text{model}}^l}$ for the whole story where N_{token} is the number of tokens/words), to the fMRI acquisition times (TR times) using Lanczos interpolation, obtaining an extracted feature of size $\mathbb{R}^{N_{\text{TR}} \times d_{\text{model}}^l}$, where N_{TR} is the number of tokens (or number of words for language models) for each story and d_{model}^l is the dimension of the l^{th} hidden layer. We constructed the feature corresponding to a given n^{th} TR ($2n$ seconds in physical time) by concatenating the representations from four previous TRs (2, 4, 6, 8 seconds before t in physical time) to get a vector of shape $\mathbb{R}^{4d_{\text{model}}^l}$ for every n^{th} TR, which we denote as $H_l^n(S_{\{t|t \leq 2n\}})$. H' denotes the additional resampling and concatenation done after applying the model, H . We used four previous time delays (2, 4, 6, 8 seconds) to account for the delay between the stimuli and brain response and to provide past stimuli information to the model.

B.4. Representations for fMRI response using PCA

To an aggregate fMRI response, $Y_{\text{org}} \in \mathbb{R}^{N_{\text{TR}} \times N_{\text{voxels}}}$, we applied PCA with 8192 maximum components along the voxel dimension using scikit-learn (Pedregosa et al., 2011), yielding an approximate projection matrix, $W \in \mathbb{R}^{N_{\text{voxels}} \times N_{8192}}$. Given N_{PCA} number of principal components to consider, we take the top N_{PCA} components to get $W_{\text{PCA}} \in \mathbb{R}^{N_{\text{voxels}} \times N_{\text{PCA}}}$, and train the encoding model to

predict the reduced dimension PCA projection of the data, $Y_{\text{PCA}} = Y_{\text{org}} W_{\text{PCA}} \in \mathbb{R}^{N_{\text{TR}} \times N_{\text{PCA}}}$. During evaluation, the trained model outputs a reduced dimension representation of the data, $\hat{Y}_{\text{PCA}}^{\text{test}} \in \mathbb{R}^{N_{\text{TR-test}} \times N_{\text{PCA}}}$, where $N_{\text{TR-test}}$ denotes the number of timepoints (TRs) in the test story. This is reconstructed back to the original voxel space by applying an inverse of the projection matrix, $\hat{Y}^{\text{test}} = \hat{Y}_{\text{PCA}}^{\text{test}} W_{\text{PCA}}^T \in \mathbb{R}^{N_{\text{TR-test}} \times N_{\text{voxels}}}$, which is later compared with the ground truth, $Y^{\text{test}} \in \mathbb{R}^{N_{\text{TR-test}} \times N_{\text{voxels}}}$.

It should be noted that due to the high dimensionality of the data, incremental PCA was used, in place of regular PCA.

B.5. Details of encoding models

The encoding model architecture is as follows:

- *Linear Regression (Linear)*: Ridge regression. Following (Antonello et al., 2024), ridge regression with bootstrapping ($n = 3$) was used to estimate the optimal regularization parameters (alphas) for each voxel. The training data was divided into chunks of length 20, with 25% used for held-out validation in each bootstrap iteration. The best alpha values were averaged across iterations, and the final model was trained on the full training dataset using these alphas.
- *Multi-Layer Perceptron (MLP)*: MLP with a single hidden layer of 256 units, applying batch normalization and dropout to prevent overfitting. The hyperbolic tangent (tanh) was used as the activation function.
- *Multi-Layer Linear (MLLinear)*: MLP but without dropout, batch normalization, and with the identity activation function.
- *Delayed Interaction MLP (DIMLP)*: MLP variant processes. Each modality through separate 256-unit hidden layers before concatenation and final linear projection.

We implemented encoding models using PyTorch. We employed the AdamW optimizer (Loshchilov, 2017) with a batch size of 128 and Mean Absolute Error (MAE) as the loss function to mitigate excessively penalizing random signal fluctuations. Our training regime consisted of 200 epochs with early stopping (patience = 10) based on validation loss, and we applied batch normalization with a momentum of 0.1. For robust evaluation, we implemented 5-fold cross-validation, averaging predictions across the five models for our final results. Hyperparameter optimization was conducted using Optuna (Akiba et al., 2019), which performed 70 trials to determine optimal values for the dropout rate (0.1 to 0.3), learning rate (10^{-5} to 10^{-1}), and weight decay (5×10^{-5} to 10^{-1}).

Ridge regression was performed using a CPU node with 96 cores (Intel(R) Xeon(R) Gold 6240R CPU @ 2.40GHz) and 512 GB of RAM. Running the audio and language models and training encoding models was done using a GPU node with 8 H100 80GB GPUs.

C. Comparison with stacked regression model of (Antonello et al., 2024)

To establish the effectiveness of our nonlinear multimodal approach, we conduct a detailed comparison with the current state-of-the-art stacked regression model (Antonello et al., 2024). Their method combines semantic and audio predictions through stacked regression followed by voxel-selection, where they decide what model to use (stacked regression or semantic linear) for each voxel based on a validation dataset. Their results are compared here and not in Table 1 due to their use of only parts of the test stories as validation, barring computation of the “Avg r^2 ” value in Table 1. For accurate comparison, we obtain and use their published model weights and features.

The evaluation protocols differ specifically for the stacked regression (SR) model: while all models (including those in Antonello et al. (2024)) primarily report performance using three test stories (Table 1), SR uniquely requires using two of these test stories for validation-based voxel selection and only using the story “wherethersmoke” for final testing.

Also, following the identification of an error in the original evaluation protocol through community feedback, we corrected the methodology for fair comparison. Note that CC_{norm} values remain consistent with Table 1 as they were originally computed using only the “wherethersmoke” story due to the unavailability of test repeats for the other two stories.

To ensure fair comparison with SR, we additionally evaluate all models using their single-story protocol in Table 2, reporting both CC_{norm} and a story-specific Avg r^2 (**single story**) metric to distinguish from our three-story evaluation. We found CC_{norm} provides more stable comparisons than r^2 in this context, as the reduced number of timepoints (251 versus 790) makes r^2 more susceptible to noisy voxels compared to CC_{norm} that accounts for these noisy voxels. This stability is reflected in the closer alignment between CC_{norm} and r^2 rankings in Table 1 compared to Table 2. Therefore, we sort Table 2 with respect to the CC_{norm} .

Also, while their approach uses LLAMA-30B’s 18th layer (denoted as semantic_A), we demonstrate competitive performance using LLAMA-7B features, consistent with our finding that encoding performance roughly plateaus beyond 7B parameters (Appendix F). For comprehensive comparison, we implement both their pre-computed validation-based voxel selection mask (“mask_A”, created using an unspecified

significance threshold) and our simpler approach (“mask”) that retains voxels showing any validation set improvement.

Table 2 demonstrates several key results about our multimodal nonlinear approach. Our multimodal MLP achieves 34.32% CC_{norm} without masking, representing a 14.4% improvement over the baseline stacked regression model, though the Avg r^2 (story) improvement is more modest at 7.7%.

Our multimodal linear encoder also outperforms stacked regression by 4.5%, supporting our hypothesis that direct concatenation enables more effective modality interaction compared to weighted averaging of unimodal predictions. The performance hierarchy (MLP > Linear > SR) suggests that both architectural choices - direct multimodal fusion and nonlinearity - contribute independently to improved predictions.

Interestingly, validation-based masking did not improve performance for either our linear or MLP models, regardless of whether using our mask or the precomputed mask_A from previous work. This suggests our models learn effective feature selection implicitly, determining when to leverage or ignore audio features for specific voxels without explicit masking. The benefit of removing masking also likely stems from our models’ ability to learn voxel-specific feature importance through direct access to input data, combined with the inherent noise in validation masks due to the limited number of timepoints.

These results demonstrate that enabling direct interaction between modalities through concatenation, combined with nonlinear processing, provides a more robust approach than previous methods relying on weighted averaging and explicit feature selection.

D. Results of more complex nonlinear models

We explored a range of more complex nonlinear models, as detailed in Table 3. Specifically, we evaluated LSTM, GRU, RNN, and Transformer architectures, each configured with a single layer. The hidden dimensions for these models were determined by experimenting with sizes of 256, 512, 768, and 1024, selecting the dimension that yielded the best performance.

All models received inputs consisting of four timepoints, consistent with the MLP model, which concatenates these timepoints. For the recurrent models (LSTM, GRU, RNN), the final predictions were generated by applying a linear projection to a weighted pooling of the outputs corresponding to the four input timepoints. In the case of the Transformer model, we utilized learnable positional embeddings along with full self-attention mechanisms, and the final prediction was obtained by linearly projecting the output of the last

Table 2. Comparing encoding performance across different models using the single test story evaluation protocol. Values show normalized correlation coefficient (CC_{norm}) and story-specific r^2 (**Avg r^2 (story)**)(distinguishing from Table 1’s three-story evaluation (**Avg r^2**)). SR refers to the previous state-of-the-art stacked regression model (Antonello et al., 2024), which combines LLM and audio predictions through weighted averaging. Two masking approaches are used: 1) “mask_A” - their pre-computed validation-based voxel selection mask, and 2) “mask” - our computed masks that retain voxels showing validation improvements. For “mask”, Linear+Mask indicates creating and applying a mask based on multimodal linear vs semantic linear performance, while MLP+Mask does the same using MLP models. semantic_A denotes features from LLAMA-30B’s 18th layer used in SR, while our models uses features from the 12th layer of LLAMA-7B. All approaches are evaluated using identical test data for fair comparison and r^2 is computed as $|r| * r$.

