

# Cortical response to naturalistic stimuli is largely predictable with deep neural networks

Meenakshi Khosla,<sup>1</sup> Gia H. Ngo,<sup>1</sup> Keith Jamison,<sup>2</sup>  
Amy Kuceyeski,<sup>2,3</sup> Mert R. Sabuncu<sup>1,2,4\*</sup>

<sup>1</sup>School of Electrical & Computer Engineering, Cornell University

<sup>2</sup>Radiology, Weill Cornell Medicine

<sup>3</sup>Brain and Mind Research Institute, Weill Cornell Medicine

<sup>4</sup>Nancy E. & Peter C. Meinig School of Biomedical Engineering, Cornell University

\*To whom correspondence should be addressed; Email: msabuncu@cornell.edu

**Naturalistic stimuli, such as movies, activate a substantial portion of the human brain, invoking a response shared across individuals. Encoding models that predict the neural response to a given stimulus can be very useful for studying brain function. However, existing neural encoding models focus on limited aspects of naturalistic stimuli, ignoring the complex and dynamic interactions of modalities in this inherently context-rich paradigm. Using movie watching data from the Human Connectome Project (HCP,  $N = 158$ ) database, we build group-level models of neural activity that incorporate several inductive biases about information processing in the brain, including hierarchical processing, assimilation over longer timescales and multi-sensory auditory-visual interactions. We demonstrate how incorporating this joint information leads to remarkable prediction performance across large areas of the cortex, well beyond the visual and auditory cortices into multi-sensory sites and frontal cortex. Furthermore, we illustrate that encoding models learn high-level concepts that generalize remarkably well to alternate task-bound paradigms. Taken together, our findings underscore the potential of neural encoding models as a powerful tool for studying brain function in ecologically valid conditions.**

## 22 Introduction

23 How are dynamic signals from multiple senses integrated in our minds to generate a coherent  
24 percept of the world? Understanding the neural basis of perception has been a longstanding  
25 goal of neuroscience. Previously, sensory perception in humans has been dominantly studied via  
26 controlled task-based paradigms that reduce computations underlying brain function into simpler,  
27 isolated components, preventing broad generalizations to new environments or tasks (*1*).  
28 Alternatively, fMRI recordings from healthy subjects during free-viewing of movies present a  
29 powerful opportunity to build ecologically-sound and generalizable models of sensory systems,  
30 known as encoding models (*2, 3, 4, 5, 6, 7*).

31 To date, however, existing works on encoding models study sensory systems individually, and  
32 often ignore the temporal context of the sensory input. In reality, the different senses are not  
33 perceived in isolation; rather, they are closely entwined through a phenomenon now well-known  
34 as multi-sensory integration (*8, 9*). For example, specific visual scenes and auditory signals  
35 occur in conjunction and this synergy in auditory-visual information can enhance perception in  
36 animals, improving object recognition and event detection as well as markedly reducing reaction  
37 times (*10*). Furthermore, our cognitive experiences unfold over time; much of the meaning we  
38 infer is from stimulation sequences rather than from instantaneous visual or auditory stimuli.  
39 This integration of information from multiple natural sensory signals over time is crucial to our  
40 cognitive experience. Yet, previous encoding methodologies have precluded the joint encoding  
41 of this rich information into a mental representation of the world.

42 Accurate group-level predictive models of whole-brain neural activity can be invaluable to the  
43 field of sensory neuroscience. These models learn to disregard the idiosyncratic signals and/or  
44 noise within each individual, while capturing only the shared response relevant to the stimuli.  
45 Naturalistic viewing engages multiple brain systems and involves several cognitive processes  
46 simultaneously, including auditory and visual processing, memory encoding and many other  
47 functions (*11*). Group-level analysis in this paradigm is enabled by the synchrony of neuronal  
48 fluctuations in large areas of the cortex across subjects (*12*). Thus far, inter-subject correlation  
49 (ISC) analysis (*12*) has been a cornerstone tool for naturalistic paradigms because of its ability  
50 to characterize the shared response across individuals. Group-level encoding models adopt an  
51 alternative approach for capturing shared response, one grounded in out-of-sample prediction  
52 and generalization (*1*). This allows them to model neural activity beyond a constrained stimulus  
53 set. However, there is a clear gap between the two mediums of analysis. While ISC analysis  
54 suggests that large areas of the cortex exhibit fluctuations that are consistent across subjects,  
55 existing neural encoding models have largely focused on predicting activity within pre-defined  
56 functional areas of the brain such as visual and auditory cortices. It is unclear how they may  
57 be scaled to develop a single predictive model for whole-brain neural responses, given that  
58 naturalistic scenes produce wide-spread cortical activations. In this paper, we aim to fill this  
59 gap: provided adequate characterization of stimuli, we hypothesize that the stable component  
60 of neural activity across a subject population, i.e., the stimulus related activity, should be pre-

61 dictable. In the present study, we aim to quantify and improve the encoding of this wide-spread  
62 stimulus-driven cortical activity using rich stimulus descriptions.

63 Brain responses in real-world conditions are highly complex and variable. Owing to their high  
64 expressive capacity, deep neural networks (DNNs) are well-suited to model the complex high-  
65 dimensional nature of neural activity in response to the multitude of signals encountered during  
66 movie-watching. Recently, DNNs optimized for image or sound recognition have emerged as  
67 powerful models of computations underlying sensory processing (4, 5, 7, 2), surpassing tradi-  
68 tional models of image or sound representation based on Gabor filters (3) and spectrotempo-  
69 ral filters (13), respectively, in higher-order processing regions. In this approach, the stimuli  
70 presented during brain activity recordings are fed as input to pre-trained neural networks and  
71 activations of individual layers are linearly transformed into predictions of neural responses in  
72 different regions of the brain. This approach affords a useful interpretation of these feature  
73 spaces as outcomes of a task-constrained optimization, shedding light on how high-level be-  
74 havioral goals, such as recognition, may constrain representations in neural systems (2). While  
75 useful, task-driven features may diverge from optimal neural representations and tuning these  
76 features to better match the latter may be both feasible and beneficial (14). This approach can  
77 help bridge the quantitative gap in explaining neural responses under realistic conditions while  
78 improving our understanding of the nature of information processing in the brain. From a purely  
79 modeling standpoint, our methodological innovations are threefold. First, we propose an end-  
80 to-end deep-learning based encoding model that extracts semantic feature maps from audio and  
81 visual recognition networks and refines them jointly to predict the evoked brain response. To  
82 this effect, we demonstrate that using different modalities concurrently leads to improvements  
83 in brain encoding. Second, we note that cognitive perception during movie-watching involves  
84 maintaining memory over time and demonstrate the suitability of recurrent neural networks  
85 (RNNs) to capture these temporal dynamics. Finally, based on existing evidence of hierarchical  
86 information processing in visual and auditory cortices (5, 7), we adopt features at multiple lev-  
87 els of abstraction rather than low level or high level stimulus characteristics alone. We embed  
88 these inductive biases about hierarchy, long-term memory and multi-modal integration into our  
89 neural architecture and demonstrate that this comprehensive deep-learning framework general-  
90 izes remarkably well to unseen data. Specifically, using fMRI recordings from a large cohort of  
91 subjects in the HCP, we build group-level encoding models that reliably predict stimuli-induced  
92 neuronal fluctuations across large parts of the cortex. As a demonstration of application, we  
93 employ these encoding models to predict neural activity in response to other task-based stimuli  
94 and report excellent transferability of these models to artificial stimuli from constrained cogni-  
95 tive paradigms. This further suggests that these encoding models are able to capture high-level  
96 mechanisms of sensory processing.

97 Approaching multi-sensory perception through the predictive lens of encoding models has sev-  
98 eral advantages. Because of their unconstrained nature, encoding models can enable data-driven  
99 exploration and catalyze new discoveries. Using six neural encoding models with different tem-  
100 poral scales and/or sensory inputs, trained only on ~36 minutes of naturalistic data per subject,

101 we can replicate findings from a large number of prior studies on sensory processing. First, by  
102 prominently highlighting the transition from short to long temporal receptive windows as we  
103 move progressively from early to high-level auditory areas, we can distinguish the cortical tem-  
104 poral hierarchy. Next, by differentiating uni-sensory cortices from multi-sensory regions such  
105 as the superior temporal sulcus and angular gyrus, we can reproduce the multi-modal architec-  
106 ture of the brain. Finally, by synthesizing neural response to arbitrary stimuli such as faces,  
107 scenes or speech, we can demonstrate the functional specialization of known brain regions for  
108 processing of these distinct categories. Altogether, our results highlight the advantages and  
109 ubiquitous applications of DNN encoding models of naturalistic stimuli.

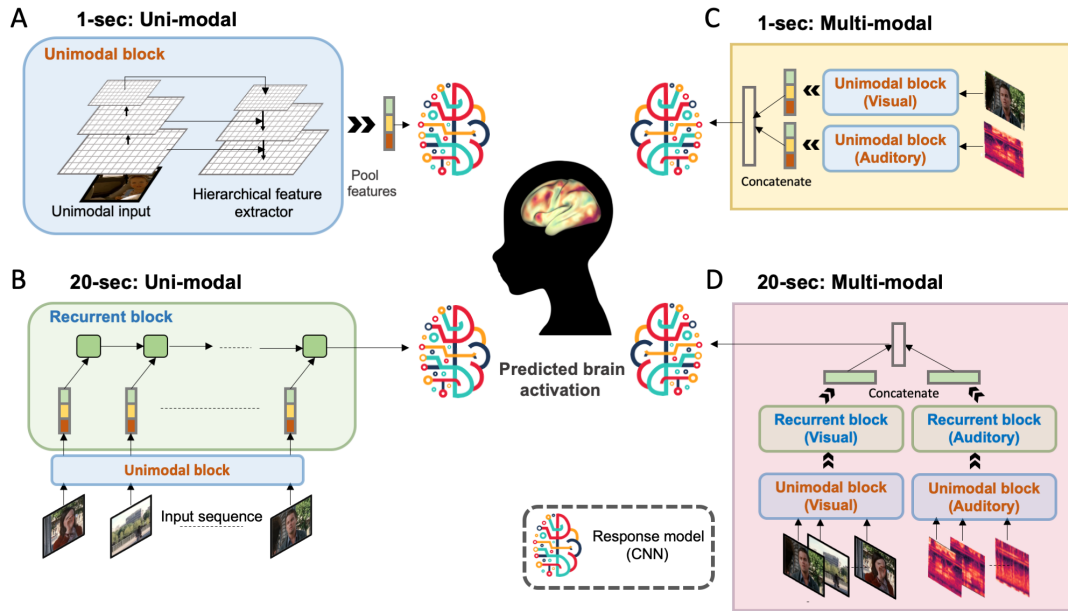
## 110 **Materials and Methods**

### 111 **Dataset**

112 We study high-resolution 7T fMRI data of 158 individuals from the Human Connectome Project  
113 movie-watching protocol comprising 4 audio-visual movie scans (15, 16). The movies repre-  
114 sent a diverse collection, ranging from short snippets of Hollywood movies to independent  
115 vimeo clips. All fMRI data was preprocessed following the HCP pipeline, which includes mo-  
116 tion and distortion correction, high-pass filtering, head motion effect regression using Friston  
117 24-parameter model, automatic removal of artifactual timeseries identified with Independent  
118 Component Analysis (ICA) as well as nonlinear registration to the MNI template space (16).  
119 Complete data acquisition and preprocessing details are described elsewhere (15, 16). Finally,  
120 whole-brain fMRI volumes of size 113x136x113 are used as the prediction target of all pro-  
121 posed encoding models. Rest periods as well as the first 20 seconds of every movie segment  
122 were discarded from all analysis, leaving ~12 minutes of audio-visual stimulation data per  
123 movie paired with the corresponding fMRI response. We estimated a hemodynamic delay of  
124 4 sec using ROI-based based encoding models, as the response latency that yields highest en-  
125 coding performance (Figure S2, see Supplementary Information for details). Thus, all proposed  
126 models are trained to use the above stimuli to predict the fMRI response 4 seconds *after* the cor-  
127 responding stimulus presentation. We train and validate our models on 3 audio-visual movies  
128 with a 9:1 split respectively and evaluate our models on the first three clips of the held-out test  
129 movie. Since the last clip in the held-out movie is repeated within the training movies, we  
130 excluded it from our analysis.

