
Multi-modal brain encoding models for multi-modal stimuli

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Despite participants engaging in single modality stimuli, such as watching images
2 or silent videos, recent work has demonstrated that multi-modal Transformer
3 models can predict visual brain activity impressively well, even with incongruent
4 modality representations. This raises the question of how accurately these multi-
5 modal models can predict brain activity when participants are engaged in multi-
6 modal stimuli. As these models grow increasingly popular, their use in studying
7 neural activity provides insights into how our brains respond to such multi-modal
8 naturalistic stimuli, i.e., where it separates and integrates information from different
9 sensory modalities. We investigate this question by using multiple unimodal
10 and two types of multi-modal models—cross-modal and jointly pretrained—to
11 determine which type of models is more relevant to fMRI brain activity when
12 participants were engaged in watching movies (videos with audio). We observe that
13 both types of multi-modal models show improved alignment in several language
14 and visual regions. This study also helps in identifying which brain regions
15 process unimodal versus multi-modal information. We further investigate the
16 impact of removal of unimodal features from multi-modal representations and
17 find that there is additional information beyond the unimodal embeddings that
18 is processed in the visual and language regions. Based on this investigation, we
19 find that while for cross-modal models, their brain alignment is partially attributed
20 to the video modality; for jointly pretrained models, it is partially attributed to
21 both the video and audio modalities. The inability of individual modalities in
22 explaining the brain alignment effectiveness of multi-modal models suggests that
23 multi-modal models capture additional information processed by all brain regions.
24 This serves as a strong motivation for the neuro-science community to investigate
25 the interpretability of these models for deepening our understanding of multi-modal
26 information processing in brain.

27 1 Introduction

28 The study of brain encoding aims at predicting the neural brain activity recordings from an input
29 stimulus representation. Recent brain encoding studies use neural models as a powerful approach to
30 better understand the information processing in the brain in response to naturalistic stimuli (Oota
31 et al., 2023a). Current encoding models are trained and tested on brain responses captured from
32 participants who are engaged in a *single stimulus modality*, using stimulus representations extracted
33 from AI systems that are pretrained on single modality, such as language (Wehbe et al., 2014; Jain &
34 Huth, 2018; Toneva & Wehbe, 2019; Caucheteux & King, 2020; Schrimpf et al., 2021; Toneva et al.,
35 2022; Aw & Toneva, 2023), vision (Yamins et al., 2014; Eickenberg et al., 2017; Schrimpf et al.,
36 2018; Wang et al., 2019) or speech (Millet et al., 2022; Vaidya et al., 2022; Tuckute et al., 2023). In
37 this paper, we build encoding models where participants are engaged with *multi-modal stimuli* (e.g.,
38 watching movies that also include audio). We explore multi-modal stimulus representations extracted

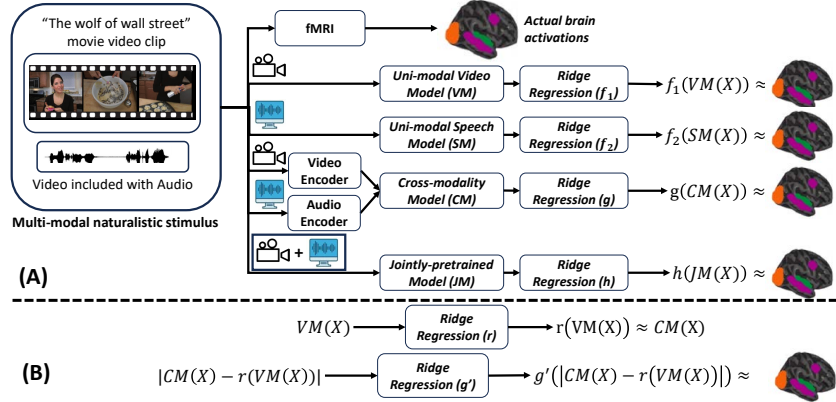


Figure 1: (A) Overview of our proposed Multi-modal Brain Encoding Pipeline. (B) Residual Analysis.

39 using Transformer (Vaswani et al., 2017) based multi-modal models. Our analysis focuses on their
 40 alignment with both uni- and multi-modal brain regions.

41 There is a growing evidence that the human brain’s ability for multi-modal processing is underpinned
 42 by synchronized cortical representations of identical concepts across various sensory modalities (Gau-
 43 thier et al., 2003; Bracci & Op de Beeck, 2023). Reflecting similar principles, the recent advances in
 44 AI systems have led to the development of multi-modal models (like CLIP (Radford et al., 2021),
 45 ImageBind (Girdhar et al., 2023), and TVLT (Tang et al., 2022)) using massive interleaved image-text
 46 data, speech-text data or video-audio-text data to represent multi-modal input. This recent progress in
 47 AI has stimulated advancements in brain encoding models (Doerig et al., 2022; Oota et al., 2022;
 48 Popham et al., 2021; Wang et al., 2022; Tang et al., 2024; Nakagi et al., 2024) that learn effectively
 49 from multiple input modalities, despite participants being engaged with single stimulus modality
 50 during experiments, e.g., watching natural scene images, or silent movie clips. However, these studies
 51 have experimented with subjects engaged with single-modality stimulus, leaving the full potential of
 52 these models in true multi-modal scenarios still unclear.

53 Using brain recordings of participants watching several popular movies included with audio (St-
 54 Laurent et al., 2023), we investigate several research questions. First, we investigate the effectiveness
 55 of multi-modal stimulus representations obtained using multi-modal models versus unimodal models
 56 for brain encoding. Multi-modal models are of two broad types: (i) cross-modal pretrained models,
 57 where first individual modality encoders are trained and then cross-modal alignment is performed, and
 58 (ii) jointly pretrained models, which involve combining data from multiple modalities and training a
 59 single joint encoder. Hence, we also investigate which of the two types (cross-modal versus joint) are
 60 better for encoding. In this work, we focus on one cross-modal (ImageBind), one jointly pretrained
 61 (TVLT), three video and two speech models. Additionally, we explore which modality representations
 62 are more brain relevant, and identify which brain regions process uni- and multi-modal information.
 63 Overall, this research utilizes various modality representations to develop encoding models based on
 64 fMRI responses within a multi-modal model framework (see Fig. 1 for workflow).

65 Using our multi-modal brain encoding approach, we examine several insights. First, we use previous
 66 neuroscience findings that have identified brain regions involved in visual, language and auditory
 67 processing, and investigate how well our model aligns with these regions when both the model and a
 68 human participant watch the same multi-modal video stimuli. Second, we expect that multi-modal
 69 models which can learn cross-modal and joint embeddings across modalities in a brain-relevant
 70 way would significantly align with these regions. However, alignment with these brain regions
 71 doesn’t necessarily mean that the model is effectively learning from multiple modalities, as unimodal
 72 models for vision or language or audio have also been shown to significantly align with these brain
 73 regions (Wehbe et al., 2014; Toneva et al., 2022; Schrimpf et al., 2021; Millet et al., 2022; Vaidya
 74 et al., 2022). To check the second aspect, we investigate this question via a direct approach, closely
 75 related to previous studies (Toneva et al., 2022; Oota et al., 2023b,c). For each modality, we analyze
 76 how the alignment between brain recordings and multi-modal model representations is affected by
 77 the elimination of information related to that particular modality from the model representation.

78 Our analysis of multi-modal brain alignment leads to several key conclusions: (1) Both cross-modal
 79 and jointly pretrained models demonstrate significantly improved brain alignment with language

80 regions (AG, PCC, PTL, and IFG) and visual regions (EVC and MT) when analyzed against unimodal
81 video data. In contrast, compared to unimodal speech-based models, all multi-modal embeddings
82 show significantly better brain alignment, except in the OV (object visual processing) region. This
83 highlights the ability of multi-modal models to capture additional information—either through
84 knowledge transfer or integration between modalities—which is crucial for multi-modal brain
85 alignment. (2) Using our residual approach, we find that the improved brain alignment in cross-
86 modal models can be partially attributed to the removal of video features alone, rather than auditory
87 features. On the other hand, the improved brain alignment in jointly pretrained models can be partially
88 attributed to the removal of both video and auditory features.

89 Overall, we make the following contributions in this paper. (1) To the best of our knowledge, this
90 study is the first to leverage both cross-modal and jointly pretrained multi-modal models to perform
91 brain alignment while subjects are engaged with multi-modal naturalistic stimuli. (2) We evaluate the
92 performance of several unimodal Transformer models (three video and two audio) and measure their
93 brain alignment. (3) Additionally, we remove unimodal features from multi-modal representations
94 to explore the impact on brain alignment before and after their removal. We will release code upon
95 publication of this paper.

96 2 Related Work

97 **Multi-modal models.** Pretrained Transformer-based models have been found to be very effective in
98 various tasks related to language (Devlin et al., 2019; Radford et al., 2019), speech (Baevski et al.,
99 2020), and images (Dosovitskiy et al., 2020). To learn associations between pairs of modalities,
100 Transformer models have been pretrained on multiple modalities, showing excellent results in multi-
101 modal tasks like visual question answering and visual common-sense reasoning. These multi-modal
102 models are pretrained in two different ways: (i) cross-modal models that integrate information
103 from multiple modalities and learn a joint encoder, such as VisualBERT (Li et al., 2019) and
104 ImageBind (Girdhar et al., 2023), and (ii) jointly pretrained models like LXMERT (Tan & Bansal,
105 2019), CLIP (Radford et al., 2021), ViLBERT (Lu et al., 2019), and TVLT (Tang et al., 2022) which
106 fuse individual modality encoders at different stages, transferring knowledge from one modality
107 to another. In this work, we investigate how the representations extracted from *cross-modal and*
108 *jointly-pretrained Transformer models* align with human brain recordings when participants engage
109 with multi-modal stimuli.