modality 1	modality 2	encoder	response	Avg r^2 (single story)	Avg CC_{norm}
semantic	audio	MLP	PCA	5.13% (+7.7%)	34.32% (+14.4%)
semantic	audio	MLP + mask	PCA	5.02% (+5.5%)	33.33% (+11.0%)
semantic	audio	DIMLP	PCA	4.93% (+3.6%)	32.59% (+8.6%)
semantic	audio	MLLinear	PCA	5.00% (+5.1%)	32.41% (+8.0%)
semantic	audio	MLP + mask _A	PCA	4.77% (+0.2%)	31.70% (+5.6%)
semantic	audio	Linear	all voxels	4.92% (+3.4%)	31.36% (+4.5%)
semantic	audio	MLP	all voxels	4.54% (-4.5%)	31.11% (+3.6%)
semantic	audio	Linear + mask	all voxels	4.90% (+2.9%)	31.09% (+3.6%)
semantic _A	audio	SR + mask _A	all voxels	4.76% (Baseline)	30.02% (Baseline)
semantic	audio	Linear	PCA	4.48% (-5.8%)	28.92% (-3.7%)
semantic	-	MLP	PCA	4.58% (-3.7%)	30.89% (+2.9%)
semantic	-	MLLinear	PCA	4.59% (-3.6%)	29.95% (-0.2%)
semantic _A	-	Linear	all voxels	4.60% (-3.3%)	29.84% (-0.6%)
semantic	-	Linear	all voxels	4.50% (-5.4%)	29.12% (-3.0%)
semantic	-	MLP	all voxels	3.97% (-16.6%)	27.45% (-8.6%)
semantic	-	Linear	PCA	4.15% (-12.8%)	26.88% (-10.4%)
audio	-	MLP	PCA	3.83% (-19.6%)	29.01% (-3.4%)
audio	-	MLP	all voxels	3.67% (-22.8%)	28.21% (-6.0%)
audio	-	MLLinear	PCA	3.66% (-23.1%)	27.50% (-8.4%)
audio	-	Linear	PCA	3.54% (-25.6%)	26.71% (-11.0%)
audio	-	Linear	all voxels	3.46% (-27.3%)	25.20% (-16.0%)

token.

Additionally, we examined the DeepMLP model, an extension of the standard MLP with two hidden layers instead of one.

Our results indicate that while the MLP with a single hidden layer outperforms linear models, introducing greater complexity—such as recurrent models or additional hidden layers—leads to overfitting and decreased performance.

E. Performance of multimodal MLP model when mixing different layers

We observe in Figure 4 that integrating the best performing layers from each modality results in the best performing multimodal model.

F. Scaling LLM and audio models does not necessarily lead to better encoders

Previous research by Antonello et al. (2024) found that increasing the size of large language models (LLMs) and audio models, such as scaling OPT from 125M to 175B parameters or Whisper from 8M to 637M parameters, enhanced encoding performance. However, performance gains plateaued for larger models like LLAMA-33B and OPT-175B, which they attributed to overfitting from larger hidden sizes.

Building on these findings, our study delves deeper into the scaling trends and offers a refined perspective on their implications for brain encoding models. For audio models, we confirm a positive correlation between model size and performance, as shown in Figure 5 (d). However, this scaling effect does not hold for language models. Specifically, LLAMA-7B, LLAMA-13B, LLAMA-33B, and LLAMA-65B exhibit comparable encoding performance, as shown in Figure 5 (b). This suggests diminishing returns beyond 7 billion parameters, a finding consistent with prior work

Table 3. Encoding performance of various nonlinear semantic encoders compared to other models. The table presents the average r^2 and normalized correlation coefficients (CC_{norm}) along with percentage changes relative to the baseline Linear model. Deep MLP refers to an MLP with two hidden layers, while MLP is an MLP with one hidden layer.

modality 1	modality 2	encoder	response	Avg r^2	Avg CC_{norm}
semantic	-	MLP	PCA	3.79% (+3.6%)	30.89% (+6.1%)
semantic	-	Linear	all voxels	3.66% (Baseline)	29.12% (Baseline)
semantic	-	LSTM	PCA	3.33% (-9.0%)	26.95% (-7.46%)
semantic	-	GRU	PCA	3.21% (-12.3%)	26.15% (-10.2%)
semantic	-	DeepMLP	PCA	3.05% (-16.7%)	27.45% (-5.73%)
semantic	-	RNN	PCA	2.99% (-18.0%)	25.42% (-12.7%)
semantic	-	Transformer	PCA	2.82% (-23.0%)	27.97% (-3.95%)

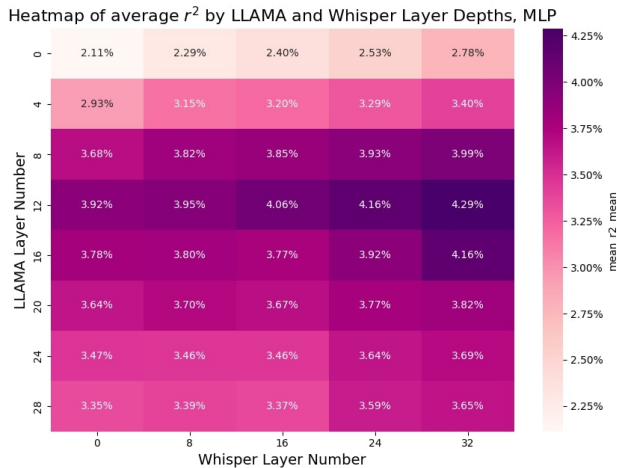


Figure 4. Heatmap showing average r^2 values for different combinations of LLAMA and Whisper layer depths using an MLP encoder. Darker colors represent higher performance, with the best results obtained when the best layers in the respective uni-modal encoding models were used.

by [Bonnasse-Gahot & Pallier \(2024\)](#), which reported performance plateaus for LLMs larger than 3 billion parameters.

We also evaluated the impact of scaling training data by examining newer versions of LLAMA and Whisper (e.g., LLAMA-1, LLAMA-2, LLAMA-3; Whisper v1, v2, v3). Despite larger datasets, newer versions did not yield significant performance improvements for either audio or semantic encoding models. This indicates that advancements in self-supervised learning (SSL) tasks, such as better next-token prediction, do not necessarily translate to more effective features for brain encoding. In essence, SSL improvements do not directly enhance brain-aligned representations.

In conclusion, our findings highlight two key points: (1) scaling language models beyond 7 billion parameters does not substantially improve encoding performance, and (2) increasing training data or using newer model versions does not enhance brain encoding feature extractors. These results

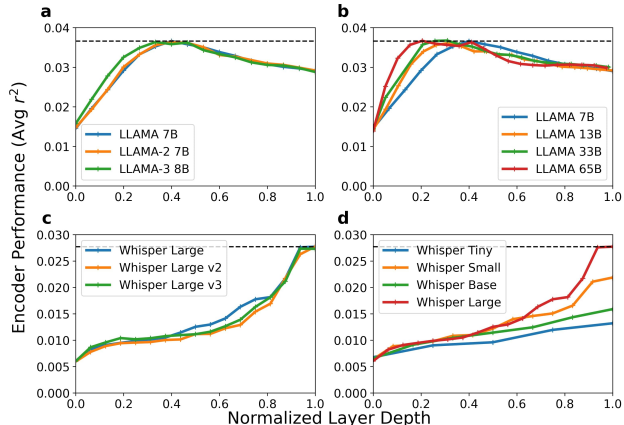


Figure 5. Encoder performance across different LLAMA and Whisper model variants, using linear regression applied to the full set of voxels. Panel (a) compares LLAMA models of various architectures (LLAMA-2 and LLAMA-3) with 7B and 8B parameters. Panel (b) presents performance across different LLAMA models of increasing sizes, from 7B to 65B. Panels (c) and (d) show the performance for different Whisper model variants, including comparisons between Whisper Large versions (c) and different model sizes (d), from Whisper Tiny to Whisper Large. Performance is measured in terms of average r^2 , plotted against normalized layer depth.

challenge the assumption that simply scaling feature extractors, as proposed by [Antonello et al. \(2024\)](#), will lead to better encoding models.

G. Context size speech models influence encoder performance

Figure 6 illustrates the impact of varying the context size (window size) of the Whisper model on encoding performance when using linear encoders, as explored in ([Oota et al., 2023](#)). The results indicate that a 16-second window size, which was used as the default throughout our study, delivers the best performance. This outcome aligns with expectations, as the selected window size is consistent with

the recommendations from (Antonello et al., 2024).

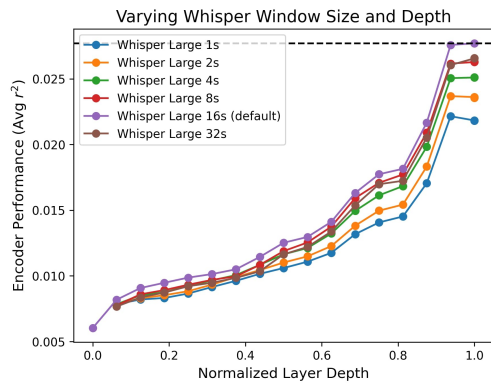


Figure 6. Encoder performance across different Whisper Large models with varying window size, using linear regression applied to the full set of voxels.

H. Performance of various encoding models using different inputs

H.1. Voxelwise r values from different encoding models and stimuli

Figures 7 represent the voxelwise correlation (r) values using various encoders and inputs for subject S1. Due to file size constraints, the plots for other subjects have been moved to the supplementary materials.

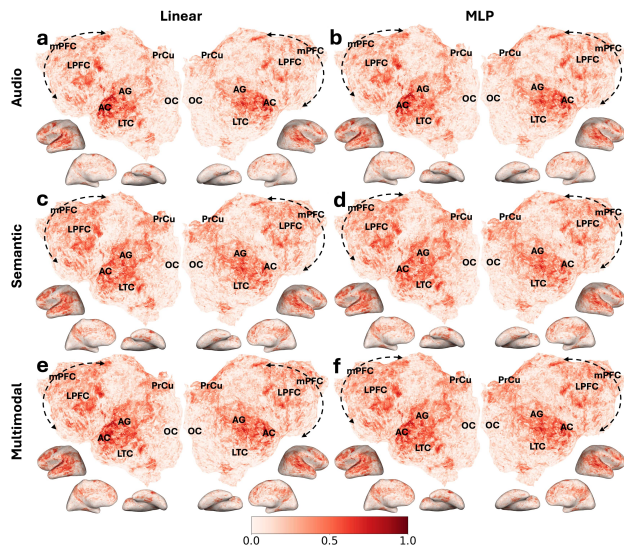


Figure 7. Voxelwise r values for Subject S1 across different input modalities and encoding models. Rows show audio-only (a,b), semantic-only (c,d), and multimodal (e,f) inputs. Columns compare Linear (left) and MLP (right) encoders. Warmer colors indicate higher prediction accuracy.