### 131 **Methodology**

132 We train six encoding models employing different facets of the complex, dynamic movie stim-  
133 ulus. These include: (1) Audio-1sec and (2) Audio-20sec models, which are trained on single  
134 audio spectrograms extracted over 1 second epochs and contiguous sequences of 20 spectro-  
135 grams spanning 20 seconds respectively; (3) Visual-1sec and (4) Visual-20sec models, trained  
136 with last frames of 1-second epochs and sequences of 20 evenly spaced frames within 20-second



**Fig. 1.** Schematic of the proposed models. (A) The short-duration (1-sec) auditory and visual models take a single image or spectrogram as input, extract multi-scale hierarchical features and feed them into a CNN-based response model to predict the whole-brain response (B) The long-duration (20-sec) uni-modal models take a sequence of images or spectrograms as input, feed their hierarchical features into a recurrent pathway and extract the last hidden state representation for the response model (C) The short-duration multi-modal model combines uni-modal features and passes them into the response model (D) The long-duration multi-modal model combines auditory and visual representations from the recurrent pathways for whole-brain prediction. Architectural details, including the feature extractor and convolutional response model are provided in Supplementary Information.

137 clips respectively; (5) Audiovisual-1sec and (6) Audiovisual-20sec models, which employ au-  
 138 dio and visual input as described above, *jointly*. All models are trained to minimize the *mean*  
 139 *squared error* between the predicted and measured whole-brain response. Figure 1 depicts the  
 140 overall methodology for training different encoding models.

## 141 Stimuli

142 **Audio** We extract mel-spectrograms over 64 frequency bands between 125-7500 Hz from  
 143 sound waveforms to represent auditory stimulus in  $\sim 1$  second epochs, following (17). The  
 144 audio spectrogram is treated as a single grayscale 96x64 image, denoted by  $x_t^a$ , for the short  
 145 duration model. For the longer-duration model, the input is simply a contiguous sequence of  
 146 20 of these gray-scale images, represented as  $s_t^a = \{x_i^a\}_{i=t-19}^t$ . This representation of audi-  
 147 tory input is also supported by strong evidence that suggests the cochlea may be providing a  
 148 spectrogram-like input to the brain for information processing (18).

149 **Visual** All videos were collected at 24 fps. We extract the last frame of every second of the  
150 video as a 720x1280x3 RGB input, denoted by  $x_t^v$ , for the 1-sec models. We emphasize that the  
151 input here is a single RGB frame and we are using the 1-sec terminology only to be consistent  
152 with the nomenclature for audio models. We further arrange the last frame of every second in  
153 a 20-second clip into a sequence of 20 images, denoted by  $s_t^v = \{x_i^v\}_{i=t-19}^t$ , to represent the  
154 continuous stream of visual stimuli. These are presented to the longer-duration Visual-20sec  
155 and Audiovisual-20sec models.

156 The inputs to the Audio-1sec, Visual-1sec, Audio-20sec, Visual-20sec, Audiovisual-1sec and  
157 Audiovisual-20sec models are thus given as  $x_t^a, x_t^v, s_t^a, s_t^v, \{x_t^a, x_t^v\}$  and  $\{s_t^a, s_t^v\}$  respectively.

### 158 **Audio-1sec and Visual-1sec models**

159 Neural encoding models comprise two components: a feature extractor, which pulls out rel-  
160 evant features,  $\mathbf{s}$ , from raw images or audio waveforms and a response model, which maps  
161 these stimuli features onto brain responses. In contrast to existing works that employ a linear  
162 response model (4, 7), we propose a CNN-based response model where stimulus features are  
163 mapped onto neural data using non-linear transformations. Previous studies have reported a  
164 cortical processing hierarchy where low-level features from early layers of a CNN-based fea-  
165 ture extractor best predict responses in early sensory areas while semantically-rich deeper layers  
166 best predict higher sensory regions (7, 5). To account for this effect, we employ a hierarchical  
167 feature extractor based on feature pyramid networks (19) that combines features from early,  
168 intermediate and later layers simultaneously. The detailed architectures of both components,  
169 including the feature extractor and convolutional response model are described in Figure S3.  
170 We employ state-of-the-art pre-trained ResNet-50 (20) and VGG-ish (17) architectures in the  
171 pyramid network to extract multi-scale features from images and audio spectrograms, respec-  
172 tively. The base architectures were selected because pre-trained weights of these networks  
173 optimized for behaviorally relevant tasks (recognition) on large datasets, namely Imagenet (21)  
174 and Youtube-8M (22), were publicly available. Resnet-50 was trained on image classification  
175 with 1000 classes, while the VGG-ish network was pre-trained on audio event recognition with  
176  $\sim 30\text{K}$  categories. Further, due to computational and memory budget, the Resnet-50 was frozen  
177 during training across all models. On the other hand, we were able to fine-tune the VGG-ish  
178 network in both the Audio and Audiovisual encoding models. We note that in contrast to im-  
179 ages, there is a clear asymmetry in the axes of a spectrogram, where the distinct meanings of  
180 time and frequency might warrant 1D convolutions over time instead of 2D convolutions over  
181 both frequency and temporal axes. However, we found the benefits of a pre-trained network to  
182 be substantial in training convergence time and hence did not explore more appropriate archi-  
183 tectures.

## 184 **Audio-20sec and Visual-20sec models**

185 Audio-20sec and Visual-20sec models employ the same feature extractor and CNN response  
186 model as their 1-second counterparts. However, here, the feature extraction step is applied on  
187 each image in a sequence of 20 frames, followed by a long short-term memory (LSTM) module  
188 to model the temporal propagation of these features. The output dimensions of the LSTM unit  
189 are set to 1024 and 512 for the visual and auditory models respectively, to ensure an equitable  
190 comparison with the corresponding 1-sec models. The last hidden state output of this LSTM  
191 unit is fed into the CNN response model with the same architecture as the 1-sec models.

## 192 **Audiovisual-1sec and Audiovisual-20sec models**

193 Meaningful comparison across different models requires the control of as many design choices  
194 as possible. To ensure fair comparisons, the Audiovisual-1sec model employs the same feature  
195 extractors as the Visual-1sec and Audio-1sec models. The only difference, here, is that the  
196 corresponding 1024-D and 512-D feature representations are concatenated before presenting  
197 to the CNN response model and the concatenated features are passed into a bottleneck layer  
198 to reduce the final feature dimensionality to the maximum among audio and visual feature  
199 dimensions, i.e., 1024, so that the multi-modal model is not equipped with a higher-dimensional  
200 feature space than the maximum among uni-modal models. We note that the response model  
201 has the same architecture across all 6 proposed models. Similarly, the Audiovisual-20sec model  
202 employs the same feature extraction scheme as the Visual-20sec and Audio-20sec models, but  
203 fuses the last hidden state output of the respective LSTM units by simple concatenation followed  
204 by a dense layer to reduce feature dimensionality to 1024 before feeding it into the response  
205 model.

## 206 **Evaluation**

207 We first evaluated the prediction accuracy of all models on the independent held-out movie by  
208 computing Pearson correlation coefficient (R) between the measured and predicted response at  
209 every voxel. Here, the ‘measured’ response refers to the group-averaged response across the  
210 same group of 158 subjects on which the models were trained. Comparison among these mod-  
211 els enables us to tease apart the sensitivity of individual voxels to input timescales and different  
212 sensory stimuli. Voxel-level correlation coefficients between the predicted and measured re-  
213 sponses were averaged to summarize the prediction accuracy of each model in relevant cortical  
214 areas (Figure 2B-F). For this region-level analysis, ROIs were derived with a comprehensive  
215 multi-modal parcellation of the human cortex (23), which was mapped onto the MNI-1.6 mm  
216 resolution template. We note that ROIs were employed only to interpret the results of the  
217 study and relate them to existing literature. We emphasize that all performance metrics re-  
218 ported henceforth are based on voxel-level correlations. It is important to note that prediction  
219 accuracy at every voxel is bounded by the proportion of non-stimulus related variance that re-  
220 flects measurement noise or other factors. We thus also show the regional level performance of

221 all models against the reliability (“noise ceiling”) of measured responses within those regions  
222 (Figure 3).