110 **Brain Encoding using Multi-modal Models.** Since human brain perceives the environment using
111 information from multiple modalities (Gauthier et al., 2003), examining the alignment between
112 language and visual representations in the brain by training encoding models on fMRI responses,
113 while extracting joint representations from multi-modal models, can offer insights into the relation-
114 ship between the two modalities. For instance, it has been shown that multi-modal models like
115 CLIP (Radford et al., 2021) better predict neural responses in the high-level visual cortex as compared
116 to previous vision-only models (Doerig et al., 2022; Wang et al., 2022). Additionally, Tang et al.
117 (2024) demonstrate the use of multi-modal models in a cross-modal experiment to assess how well
118 the language encoding models can predict movie-fMRI responses and how well the vision encoding
119 models can predict narrative story-fMRI. Nakagi et al. (2024) analyzed fMRI related to video content
120 viewing and found distinct brain regions associated with different semantic levels, highlighting the
121 significance of modeling various levels of semantic content simultaneously. However, these studies
122 have experimented with subjects engaged with single-modality stimulus, leaving the full potential of
123 these models in true multi-modal scenarios still unclear. Recently, Dong & Toneva (2023) interpreted
124 the effectiveness of pretrained versus finetuned multi-modal video transformer using video+text
125 stimuli-based brain activity. However, they did not perform any cross-modal vs jointly-pretrained
126 model analysis or analysis of multi-modal versus unimodal models, leaving it unclear which type
127 of multi-modal models perform best for brain activity prediction. Further, unlike them, we study
128 video+audio stimuli, and perform comprehensive residual analysis.

129 3 Dataset Curation

130 **Brain Imaging Dataset.** We experiment with a multi-modal naturalistic fMRI dataset, Movie10 (St-
131 Laurent et al., 2023) obtained from the Courtois NeuroMod databank. This dataset was collected
132 while six human subjects passively watched four different movies: *The Bourne supremacy* (~100

133 mins), *The wolf of wall street* (~170 mins), *Hidden figures* (~120 mins) and *Life* (~50 mins). Among
134 these, *Hidden figures* and *Life* are repeated twice, with the repeats used for testing and the remaining
135 movies for training. In this work, we use *Life* movie for testing where we average the two repetitions
136 to reduce noise in brain data. This dataset is one of the largest publicly available multi-modal fMRI
137 dataset in terms of number of samples per participant. It includes 4024 TRs (Time Repetitions) for
138 *The Bourne supremacy*, 6898 TRs for *The wolf of wall street* used in train and 2028 TRs for *Life* in
139 test. The fMRI data is collected every 1.49 seconds (= 1TR).

140 The dataset is already preprocessed and projected onto the surface space (“fsaverage6”). We use the
141 multi-modal parcellation of the human cerebral cortex based on the Glasser Atlas (which consists
142 of 180 regions of interest in each hemisphere) to report the ROI (region of interest) analysis for the
143 brain maps (Glasser et al., 2016). This includes four visual processing regions (early visual (EV),
144 object-related areas (LO), face-related areas (OFA) and scene-related areas (PPA)), one early auditory
145 area (AC), and eight language-relevant regions, encompassing broader language regions: angular
146 gyrus (AG), anterior temporal lobe (ATL), posterior temporal lobe (PTL), inferior frontal gyrus (IFG),
147 inferior frontal gyrus orbital (IFGOrb), middle frontal gyrus (MFG), posterior cingulate cortex (PCC)
148 and dorsal medium prefrontal cortex (dmPFC), based on the Fedorenko lab’s language parcels (Milton
149 et al., 2021; Desai et al., 2022). We list the detailed sub-ROIs of these ROIs in Appendix B.

150 **Estimating dataset cross-subject prediction accuracy.** To account for the intrinsic noise in
151 biological measurements, we adapt Schrimpf et al. (2021)’s method to estimate the cross-subject
152 prediction accuracy for a model’s performance for the Movie10 fMRI datasets. By subsampling
153 fMRI datasets from 6 participants, we generate all possible combinations of s participants ($s \in [2,6]$)
154 for watching movies, and use a voxel-wise encoding model (see Sec. 5) to predict one participant’s
155 response from others. Note that the estimated cross-subject prediction accuracy is based on the
156 assumption of a perfect model, which might differ from real-world scenarios, yet offers valuable
157 insights into model’s performance. We estimate cross-subject prediction accuracy in three settings:
158 (i) training with *The Bourne supremacy* and testing with *Life* data, (ii) training with *The wolf of wall*
159 *street* and testing with *Life* data, and (iii) training with both *The Bourne supremacy* and *The wolf*
160 *of wall street* and testing with *Life* data. We present the average cross-subject prediction accuracy
161 across voxels for the *Movie10 fMRI* dataset and across the three settings in Appendix A.

162 4 Methodology

163 4.1 Multi-modal models

164 To analyse how human brain process information while engaged in multi-modal stimuli, we use recent
165 popular deep learning models to explore multiple modalities information and build the encoding
166 models in two different ways: “cross-modality pretraining” and “joint pretraining”.

167 **Cross-modality Pretrained Multi-modal Models.** Cross-modality representations involve transfer-
168 ring information or learning from one modality to another. For example, in a cross-modal learning
169 scenario, text descriptions can be used to improve the accuracy of image/video recognition tasks.
170 This approach is often used in scenarios where one modality might have limited data or less direct
171 relevance but can be informed by another modality.

172 Recently, a cross-modal model called ImageBind (IB) (Girdhar et al., 2023) has shown immense
173 promise in binding data from six modalities at once, without the need for explicit supervision.
174 ImageBind model uses separate encoders for each individual modality and learns a single shared
175 representation space by leveraging multiple types of image-paired data. ImageBind consists of 12
176 layers and outputs a 1024 dimensional representation for each modality.

177 **Jointly Pretrained Multi-modal Models.** Jointly pretrained multi-modal model representations,
178 on the other hand, involve combining data from multiple modalities to build a more comprehensive
179 joint understanding to improve decision-making processes. The system processes these diverse inputs
180 concurrently to make more informed and robust decisions.

181 TVLT (Zellers et al., 2022) is an end-to-end Text-less Vision-Language multi-modal Transformer
182 model for learning joint representations of video and speech from YouTube videos. This joint encoder
183 model consists of a 12-layer encoder (hidden size 768) and uses masked autoencoding objective for
184 both videos and speech. Given the video-speech pairs, the TVLT model provides 768 dimensional
185 representations for each modality across 12 layers.

186 **Extraction of multi-modal features.** To extract video and audio embedding representations from
187 multi-modal models for the brain encoding task, we input video and audio pairs at each TR and
188 obtain the aligned embeddings for the two modalities. Here, we first segment the input video and
189 audio into clips corresponding to 1.49 seconds, which matches the fMRI image rate. For both the
190 models, ImageBind and TVLT, we use the pretrained Transformer weights. ImageBind generates
191 an embedding for each modality (IB video and IB audio) in an aligned space. We concatenate these
192 embeddings to create what we refer to as IB concat embeddings. On the other hand, TVLT provides
193 a joint embedding across all modalities at each layer. Only for the last layer, TVLT provides an
194 embedding for each modality.

195 4.2 Unimodal Models

196 To investigate the effectiveness of multi-modal representations in comparison to representations for
197 individual modalities, we use the following methods to obtain embeddings for individual modalities.

198 **Video-based models.** To extract representations of the video stimulus, we use three popular pretrained
199 Transformer video-based models from Huggingface (Wolf et al., 2020): (1) Vision Transformer Base
200 (ViT-B) (Dosovitskiy et al., 2020), (2) Video Masked Autoencoders (VideoMAE) (Tong et al., 2022)
201 and (3) Video Vision Transformer (ViViT) (Arnab et al., 2021). Details of each model are reported in
202 Table 1 in Appendix.

203 **Speech-based models.** Similar to video-based models, we use two popular pretrained Transformer
204 speech-based models from Huggingface: (1) Wav2Vec2.0 (Baevski et al., 2020) and (2) AST (Baade
205 et al., 2022). Details of each model are reported in Table 1 in Appendix.

206 **Extraction of video features.** ViT-B (Dosovitskiy et al., 2020), the underlying video encoder model
207 for ImageBind is used for extracting representations for all frames in each TR for every video. To
208 extract embedding at each TR, we average all frame embeddings and obtain the corresponding video
209 representation. For VideoMAE and ViViT, we directly obtain the video embeddings for each TR. All
210 3 models provide 768 dimensional representations and all of them are 12-layer Transformer encoders.

211 **Extraction of speech features.** To explore whether speech models incorporate linguistic information,
212 we extract representations beyond 1.49 secs, i.e., we considered context window of 16 secs with
213 stride of 100 msec and considered the last token as the representative for each context window. The
214 pretrained speech-based models output token representations at different layers. Both Wav2Vec2.0
215 and AST models provide 768 dimensional representations and all of them are 12-layer Transformer
216 encoders. Finally, we align these representations with the fMRI data acquisition rate by downsampling
217 the stimulus features with a 3-lobed Lanczos filter, thus producing chunk-embeddings for each TR.

218 5 Experimental Setup

219 **Encoding Model.** We train bootstrap ridge regression based voxel-wise encoding models (Deniz
220 et al., 2019) to predict the fMRI brain activity associated with the stimulus representations obtained
221 from the individual modalities (speech and video) and multi-modal embeddings from cross-modal and
222 jointly pretrained multi-modal models. For each subject, we account for the delay in the hemodynamic
223 response by modeling hemodynamic response function using a finite response filter (FIR) per voxel
224 with 5 temporal delays (TRs) corresponding to ~ 7.5 seconds (Huth et al., 2022). Formally, at each
225 time step t , we encode the stimuli as $X_t \in \mathbb{R}^D$ and brain region voxels $Y_t \in \mathbb{R}^V$, where D denotes
226 the dimension of the concatenation of delayed 5 TRs, and V denotes the number of voxels. Overall,
227 with N such TRs, we obtain N training examples.

228 **Train-test Setup.** We build encoding models in three settings: (i) We used all data samples from
229 10 training sessions of the *The Bourne supremacy* movie for training and tested generalization on
230 samples from the test sessions (5 sessions) of the *Life* movie. (ii) We used data from 17 training
231 sessions of the *The wolf of wall street* movie for training, with the *Life* movie used for testing. (iii)
232 We combined data from the *The Bourne supremacy* and *The wolf of wall street* movies for training,
233 and tested on the *Life* movie.

234 **Removal of a single modality features from multi-modal representations.** To remove features
235 for a particular modality m from multi-modal model representations, we rely on a simple method
236 proposed previously by Toneva et al. (2022) and Oota et al. (2023b), in which the linear contribution

237 of the features to the multi-modal model activations is removed via ridge regression. Specifically, for
238 this ridge regression the feature vector corresponding to modality m is considered as input and the
239 multi-modal representations are the target. We compute the residuals by subtracting the predicted
240 multi-modal feature representations from the actual multi-modal features resulting in the (linear)
241 removal of feature vector for modality m from the pretrained multi-modal embeddings. Because
242 the brain prediction method is also a linear function, this linear removal limits the contribution of
243 features for modality m to the eventual brain alignment. See Fig. 1(B).