H.2. ROI-wise r values from different encoding models and stimuli

Figure 8 shows the r value for different encoding models and stimuli averaged across subjects.

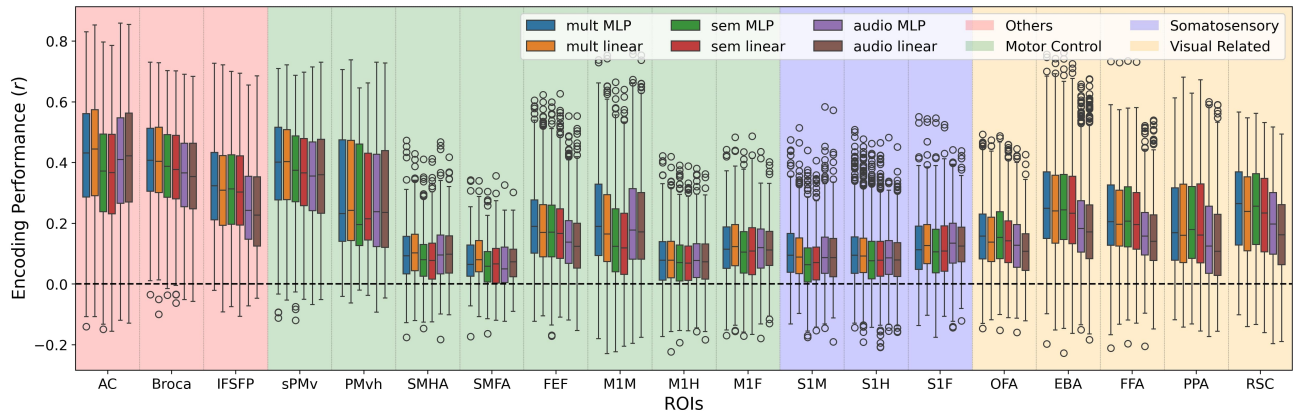


Figure 8. Box plot showing r across different regions of interest (ROIs), where the r values are aggregated over all subjects. *multi* refers to multimodal, and *sem* refers to semantic encoders. ROIs are grouped and color-coded by their functions.

I. Improvements from nonlinearity

I.1. Layer-wise performance increases from MLP

Figure 9 shows that MLP improves encoding performance for both language and audio models, regardless of what layer is used for the MLP encoding model.

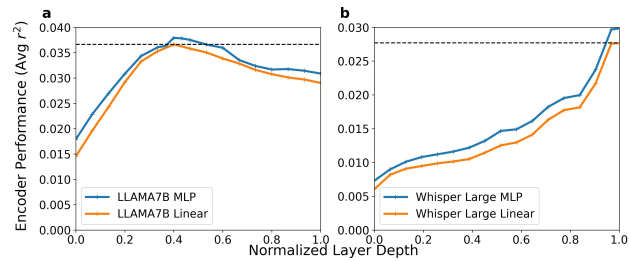


Figure 9. Average voxel-wise r^2 values, computed as the mean across three subjects, for each layer of the (a) language (LLAMA7B) and (b) audio (Whisper Large) models. Comparisons are shown between the MLP and linear encoders, and dashed black lines indicate the best performance for linear encoders

I.2. Voxelwise improvements from MLP (r analysis)

Figures 10, 11, and 12 each represent the performance improvements in voxelwise correlation values for semantic, audio, and multimodal inputs, respectively, for each subject.

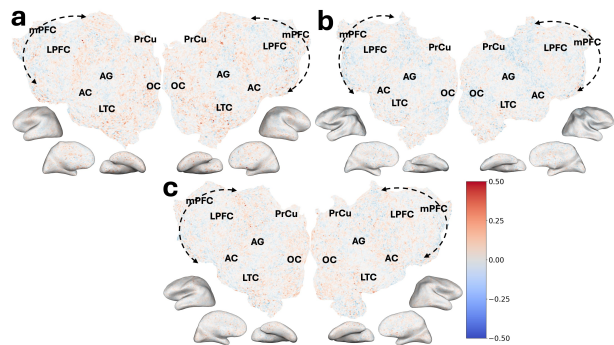


Figure 10. Encoding model performance improvements. (a-c) Vox- elwise Δr (MLP performance minus linear performance) for semantic input for subjects S1, S2, S3, respectively. Positive values indicate MLP outperformance.

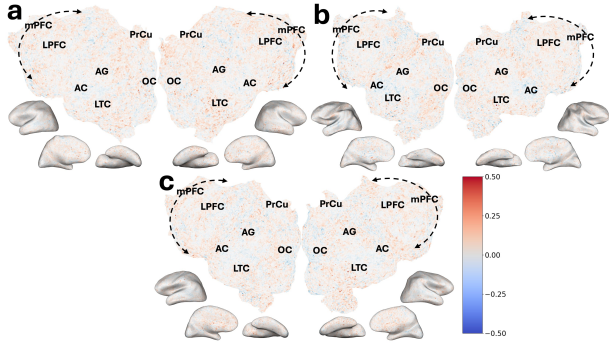


Figure 11. Encoding model performance improvements. (a-c) Vox-
elwise Δr (MLP performance minus linear performance) for audio
input for subjects S1, S2, S3, respectively. Positive values indicate
MLP outperformance.

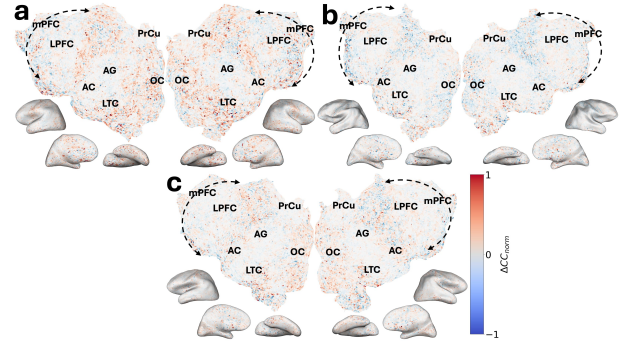


Figure 13. Encoding model performance improvements. (a-c) Vox-
elwise ΔCC_{norm} (MLP performance minus linear performance)
for semantic input for subjects S1, S2, S3, respectively. Positive
values indicate MLP outperformance.

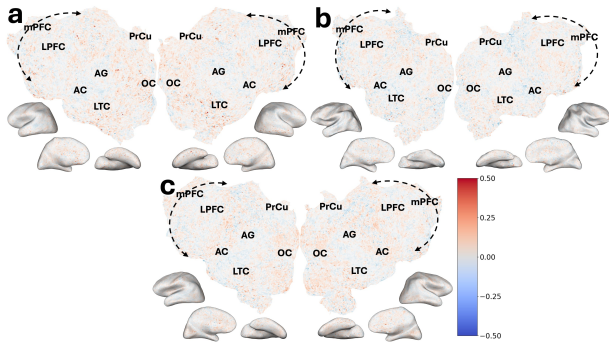


Figure 12. Encoding model performance improvements. (a-c) Vox-
elwise Δr (MLP performance minus linear performance) for mul-
timodal input for subjects S1, S2, S3, respectively. Positive values
indicate MLP outperformance.

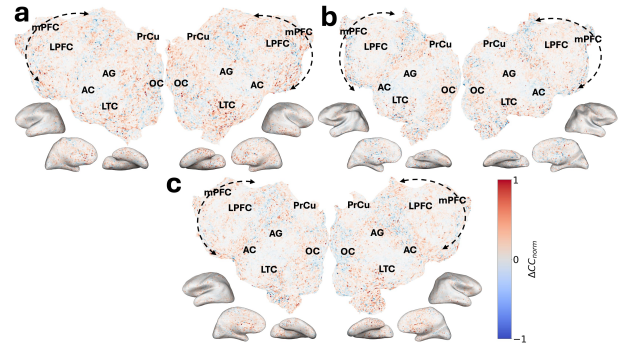


Figure 14. Encoding model performance improvements. (a-c) Vox-
elwise ΔCC_{norm} (MLP performance minus linear performance)
for audio input for subjects S1, S2, S3, respectively. Positive values
indicate MLP outperformance.

I.3. Voxelwise improvements from MLP (CC_{norm} analysis)

Figures 14, 13, and 15 each represent the performance im-
provements in voxelwise CC_{norm} values for semantic, au-
dio, and multimodal inputs, respectively, for each subject.
The improvements are more pronounced with CC_{norm} com-
pared to r as noise is taken into account.

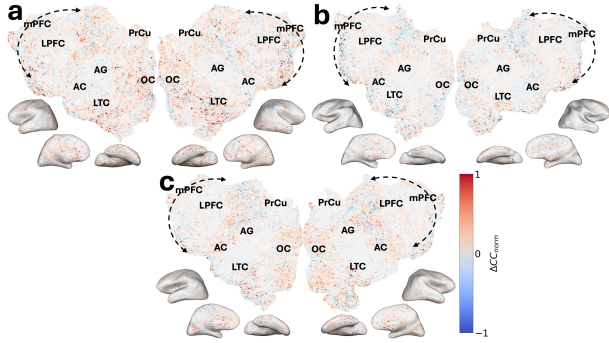


Figure 15. Encoding model performance improvements. (a-c) Voxewise ΔCC_{norm} (MLP performance minus linear performance) for multimodal input for subjects S1, S2, S3, respectively. Positive values indicate MLP outperformance.