223 *Noise ceiling estimation:*

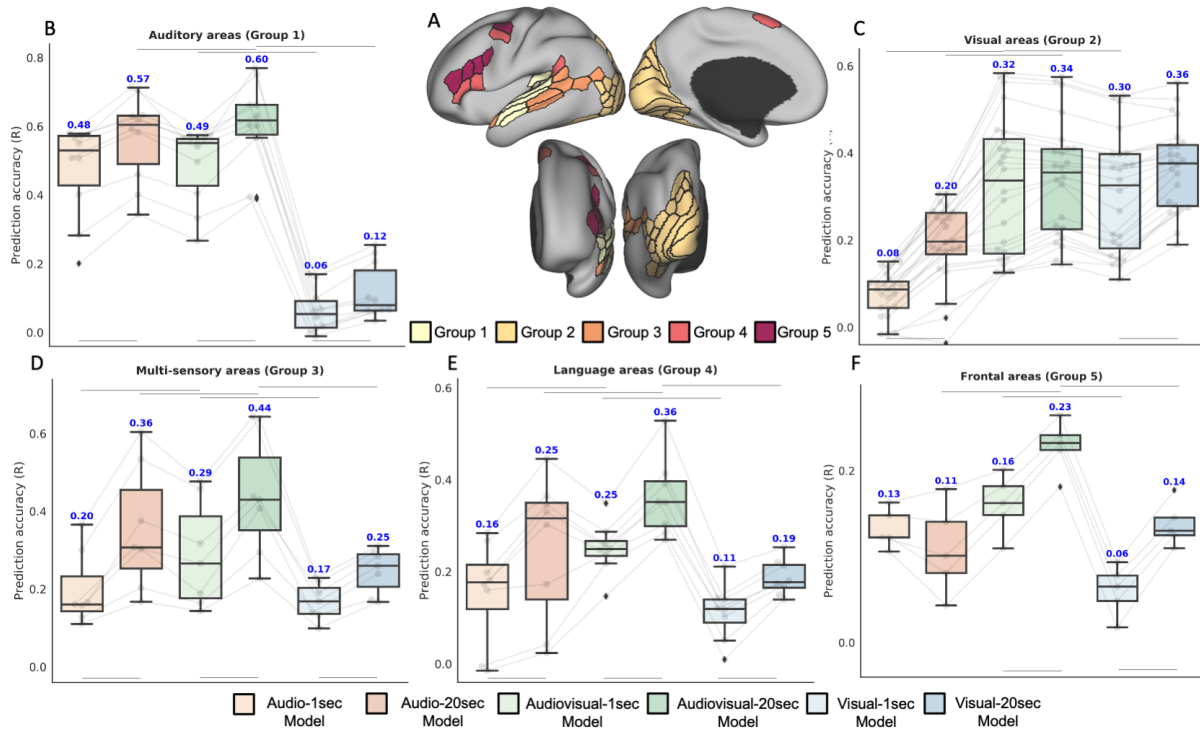
224 The reliability of the group-averaged response at each voxel is estimated from a short 84 second  
225 clip that was repeatedly presented at the end of all movie sessions. We compute an effective up-  
226 per bound on our performance metric, i.e., the correlation coefficient, as the correlation between  
227 the measured fMRI response (group-mean) during different runs. We repeat this process 6 times  
228 (choosing pairs from 4 repeat measurements) to get a mean noise ceiling estimate per voxel, as  
229 shown in Figure 3D. We divide the voxel-level prediction accuracy (R) by this noise ceiling to  
230 get noise-normalized prediction accuracy of all models in left panels of Figure 3A-C. We note  
231 that this noise ceiling is computed on the repeated video clip, which is distinct from the test  
232 movie on which the model performance metrics are computed. Direct comparison against this  
233 noise ceiling can be sub-optimal, especially if the properties of the group-averaged response  
234 vary drastically across the two stimulus conditions. We address this limitation during model  
235 evaluation against data from a held-out independent group of subjects by computing a more  
236 suitable upper bound, which is achievable by a group-level encoding model (Figure S7, see  
237 Supplementary Information for more details). As we demonstrate in the results (Figure S7, S8),  
238 the trend and spatial distribution of model performance against noise ceiling remains unchanged  
239 across the model evaluation and noise ceiling estimation method.

## 240 Results

### 241 **Multi-sensory inputs and longer time-scales lead to the best encoding performance with** 242 **significant correlations across a large proportion of the stimulus-driven cortex**

243 To gain quantitative insight into the influence of temporal history and multi-sensory inputs  
244 on encoding performance across the brain, we computed the mean prediction accuracy in five  
245 groups of regions defined as per the HCP MMP parcellation (23), namely, (1) auditory regions  
246 comprising both early and association areas, (2) early visual and visual association regions, (3)  
247 known multi-sensory sites and regions forming a bridge between higher auditory and higher  
248 visual areas, (4) language-associated regions, and (5) frontal cortical areas. As our research  
249 concerns stimulus-driven processing, only ROIs belonging to the “stimulus-driven” cortex were  
250 included in the above groups (Table S2, see Supplementary Information for the definition of  
251 “stimulus-driven” cortex). Groups 1 and 2, which are associated with a single modality (audi-  
252 tory or visual) do not show any marked improvement from audio-visual multi-sensory inputs  
253 and are best predicted by features of their respective sensory stimulus (Figure 2B,C). The per-  
254 formance boost with multi-sensory inputs is more pronounced in groups 3, 4 and 5 which are  
255 not preferentially associated with a single modality, but are involved in higher-order processing  
256 of sensory stimuli (Figure 2D-F). Further, temporal history of the stimulus yields consistent  
257 improvement in prediction performance in almost all groups of regions, albeit to different ex-





**Fig. 2.** Regional predictive accuracy for the test movie. (B)-(F) depict quantitative evaluation metrics for all the proposed models across major groups of regions as identified in the HCP MMP parcellation (A). Predictive accuracy of all models is summarized across (B) auditory, (C) visual, (D) multi-sensory, (E) language and (F) frontal areas. Box plots depict quartiles and swarmplots depict mean prediction accuracy of every ROI in the group. For language areas (Group 4), left and right hemisphere ROIs are shown as separate points in the swarmplot because of marked differences in prediction accuracy. Statistical significance tests (results indicated with horizontal bars) are performed to compare 1-sec and 20-sec models of the same modality (3 comparisons) or uni-modal against multi-modal models of the same duration (4 comparisons) using paired t-test ( $p$ -value  $< 0.05$ , Bonferroni corrected) on mean prediction accuracy within ROIs of each group.

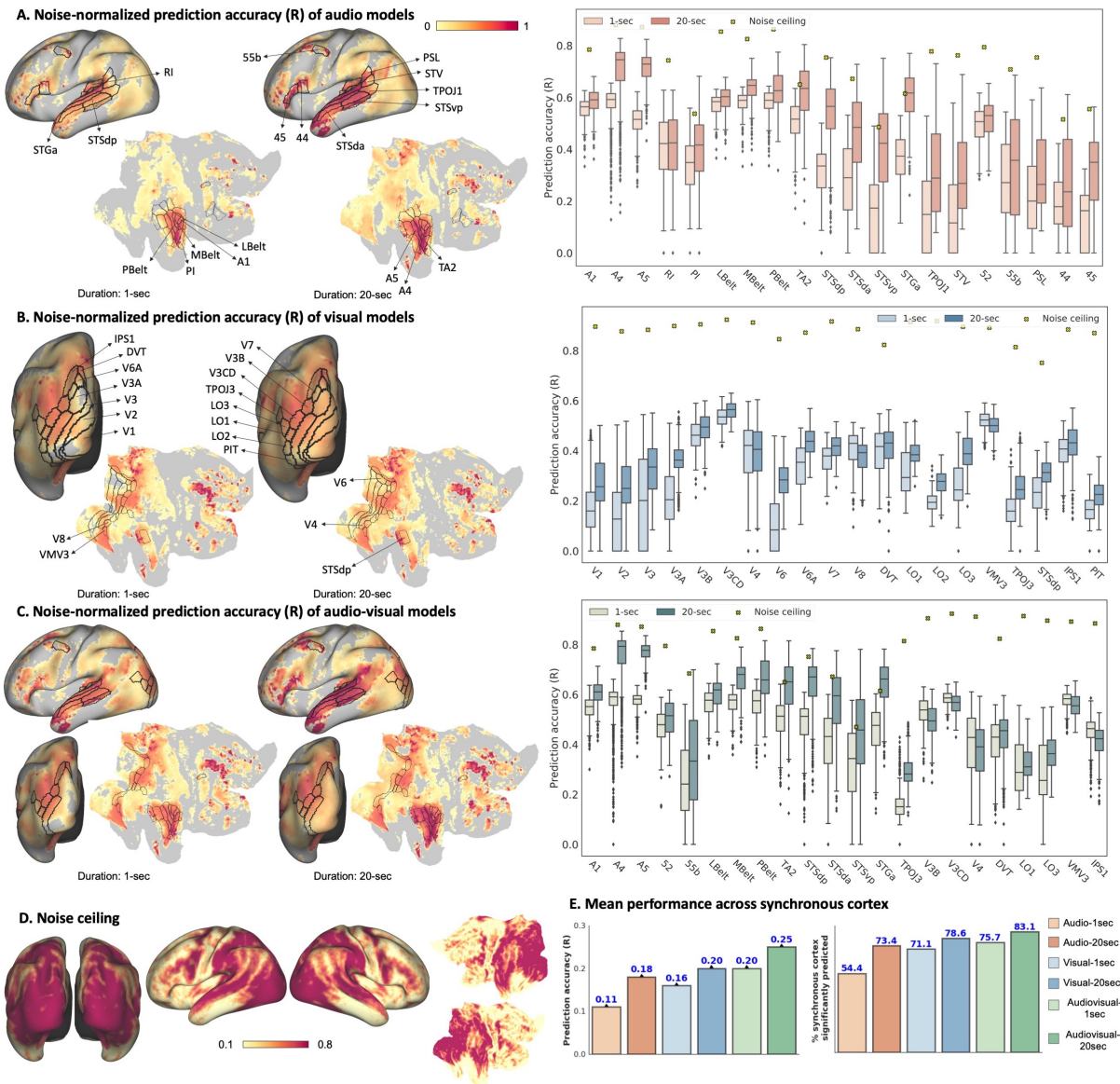
258 tents. Improvements in groups 3, 4 and 5 agree well with the idea that higher-order sensory  
259 processing as well as cognitive and perceptual processes, such as attention and working mem-  
260 ory, are hinged upon the history of sensory stimuli; therefore, accumulated information benefits  
261 response prediction in regions recruited for these functions. Further, both auditory and visual  
262 association cortices are known to contain regions that are responsive to sensory information ac-  
263 cumulated over the order of seconds (24). This potentially explains the significant improvement  
264 observed for long-timescale encoding models compared to their short-timescale counterparts in  
265 these sensory cortices (Figure 4). Together, the Audiovisual-20sec model integrating audio-  
266 visual multi-sensory information over longer time-scales yields maximum prediction accuracy  
267 (R) and highest percentage (~ 83 percent) of significantly predicted voxels across the stimulus-  
268 driven cortex (Figure 3E), suggesting that the Audiovisual-20sec model can adequately capture  
269 complementary features of each additional facet (multi-sensory stimuli / temporal information)  
270 of the sensory environment.