244 **Evaluation Metrics.** We evaluate our models using Pearson Correlation (PC) which is a standard
245 metric for evaluating brain alignment (Jain & Huth, 2018; Schrimpf et al., 2021; Goldstein et al.,
246 2022). Let TR be the number of time repetitions in the test set. Let $Y = \{Y_i\}_{i=1}^{TR}$ and $\hat{Y} = \{\hat{Y}_i\}_{i=1}^{TR}$
247 denote the actual and predicted value vectors for a single voxel. Thus, Y and $\hat{Y} \in \mathbb{R}^{TR}$. We use
248 Pearson Correlation (PC) which is computed as $\text{corr}(Y, \hat{Y})$ where corr is the correlation function.

249 The final measure of a model’s performance is obtained by calculating Pearson’s correlation between
250 the model’s predictions and neural recordings. This correlation is then divided by the estimated
251 cross-subject prediction accuracy and averaged across voxels, regions, and participants, resulting in
252 a standardized measure of performance referred to as normalized brain alignment. For calculating
253 normalized alignment, we select the voxels whose cross-subject prediction accuracy is ≥ 0.05 .

254 **Implementation Details for Reproducibility.** All experiments were conducted on a machine with
255 1 NVIDIA GeForce-GTX GPU with 16GB GPU RAM. We used bootstrap ridge-regression with
256 the following parameters: MSE loss function; L2-decay (λ) varied from 10^1 to 10^3 ; the best λ was
257 chosen by tuning on validation data that comprised a randomly chosen 10% subset from the train set
258 used only for hyper-parameter tuning.

259 **Statistical Significance.** To determine if normalized predictivity scores significantly higher than
260 chance, we run a permutation test using blocks of 10 contiguous fMRI TRs (considering the slowness
261 of hemodynamic response) rather than individual TRs. By permuting predictions 5000 times, we
262 create an empirical distribution for chance performance, from which we estimate the p-value of
263 the actual performance. To estimate the statistical significance of performance differences, such
264 as between the model’s predictions and chance or residual predictions and chance, we utilized the
265 Wilcoxon signed-rank test (Conover, 1999), applying it to the mean normalized predictivity for the
266 participants. In all cases, we denote significant differences ($p \leq 0.05$) with a * or \wedge .

267 6 Results

268 6.1 How effective are multi-modal representations obtained from multi-modal models?

269 In Fig. 2, we present the average normalized brain alignment scores for both multi-modal and
270 individual modality features. Specifically, we show the normalized brain alignment for cross-modality
271 (ImageBind), jointly pretrained multi-modal (TVLT), and the average from individual video and
272 speech models. The results are shown for whole brain, and also for average across language and
273 visual ROIs. Results for individual ROIs are in Fig. 3.

274 **Baseline comparison.** To compare the brain predictivity of multi-modal and unimodal models against
275 baseline performance, we employ randomly generated vector embeddings to predict brain activity as
276 baseline. We observe that the brain alignment from a random vector is significantly lower than that
277 of both multi-modal and unimodal models across the whole brain and language-visual processing
278 regions. This shows that the representations from these multi-modal models are significant enough
279 for learning non-trivial alignment with the fMRI recordings of multi-modal stimuli.

280 **Cross-modal vs. Jointly pretrained multi-modal models vs. Unimodal Models.** Fig. 2(left)
281 displays results for whole brain analysis, where the IB Concat bar plot corresponds to results for
282 representations from a cross-modal model, while TVLT Joint bar plot corresponds to results for
283 representations from a jointly pretrained multi-modal model. From Fig. 2(left), we make the following
284 observations: (i) At the whole brain level, the Wilcoxon signed-rank test shows that the differences in
285 embeddings from the IB Concat and TVLT models are not statistically significant. (ii) The multi-
286 modal embeddings show improved brain alignment compared to unimodal models. Specifically,
287 cross-modal embeddings are significantly better than both unimodal video and speech models, while
288 jointly pretrained embeddings are significantly better than speech models. This implies that cross-

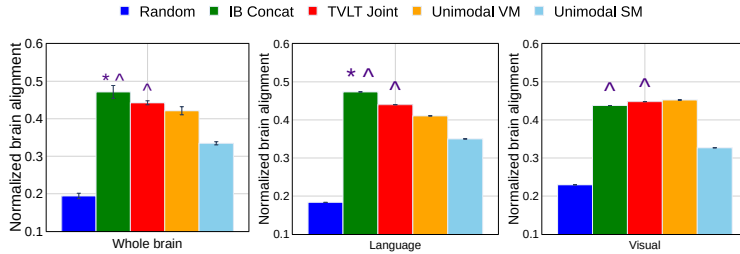


Figure 2: Average normalized brain alignment for both multi-modal and individual modality features across whole brain, language, and visual regions. Error bars indicate the standard error of the mean across participants. * indicates cases where multi-modal embeddings are significantly better than unimodal video models (VM), i.e., $p \leq 0.05$. ^, indicates cases where multi-modal embeddings are significantly better than unimodal speech models (SM), i.e., $p \leq 0.05$.

289 modal embeddings contain additional information beyond the two modalities, while embeddings
 290 from a jointly pretrained model do not provide extra information beyond unimodal visual information
 291 but do contain additional information beyond unimodal speech.

292 We also present average results across language and visual regions in Figs. 2 (middle), and 2(right),
 293 respectively. The Wilcoxon signed-rank test shows that the differences in embeddings from the IB
 294 Concat and TVLT models are not statistically significant when averaged across language and visual
 295 regions. Similar to whole brain performance, in the language regions, cross-modal embeddings
 296 are significantly better than both unimodal video and speech models, while jointly pretrained em-
 297 beddings are significantly better than unimodal speech models. In contrast, for the visual regions,
 298 the normalized brain alignment of cross-modal and jointly pretrained embeddings is similar to the
 299 performance of unimodal video models. This implies that when we average across visual regions,
 300 there is no additional information beyond unimodal video features. However, when compared to
 301 unimodal speech features, both multi-modal embeddings show significant improvement.

302 Since we didn't observe any significant difference at the whole brain level and when averaged across
 303 language and visual regions, between cross-modal and jointly pretrained multi-modal models, we
 304 attempt to seek if there any any differences when we pay a closer look at the individual ROIs. We
 305 present results for language and visual regions such as Angular gyrus (AG), the posterior temporal
 306 lobe (PTL), and the inferior frontal gyrus (IFG) in Fig. 3. Additionally, we cover visual regions
 307 like early visual cortex (EVC), scene visual areas (SV) and middle temporal gyrus (MT), as well
 308 as early auditory cortex (AC). In this figure, we also report the average normalized brain alignment
 309 of each modality obtained from multi-modal models. Unlike the whole brain analysis, we observe
 310 some differences between cross-modal and jointly pretrained models in several language and visual
 311 ROIs. Results for other ROIs are in Fig. 7 in Appendix. Our observations are as follows: (i) Cross-
 312 modal IB Concat embeddings are significantly better than TVLT Joint embeddings in semantic
 313 regions such as AG and PCC, as well as the multi-modal processing region MT. (ii) Conversely,
 314 TVLT Joint embeddings are significantly better than IB Concat embeddings in dmPFC regions.
 315 While considering both joint and each modality embeddings from multi-modal models, we make the
 316 following observations from Fig. 3: (1) Cross-modal IB video embeddings exhibit improved brain
 317 alignment compared to unimodal video in the AG and MT regions with the exceptions of the PTL
 318 and AC regions. But this is not the case for IB audio vs unimodal audio. This suggests that video
 319 modality information is more relevant and beneficial in the brain for IB Concat embeddings from
 320 cross-modality models. (2) TVLT video embeddings show improved brain alignment in the AG, PTL,
 321 PCC, dmPFC and EVC regions, with other regions displaying similar normalized brain alignment
 322 unimodal video embeddings. (3) Consistent with the cross-modality models, in jointly pretrained
 323 TVLT models, TVLT video embeddings significantly outperform TVLT audio embeddings, except in
 324 PTL region. These observations indicate that video information is advantageous for both cross-modal
 325 and jointly pretrained models, whereas audio embeddings mainly benefit the PTL region.

326 6.2 Which brain regions process uni- and multi-modal information?

327 From Fig. 3, we observe that multi-modal video embeddings exhibit improved brain alignment not
 328 only in the whole brain but also in various language, visual and multi-modal regions. For instance,
 329 the cross-modal IB Concat embeddings demonstrate superior brain alignment compared to unimodal
 330 video-based models in areas such as the AG, PTL, IFG, and PCC. Moreover, TVLT-joint embeddings

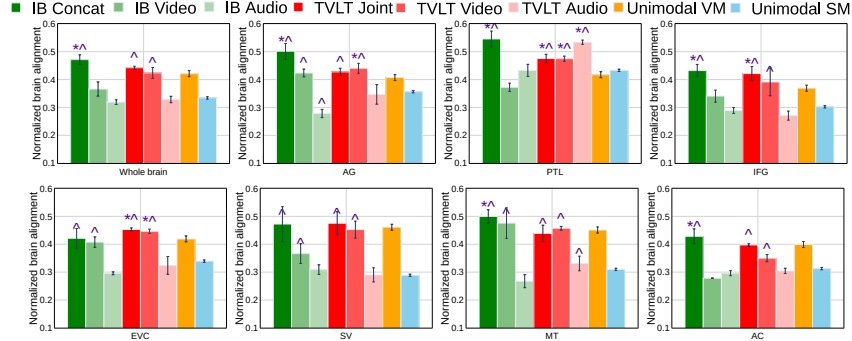


Figure 3: Average normalized brain alignment for video and audio modalities from multi-modal and individual modality features across whole brain and several ROIs of language (AG, PTL and IFG), visual (EVC, SV and MT) and auditory cortex (AC). Error bars indicate the standard error of the mean across participants. * indicates cases where multi-modal embeddings are significantly better than unimodal video models (VM), i.e., $p \leq 0.05$. ^ indicates cases where multi-modal embeddings are significantly better than unimodal speech models (SM), i.e., $p \leq 0.05$.