I.4. Better spatio-temporal compartmentalization of brain function

To compare the performance between Whisper and LLAMA models, we define the Relative Error Difference (RED) for each voxel v at time t as:

$$\text{RED}(v, t) = |f_{\text{semantic}}(v, t) - y(v, t)| - |f_{\text{audio}}(v, t) - y(v, t)|$$

where $f_{\text{semantic}}(v, t)$ is the prediction from the semantic encoding model for voxel v at time t , $f_{\text{audio}}(v, t)$ is the prediction from the audio encoding model for voxel v at time t , and $y(v, t)$ represents the true value at voxel v and time t . A positive RED value indicates that the audio model outperforms the semantic model at that specific voxel and time, while a negative value indicates that the semantic model performs better.

In this analysis, we computed the RED between Whisper and LLAMA models for each voxel v at a given time t . For each region of interest (ROI), the average RED is calculated as:

$$\text{RED}_{\text{ROI}}(t) = \frac{1}{N} \sum_{v \in \text{ROI}} \text{RED}(v, t)$$

Where N is the number of voxels in the ROI. The correlation matrices were then computed over these ROI time series for both linear and nonlinear (MLP) encoders (Figure 16 (b, c)). A high correlation between two ROIs indicates that their semantic/audio processing temporal dynamics are similar over time.

For comparison, functional connectivity (FC) was also computed using the average fMRI signal for each voxel (Figure 16 a). Hierarchical clustering was then performed on the

correlation matrices, producing the dendrograms in panels (d-f).

As shown in Figure 16, panel (d) does not exhibit meaningful compartmentalization, indicating that the ROIs are not functionally clustered based on FC. However, the correlation matrices derived from RED (panels b, c) demonstrate clear block-diagonal structures, suggesting better functional compartmentalization. The dendrograms in panels (e, f) show that the ROIs cluster according to their functional roles, where the somatosensory and motor areas, visual areas, and auditory areas are grouped (even lower levels are grouped well (M1H/S1H, M1M/S1M, M1F/S1F, SMHA/SMFA, Broca/SPMv are grouped)) with nonlinear (MLP) models (f) achieving more accurate clustering than linear models (e). Specifically, panel (e) incorrectly clusters SMFA with S1M and M1M, whereas panel (f) correctly clusters SMHA and SMFA together before clustering them with other sensory and motor-related regions.

This study presents a novel approach, as it is the first to use fMRI language encoding models to group ROIs based not only on spatial dynamics but also on their temporal processing dynamics. Traditionally, voxel-wise functional classification or grouping has been the norm in fMRI analysis, focusing solely on spatial relationships. However, here with the help of fMRI encoders, we incorporate both spatial and temporal information, allowing for a more comprehensive view of brain function, especially in the context of semantic and auditory encoding.

In summary, using nonlinear (MLP) models leads to better functional compartmentalization. In fact, modularity Q values further confirm this: FC (a) scored 0.068, linear encoders (b) scored 0.145, and nonlinear encoders (c) scored 0.155, highlighting the improved functional clustering achieved with better encoders.

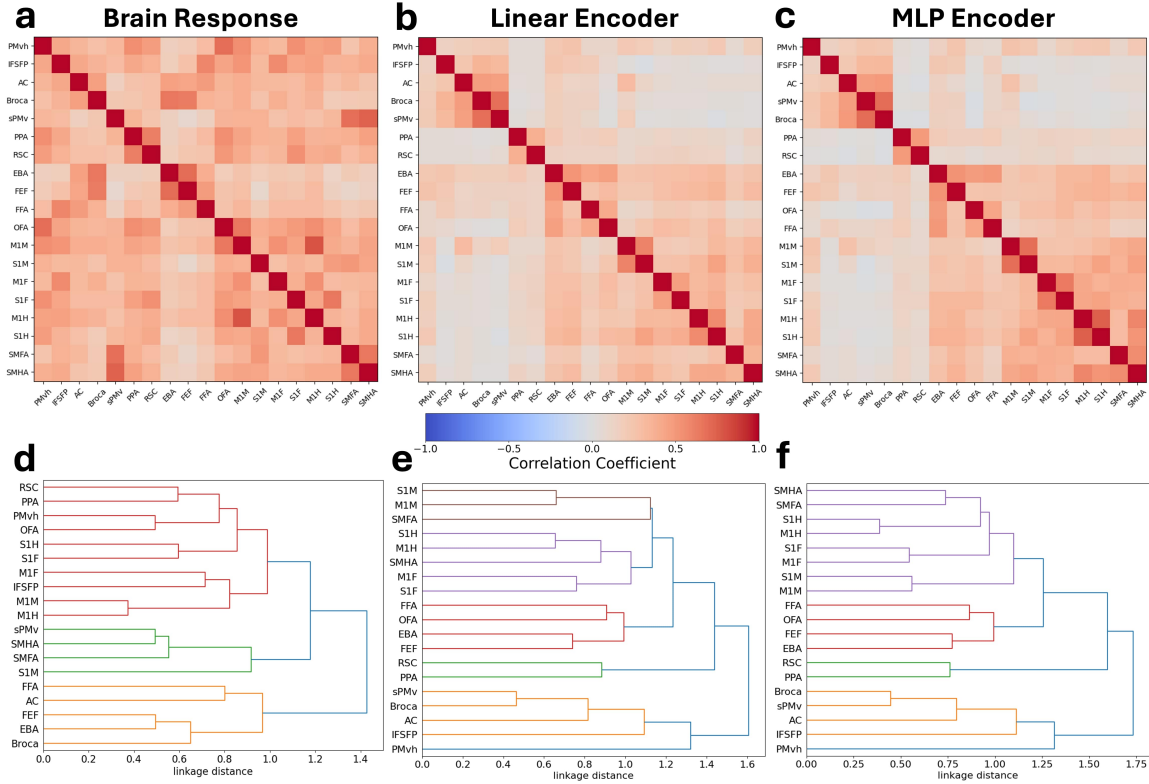


Figure 16. Spatio-temporal clustering based on Relative Error Difference (RED) between semantic and audio encoding models. Panels (a-c) display correlation matrices representing the temporal relationships between regions of interest (ROIs). For consistency, all the ROIs in (a,b,c) are ordered according to the most optimal ordering for (c). Panel (a) shows the functional connectivity (FC) matrix, calculated from the average fMRI signals. Panel (b) presents the correlation matrix from Relative Error Difference between Whisper and LLAMA using linear encoders, while panel (c) uses nonlinear (MLP) encoders, showing better functional compartmentalization with stronger block-diagonal structures. Panels (d-f) depict hierarchical clustering dendrograms derived from the correlation matrices in panels (a-c). Panel (d), based on FC, shows no clear compartmentalization of ROIs. Panel (e), based on linear encoders, show almost perfect functional clustering, though with inaccuracies (e.g., SMFA clustered with S1M/M1M). Panel (f), based on nonlinear (MLP) encoders, achieves better functional clustering, correctly grouping motor-related regions. The modularity Q values confirm this improvement: FC (a) scored 0.068, linear encoders (b) scored 0.145, and nonlinear encoders (c) scored 0.155, highlighting the advantage of nonlinear encoders for functional organization.

J. Improvements from multimodality

J.1. Voxelwise improvements from multimodality (r analysis)

This section shows the subject-wise plots of voxelwise Δr between multimodal linear/MLP and semantic/audio linear models (Figure 17, Figure 18). We observe consistent patterns of improvement when using multimodal models.

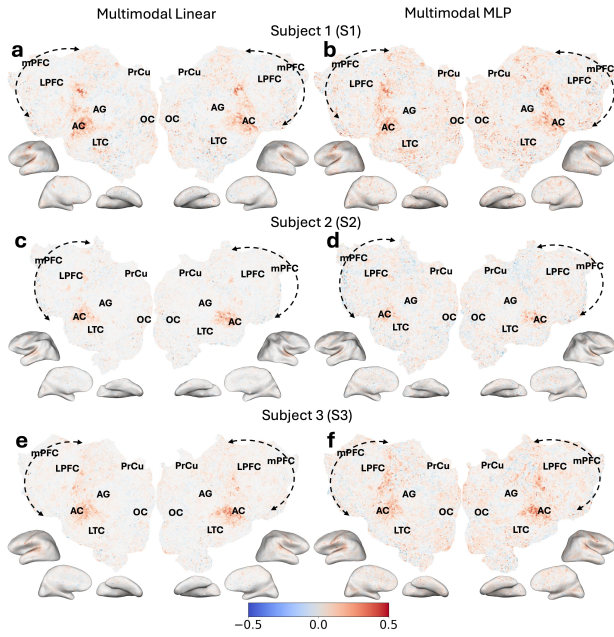


Figure 17. Subject-wise voxelwise Δr plots of multimodal models compared to semantic models. Panels (a-f) display voxelwise Δr values comparing multimodal and unimodal models across three subjects. Panels a, c, e show the difference between multimodal linear and semantic linear models, while panels b, d, f compare multimodal MLP and semantic linear models. Each row represents a different subject: Subject 1 (S1) in panels a-b, Subject 2 (S2) in panels c-d, and Subject 3 (S3) in panels e-f. Warmer colors indicate regions where the multimodal models outperform the unimodal linear models in prediction accuracy. The spatial patterns highlight enhanced encoding performance in key areas associated with semantic and auditory processing, such as the medial prefrontal cortex (mPFC), angular gyrus (AG), precuneus (PrCu), and lateral temporal cortex (LTC), emphasizing the benefits of multimodal models in capturing complex brain activity.