### 271 **Longer time-scales improve encoding performance, particularly in higher order auditory** 272 **areas**

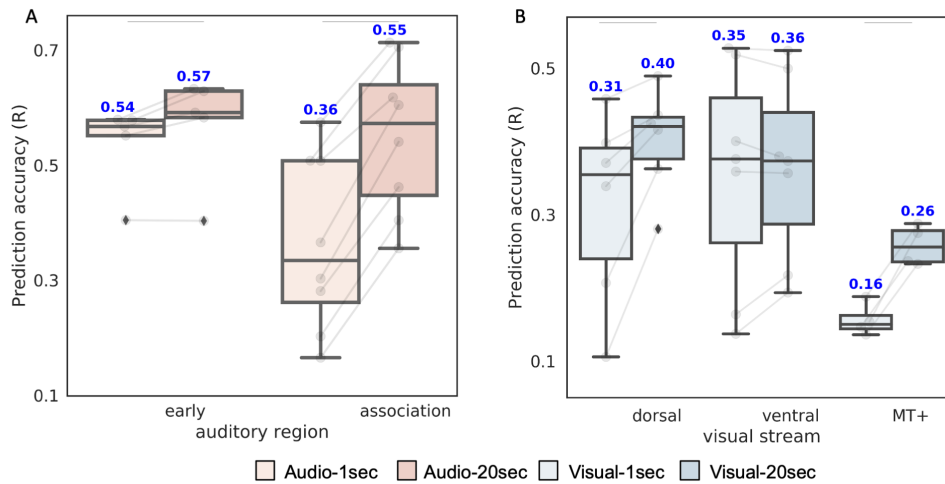
273 As a movie unfolds over time, the dynamic stream of multi-modal stimuli continuously up-  
274 dates our neural codes. Evidence from neuroimaging experiments suggests that different brain  
275 regions integrate information at different timescales; a cortical temporal hierarchy is reported  
276 for auditory perception where early auditory areas encode short timescale events while higher  
277 association areas process information over longer spans (25). This temporal gradient of audi-  
278 tory processing is well-replicated within our study. Comparison of 1-sec and 20-sec models  
279 allows us to distinguish brain regions that process information at shorter timescales from those  
280 that rely on longer dynamics. There is a negligible contribution of longer timescale inputs  
281 on prediction correlations in regions within early auditory cortex, such as A1, LBelt, PBelt,  
282 MBelt and Restro-insular cortex (RI) (Figure 3A, 4A), in line with previous reports suggesting  
283 short temporal receptive windows (TRWs) of early sensory regions (25). Shorter integration  
284 windows are in agreement with the notion that these regions facilitate rapid processing of the  
285 instantaneous incoming auditory input. In contrast, response in voxels within auditory associ-  
286 ation ROIs lying mainly in the superior temporal sulcus or along the temporal gyrus (A4, A5,  
287 STSda, STSva, STSdp, STSvp, STGa, TA2) is seen to be much better predicted with longer  
288 time-scales (Figure 3A, 4A). Cumulatively across association ROIs, Audio-20sec model yields  
289 a highly significant improvement in prediction accuracy (~50%) over the Audio-1sec model, in  
290 comparison to a marginal improvement (~5%) across early auditory ROIs.

### 291 **Longer time-scales lead to significantly better predictions in the dorsal visual stream and** 292 **MT+ complex**

293 The distinct association of dorsal visual stream with spatial localization and action-oriented  
294 behaviors and ventral visual stream with object identification is well documented in the liter-

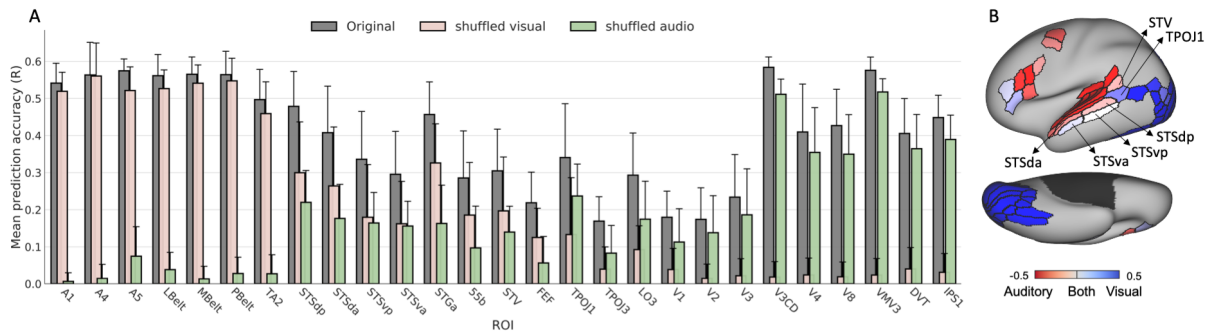


**Fig. 3.** Predictive accuracy of uni-modal (A,B) and multi-modal (C) models over the whole brain in the test movie. Colors on the brain surface indicate the Pearson's correlation coefficient between the predicted timeseries at each voxel and the true voxel's timeseries normalized by the noise ceiling (D) computed on repeated validation clips. Only significantly predicted voxels ( $p$ -value  $< 0.5$ , FDR corrected) are colored. ROI box plots depict the un-normalized correlation coefficients between the predicted and measured response of voxels in each ROI and the respective noise ceiling for the mean. (E) shows the percentage of voxels in stimulus-driven cortex that are significantly predicted by each model and mean prediction accuracy across the stimulus-driven cortex.



**Fig. 4.** Influence of temporal history on encoding performance. (A) Mean predictive performance of Audio-1sec and Audio-20sec models in early auditory and association auditory cortex ROIs. A major boost in encoding performance is seen across auditory association regions with the 20-sec model. (B) Mean predictive performance of Visual-1sec and Visual-20sec models across ROIs in the dorsal, ventral and MT+ regions. Dorsal stream and MT+ ROIs exhibit a significant improvement with Visual-20sec model but no effect is observed for the ventral stream. Boxplots are overlaid on top of the beeswarm plot to depict quartiles. Horizontal bars indicate significant differences between models in the mean prediction accuracy within ROIs of each stream using paired t-test (p-value < 0.05).

295 ature (26). Another specialized visual area is the medial temporal complex (MT+), which has  
 296 been shown to play a central role in motion processing. The functional division between these  
 297 streams thus suggests a stronger influence of temporal dynamics on responses along the dorsal  
 298 pathway and MT+ regions. To test this hypothesis, we contrast the encoding performance of  
 299 Visual-1sec and Visual-20sec models across the three groups by averaging voxel-wise correla-  
 300 tions in their constituent ROIs. In accordance with the dorsal/ventral/MT+ stream definition in  
 301 the HCP MMP parcellation, we use the following ROIs for analysis: (a) dorsal: V3A, V3B, V6,  
 302 V6A, V7, IPS1 (b) ventral: V8, Ventral Visual Complex (VVC), PIT complex, Fusiform Face  
 303 Complex (FFC) and Ventro-medial Visual areas 1,2 and 3 (c) MT+: MT, MST, V4t, FST. Figure  
 304 4B demonstrates the distribution of mean correlations over these ROIs for different models and  
 305 streams. Our findings suggest that temporal history, as captured by the Visual-20sec model,  
 306 can be remarkably beneficial to response prediction across the dorsal visual stream ( 30% im-  
 307 provement over Visual-1sec model) and the MT+ complex ( 62% improvement over Visual-1sec  
 308 model), in agreement with our *a priori* hypothesis . Further, in our experiments, no marked im-  
 309 provement was observed for the ventral visual stream, indicating a non-significant influence of  
 310 temporal dynamics on these regions.



**Fig. 5.** Sensitivity of ROIs to different sensory inputs. (A) Predictive accuracy (R) of audiovisual encoding model with and without input distortions, (B) Sensory sensitivity index of different brain regions as determined using performance metrics under input distortion (see Supplementary Information for details). Regions dominated by a single modality are shown in darker colors, whereas light-colored regions are better predicted by a combination of auditory and visual information. Red indicates auditory-dominant regions whereas blue indicates visual dominance.

311 **Auditory and visual stimuli features jointly approach the noise ceiling in multi-sensory**  
 312 **areas**

313 Examining prediction accuracy against response reliability allows us to quantify how far we are  
 314 from explaining predictable neural activity. A high fraction of the stimulus-driven cortex (~  
 315 83%) is predictable with a longer timescale input and joint audiovisual features. Notably, areas  
 316 extending anteriorly and posteriorly from the primary auditory cortex such as the posterior STS,  
 317 STGa and TA2 achieve prediction correlations close to the noise ceiling with the Audiovisual-  
 318 20 sec model (Figure 3C), suggesting that DNN representations are remarkably suited to encode  
 319 their response.

320 Interestingly, performance in auditory regions is much closer to the noise ceiling than visual  
 321 regions. Understanding audition and vision in the same space further allows us to appreciate  
 322 the differences between these modalities. While this may suggest that audition is perhaps a  
 323 simpler modality to model, the differences could also result from a bias of the dataset. A more  
 324 diverse sampling of acoustic stimuli in the training set could allow the model to generalize better  
 325 in auditory regions. Furthermore, in contrast to auditory stimulation where all subjects hear the  
 326 same sounds, visual stimulation can elicit highly varied responses dependent on gaze location.  
 327 This variability could plausibly make group-level visual encoding a more difficult task.

328 **Joint encoding models tease apart the modal sensitivity of voxels throughout the sensory**  
 329 **cortex**

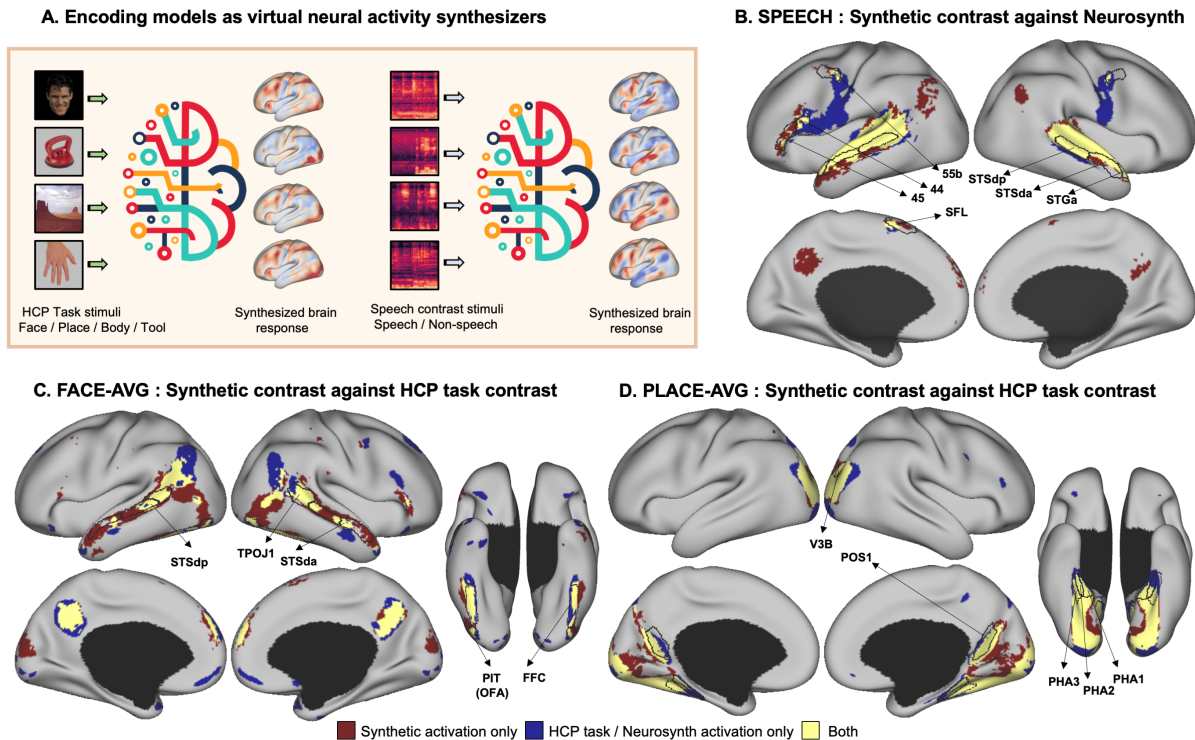
330 Neural patterns evoked by movies are not simply a conjunction of activations in modality-  
 331 specific cortices by their respective uni-sensory inputs; rather, there are known cross-modal  
 332 influences as well as regions that receive afferents from multiple senses (27). Can we interro-