331 show notable enhancements in the AG, PTL, IFG, PCC, dmPFC and EVC regions. In contrast,
 332 compared to unimodal speech-based models, all multi-modal embeddings display significantly better
 333 brain alignment, except the OV (object visual processing) region. Overall, this observation suggests
 334 that integrating multiple modalities leads to transferring information from one modality to another,
 335 resulting in improved brain predictability. Based on these, it can be inferred that these multi-modal
 336 models can indeed learn multi-modal linkages that are relevant to the brain.

337 When subjects engage with multi-modality stimuli, we observe that multi-modal embeddings show
 338 improvements in semantic regions such as the AG, PCC and dmPFC, and syntactic regions such as
 339 the PTL and IFG. Overall, we find that multi-modal information is processed in only a few regions.
 340 Furthermore, several regions, including the SV (scene visual area), EVC (early visual cortex), ATL
 341 (anterior temporal lobe), IFGOrb, MFG, and dmPFC, exhibit similar brain alignment with both
 342 unimodal and multi-modal embeddings.

343 6.3 How is the brain alignment of multi-modal features affected by the elimination of a 344 particular modality?

345 To understand the contribution of each modality to the multi-modal brain alignment for multi-modal
 346 naturalistic stimulus, we perform residual analyses by removing the unimodality features from
 347 multi-modal joint representations as well as multi-modal video or audio representations from joint
 348 representations and measure the differences in brain alignment before and after removal modality-
 349 specific features. Fig. 4 displays the normalized brain alignment for language (AG) and visual regions
 350 (MT). We note a decrease in brain alignment for both the AG and MT regions following the removal
 351 of video embeddings from cross-modality models, whereas the removal of audio embeddings does
 352 not affect the brain alignment. On the other hand, for jointly pretrained models, removal of both
 353 video and audio embeddings partially impacts the brain alignment. We observe similar findings for
 354 language ROIs such as PTL, MFG, ATL, PCC and visual regions EVC, OV and FV, as shown in
 355 Figs. 9 and 10 in Appendix. These results suggest that there is additional information beyond the
 356 unimodal embeddings considered in this study that is processed in the visual and language regions.

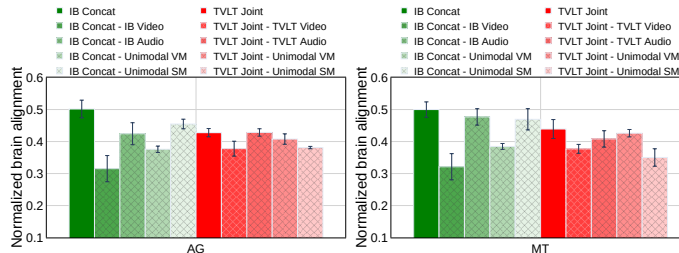


Figure 4: Residual analysis: Average normalized brain alignment was computed across participants before and after removal of video and audio embeddings from both jointly pretrained and cross-modality models. Error bars indicate the standard error of the mean across participants.

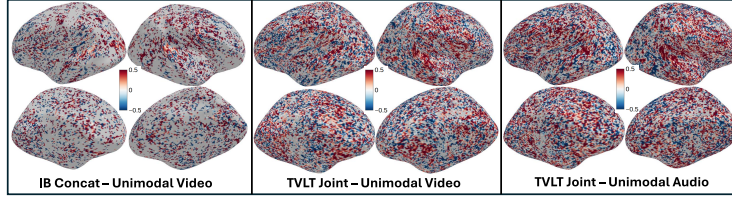


Figure 5: Percentage decrease of brain alignment after removal of (left) Unimodal VM embeddings from IB-Concat (middle) Unimodal VM embeddings from jointly pretrained TVLT, and (right) Unimodal SM embeddings from TVLT Joint. Colorbar indicates the percentage of decrease where red denotes higher and white denotes zero.

357 **Qualitative analysis.** We compute the percentage decrease in alignment for each voxel following the
 358 removal of unimodal video embeddings from the IB Concat (cross-modality) and the TVLT Joint
 359 (jointly pretrained model), with projections onto the brain surface averaged across participants, as
 360 depicted in Fig. 5. The colorbar shows the percentage decrease in brain alignment, where red voxels
 361 indicate a higher percentage decrease and white voxels indicate areas where unimodal video features
 362 do not contribute any shared information within the multi-modal context. We observe that removal of
 363 unimodal video features leads to a significant drop (40-50%) in performance in the visual regions for
 364 IB Concat, and in language regions (PTL & MFG) for TVLT Joint.

365 7 Discussion

366 Using multi-modal model representations, including both cross-modal and jointly pretrained types,
 367 we evaluated how these representations can predict fMRI brain activity when participants are
 368 engaged in multi-modal naturalistic stimuli. Further, we compared both multi-modal and unimodal
 369 representations and observed their alignment with both unimodal and multi-modal brain regions.
 370 This is achieved by removing information related to unimodal stimulus features (audio and video)
 371 and observing how this perturbation affects the alignment with fMRI brain recordings acquired while
 372 participants are engaged in watching multi-modal naturalistic movies.

373 Our analysis of multi-modal brain alignment yields several important conclusions: (1) The improved
 374 brain alignment of the multi-modal models over unimodal models, across several language, visual, and
 375 auditory regions is only partially attributable to the video and audio stimulus features presented to the
 376 model. A deeper understanding of these models is required to shed light on the underlying information
 377 processing of both unimodal and multi-modal information. (2) Cross-modal representations have
 378 significantly improved brain alignment in language regions such as AG, PCC and PTL. This variance
 379 can be partially attributed to the removal of video features alone, rather than auditory features. (3)
 380 Video embeddings from multi-modal models exhibit higher brain alignment than audio embeddings,
 381 except in the PTL and AC regions. This suggests that audio-based models may encode weaker brain-
 382 relevant semantics. (4) Both cross-modal and jointly pretrained models demonstrate significantly
 383 improved brain alignment with language regions (AG, PCC, PTL and IFG) compared to visual regions
 384 when analyzed against unimodal video data. In contrast, when compared to unimodal audio-based
 385 models, all multi-modal embeddings display significantly better brain alignment, with the exception
 386 of the OV region. This underscores the capability of multi-modal models to capture additional
 387 information—either through knowledge transfer or integration between modalities—crucial for
 388 multi-modal brain alignment.

389 **Limitations.** The low alignment scores clearly show that despite the increasing popularity of multi-
 390 modal models in tackling complex tasks such as visual question answering, we are still far from
 391 developing a model that fully encapsulates the complete information processing steps involved
 392 in handling multi-modal naturalistic information in the brain. In the future, by fine-tuning these
 393 multi-modal models on specific tasks such as generating captions for videos, we can better leverage
 394 their alignment strengths. This approach will allow us to explore task-level brain alignment of three
 395 modalities—video, audio, and text—more effectively. Further, multi-modal large language models
 396 (MLLMs) (Zhang et al., 2023; Ataallah et al., 2024; Wu et al., 2023) that align visual features from
 397 video frames into the LLM embedding space via a trainable linear projection layer, offer promise for
 398 enhanced multi-modal capabilities. We would further extend this work by comparing the region-wise
 399 brain alignment performance of these multi-modal LLM models with existing approaches.

400 References

- 401 Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid.
402 Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on*
403 *computer vision*, pp. 6836–6846, 2021.
- 404 Kirolos Ataallah, Xiaoqian Shen, Eslam Abdelrahman, Essam Sleiman, Deyao Zhu, Jian Ding, and
405 Mohamed Elhoseiny. Minigt4-video: Advancing multimodal llms for video understanding with
406 interleaved visual-textual tokens. *arXiv preprint arXiv:2404.03413*, 2024.
- 407 Khai Loong Aw and Mariya Toneva. Training language models to summarize narratives improves
408 brain alignment. In *The Eleventh International Conference on Learning Representations*, 2023.
- 409 Alan Baade, Puyuan Peng, and David Harwath. Mae-ast: Masked autoencoding audio spectrogram
410 transformer. *Interspeech 2022*, 2022.
- 411 Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework
412 for self-supervised learning of speech representations. *Advances in neural information processing*
413 *systems*, 2020.
- 414 Cordell M Baker, Joshua D Burks, Robert G Briggs, Andrew K Conner, Chad A Glenn, Kathleen N
415 Taylor, Goksel Sali, Tressie M McCoy, James D Battiste, Daniel L O’Donoghue, et al. A connec-
416 tomic atlas of the human cerebrum—chapter 7: the lateral parietal lobe. *Operative Neurosurgery*,
417 15(suppl_1):S295–S349, 2018.
- 418 Stefania Bracci and Hans P Op de Beeck. Understanding human object vision: a picture is worth a
419 thousand representations. *Annual review of psychology*, 74:113–135, 2023.
- 420 Charlotte Caucheteux and Jean-Rémi King. Language processing in brains and deep neural networks:
421 computational convergence and its limits. *BioRxiv*, 2020.
- 422 William Jay Conover. *Practical nonparametric statistics*, volume 350. john wiley & sons, 1999.
- 423 Fatma Deniz, Anwar O Nunez-Elizalde, Alexander G Huth, and Jack L Gallant. The representation
424 of semantic information across human cerebral cortex during listening versus reading is invariant
425 to stimulus modality. *Journal of Neuroscience*, 2019.
- 426 Rutvik Desai, Usha Tadimeti, and Nicholas Riccardi. Proper and common names in the semantic
427 system, 2022.
- 428 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
429 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*
430 *the North American Chapter of the Association for Computational Linguistics: Human Language*
431 *Technologies, Volume 1 (Long and Short Papers)*, 2019.
- 432 Adrien Doerig, Tim C Kietzmann, Emily Allen, Yihan Wu, Thomas Naselaris, Kendrick Kay,
433 and Ian Charest. Semantic scene descriptions as an objective of human vision. *arXiv preprint*
434 *arXiv:2209.11737*, 2022.
- 435 Dota Tianai Dong and Mariya Toneva. Vision-language integration in multimodal video transformers
436 (partially) aligns with the brain. *arXiv preprint arXiv:2311.07766*, 2023.
- 437 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
438 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image
439 is worth 16x16 words: Transformers for image recognition at scale. In *International Conference*
440 *on Learning Representations*, 2020.
- 441 Michael Eickenberg, Alexandre Gramfort, Gaël Varoquaux, and Bertrand Thirion. Seeing it all:
442 Convolutional network layers map the function of the human visual system. *NeuroImage*, 152:
443 184–194, 2017.
- 444 Isabel Gauthier, Thomas W James, Kim M Curby, and Michael J Tarr. The influence of conceptual
445 knowledge on visual discrimination. *Cognitive Neuropsychology*, 20(3-6):507–523, 2003.