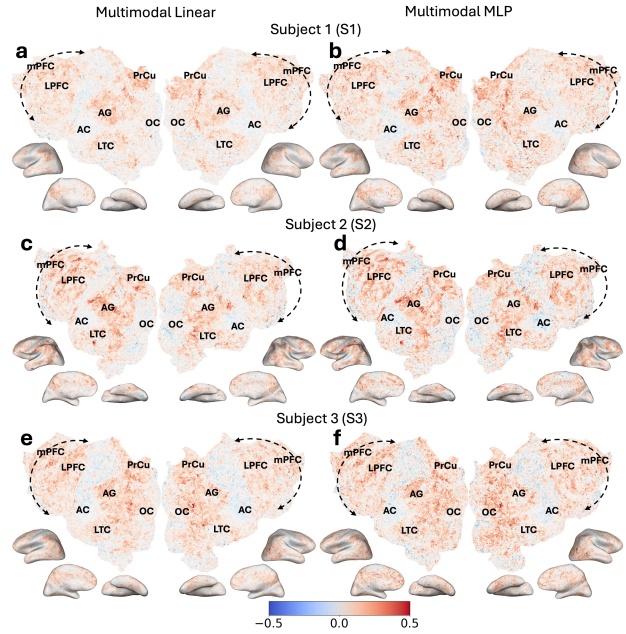


Figure 18. Subject-wise voxelwise Δr plots of multimodal models compared to audio models. Panels (a-f) display voxelwise Δr values comparing multimodal and unimodal models across three subjects. Panels a, c, e show the difference between multimodal linear and audio linear models, while panels b, d, f compare multimodal MLP and audio linear models. Each row represents a different subject: Subject 1 (S1) in panels a-b, Subject 2 (S2) in panels c-d, and Subject 3 (S3) in panels e-f. Warmer colors indicate regions where the multimodal models outperform the unimodal linear models in prediction accuracy.

J.2. Voxelwise improvements from multimodality (CC_{norm} analysis)

This section shows the subject-wise plots of voxelwise ΔCC_{norm} between multimodal linear/MLP and semantic/audio linear models (Figure 20, Figure 20). We observe consistent patterns of improvement when using multimodal models. The improvements are more noticeable with CC_{norm} compared to r as noise is taken into account.

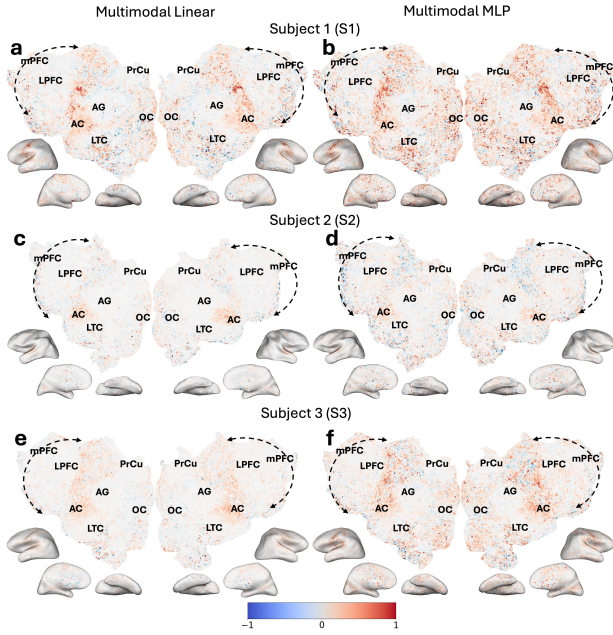


Figure 19. Subject-wise voxelwise ΔCC_{norm} plots of multimodal models compared to semantic models. Panels (a-f) display voxelwise ΔCC_{norm} values comparing multimodal and unimodal models across three subjects. Panels a, c, e show the difference between multimodal linear and semantic linear models, while panels b, d, f compare multimodal MLP and semantic linear models. Each row represents a different subject: Subject 1 (S1) in panels a-b, Subject 2 (S2) in panels c-d, and Subject 3 (S3) in panels e-f. Warmer colors indicate regions where the multimodal models outperform the unimodal linear models in prediction accuracy.

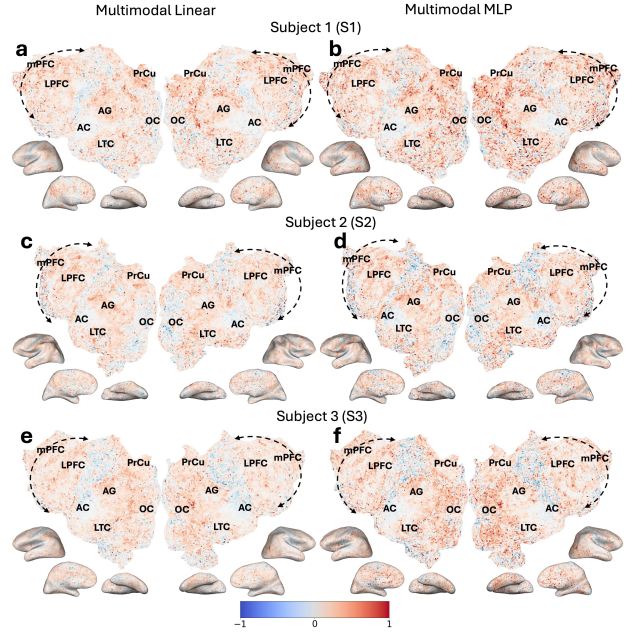


Figure 20. Subject-wise voxelwise ΔCC_{norm} plots of multimodal models compared to audio models. Panels (a-f) display voxelwise ΔCC_{norm} values comparing multimodal and unimodal models across three subjects. Panels a, c, e show the difference between multimodal linear and audio linear models, while panels b, d, f compare multimodal MLP and audio linear models. Each row represents a different subject: Subject 1 (S1) in panels a-b, Subject 2 (S2) in panels c-d, and Subject 3 (S3) in panels e-f. Warmer colors indicate regions where the multimodal models outperform the unimodal linear models in prediction accuracy.

J.3. ROI predictions improvements from multimodality

This section shows the ROI-wise improvements from using multimodal models (Figure 21)

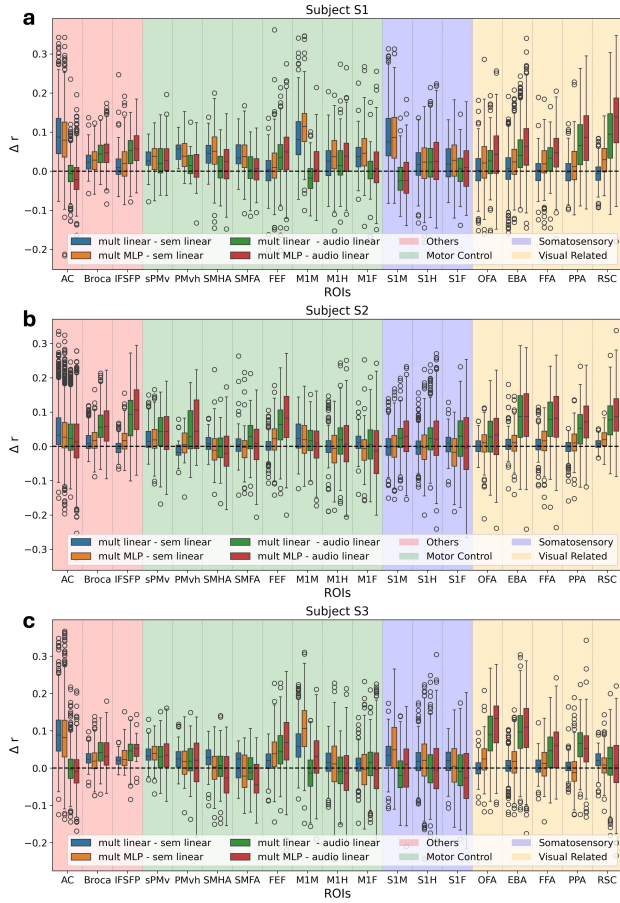


Figure 21. Subject-wise boxplots of performance differences (Δr) across different ROIs. The comparisons are made between different stimuli and encoding models: multimodal linear and multimodal MLP (mult MLP) models are compared against semantic (sem) and audio linear models. The ROIs are grouped into functional categories.

K. Improvements from nonlinearity and multimodality

K.1. Voxelwise improvements from DIMLP, and additional improvements from MLP (r analysis)

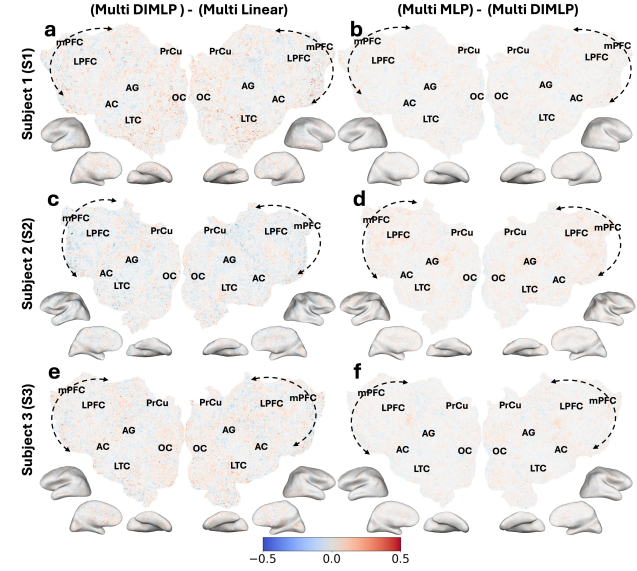


Figure 22. Nonlinearity Enhances Multimodal fMRI Predictions. Panels (a, c, e) show the voxelwise Δr values (DIMLP minus linear model), illustrating the improvements achieved through nonlinear processing within each modality, while largely limiting cross-modal interactions. Panels (b, d, f) display voxelwise Δr values (Multi MLP minus Multi DIMLP), highlighting the additional benefits of allowing nonlinear interactions between modalities (“Multi” denotes Multimodal). Each row represents the same subject: Subject 1 (S1) in panels a-b, Subject 2 (S2) in panels c-d, and Subject 3 (S3) in panels e-f. Warmer colors indicate regions where the nonlinear models outperform linear models.