333 gate a joint encoding model to reveal the individual contribution of auditory and visual features  
334 in encoding response across different brain regions? To address this question, we shuffled  
335 inputs of either modality along the temporal axis during inference. We measured test perfor-  
336 mance of the trained audio-visual model on predictions generated by shuffling inputs of one  
337 modality while keeping the other one intact. This distortion at test time allows us to identify  
338 areas that are preferentially associated with either visual or auditory modality. We hypothe-  
339 sized that regions encoding multi-sensory information will incur loss in prediction accuracy  
340 upon distortion of both auditory and visual information. Further, uni-sensory regions will likely  
341 be adversely affected by distortion of either auditory or visual information but not both. To  
342 test this hypothesis, we further developed a sensory-sensitivity index that directly reflects the  
343 sensitivity of individual brain regions to information about auditory or visual stimuli (see Sup-  
344 plementary Information for details). For this examination, we utilized the Audiovisual-1sec  
345 model to avoid potential confounds associated with temporal history, although analysis of the  
346 Audiovisual-20sec model showed similar results. Figure 5 demonstrates the result of this analy-  
347 sis on sensory-specific regions as well as regions known for their involvement in multi-sensory  
348 integration. The benefit from (non-distorted) multi-sensory inputs to the prediction correlations  
349 of the Audio-visual model is most remarkably seen in posterior STS, STGa and sensory-bridge  
350 regions such as the temporal-parietal-occipital junction (TPOJ1-3) and superior temporal visual  
351 (STV) area. Another region that seems to be employing features of both modalities, albeit to a  
352 lesser extent, is the frontal eye field (FEF), whose recruitment in audiovisual attention is well  
353 studied (28).

354 Classically, multi-sensory integration hubs are identified as regions that show enhanced activ-  
355 ity in response to multi-sensory stimulation as opposed to presentation of either uni-sensory  
356 stimuli based on some statistical criteria (29). Accordingly, the posterior STS is consistently  
357 described as a multi-sensory convergence site for audio-visual stimuli (27, 30, 29, 9). Its role  
358 in audiovisual linguistic integration has also been well-studied in the literature (28). Other  
359 multi-sensory integration sites reported extensively in prior literature include the temperopari-  
360 etal junction (9, 27, 28) and superior temporal angular gyrus (31). Our findings above lend strong  
361 support for the multi-sensory nature of all these regions.

## 362 **Encoding models as virtual neural activity synthesizers**

363 Next, we sought to characterize whether encoding models can generalize to novel task paradigms.  
364 By predicting neural activity for different visual categories from the category-specific represen-  
365 tation task within the HCP Working Memory (WM) paradigm, we generated synthetic func-  
366 tional localizers for the two most common visual classes: faces and places. Specifically, we  
367 predict brain response to visual stimuli, comprising faces, places, tools and body parts, from  
368 the HCP task battery. We use the *predicted* response to synthesize contrasts (FACES-AVG and  
369 PLACES-AVG) by computing the difference between mean activations predicted for the cate-  
370 gory of interest (*faces* or *places* respectively) and the average mean activations of all categories  
371 at each voxel (Figure 6). The predicted and measured contrasts are thresholded to keep top 5%



**Fig. 6.** Encoding models as virtual brain activity synthesizers. (A) Synthetic contrasts are generated from trained encoding models by contrasting their “synthesized” (i.e., predicted) response to different stimulus types. (B) Comparison of the synthesized contrast for ‘speech’ against the speech association template on *neurosynth*. (C-D) compare the synthesized contrasts for ‘faces’ and ‘places’ against the corresponding contrasts derived from HCP tfMRI experiments.

372 of the voxels.

373 We observe a notable overlap between the synthetic and measured group-level contrasts. Fur-  
374 ther, our findings are consistent with the well-known cortical specificity of neuronal activations  
375 for processing of *faces* and *places*. Both the synthetic and measured *faces* contrasts are con-  
376 sistent with previously identified regions for face-specific processing, including the fusiform  
377 face area (corresponds to fusiform face complex (FFC) in Figure 6), the occipital face area in  
378 lateral occipital cortex (overlaps with the PIT complex in HCP MMP parcellation), and regions  
379 within temporo-parieto-occipital junction and STS (32, 33). Among these, the *selective* role of  
380 the Fusiform Face Area in face processing has been most consistently and robustly established.  
381 Another region known to respond more strongly to faces than other object categories, namely  
382 posterior STS, has been previously implicated in processing of facial emotions (32).

383 Similarly, both synthetic and measured *places* contrasts highlight cortical regions thought to be  
384 prominent in selective processing of visual scenes. These include the parahippocampal areas  
385 (PHA1-3), retrosplenial cortex (POS1 in HCP MMP parcellation) and the transverse occipital  
386 sulcus (TOS), which comprises the occipital place area (OPA) (34).

387 Cortical areas related to speech processing are similarly discovered using our models by con-  
388 trasting activations predicted for speech stimuli against non-speech stimuli such as environmen-  
389 tal sounds (Figure 6B, see Supplementary Information for more details). The synthetic contrast  
390 shows increased activation in language-related areas of the HCP MMP parcellation such as 55b,  
391 44 and the superior frontal language (SFL) area with left-lateralization, in accordance with pre-  
392 vious language fMRI studies (35). In addition, areas tuned for voice processing in STS (36) are  
393 also highlighted. The synthetic map also shows highest correlation with ‘speech’ on *neurosynth*  
394 term-based meta-analysis (37) and overlaps considerably with the speech association template  
395 on the platform.

### 396 **Additional analyses**

397 In prior studies, neural response prediction is done via regularized regression, where the signal  
398 at each voxel is modeled as a weighted sum of stimulus features with appropriate regulariza-  
399 tion on the regression weights. Following earlier works, we also train  $l_2$ -regularized regression  
400 models using features derived from hierarchical convolutional networks trained on image or  
401 sound recognition such as those used in the proposed models, as well as semantic categories  
402 features labelled using the WordNet semantic taxonomy similar to (38). The latter are typically  
403 used for mapping the semantic tuning of individual voxels across the cortex. Our models con-  
404 sistently outperform the baselines, further illustrating the benefits of the proposed methodology  
405 (Figure S4(A)-(C), see Supplementary Information for more details). Additionally, we also per-  
406 formed ablation studies to understand the influence of different network components, namely  
407 the “non-linear” response model as well as the “hierarchical” feature extractor on model predic-  
408 tion performance and found that both components improve performance, although their relative  
409 contribution is stronger in visual encoding models than auditory models (Figure S4D, see Sup-



410plementary Information for more details). The superior predictive performance of our models  
411in comparison to the classical approach along with our ablation studies suggest that an inter-  
412play of end-to-end optimization with a non-linear response model can jointly afford improved  
413generalization performance.

414To test the generality of the models beyond the subject population they were trained on, we fur-  
415ther compared the predictions of all models against the group-averaged response of a held-out  
416group within HCP comprising 20 novel subjects distinct from the 158 individuals used in the  
417training set, on the same independent held-out movie. The noise ceiling for this group was com-  
418puted as the correlation coefficient between the mean measured response for the *independent*  
419test movie across all 158 subjects in the training set and the group-averaged response computed  
420over the 20 new subjects. This metric captures the response component shared across inde-  
421pendent groups of subjects and thus reflects the true upper bound achievable by a group-level  
422encoding model. As shown in Figure S7 (see Supplementary Information for more details), the  
423models can accurately predict neural responses as measured with respect to the group mean of  
424the *held-out* subjects, with the Audiovisual-20sec model performance even approaching noise  
425ceiling in some regions, particularly the higher-order auditory association regions and multi-  
426sensory sites such as the posterior STS. Importantly, the predictivities across the cortical sur-  
427face are consistent with the performance metrics reported for the training subject population in  
428Figure 3. Finally, by comparing model predictions against neural responses at the single subject  
429level for subjects from the held-out group, we further demonstrate that the Audiovisual-20sec  
430model can also successfully capture the response component that individual subjects share with  
431the population (Figure S9, see Supplementary Information for details).

## 432 Discussion

433Free viewing of dynamic audio-visual movies enables an ecologically valid analysis of a col-  
434lective set of functional processes at once, including temporal assimilation and audio-visual  
435integration in addition to momentary sensory-specific processing. Perception, under such stim-  
436ulation, thus recruits sensory systems as well as areas subserving more sophisticated cogni-  
437tive processing. Building quantitatively accurate models of neural response across widespread  
438cortical regions to such real-life, continuous stimuli thus requires an integrated modelling of  
439these disparate computations on sensory inputs. In this paper, we have presented six deep neu-  
440ral network based encoding models with varying sensory and temporal information about the  
441audio-visual stimulus. Subsequently, we queried the role of input history and different sen-  
442sory information on prediction performance across individual regions of the cortex. We have  
443shown that exploiting the richness of the stimulus along the time axis and sensory modality  
444substantially increases the predictive accuracy of neural responses throughout the cortex, so  
445far as approaching the noise ceiling for voxels in some known multi-sensory sites, such as the  
446posterior STS (27, 30, 29, 9).

447 Auditory and visual scenes are the principal input modalities to the brain during naturalistic  
448 viewing. Yet, existing encoding models ignore their interactions. We employ a common strat-  
449 egy in multi-modal machine learning settings, namely feature fusion, to jointly model auditory  
450 and visual signals from the environment. We find that minimizing the prediction error is a  
451 useful guiding principle to learn useful joint representations from an audio-visual stimulation  
452 sequence and demonstrate that models that consume multi-modal signals concurrently, namely,  
453 Audiovisual-1sec and Audiovisual-20sec, can not only predict the respective uni-modal cortices  
454 slightly better but also lead to remarkable improvements in predicting response of multi-sensory  
455 and frontal brain regions (Figure 2). Further, we show that multi-modal neural encoding models  
456 not only boost performance in large areas of the cortex relative to their uni-modal counterparts  
457 (Figure 2,3E), but also shed light on how neural resources are spatially distributed across the  
458 cortex for dynamic multi-sensory perception (Figure 5). The predictivity of different sensory  
459 inputs for neural response, as evaluated on independent held-out data, can facilitate reverse in-  
460 ference by identifying the sensory-associations of different brain regions, providing clues into  
461 the multi-sensory architecture of the cortex. By comparative analysis of predictive performance  
462 in different regions across models (Figure 2) as well as perturbation analysis within the multi-  
463 modal model (Figure 5), we identify a number of regions that are consistently sensitive to both  
464 auditory and visual information, most notably the superior temporal sulcus and some frontal  
465 regions. Regions within inferior frontal cortex, have been implicated in the processing of visual  
466 speech, guiding sensory inferences about the likely common cause of multi-modal auditory and  
467 visual signals, as well as resolving sensory conflicts (39). Prior research has also implicated  
468 an extensive network of inferior frontal and premotor regions in comprehending audiovisual  
469 speech, suggesting that they bind information from both modalities (40). While unveiling the  
470 causal sequence of events for a mechanistic understanding of multi-sensory perception is not  
471 possible with the proposed approach, our findings align well with commonly held theories of  
472 sensory fusion which suggest that uni-sensory signals are initially processed in segregated re-  
473 gions and eventually fused in regions within superior temporal lobe, occipital-temporal junction  
474 and frontal areas (27). This proposition is corroborated by our experiments as response predic-  
475 tion in these regions is best achieved by a combination of both sensory inputs (Figure 3,5).