- 446 Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand
447 Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the*
448 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15180–15190, 2023.
- 449 Matthew F Glasser, Timothy S Coalson, Emma C Robinson, Carl D Hacker, John Harwell, Essa
450 Yacoub, Kamil Ugurbil, Jesper Andersson, Christian F Beckmann, Mark Jenkinson, et al. A
451 multi-modal parcellation of human cerebral cortex. *Nature*, 536(7615):171–178, 2016.
- 452 Ariel Goldstein, Zaid Zada, Eliav Buchnik, Mariano Schain, Amy Price, Bobbi Aubrey, Samuel A
453 Nastase, Amir Feder, Dotan Emanuel, Alon Cohen, et al. Shared computational principles for
454 language processing in humans and deep language models. *Nature neuroscience*, 25(3):369–380,
455 2022.
- 456 Alexander G Huth, Shinji Nishimoto, An T Vu, and T Dupre La Tour. Gallant lab natural short clips
457 3t fmri data. *50 GiB*, 2022.
- 458 Shailee Jain and Alexander G Huth. Incorporating context into language encoding models for fmri.
459 In *NIPS*, pp. 6629–6638, 2018.
- 460 Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple
461 and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- 462 Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: pretraining task-agnostic visiolinguistic
463 representations for vision-and-language tasks. In *Proceedings of the 33rd International Conference*
464 *on Neural Information Processing Systems*, pp. 13–23, 2019.
- 465 Juliette Millet, Charlotte Caucheteux, Yves Boubenec, Alexandre Gramfort, Ewan Dunbar, Christophe
466 Pallier, Jean-Remi King, et al. Toward a realistic model of speech processing in the brain with
467 self-supervised learning. *Advances in Neural Information Processing Systems*, 35:33428–33443,
468 2022.
- 469 Camille K Milton, Vukshitha Dhanaraj, Isabella M Young, Hugh M Taylor, Peter J Nicholas, Robert G
470 Briggs, Michael Y Bai, Rannulu D Fonseka, Jorge Hormovas, Yueh-Hsin Lin, et al. Parcellation-
471 based anatomic model of the semantic network. *Brain and behavior*, 11(4):e02065, 2021.
- 472 Yuko Nakagi, Takuya Matsuyama, Naoko Koide-Majima, Hiroto Yamaguchi, Rieko Kubo, Shinji
473 Nishimoto, and Yu Takagi. The brain tells a story: Unveiling distinct representations of semantic
474 content in speech, objects, and stories in the human brain with large language models. *bioRxiv*, pp.
475 2024–02, 2024.
- 476 Subba Reddy Oota, Jashn Arora, Vijay Rowtula, Manish Gupta, and Raju S Bapi. Visio-linguistic
477 brain encoding. In *COLING*, pp. 116–133, 2022.
- 478 Subba Reddy Oota, Manish Gupta, Raju S Bapi, Gael Jobard, Frédéric Alexandre, and Xavier Hinaut.
479 Deep neural networks and brain alignment: Brain encoding and decoding (survey). *arXiv preprint*
480 *arXiv:2307.10246*, 2023a.
- 481 Subba Reddy Oota, Manish Gupta, and Mariya Toneva. Joint processing of linguistic properties in
482 brains and language models. *NeurIPS*, 2023b.
- 483 Subba Reddy Oota, Agarwal Veeral, Marreddy Mounika, Gupta Manish, and Raju Surampudi Bapi.
484 Speech taskonomy: Which speech tasks are the most predictive of fmri brain activity? In *24th*
485 *INTERSPEECH Conference*, 2023c.
- 486 Sara F Popham, Alexander G Huth, Natalia Y Bilenko, Fatma Deniz, James S Gao, Anwar O Nunez-
487 Elizalde, and Jack L Gallant. Visual and linguistic semantic representations are aligned at the
488 border of human visual cortex. *Nature neuroscience*, 24(11):1628–1636, 2021.
- 489 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
490 models are unsupervised multitask learners. *OpenAI*, 2019.
- 491 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
492 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
493 models from natural language supervision. *Image*, 2:T2, 2021.

- 494 Martin Schrimpf, Jonas Kubilius, Ha Hong, Najib J Majaj, Rishi Rajalingham, Elias B Issa, Kohitij
495 Kar, Pouya Bashivan, Jonathan Prescott-Roy, Franziska Geiger, et al. Brain-score: Which artificial
496 neural network for object recognition is most brain-like? *BioRxiv*, pp. 407007, 2018.
- 497 Martin Schrimpf, Idan Asher Blank, Greta Tuckute, Carina Kauf, Eghbal A Hosseini, Nancy Kan-
498 wisher, Joshua B Tenenbaum, and Evelina Fedorenko. The neural architecture of language:
499 Integrative modeling converges on predictive processing. *Proceedings of the National Academy of*
500 *Sciences*, 2021.
- 501 Marie St-Laurent, Basile Pinsard, Oliver Contier, Katja Seeliger, Valentina Borghesani, Julie Boyle,
502 Pierre Bellec, and Martin Hebart. cneuromod-things: a large-scale fmri dataset for task-and
503 data-driven assessment of object representation and visual memory recognition in the human brain.
504 *Journal of Vision*, 23(9):5424–5424, 2023.
- 505 Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transform-
506 ers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*
507 *and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*,
508 pp. 5100–5111, 2019.
- 509 Jerry Tang, Meng Du, Vy Vo, Vasudev Lal, and Alexander Huth. Brain encoding models based on
510 multimodal transformers can transfer across language and vision. *Advances in Neural Information*
511 *Processing Systems*, 36, 2024.
- 512 Zineng Tang, Jaemin Cho, Yixin Nie, and Mohit Bansal. Tvlt: Textless vision-language transformer.
513 *Advances in Neural Information Processing Systems*, 35:9617–9632, 2022.
- 514 Mariya Toneva and Leila Wehbe. Interpreting and improving natural-language processing (in
515 machines) with natural language-processing (in the brain). *Advances in Neural Information*
516 *Processing Systems*, 32, 2019.
- 517 Mariya Toneva, Tom M Mitchell, and Leila Wehbe. Combining computational controls with natural
518 text reveals aspects of meaning composition. *Nature Computational Science*, 2(11):745–757, 2022.
- 519 Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-
520 efficient learners for self-supervised video pre-training. *Advances in neural information processing*
521 *systems*, 35:10078–10093, 2022.
- 522 Greta Tuckute, Jenelle Feather, Dana Boebinger, and Josh H McDermott. Many but not all deep
523 neural network audio models capture brain responses and exhibit correspondence between model
524 stages and brain regions. *Plos Biology*, 21(12):e3002366, 2023.
- 525 Aditya Vaidya, Shailee Jain, and Alexander Huth. Self-supervised models of audio effectively
526 explain human cortical responses to speech. In *International Conference on Machine Learning*, pp.
527 21927–21944. PMLR, 2022.
- 528 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
529 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
530 *systems*, 30, 2017.
- 531 Aria Wang, Michael Tarr, and Leila Wehbe. Neural taskonomy: Inferring the similarity of task-derived
532 representations from brain activity. *NeurIPS*, 32:15501–15511, 2019.
- 533 Aria Y Wang, Kendrick Kay, Thomas Naselaris, Michael J Tarr, and Leila Wehbe. Natural language
534 supervision with a large and diverse dataset builds better models of human high-level visual cortex.
535 *BioRxiv*, pp. 2022–09, 2022.
- 536 Leila Wehbe, Brian Murphy, Partha Talukdar, Alona Fyshe, Aaditya Ramdas, and Tom Mitchell.
537 Simultaneously uncovering the patterns of brain regions involved in different story reading subpro-
538 cesses. *PloS one*, 11, 2014.
- 539 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
540 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art
541 natural language processing. In *Proceedings of the 2020 conference on empirical methods in*
542 *natural language processing: system demonstrations*, pp. 38–45, 2020.

- 543 Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal
544 llm. *arXiv preprint arXiv:2309.05519*, 2023.
- 545 Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J
546 DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual
547 cortex. *PNAS*, 111(23):8619–8624, 2014.
- 548 Rowan Zellers, Jiasen Lu, Ximing Lu, Youngjae Yu, Yanpeng Zhao, Mohammadreza Salehi, Aditya
549 Kusupati, Jack Hessel, Ali Farhadi, and Yejin Choi. Merlot reserve: Neural script knowledge
550 through vision and language and sound. In *Proceedings of the IEEE/CVF Conference on Computer
551 Vision and Pattern Recognition*, pp. 16375–16387, 2022.
- 552 Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language
553 model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023.

554 **NeurIPS Paper Checklist**

555 **1. Claims**

556 Question: Do the main claims made in the abstract and introduction accurately reflect the
557 paper's contributions and scope?

558 Answer: [Yes]

559 Justification: We have ensured that the main claims made in the abstract and introduction
560 are directly correlating to the research findings and the methods we have employed.

561 Guidelines:

- 562 • The answer NA means that the abstract and introduction do not include the claims
563 made in the paper.
- 564 • The abstract and/or introduction should clearly state the claims made, including the
565 contributions made in the paper and important assumptions and limitations. A No or
566 NA answer to this question will not be perceived well by the reviewers.
- 567 • The claims made should match theoretical and experimental results, and reflect how
568 much the results can be expected to generalize to other settings.
- 569 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
570 are not attained by the paper.

571 **2. Limitations**

572 Question: Does the paper discuss the limitations of the work performed by the authors?

573 Answer: [Yes]

574 Justification: The paper discusses the main limitations of the work performed by the authors
575 in the discussion section.