K.2. Voxelwise improvements from DIMLP, and additional improvements from MLP (CC_{norm} analysis)

Figure 23 shows the voxel-wise performance improvements in voxelwise CC_{norm} values when incorporating nonlinear interactions. The improvements are more pronounced with CC_{norm} compared to r as noise is taken into account.

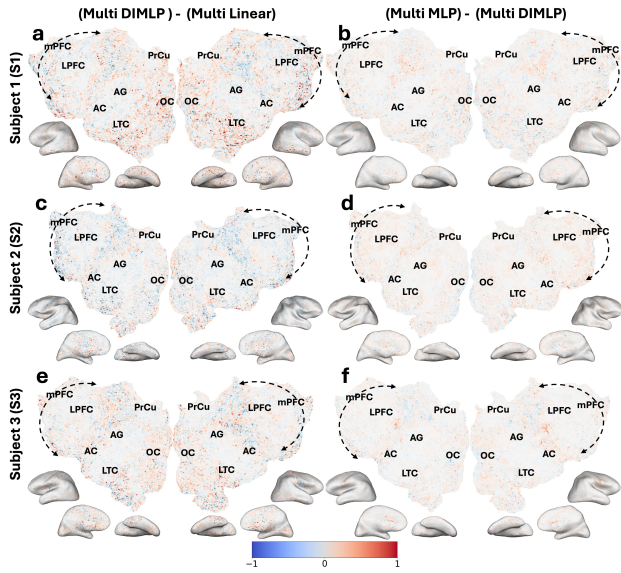


Figure 23. Nonlinearity Enhances Multimodal fMRI Predictions. Panels (a, c, e) show the voxelwise ΔCC_{norm} values (DIMLP minus linear model), illustrating the improvements achieved through nonlinear processing within each modality, while largely limiting cross-modal interactions. Panels (b, d, f) display voxelwise ΔCC_{norm} values (Multi MLP minus Multi DIMLP), highlighting the additional benefits of allowing nonlinear interactions between modalities (“Multi” denotes Multimodal). Each row represents the same subject: Subject 1 (S1) in panels a-b, Subject 2 (S2) in panels c-d, and Subject 3 (S3) in panels e-f. Warmer colors indicate regions where the nonlinear models outperform linear models.

K.3. ROI-wise improvements of multimodal DIMLP and MLP from multimodal linear model

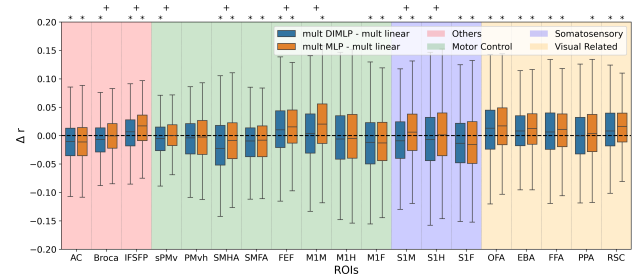


Figure 24. Box plot showing Δr across ROIs, where the Δr values are aggregated over all subjects. *multi* refers to multimodal, and *sem* refers to semantic encoders, and *DIMLP* refers to Delayed Interaction MLP, where only a *linear* interaction between modalities is allowed. The ROIs are color-coded by function. Regions where $\Delta r > 0$ with a p-value less than 0.05 are indicated by * symbols. Additionally, + symbols denote ROIs where there is a statistically significant difference (p-value < 0.05) between the two models based on a pairwise t-test. Voxelwise and ROI-wise plots for each subjects can be found in Figure 22 (Appendix), and Figure 25 (Appendix), respectively.

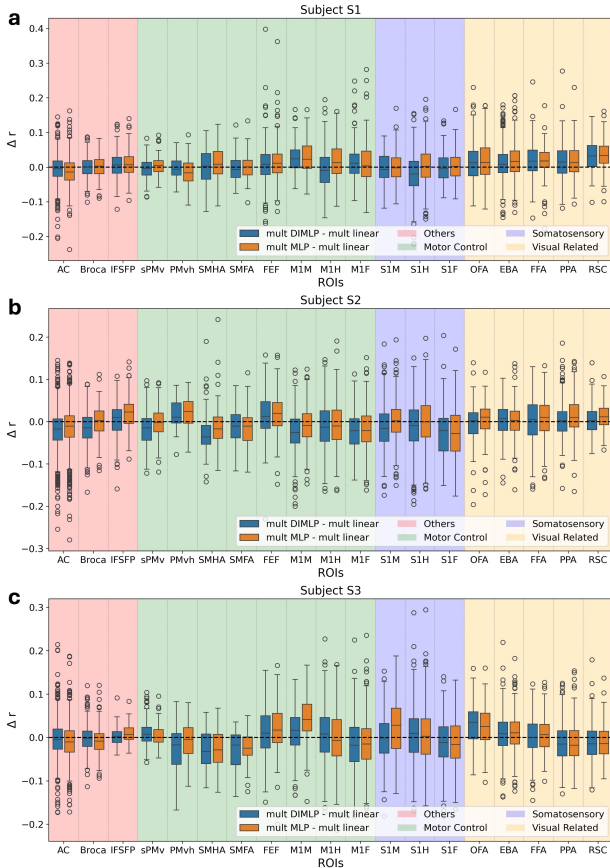


Figure 25. Subject-wise boxplots of voxel-wise differences (Δr) across different ROIs. The comparisons are made between different encoding models: multimodal MLP and multimodal DIMLP models are compared against multimodal linear models. The ROIs are grouped into functional categories.

L. Variance partitioning analysis

To quantify the unique contributions of different feature spaces in our nonlinear multimodal encoding models, we employed a variance partitioning analysis similar to (de Heer et al., 2017). This approach allowed us to determine how much variance could be uniquely explained by each feature versus that explained by a multiple features. We estimated both the fraction of variance explained by each feature space individually and the fraction that might be equally well explained by combinations of feature spaces.

We show our variance partitioning analysis results in three complementary ways: 1) voxel-wise variance partition results (Appendix L.2), 2) voxel-wise plots showing the largest variance partition for each voxel (Appendix L.3), and 3) ROI-wise Venn diagrams illustrating the distribution of variance explained across different brain regions (Appendix L.4).

For this analysis, we fit models with all possible combinations of feature spaces: two single-feature models (audio and semantic), one model combining both features (semantic-audio), and examined the distribution of variance explained within brain regions. This allowed us to decompose the total explained variance into three components: variance uniquely explained by audio features, variance uniquely explained by semantic features, and variance jointly explained by both feature spaces.

L.1. Summary of variance partitioning results

Looking at the results of Appendix L.2, we observe that joint variance dominates across most cortical regions, contrasting with (de Heer et al., 2017) where semantic only features showed greater dominance. This difference likely stems from our feature choices - whereas (de Heer et al., 2017) used spectral and articulatory features that primarily contained information relevant mostly only to auditory cortex, our use of Whisper features provides richer auditory representations that enable better predictions beyond traditional auditory regions. This finding aligns with our earlier argument (Section 3.2.2) that multiple modalities jointly contribute to neural computations across the cortex rather than having one modality dominate.

The dominance pattern of joint variance is consistent both within and near AC, with a notable exception in early auditory regions where audio features show unique contributions. This hierarchical organization suggests that while early AC predominantly processes pure acoustic information, later AC regions integrate both semantic and auditory features for higher-level speech processing. The unique contribution of audio features in early AC is noteworthy as it suggests preservation of modality-specific processing at early sensory stages despite using rich Whisper features.

Also, Appendix L.3 reveals distinct spatial patterns in feature representation across cortical regions. The prefrontal cortex exhibits mixed dominance patterns, showing both joint semantic-audio representation and semantic-only areas. While early auditory cortex shows expected unique audio contributions, we also observe audio-specific representation in motor-sensory mouth areas (M1M, S1M), though this pattern varies across subjects.

The ROI-wise analysis in Appendix L.4 reveals that joint semantic-audio features dominate cortical representation, accounting for approximately 65% of significantly predicted voxels across the entire cortex. Core language-processing regions (AC, Broca’s area, sPMv) show particularly strong joint representation (around 80 to 90%), supporting our hypothesis that speech comprehension relies on integrated multimodal processing. This integration is consistently observed across subjects, though some ROIs (e.g., PMvh in Subject S2 with only 14 voxels) have insufficient data for

reliable interpretation. The transition from linear to MLP encoders increases the total number of significantly predicted voxels while maintaining similar representation patterns, indicating that nonlinear encoding primarily enhances prediction accuracy rather than fundamentally altering feature representation structure.

L.2. Variance partitioning of various models

Due to file size constraints, we only show the voxel-wise variance partitioning result of subject S1 using linear encoders, Figure 26. The rest have been moved to the supplementary material.

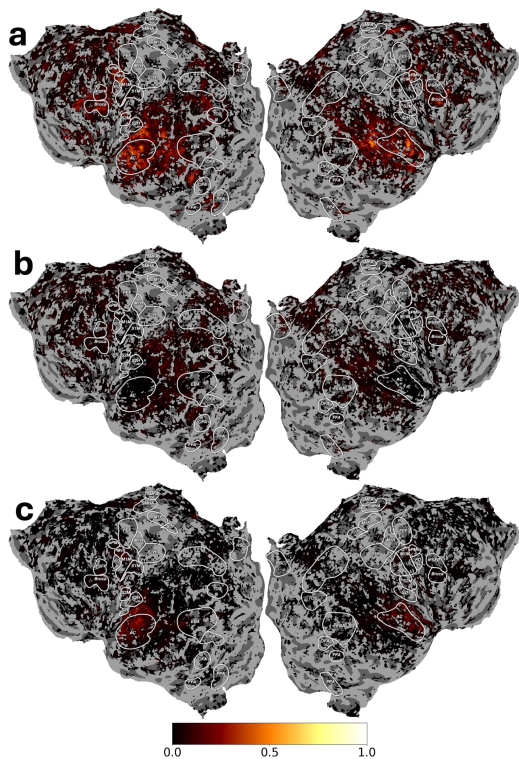


Figure 26. Voxelwise variance partitioning analysis showing the contributions of different feature types to prediction accuracy for a subject S1 using linear models. The flatmaps display (a) variance jointly explained by audio and semantic features, (b) variance uniquely explained by semantic features, and (c) variance uniquely explained by audio features. Values shown are normalized correlations (CC_{norm}) for voxels where the joint model achieved significant prediction ($q(\text{FDR}) < 0.01$).