476 A linear response model with pre-trained and non-trainable feature extractors, while simple and  
477 interpretable, imposes a strong constraint on the feature-response relationship. The underlying  
478 assumption is that neural networks optimized for performance on behaviorally relevant tasks,  
479 are mappable to neural data with a linear transform. We designed a flexible model, capable  
480 of capturing complex non-linear transformations from stimulus feature space to neural space,  
481 leading to more quantitatively accurate models that are better aligned with sensory systems.  
482 Even better accounts of cortical responses are then obtained by interlacing dynamic, multi-  
483 modal representation learning with whole-brain activation regression in an end-to-end fashion.  
484 Using these rich stimulus descriptions, we demonstrated a widespread predictability map across  
485 the cortex, that covers a large portion (~83%) of the stimulus-driven cortex (Figure 3C,E), in-  
486 cluding association and some frontal regions. While inter-subject correlations in these regions

487 are frequently reported (12, 41), suggesting their involvement in stimulus-driven processing,  
488 response predictability in these areas had remained elusive so far. Further, the cortical predic-  
489 tivity is maintained even as we compare model predictions against neural responses of held-out  
490 subjects (Figure S7 and S9), suggesting that the proposed models are capable of successfully  
491 capturing the “shared” or stimulus-driven response component. These results provide com-  
492 pelling evidence that deep neural networks trained end-to-end can learn to capture the complex  
493 computations underlying sensory perception of real-life, continuous stimuli.

494 We further demonstrated that encoding models can form an alternative framework for prob-  
495 ing the time-scales of different brain regions. While primary auditory and auditory belt cor-  
496 tex (comprising A1, PBelt, LBelt, Mbelt) as well as the ventral visual stream benefit only  
497 marginally from temporal information, there is a remarkable improvement in prediction per-  
498 formance in auditory and visual association and pre-frontal cortices, most notably in superior  
499 temporal lobe, visuomotor regions within the dorsal stream such as V6A, temporal parietal oc-  
500 cipital junction and inferior frontal regions. The improvement in prediction performance with  
501 the 20-second input is consistently seen for both uni-modal and multi-modal models. It is im-  
502 portant to acknowledge that directly comparing the prediction accuracies of static (1-sec) and  
503 recurrent (20-sec) models to infer processing timescales of different brain regions has its limita-  
504 tions. First, this analysis can be confounded by the slow hemodynamic response as performance  
505 improvement may be driven in part by the slow and/or spatially varying dynamics. Based on  
506 our analysis with ROI-level encoding models, the latter seems like a less plausible explanation  
507 (Figure S2, see Supplementary Information for details). Further, we performed additional anal-  
508 yses to understand the relationship between performance improvement in individual voxels and  
509 their autocorrelation properties and found a strong correspondence between the two, suggesting  
510 that the distribution of performance improvement across the cortex broadly agrees well with  
511 processing timescales (Figure S5, see Supplementary Information for details).

512 Predictions from long-timescale models are based on temporal history as provided in stimulus  
513 sequences, and not just the instantaneous input. Modeling dynamics within these sequences  
514 appropriately is crucial to probe effects of temporal accumulation. RNNs have internal memo-  
515 ries that capture long-term temporal dependencies relevant for the prediction task, in this case,  
516 encoding brain response, while discarding task-irrelevant content. We compare this modeling  
517 choice against a regularized regression approach on stimulus features concatenated within T-  
518 second clips, with T ranging between 1 and 20 (Figure S4, see Supplementary Information for  
519 details). The inferior performance compared to our proposed models as well as a non-increasing  
520 performance trend against T for these linear models indicates that accumulation of temporal  
521 information by simply concatenating stimulus features over longer temporal windows is insuf-  
522 ficient; rather, models that can efficiently store and access information over longer spans, such  
523 as RNNs with sophisticated gating mechanisms, are much more suitable for modeling neural  
524 computations that unfold over time. Since activations of units within RNNs depend not only  
525 on the incoming stimulus, but also on the “current” state of the network as influenced by past  
526 stimuli, they are capable of holding short-term events into memory. Adding the RNN module

527 can thus be viewed as augmenting the encoding models with working memory.

528 Investigating timescales of representations across brain regions by understanding the influence  
529 of contextual representations on language processing in the brain, as captured by LSTM lan-  
530 guage models for instance, has become a major research focus recently (42). In these language  
531 encoding models for fMRI, past context has been shown to be beneficial in neural response pre-  
532 diction, surpassing word embedding models. However, models that explain neural responses  
533 under dynamic natural vision while exploiting the rich temporal context have not yet been rig-  
534 orously explored with human fMRI datasets. In a previous study with awake mice, recurrent  
535 processing was shown to be useful in modelling the spiking activity of V1 neurons in response  
536 to natural videos (43). In dynamic continuous visual stimulation fMRI paradigms, a common  
537 practice is to concatenate multiple delayed copies of the stimulus to model the hemodynamic  
538 response function as a linear finite impulse response (FIR) function (38). However, since the  
539 feature dimensionality scales linearly with time-steps, this approach is limited to HRF mod-  
540 eling and is not feasible to capture longer dynamics of the order of tens of seconds. Another  
541 approach is to employ features from neural networks trained on video tasks, such as action  
542 recognition (6). However, these encoding models are constrained to capture one aspect of dy-  
543 namic visual scenes and are likely useful to predict neural responses in highly localized brain  
544 regions. Most studies in visual encoding remain limited to static stimuli and evoked responses  
545 in relatively small cortical populations.

546 Our brain has evolved to process ‘natural’ images and sounds. In fact, recent evidence has  
547 shown that sensory systems are intrinsically more attuned to features of naturalistic stimuli  
548 and such stimuli can induce stronger neural responses than task-based stimuli (44). Here, we  
549 demonstrate that encoding models trained with naturalistic data are not limited to modeling  
550 responses of their constrained stimuli set. Instead, by learning high-level concepts of sensory  
551 processing, these models can also generalize to out-of-domain data and replicate results of al-  
552 ternate task-bound paradigms. While our models were trained on complex and cluttered movie  
553 scenes, we tested their ability to predict response to relatively simple stimuli from HCP task bat-  
554 tery, such as faces and scenes (Figure 6). The remarkable similarity between the predicted and  
555 measured contrasts in all cases suggests that ‘synthetic’ brain voxels, predicted by the trained  
556 DNNs, correspond well with the target voxels they were trained to model. We thus provide  
557 evidence that these encoding models are encapsulating stimulus-to-brain relationships extending  
558 beyond the experimental world they were trained in. On the other hand, classical fMRI experi-  
559 ments, for instance task contrasts, don’t generalize outside the experimental circumstance they  
560 were based on. This preliminary evidence suggests that encoding models can serve as promis-  
561 ing alternatives for circumventing the use of contrast conditions to study hypotheses regarding  
562 the functional specialization of different brain regions. Embedded knowledge within these de-  
563 scriptive models of the brain, could also be harnessed in other applications, such as independent  
564 neural population control by optimally synthesizing stimuli to elicit a desired neural activation  
565 pattern (45).

566 With purely data-driven exploration of fMRI recordings under a hypothesis-free naturalistic  
567 experiment, our models replicate the results of previous neuroimaging studies operating under  
568 controlled task-based regimes. Our analysis lends support to existing theories of perception  
569 which suggest that primary sensory cortices build representations at short timescales and lead  
570 up to multi-modal representations in posterior portions of STS (25). Encoding performance  
571 in these regions is consistently improved with longer timescales as well as multi-sensory in-  
572 formation. We reasoned that regions that are sensitive to multi-modal signals and/or longer  
573 stimulus dynamics could be distinguished by interrogating the performance of these models  
574 on unseen data. To date, encoding models have been rarely used in this manner to assess in-  
575 tegration timescales or sensory-sensitivity of different brain regions. Classically, processing  
576 timescales have been probed using various empirical strategies, for example, by observing ac-  
577 tivity decay over brief stimulus presentations or by comparing auto-correlation characteristics  
578 of resting-state and stimulus-evoked activity (46). Further, multi-sensory regions are identified  
579 via carefully-constructed experiments with uni-modal and multi-modal stimulus presentations,  
580 followed by analysis of interaction effects using statistical approaches (27). Here, we suggest  
581 that encoding models can form an alternate framework to reveal clues into these functional prop-  
582 erties that can be rigorously validated with future investigation. As with interpreting the results  
583 of any predictive model, one should, however, proceed with caution. Sounds are generated by  
584 events; this implies that sound representations implicitly convey information about actions that  
585 generated them. Similarly, visual imagery provides clues into auditory characteristics, such as  
586 the presence of absence of speech. Thus, it is difficult to completely disentangle the individual  
587 contributions of auditory and visual features to prediction performance across cortical regions.  
588 Similarly, longer time-scale inputs can lead to a more robust estimate of the momentary sen-  
589 sory signal, potentially confounding the interpretations of TRWs. Here, we contend that these  
590 models can, nonetheless, serve as powerful hypothesis generation tools.

591 The methodological innovations in this study must also be considered in light of their limi-  
592 tations. Due to high dimensionality of features in early layers of the ResNet architecture for  
593 high-dimensional visual inputs, we employ pooling operations on these feature maps. Thus,  
594 low-level visual features, such as orientations, are compromised. The consequent unfavorable  
595 outcome is a low predictive performance in V1. Further, since different subjects can focus on  
596 different parts of the stimulus, group-level models can also blur out the precise object orienta-  
597 tion information. This is particularly relevant for complex naturalistic stimuli such as movies.  
598 In the future, incorporating eye gaze data into these models can be an interesting exploration.  
599 Furthermore, due to computational constraints, the proposed model is only able to examine the  
600 effects of stimuli up to 20 seconds in the past. However, previous research with naturalistic  
601 stimuli has shown that some brain regions maintain memory of the order of minutes during  
602 naturalistic viewing (47). Existing evidence also suggests that neural activity is structured into  
603 semantically meaningful and coherent events (25). Capturing long-range context in encoding  
604 models can be a challenging, yet fruitful endeavour yielding potentially novel insights into  
605 memory formation.