576 Guidelines:

- 577 • The answer NA means that the paper has no limitation while the answer No means that
578 the paper has limitations, but those are not discussed in the paper.
- 579 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 580 • The paper should point out any strong assumptions and how robust the results are to
581 violations of these assumptions (e.g., independence assumptions, noiseless settings,
582 model well-specification, asymptotic approximations only holding locally). The authors
583 should reflect on how these assumptions might be violated in practice and what the
584 implications would be.
- 585 • The authors should reflect on the scope of the claims made, e.g., if the approach was
586 only tested on a few datasets or with a few runs. In general, empirical results often
587 depend on implicit assumptions, which should be articulated.
- 588 • The authors should reflect on the factors that influence the performance of the approach.
589 For example, a facial recognition algorithm may perform poorly when image resolution
590 is low or images are taken in low lighting. Or a speech-to-text system might not be
591 used reliably to provide closed captions for online lectures because it fails to handle
592 technical jargon.
- 593 • The authors should discuss the computational efficiency of the proposed algorithms
594 and how they scale with dataset size.
- 595 • If applicable, the authors should discuss possible limitations of their approach to
596 address problems of privacy and fairness.
- 597 • While the authors might fear that complete honesty about limitations might be used by
598 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
599 limitations that aren't acknowledged in the paper. The authors should use their best
600 judgment and recognize that individual actions in favor of transparency play an impor-
601 tant role in developing norms that preserve the integrity of the community. Reviewers
602 will be specifically instructed to not penalize honesty concerning limitations.

603 **3. Theory Assumptions and Proofs**

604 Question: For each theoretical result, does the paper provide the full set of assumptions and
605 a complete (and correct) proof?

606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658

Answer: [NA]

Justification: Our paper does not require any explicit theorems and proofs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper has delineated all the information related to the experimental setup in the experimental setup section.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

659 Question: Does the paper provide open access to the data and code, with sufficient instruc-
660 tions to faithfully reproduce the main experimental results, as described in supplemental
661 material?

662 Answer: [NA]

663 Justification: We will release the code upon acceptance. The dataset is publicly available
664 through a licence.

665 Guidelines:

- 666 • The answer NA means that paper does not include experiments requiring code.
- 667 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
668 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 669 • While we encourage the release of code and data, we understand that this might not be
670 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
671 including code, unless this is central to the contribution (e.g., for a new open-source
672 benchmark).
- 673 • The instructions should contain the exact command and environment needed to run to
674 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
675 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 676 • The authors should provide instructions on data access and preparation, including how
677 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 678 • The authors should provide scripts to reproduce all experimental results for the new
679 proposed method and baselines. If only a subset of experiments are reproducible, they
680 should state which ones are omitted from the script and why.
- 681 • At submission time, to preserve anonymity, the authors should release anonymized
682 versions (if applicable).
- 683 • Providing as much information as possible in supplemental material (appended to the
684 paper) is recommended, but including URLs to data and code is permitted.

685 6. Experimental Setting/Details

686 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
687 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
688 results?

689 Answer: [Yes]

690 Justification: We provide all the training and test details in the experimental setup.

691 Guidelines:

- 692 • The answer NA means that the paper does not include experiments.
- 693 • The experimental setting should be presented in the core of the paper to a level of detail
694 that is necessary to appreciate the results and make sense of them.
- 695 • The full details can be provided either with the code, in appendix, or as supplemental
696 material.

697 7. Experiment Statistical Significance

698 Question: Does the paper report error bars suitably and correctly defined or other appropriate
699 information about the statistical significance of the experiments?

700 Answer: [Yes]

701 Justification: We conducted our experiments multiple times across 6 participants and took
702 the average results. We also include error bars in the plots.

703 Guidelines:

- 704 • The answer NA means that the paper does not include experiments.
- 705 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
706 dence intervals, or statistical significance tests, at least for the experiments that support
707 the main claims of the paper.
- 708 • The factors of variability that the error bars are capturing should be clearly stated (for
709 example, train/test split, initialization, random drawing of some parameter, or overall
710 run with given experimental conditions).

- 711 • The method for calculating the error bars should be explained (closed form formula,
712 call to a library function, bootstrap, etc.)
- 713 • The assumptions made should be given (e.g., Normally distributed errors).
- 714 • It should be clear whether the error bar is the standard deviation or the standard error
715 of the mean.
- 716 • It is OK to report 1-sigma error bars, but one should state it. The authors should
717 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
718 of Normality of errors is not verified.
- 719 • For asymmetric distributions, the authors should be careful not to show in tables or
720 figures symmetric error bars that would yield results that are out of range (e.g. negative
721 error rates).
- 722 • If error bars are reported in tables or plots, The authors should explain in the text how
723 they were calculated and reference the corresponding figures or tables in the text.

724 8. Experiments Compute Resources

725 Question: For each experiment, does the paper provide sufficient information on the com-
726 puter resources (type of compute workers, memory, time of execution) needed to reproduce
727 the experiments?

728 Answer: [Yes]

729 Justification: We have included the specifications of the hardware and software environments
730 to ensure the reproducibility of our results.

731 Guidelines:

- 732 • The answer NA means that the paper does not include experiments.
- 733 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
734 or cloud provider, including relevant memory and storage.
- 735 • The paper should provide the amount of compute required for each of the individual
736 experimental runs as well as estimate the total compute.
- 737 • The paper should disclose whether the full research project required more compute
738 than the experiments reported in the paper (e.g., preliminary or failed experiments that
739 didn't make it into the paper).

740 9. Code Of Ethics

741 Question: Does the research conducted in the paper conform, in every respect, with the
742 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

743 Answer: [Yes]

744 Justification: The research conducted in this paper fully conforms with the NeurIPS Code of
745 Ethics in every respect.

746 Guidelines:

- 747 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 748 • If the authors answer No, they should explain the special circumstances that require a
749 deviation from the Code of Ethics.
- 750 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
751 eration due to laws or regulations in their jurisdiction).

752 10. Broader Impacts

753 Question: Does the paper discuss both potential positive societal impacts and negative
754 societal impacts of the work performed?

755 Answer: [Yes]

756 Justification: The paper explores how the advancements and applications of our findings
757 could benefit society in terms of computational neuroscience research by specifically inves-
758 tigating the effectiveness of the current state of the art multimodal models in encoding brain
759 activity.

760 Guidelines:

- 761 • The answer NA means that there is no societal impact of the work performed.

- 762 • If the authors answer NA or No, they should explain why their work has no societal
763 impact or why the paper does not address societal impact.
- 764 • Examples of negative societal impacts include potential malicious or unintended uses
765 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
766 (e.g., deployment of technologies that could make decisions that unfairly impact specific
767 groups), privacy considerations, and security considerations.
- 768 • The conference expects that many papers will be foundational research and not tied
769 to particular applications, let alone deployments. However, if there is a direct path to
770 any negative applications, the authors should point it out. For example, it is legitimate
771 to point out that an improvement in the quality of generative models could be used to
772 generate deepfakes for disinformation. On the other hand, it is not needed to point out
773 that a generic algorithm for optimizing neural networks could enable people to train
774 models that generate Deepfakes faster.
- 775 • The authors should consider possible harms that could arise when the technology is
776 being used as intended and functioning correctly, harms that could arise when the
777 technology is being used as intended but gives incorrect results, and harms following
778 from (intentional or unintentional) misuse of the technology.
- 779 • If there are negative societal impacts, the authors could also discuss possible mitigation
780 strategies (e.g., gated release of models, providing defenses in addition to attacks,
781 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
782 feedback over time, improving the efficiency and accessibility of ML).

783 11. Safeguards

784 Question: Does the paper describe safeguards that have been put in place for responsible
785 release of data or models that have a high risk for misuse (e.g., pretrained language models,
786 image generators, or scraped datasets)?

787 Answer: [NA]

788 Justification: Our research does not pose any risks for misuse.

789 Guidelines:

- 790 • The answer NA means that the paper poses no such risks.
- 791 • Released models that have a high risk for misuse or dual-use should be released with
792 necessary safeguards to allow for controlled use of the model, for example by requiring
793 that users adhere to usage guidelines or restrictions to access the model or implementing
794 safety filters.
- 795 • Datasets that have been scraped from the Internet could pose safety risks. The authors
796 should describe how they avoided releasing unsafe images.
- 797 • We recognize that providing effective safeguards is challenging, and many papers do
798 not require this, but we encourage authors to take this into account and make a best
799 faith effort.

800 12. Licenses for existing assets

801 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
802 the paper, properly credited and are the license and terms of use explicitly mentioned and
803 properly respected?

804 Answer: [Yes]

805 Justification: We have explicitly cited the datasets, code and models used.

806 Guidelines:

- 807 • The answer NA means that the paper does not use existing assets.
- 808 • The authors should cite the original paper that produced the code package or dataset.
- 809 • The authors should state which version of the asset is used and, if possible, include a
810 URL.
- 811 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 812 • For scraped data from a particular source (e.g., website), the copyright and terms of
813 service of that source should be provided.

- 814
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 815
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 816
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.
- 817
- 818
- 819
- 820
- 821

822 13. **New Assets**

823 Question: Are new assets introduced in the paper well documented and is the documentation
824 provided alongside the assets?

825 Answer: [NA]

826 Justification: We will try to opensource the code and provide complete documentation for
827 our assets upon acceptance.

828 Guidelines:

- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 829
- 830
- 831
- 832
- 833
- 834
- 835
- 836

837 14. **Crowdsourcing and Research with Human Subjects**

838 Question: For crowdsourcing experiments and research with human subjects, does the paper
839 include the full text of instructions given to participants and screenshots, if applicable, as
840 well as details about compensation (if any)?

841 Answer: [NA]

842 Justification: We use publicly available fMRI dataset and do not collect any new data.

843 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 844
- 845
- 846
- 847
- 848
- 849
- 850
- 851

852 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 853 Subjects**

854 Question: Does the paper describe potential risks incurred by study participants, whether
855 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
856 approvals (or an equivalent approval/review based on the requirements of your country or
857 institution) were obtained?

858 Answer: [NA]

859 Justification: We use publicly available fMRI dataset and do not collect any new data.

860 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- 861
- 862
- 863
- 864
- 865

866
867
868
869
870

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

871 A Cross-subject prediction accuracy

872 We estimate cross-subject prediction accuracy in three settings: (i) training with *The Bourne supremacy* and testing with *Life* data, (ii) training with *The wolf of wall street* and testing with *Life* data, and (iii) training with both *The Bourne supremacy* and *The wolf of wall street* and testing with *Life* data. We present the average cross-subject prediction accuracy across voxels for the *Movie10 fMRI* dataset and across the three settings in Fig. 6.