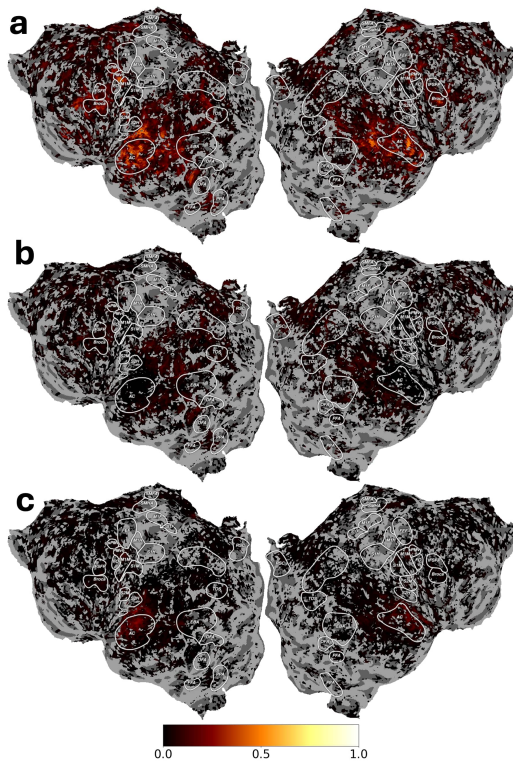


Figure 27. Voxelwise variance partitioning analysis showing the contributions of different feature types to prediction accuracy for a subject S1 using MLP models. The flatmaps display (a) variance jointly explained by audio and semantic features, (b) variance uniquely explained by semantic features, and (c) variance uniquely explained by audio features. Values shown are normalized correlations (CC_{norm}) for voxels where the joint model achieved significant prediction ($q(\text{FDR}) < 0.01$).

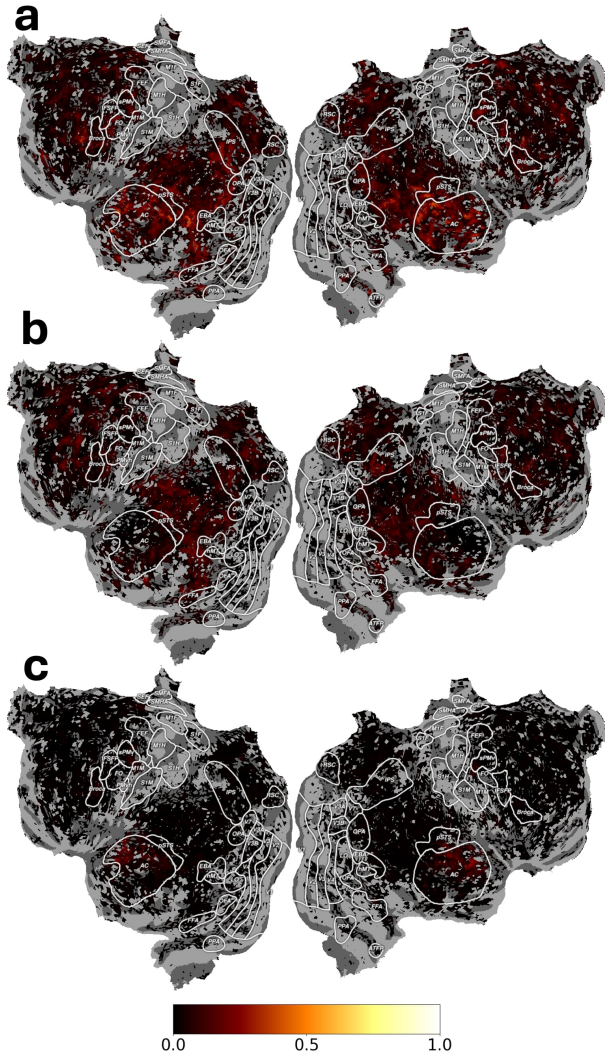


Figure 28. Voxelwise variance partitioning analysis showing the contributions of different feature types to prediction accuracy for a subject S2 using linear models. The flatmaps display (a) variance jointly explained by audio and semantic features, (b) variance uniquely explained by semantic features, and (c) variance uniquely explained by audio features. Values shown are normalized correlations (CC_{norm}) for voxels where the joint model achieved significant prediction ($q(\text{FDR}) < 0.01$).

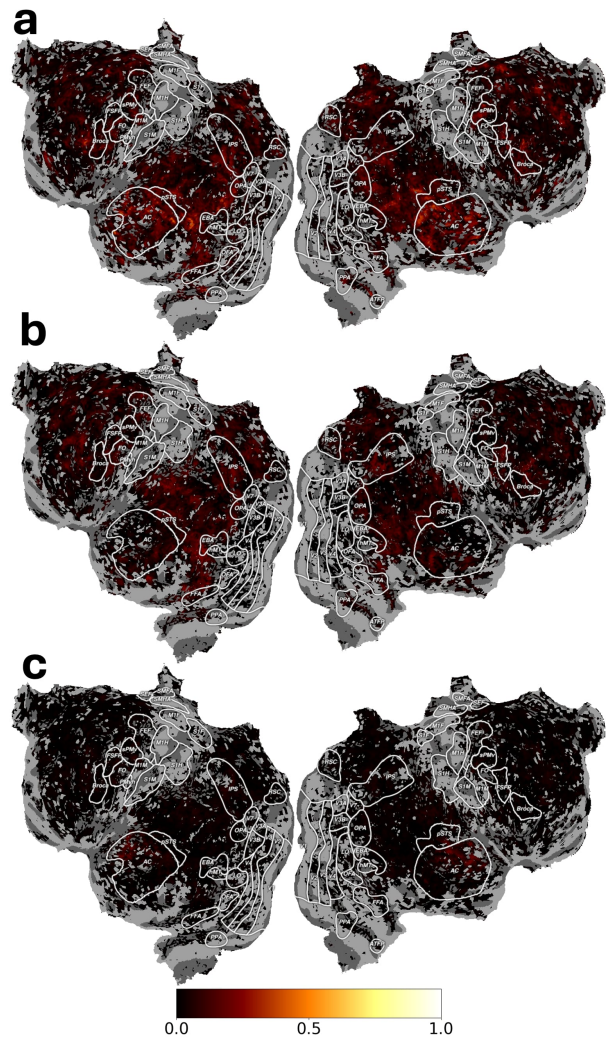


Figure 29. Voxelwise variance partitioning analysis showing the contributions of different feature types to prediction accuracy for a subject S2 using MLP models. The flatmaps display (a) variance jointly explained by audio and semantic features, (b) variance uniquely explained by semantic features, and (c) variance uniquely explained by audio features. Values shown are normalized correlations (CC_{norm}) for voxels where the joint model achieved significant prediction ($q(\text{FDR}) < 0.01$).

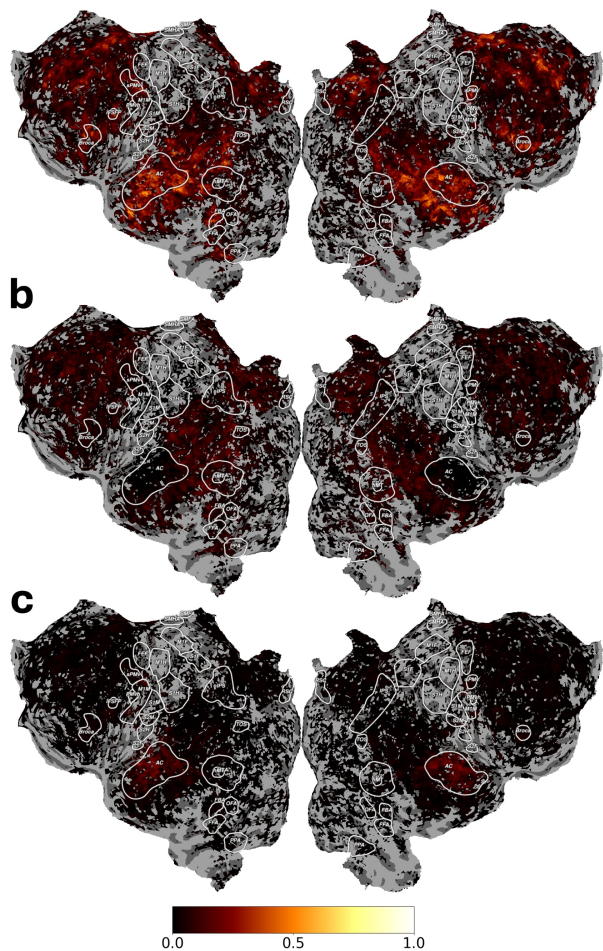


Figure 30. Voxelwise variance partitioning analysis showing the contributions of different feature types to prediction accuracy for a subject S3 using linear models. The flatmaps display (a) variance jointly explained by audio and semantic features, (b) variance uniquely explained by semantic features, and (c) variance uniquely explained by audio features. Values shown are normalized correlations (CC_{norm}) for voxels where the joint model achieved significant prediction ($q(\text{FDR}) < 0.01$).