606 There are also inherent differences between proposed neural network models and biological  
607 networks. DNNs fail to capture known properties of biological networks such as local recur-  
608 rence, however, they have been found to be useful for modelling neural activity across different  
609 sensory systems. At present, feed-forward DNNs trained on recognition tasks constitute the  
610 best predictors of sensory cortical activations in both humans and non-human primates (2).  
611 In light of this observation, a recent study proposed that very deep feed-forward only CNNs  
612 (for example, ResNet-50 as employed in this study for visual feature extraction) might im-  
613 plicitly be approximating ‘unrolled’ versions of recurrent computations of the ventral visual  
614 stream (48). Object recognition studies on non-human primates have also hinted at a functional  
615 correspondence between recurrence and deep non-linear transformations (49). Although the  
616 functional significance of intra-regional recurrent circuits in core object recognition is still un-  
617 der debate, mounting evidence suggests they may be subserving recognition under challenging  
618 conditions (49, 50). Thus, investigation of more neurobiologically plausible models of the cor-  
619 tex that innately model intra-regional recurrent computations should be explored in the future,  
620 especially in relation to their role in visual recognition.

## 621 **Concluding remarks**

622 Comprehensive descriptive models of the brain need comprehensive accounts of the stimulus.  
623 Using a novel group-level encoding framework, we showed that ‘reliable’ cortical responses  
624 to naturalistic stimuli can be accurately predicted across large areas of the cortex using multi-  
625 sensory information over longer time-scales. Since our models were trained on a large-scale,  
626 multi-subject and open-source dataset, we believe these results could provide an important point  
627 of reference against which encoding models for naturalistic stimuli can be assayed in the future.  
628 The continued interplay of artificial neural networks and neuroscience can pave the way for  
629 several exciting discoveries, bringing us one step closer to understanding the neural code of  
630 perception under realistic conditions.

## 631 **H2: Supplementary Materials**

632 Supplementary Text  
633 Figs. S1 to S9  
634 Tables S1 to S2  
635 References (51-53)

## 636 **References**

- 637 1. G. Varoquaux, R. A. Poldrack, Predictive models avoid excessive reductionism in cognitive  
638 neuroimaging. *Curr. Opin. Neurobiol.* **55**, 1–6 (2019).

- 639 2. D. L. Yamins, H. Hong, C. F. Cadieu, E. A. Solomon, D. Seibert, J. J. DiCarlo,  
640 Performance-optimized hierarchical models predict neural responses in higher visual cortex.  
641 *Proc. Natl. Acad. Sci. U.S.A.* **111**, 8619–8624 (2014).
- 642 3. K. N. Kay, T. Naselaris, R. J. Prenger, J. L. Gallant, Identifying natural images from human  
643 brain activity. *Nature* **452**, 352–355 (2008).
- 644 4. H. Wen, J. Shi, Y. Zhang, K. H. Lu, J. Cao, Z. Liu, Neural encoding and decoding with  
645 deep learning for dynamic natural vision. *Cerebral Cortex* **28**, 4136–4160 (2018).
- 646 5. U. Güçlü, M. A. van Gerven, Deep Neural Networks Reveal a Gradient in the Complexity  
647 of Neural Representations across the Ventral Stream. *J. Neurosci.* **35**, 10005–10014 (2015).
- 648 6. U. Güçlü, M. A. van Gerven, Increasingly complex representations of natural movies across  
649 the dorsal stream are shared between subjects. *NeuroImage* **145**, 329–336 (2017).
- 650 7. A. J. Kell, D. L. Yamins, E. N. Shook, S. V. Norman-Haignere, J. H. McDermott, A  
651 Task-Optimized Neural Network Replicates Human Auditory Behavior, Predicts Brain Responses,  
652 and Reveals a Cortical Processing Hierarchy. *Neuron* **98**, 630–644.e16 (2018).
- 653 8. A. J. King, G. A. Calvert, Multisensory integration: perceptual grouping by eye and ear.  
654 *Curr. Biol.* **11**, R322–325 (2001).
- 655 9. J. Driver, T. Noesselt, Multisensory interplay reveals crossmodal influences on ‘sensory-  
656 specific’ brain regions, neural responses, and judgments. *Neuron* **57**, 11–23 (2008).
- 657 10. J. Miller, Divided attention: evidence for coactivation with redundant signals. *Cogn Psychol* **14**,  
658 247–279 (1982).
- 659 11. S. Sonkusare, M. Breakspear, C. Guo, Naturalistic Stimuli in Neuroscience: Critically Ac-  
660 claimed. *Trends Cogn. Sci. (Regul. Ed.)* **23**, 699–714 (2019).
- 661 12. U. Hasson, Y. Nir, I. Levy, G. Fuhrmann, R. Malach, Intersubject synchronization of cortical  
662 activity during natural vision. *Science* **303**, 1634–1640 (2004).
- 663 13. M. Schönwiesner, R. J. Zatorre, Spectro-temporal modulation transfer function of single  
664 voxels in the human auditory cortex measured with high-resolution fmri. *Proceedings of*  
665 *the National Academy of Sciences* **106** **34**, 14611–6 (2009).
- 666 14. D. Schwartz, M. Toneva, L. Wehbe, Inducing brain-relevant bias in natural language pro-  
667 cessing models. *NeurIPS* (2019).
- 668 15. M. F. Glasser, S. N. Sotiropoulos, J. A. Wilson, T. S. Coalson, B. Fischl, J. L. Andersson,  
669 J. Xu, S. Jbabdi, M. Webster, J. R. Polimeni, D. C. Van Essen, M. Jenkinson, The minimal  
670 preprocessing pipelines for the Human Connectome Project. *Neuroimage* **80**, 105–124  
671 (2013).

- 672 16. A. T Vu, K. Jamison, M. F. Glasser, S. M. Smith, T. Coalson, S. Moeller, E. J. Auerbach,  
673 K. Ugurbil, E. Yacoub, Tradeoffs in pushing the spatial resolution of fMRI for the 7T  
674 Human Connectome Project. *Neuroimage* **154**, 23–32 (2017).
- 675 17. S. Hershey, S. Chaudhuri, D. P. W. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal,  
676 D. Platt, R. A. Saurous, B. Seybold, M. Slaney, R. J. Weiss, K. W. Wilson, CNN architec-  
677 tures for large-scale audio classification. *2017 IEEE International Conference on Acoustics,*  
678 *Speech and Signal Processing (ICASSP)* pp. 131–135 (2016).
- 679 18. A. S. Bregman, Auditory scene analysis. *MIT press* (2001).
- 680 19. T.-Y. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, S. J. Belongie, Feature pyramid  
681 networks for object detection. *2017 IEEE Conference on Computer Vision and Pattern*  
682 *Recognition (CVPR)* pp. 936–944 (2016).
- 683 20. K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition. *2016 IEEE*  
684 *Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 770–778 (2015).
- 685 21. J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, F.-F. Li, Imagenet: A large-scale hierarchical  
686 image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition* pp.  
687 248–255 (2009).
- 688 22. S. Abu-El-Haija, N. Kothari, J. Lee, A. Natsev, G. Toderici, B. Varadarajan, S. Vi-  
689 jayanarasimhan, Youtube-8M: A large-scale video classification benchmark. *ArXiv*  
690 **abs/1609.08675** (2016).
- 691 23. M. F. Glasser, T. S. Coalson, E. C. Robinson, C. D. Hacker, J. Harwell, E. Yacoub, K. Ugur-  
692 bil, J. Andersson, C. F. Beckmann, M. Jenkinson, S. M. Smith, D. C. Van Essen, A multi-  
693 modal parcellation of human cerebral cortex. *Nature* **536**, 171–178 (2016).
- 694 24. U. Hasson, E. Yang, I. Vallines, D. J. Heeger, N. Rubin, A hierarchy of temporal receptive  
695 windows in human cortex. *J. Neurosci.* **28**, 2539–2550 (2008).
- 696 25. C. Baldassano, J. Chen, A. Zadbood, J. W. Pillow, U. Hasson, K. A. Norman, Discovering  
697 Event Structure in Continuous Narrative Perception and Memory. *Neuron* **95**, 709–721  
698 (2017).
- 699 26. M. A. Goodale, A. D. Milner, Separate visual pathways for perception and action. *Trends*  
700 *Neurosci.* **15**, 20–25 (1992).
- 701 27. G. A. Calvert, Crossmodal processing in the human brain: insights from functional neu-  
702 roimaging studies. *Cereb. Cortex* **11**, 1110–1123 (2001).
- 703 28. T. Raij, K. Uutela, R. Hari, Audiovisual integration of letters in the human brain. *Neuron*  
704 **28**, 617–625 (2000).



- 705 29. M. S. Beauchamp, Statistical criteria in fMRI studies of multisensory integration. *Neuroinformatics* **3**, 93–113 (2005).  
706
- 707 30. M. S. Beauchamp, B. D. Argall, J. Bodurka, J. H. Duyn, A. Martin, Unraveling multisensory integration: patchy organization within human STS multisensory cortex. *Nat. Neurosci.* **7**, 1190–1192 (2004).  
708  
709
- 710 31. G. A. Calvert, E. T. Bullmore, M. J. Brammer, R. Campbell, S. C. Williams, P. K. McGuire, P. W. Woodruff, S. D. Iversen, A. S. David, Activation of auditory cortex during silent lipreading. *Science* **276**, 593–596 (1997).  
711  
712
- 713 32. N. Kanwisher, G. Yovel, The fusiform face area: a cortical region specialized for the perception of faces. *Philos. Trans. R. Soc. Lond., B, Biol. Sci.* **361**, 2109–2128 (2006).  
714
- 715 33. I. Tavor, M. Yablonski, A. Mezer, S. Rom, Y. Assaf, G. Yovel, Separate parts of occipitotemporal white matter fibers are associated with recognition of faces and places. *Neuroimage* **86**, 123–130 (2014).  
716  
717
- 718 34. S. Nasr, N. Liu, K. J. Devaney, X. Yue, R. Rajimehr, L. G. Ungerleider, R. B. Tootell, Scene-selective cortical regions in human and nonhuman primates. *J. Neurosci.* **31**, 13771–13785 (2011).  
719  
720
- 721 35. J. A. Frost, J. R. Binder, J. A. Springer, T. A. Hammeke, P. S. Bellgowan, S. M. Rao, R. W. Cox, Language processing is strongly left lateralized in both sexes. Evidence from functional MRI. *Brain* **122** ( Pt 2), 199–208 (1999).  
722  
723
- 724 36. P. Belin, R. J. Zatorre, P. Lafaille, P. Ahad, B. Pike, Voice-selective areas in human auditory cortex. *Nature* **403**, 309–312 (2000).  
725
- 726 37. T. Yarkoni, R. A. Poldrack, T. E. Nichols, D. C. Van Essen, T. D. Wager, Large-scale automated synthesis of human functional neuroimaging data. *Nat. Methods* **8**, 665–670 (2011).  
727  
728
- 729 38. A. G. Huth, S. Nishimoto, A. Vu, J. Gallant, A continuous semantic space describes the representation of thousands of object and action categories across the human brain. *Neuron* **76**, 1210–1224 (2012).  
730  
731
- 732 39. Y. Cao, C. Summerfield, H. Park, B. L. Giordano, C. Kayser, Causal Inference in the Multisensory Brain. *Neuron* **102**, 1076–1087 (2019).  
733
- 734 40. S. M. Wilson, I. Molnar-Szakacs, M. Iacoboni, Beyond superior temporal cortex: intersubject correlations in narrative speech comprehension. *Cereb. Cortex* **18**, 230–242 (2008).  
735
- 736 41. I. P. Jääskeläinen, K. Koskentalo, M. H. Balk, T. Autti, J. Kauramäki, C. Pomren, M. Sams, Inter-subject synchronization of prefrontal cortex hemodynamic activity during natural viewing. *The Open Neuroimaging Journal* **2**, 14 - 19 (2008).  
737  
738