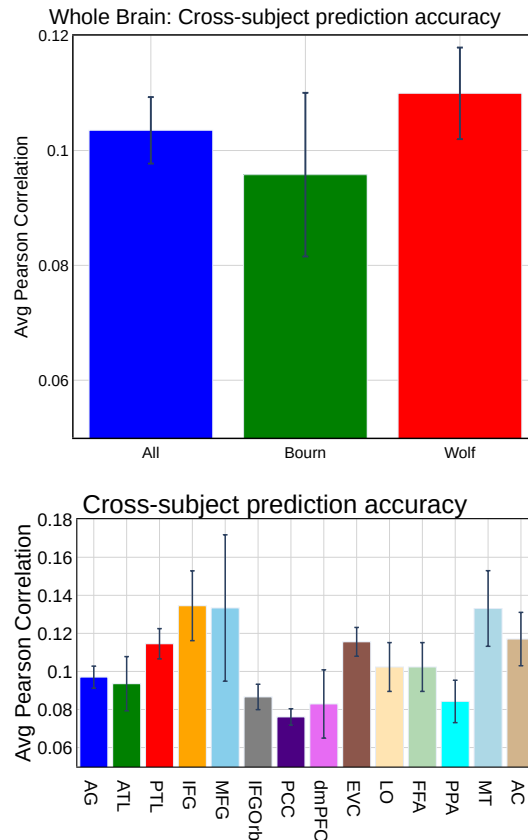


Figure 6: Cross-subject prediction accuracy: (top) across whole brain, (bottom) across language, visual and auditory regions.

877 B Detailed sub-ROIs of language, visual and auditory regions

878 The data covers seven brain regions of interest (ROIs) in the human brain with the following sub-
879 divisions: (i) early visual (EV: V1, V2, V3, V3B, and V4); (ii) object-related areas (LO1 and LO2);
880 (iii) face-related areas (OFA), (iv) scene-related areas (PPA), (v) middle temporal (MT: MT, MST,
881 LO3, FST and V3CD), (vi) late language regions, encompassing broader language regions: angular
882 gyrus (AG: PFm, PGs, PGi, TPOJ2, TPOJ3), lateral temporal cortex (LTC: STSda, STSva, STGa,
883 TE1a, TE2a, TGv, TGd, A5, STSdp, STSvp, PSL, STV, TPOJ1), inferior frontal gyrus (IFG: 44, 45,
884 IFJa, IFSp) and middle frontal gyrus (MFG: 55b) (Baker et al., 2018; Milton et al., 2021; Desai et al.,
885 2022).

886 **C Details of pretrained Transformer models**

Table 1: Pretrained Transformer-based Encoder Models. All models have 12 layers.

Model Name	Pretraining
Cross-modal Pretrained (ImageBind) Jointly Pretrained (TVLT)	Video & Audio Video & Audio
ViT-B VideoMAE ViViT	Image Video Video
Wav2Vec2.0-base AST	Speech Speech

887 **D Effectiveness of multi-modal vs unimodal representations for various brain**
888 **regions**

889 We now present the results for per unimodal video model and per speech model in Fig. 8. Similar to
890 the average results of unimodal video and speech models, we observe that multi-modal models exhibit
891 better normalized brain alignment than individual unimodal video and speech models across language
892 and visual regions. Among unimodal speech models, the AST model shows better normalized brain
893 alignment than the Wav2vec2.0 model. Among unimodal video models, each unimodal video model
894 displays notably consistent performance across regions.

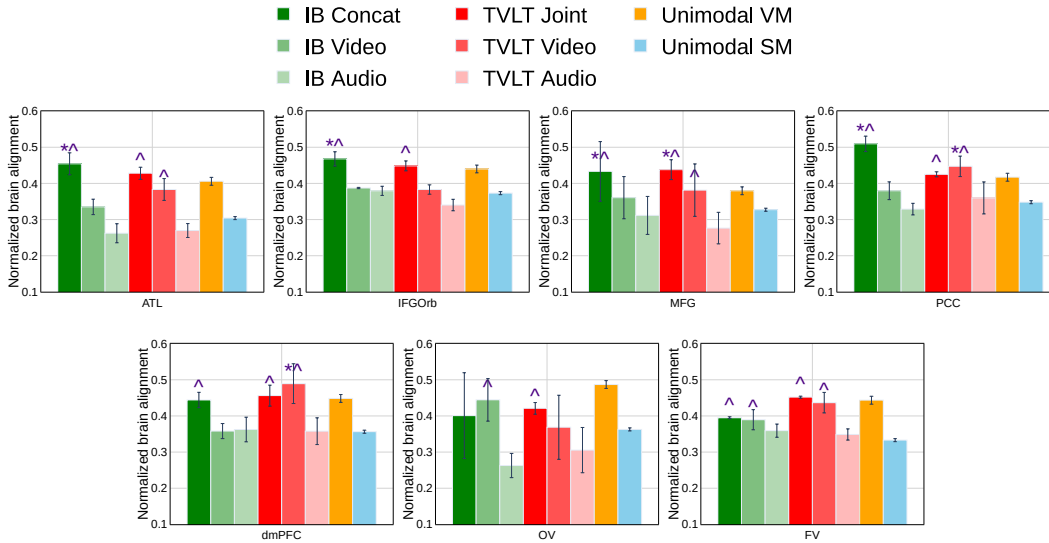


Figure 7: Average normalized brain alignment for per video and audio modalities from multi-modal and individual modality features across whole brain and several ROIs of language (ATL, IFGOrb, MFG, PCC, dmPFC) and visual (OV, FV). Error bars indicate the standard error of the mean across participants.

895 **E How is the brain alignment of multi-modal features affected by the**
896 **elimination of a particular modality?**

897 To understand the contribution of each modality to the multi-modal brain alignment for multi-modal
898 naturalistic stimulus, we perform residual analyses by removing the unimodality features from
899 multi-modal joint representations as well as multi-modal video or audio representations from joint
900 representations and measure the differences in brain alignment before and after removal modality-
901 specific features. Figs. 9 and 10 display the normalized brain alignment for language ROIs such as
902 PTL, MFG, ATL, PCC and visual regions EVC, OV and FV. We note a decrease in brain alignment
903 for these regions following the removal of video embeddings from cross-modality models, whereas
904 the removal of audio embeddings does not affect the brain alignment. On the other hand, for jointly
905 pretrained models, removal of both video and audio embeddings partially impacts the brain alignment.

906 **F Layerwise brain alignment**

907 We now plot the layer-wise normalized brain alignment for the Unimodal models and TVLT joint
908 model, as shown in Fig. 11. Observation from Fig. 11 indicates a consistent drop in performance from
909 early to lower layers, specifically for both TVLT joint and unimodal video models. The key finding
910 here is that our results that TVLT joint embeddings showcase improved brain alignment across all the
911 layers compared to unimodal video and speech embeddings.

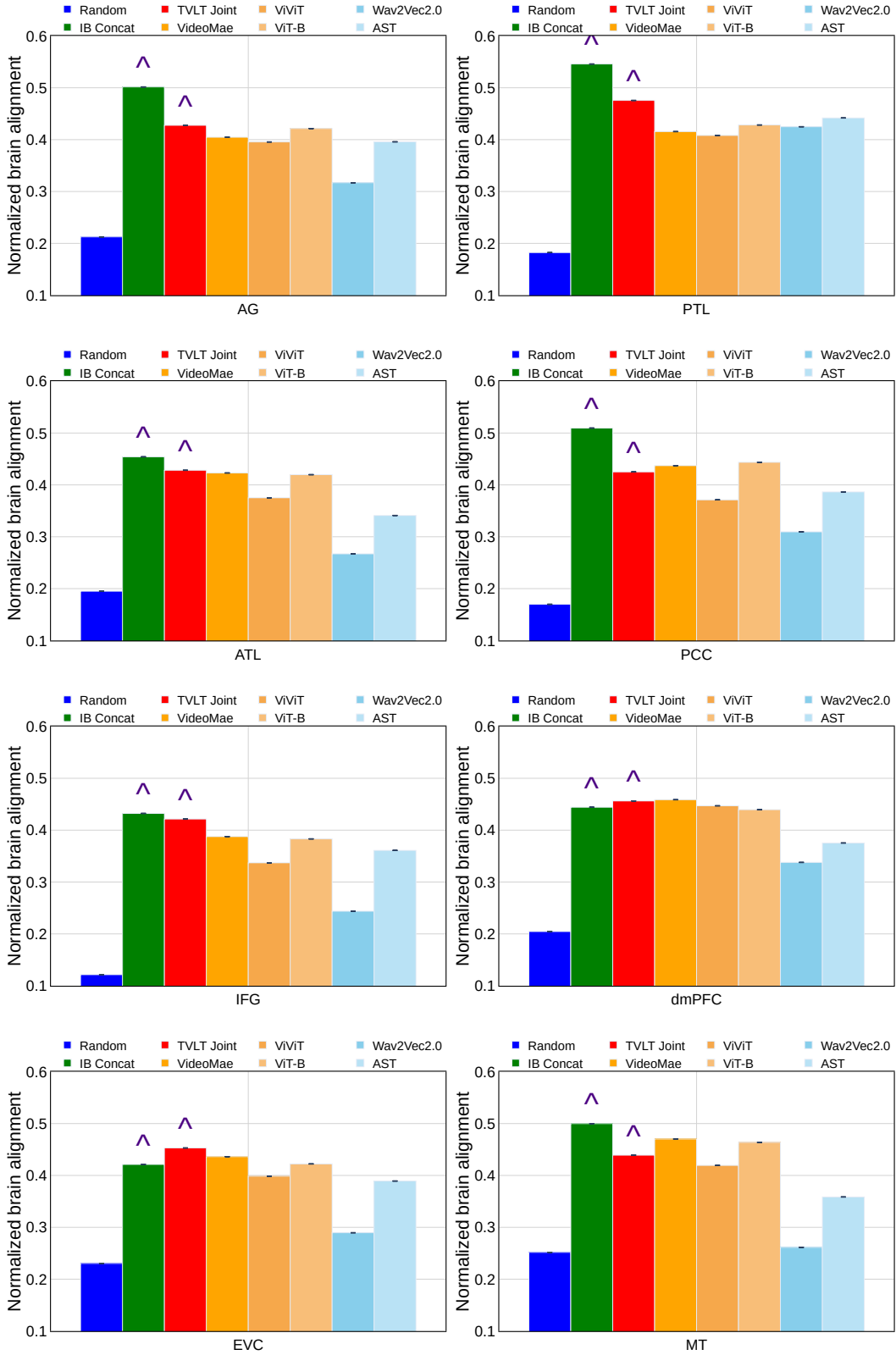


Figure 8: Average normalized brain alignment for video and audio modalities from multi-modal and individual modality features across whole brain and several ROIs of language (ATL, ATL, PTL, IFG, PCC, dmPFC) and visual (EVC, MT). Error bars indicate the standard error of the mean across participants.