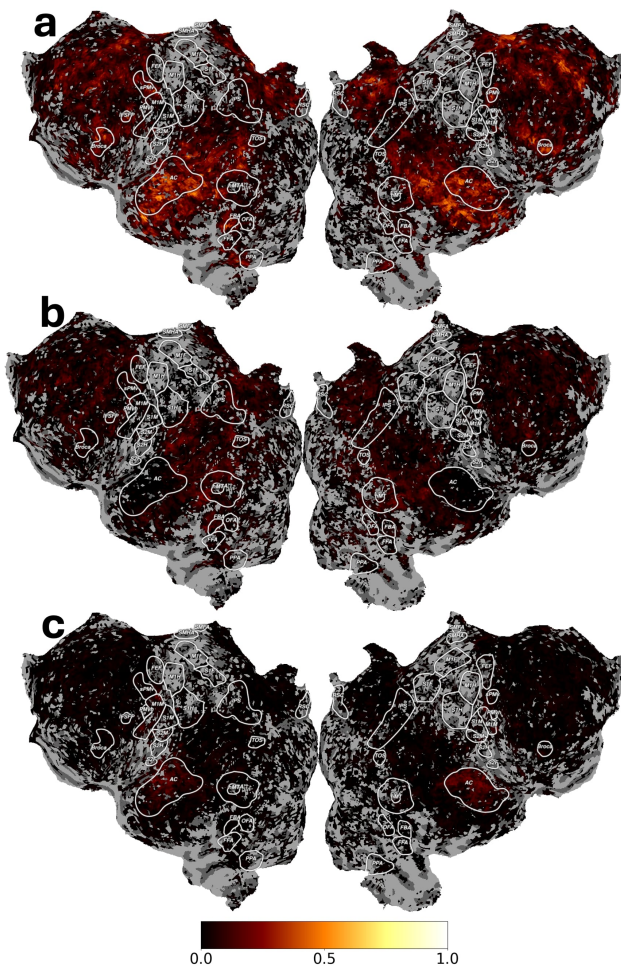


Figure 31. Voxelwise variance partitioning analysis showing the contributions of different feature types to prediction accuracy for a subject S3 using MLP models. The flatmaps display (a) variance jointly explained by audio and semantic features, (b) variance uniquely explained by semantic features, and (c) variance uniquely explained by audio features. Values shown are normalized correlations (CC_{norm}) for voxels where the joint model achieved significant prediction ($q(\text{FDR}) < 0.01$).

L.3. Largest variance partitioning for each voxel

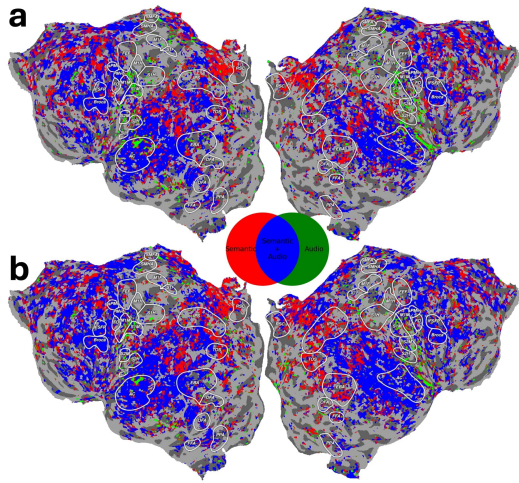


Figure 32. Voxelwise analysis showing the largest variance explained by each feature type for all significantly predicted voxels ($q(\text{FDR}) < 0.01$) for subject S1. The flatmaps display which feature partition (semantic in red, audio in green, or their combination in blue) best explains the variance in each cortical voxel using (a) linear and (b) MLP encoders, with outlined regions indicating key functional areas.

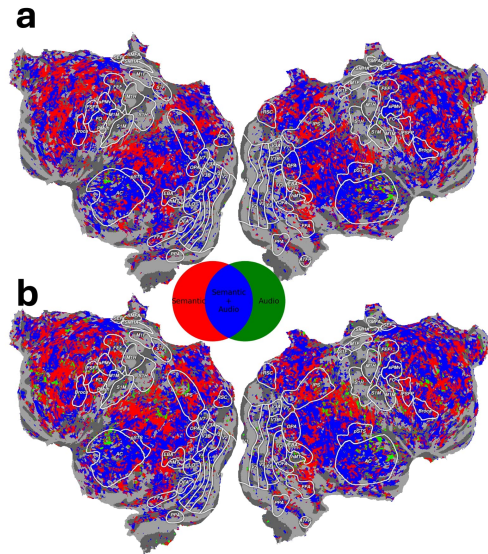


Figure 33. Same as Figure 32, but for subject S2

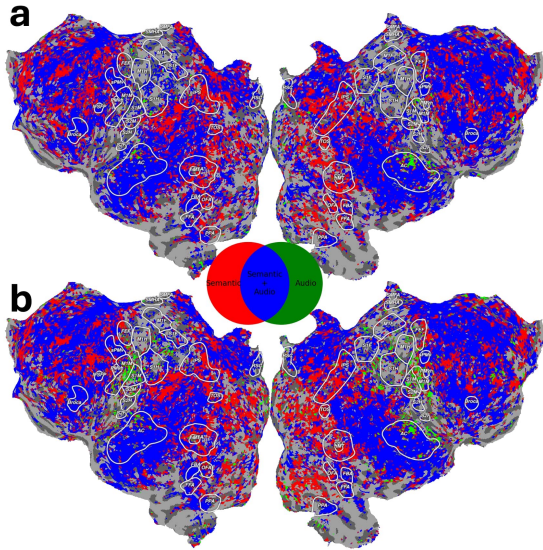


Figure 34. Same as Figure 32, but for subject S3

L.4. Variance partitioning Venn diagram

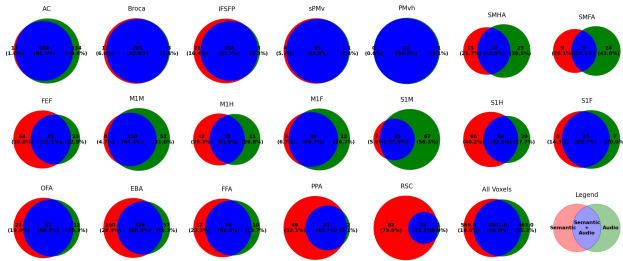


Figure 35. Venn diagrams showing the distribution of explained variance across different brain regions of interest (ROIs) for subject S1, using linear encoder. Each diagram displays the unique and shared variance explained by semantic features (red), audio features (green), and their overlap (blue). Values indicate the number of significantly predicted voxels and their percentages. Only the voxels that was predicted statistically significantly ($q(\text{FDR}) < 0.01$) was used in the analysis

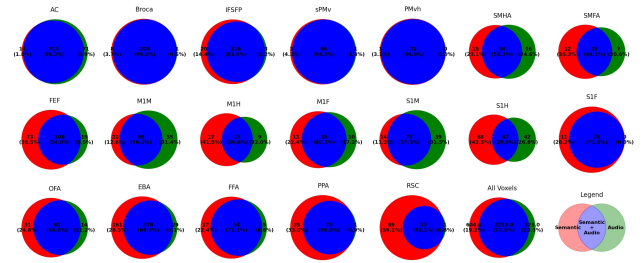


Figure 36. Venn diagrams showing the distribution of explained variance across different brain regions of interest (ROIs) for subject S1, using MLP encoder. Refer to Fig 35 for more detail.

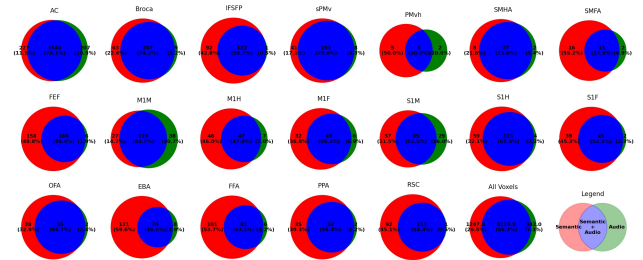


Figure 37. Venn diagrams showing the distribution of explained variance across different brain regions of interest (ROIs) for subject S2, using linear encoder. Refer to Fig 35 for more detail.

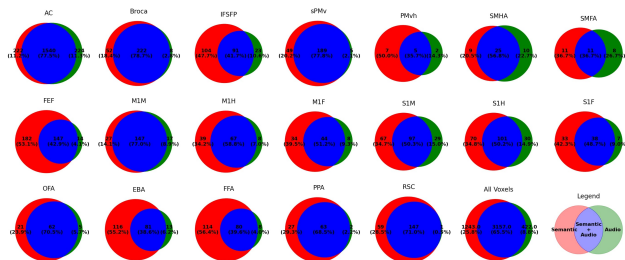


Figure 38. Venn diagrams showing the distribution of explained variance across different brain regions of interest (ROIs) for subject S2, using MLP encoder. Refer to Fig 35 for more detail.

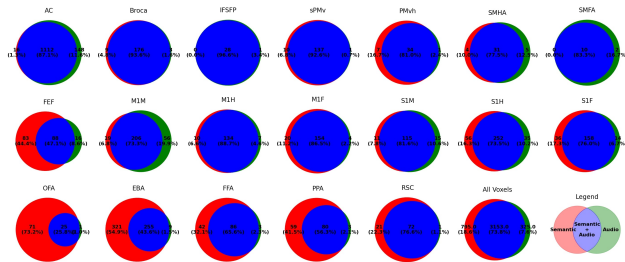


Figure 39. Venn diagrams showing the distribution of explained variance across different brain regions of interest (ROIs) for subject S3, using linear encoder. Refer to Fig 35 for more detail.

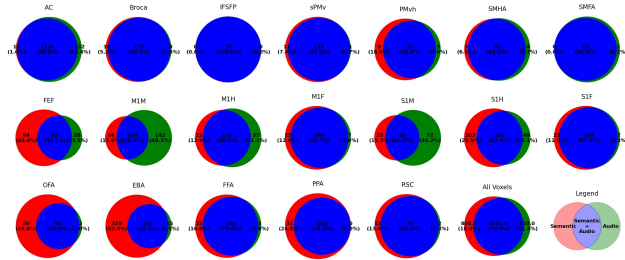


Figure 40. Venn diagrams showing the distribution of explained variance across different brain regions of interest (ROIs) for subject S3, using MLP encoder. Refer to Fig 35 for more detail.