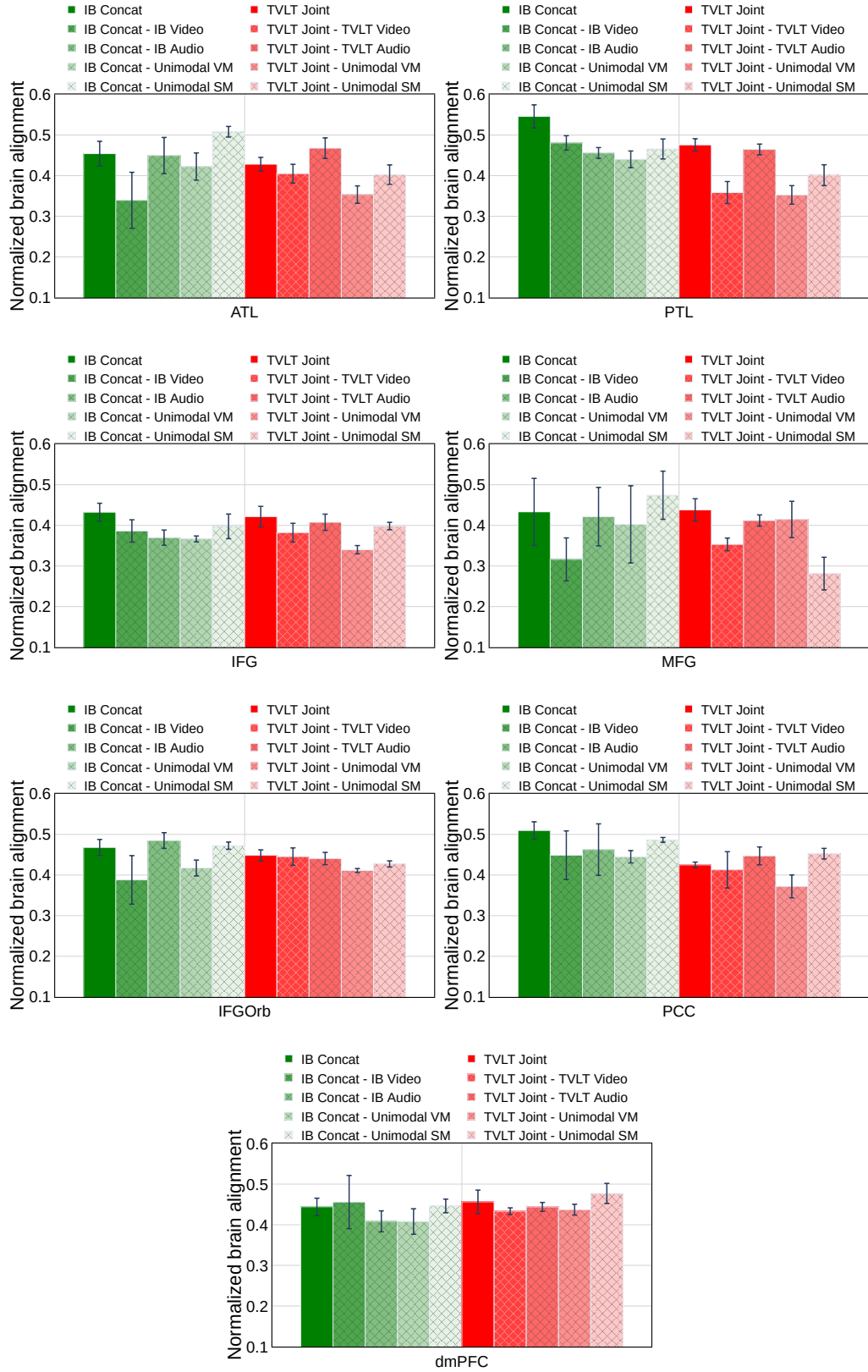


Figure 9: Residual analysis for ATL, PTL, IFG, MFG, IFGOrb, PCC and dmPFC regions: Average normalized brain alignment was computed across participants before and after removal of video and audio embeddings from both jointly pretrained and cross-modality models. Error bars indicate the standard error of the mean across participants.

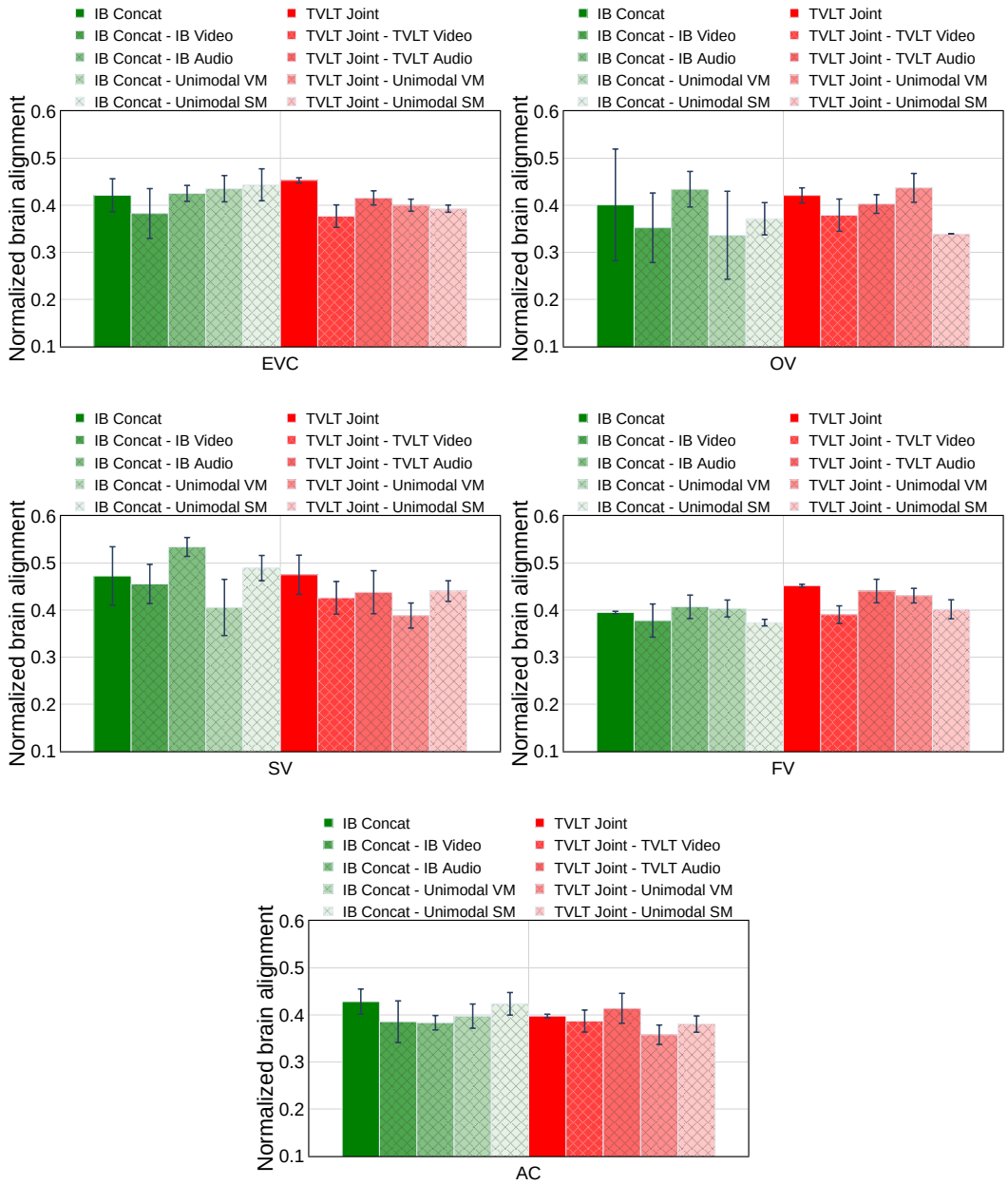


Figure 10: Residual analysis for EVC, OV, SV, FV and AC regions: Average normalized brain alignment was computed across participants before and after removal of video and audio embeddings from both jointly pretrained and cross-modality models. Error bars indicate the standard error of the mean across participants.

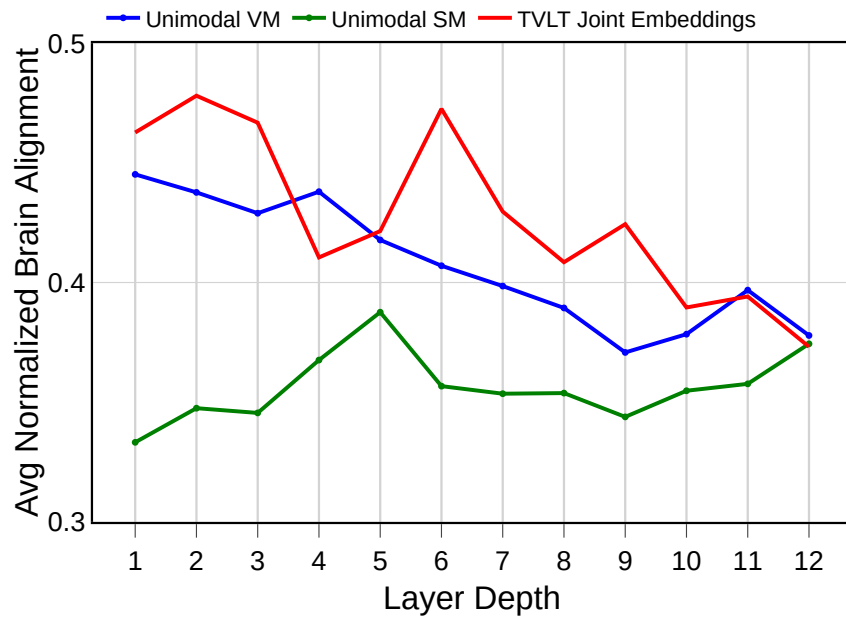


Figure 11: Normalized brain alignment across layers for multi-modal model (TVLT joint embeddings) and unimodal video and speech models.