- 739 42. S. Jain, A. Huth, Incorporating context into language encoding models for fMRI. *NIPS*  
740 (2018).
- 741 43. F. H. Sinz, A. S. Ecker, P. G. Fahey, E. Y. Walker, E. Cobos, E. Froudarakis, D. Yatsenko,  
742 X. Pitkow, J. Reimer, A. S. Tolias, Stimulus domain transfer in recurrent models for large  
743 scale cortical population prediction on video. *bioRxiv* (2018).
- 744 44. J. Schultz, K. S. Pilz, Natural facial motion enhances cortical responses to faces. *Exp Brain*  
745 *Res* **194**, 465–475 (2009).
- 746 45. P. Bashivan, K. Kar, J. DiCarlo, Neural population control via deep image synthesis. *Sci-*  
747 *ence* **364** (2019).
- 748 46. J. Chen, U. Hasson, C. J. Honey, Processing Timescales as an Organizing Principle for  
749 Primate Cortex. *Neuron* **88**, 244–246 (2015).
- 750 47. U. Hasson, J. Chen, C. J. Honey, Hierarchical process memory: memory as an integral  
751 component of information processing. *Trends Cogn. Sci. (Regul. Ed.)* **19**, 304–313 (2015).
- 752 48. Q. Liao, T. A. Poggio, Bridging the gaps between residual learning, recurrent neural net-  
753 works and visual cortex. *ArXiv* **abs/1604.03640** (2016).
- 754 49. K. Kar, J. Kubilius, K. Schmidt, E. B. Issa, J. J. DiCarlo, Evidence that recurrent circuits are  
755 critical to the ventral stream’s execution of core object recognition behavior. *Nat. Neurosci.*  
756 **22**, 974–983 (2019).
- 757 50. D. Wyatte, D. J. Jilk, R. C. O’Reilly, Early recurrent feedback facilitates visual object  
758 recognition under challenging conditions. *Front Psychol* **5**, 674 (2014).
- 759 51. Y. Benjamini, Y. Hochberg, Controlling the false discovery rate: a practical and powerful  
760 approach to multiple testing. *J. R. Stat. Soc. B.* **57**, 289-300 (1995).
- 761 52. A. Nagrani, J. S. Chung, W. Xie, A. Zisserman, Voxceleb: Large-scale speaker verification  
762 in the wild. *Comput. Speech Lang.* **60** (2020).
- 763 53. K. J. Piczak, Esc: Dataset for environmental sound classification. *MM* (2015).

## 764 **Acknowledgments**

765 This work was supported by NIH grants R01LM012719 (MS), R01AG053949 (MS), R21NS10463401  
766 (AK), R01NS10264601A1 (AK), the NSF NeuroNex grant 1707312 (MS), the NSF CAREER  
767 1748377 grant (MS) and Anna-Maria and Stephen Kellen Foundation Junior Faculty Fellowship  
768 (AK).

## 769 **Data and Software availability**

770 All experiments in this study are based on the Human Connectome Project movie-watching  
771 database. The dataset is publicly available for download through the ConnectomeDB software  
772 (<https://db.humanconnectome.org/>). Throughout this study, we utilized 7T fMRI data from the  
773 ‘Movie Task fMRI 1.6mm/59k FIX-Denoised’ package within HCP. The network implemen-  
774 tation, analysis codes as well as trained model weights will be made available on the project  
775 Github page.

## 776 **Supplementary materials**

### 777 **HCP Movies**

778 Table S1 summarizes the HCP movie-watching dataset split used for training and evaluating all models.

**Table S1.** HCP dataset split

Movie	Split	Stimulus-response pairs per subject
7T_MOVIE1_CC1.v2	Training/Validation	652
7T_MOVIE2_HO1.v2	Training/Validation	716
7T_MOVIE3_CC2.v2	Training/Validation	669
7T_MOVIE4_HO2.v2	Testing	699

779

### 780 **Region of Interest (ROI) selection**

781 ROIs were selected for each analysis based on the descriptions provided in the neuroanatomical  
782 supplementary results of the HCP MMP parcellation (23) and an extensive literature review.  
783 For Figure 2 in the main text and Figure S8, ROIs were thus assigned to groups 1-5 according  
784 to Table S2).

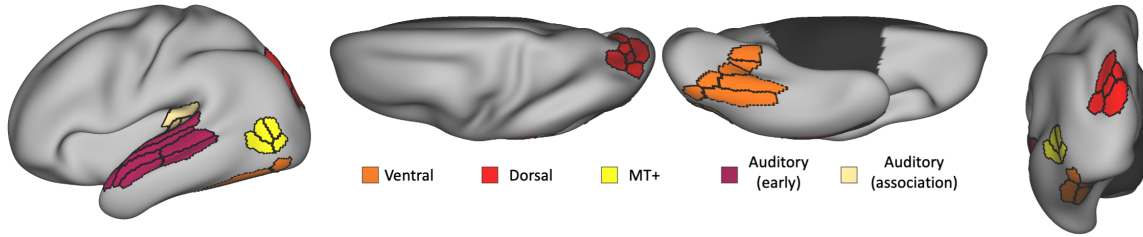
**Table S2.** ROI categorization

Group	ROIs
1. Auditory	A1, LBelt, PBelt, MBelt, RI, STSda, STSva, A4, A5, TA2
2. Visual	V1, V2, V3, V3A, V3B, V3CD, V4, V4t, V6, V6A, V7, V8, DVT, LO1-3, PIT, FFC, VMV1-3, IPS1, MT, VVC
3. Multi-sensory + sensory bridges	STSdp, STSvp, STGa, STV, TPOJ1-3
4. Language	55b, SFL, PSL, 44, 45
5. Frontal	IFSa, IFSp, IFJa, IFJp, FEF

785 Dorsal and ventral visual stream ROIs as well as early and association auditory cortex ROIs  
786 in Figure 4 (main text) were derived from the explicit stream segregation and categorization  
787 described in the HCP MMP parcellation (23) and are defined here for quick reference.

- 788 • Dorsal: V3A, V3B, V6, V6A, V7, IPS1
- 789 • Ventral: V8, VVC, PIT, FFC, VMV1-3
- 790 • MT+: MT, MST, V4t, FST
- 791 • Early auditory: A1, PBelt, MBelt, RBelt, RI
- 792 • Association auditory: A4, A5, TA2, STGa, STSdp, STSda, STSvp, STSva

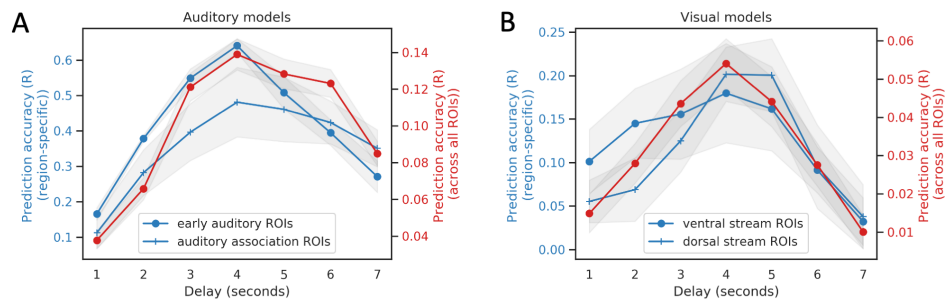
793 All ROIs are shown in Figure S1



**Fig. S1.** Group segregation from the HCP MMP parcellation.

## 794 Estimating BOLD response delay

795 BOLD response delay was estimated using ROI-level encoding models due to their faster iter-  
 796 ation times in comparison to voxel-wise encoding. The input to these models was the prepro-  
 797 cessed stimuli as described for voxel-wise encoding with the same train-validation-test split,  
 798 and the output was the evoked ROI-level fMRI response at different lags (1-7 seconds) from  
 799 the stimulus. Thus, the output is a 360-D vector corresponding to the mean fMRI response in  
 800 each ROI of the HCP MMP parcellation. The feature extractors were identical to those in the  
 801 proposed voxel-wise auditory and visual models. However, instead of a convolutional response  
 802 model, here, the response model comprised two fully connected layers with output dimensions  
 803 of 512 and 360 with an exponential linear unit and linear activation respectively. All models  
 804 were trained for 20 epochs with a batch size of 4 and a learning rate of  $1e-4$ . Validation curves  
 were monitored to ensure convergence. Prediction accuracy of each model was computed as



**Fig. S2.** ROI-based encoding performance for estimating delay. (A) depicts the estimated mean and standard error of the prediction accuracy (R) across various delays (1-7s) within the early auditory and association auditory group (blue) as well as across all ROIs (red), as obtained using the single epoch (1s) auditory model. (B) depicts the estimated mean and standard error of the prediction accuracy (R) for various delays (1-7s) within the primary and dorsal visual streams (blue) as well as across all ROIs (red), as obtained using the single frame visual model. Gray regions depict the standard error in estimating mean across ROIs within each group. ROI categorization is described in the sub-section on ROI selection.

805 the mean Pearson's correlation coefficient between the predicted and measured response across  
 806 all ROIs, in the held-out movie dataset. Based on Figure S2, we estimated a response delay  
 807 of 4 seconds, as this lag yielded the maximum prediction accuracy across all ROIs for both  
 808

809 auditory and visual ROI-level models. Further, even while restricting the prediction accuracy  
810 (R) to ROIs within different cortical areas (such as the early/association auditory areas or the  
811 dorsal/ventral visual stream), the optimal lag was consistently 4 seconds, suggesting that the  
812 difference in performance of 1-sec and 20-sec models in these regions (Figure 4) is not largely  
813 driven by differences in the hemodynamic response function (HRF).

## 814 **Defining the stimulus-driven or “synchronous” cortex**

815 We isolated voxels involved in stimulus-driven processing, termed “synchronous” or “stimulus-  
816 driven” voxels, by computing mean inter-group correlations over all training movies. Inter-  
817 group correlations were computed by splitting the entire group of subjects into two halves and  
818 computing correlations between the mean response time-course of each half (comprising 79  
819 subjects) at every voxel. We employed a liberal threshold of 0.15 for this correlation value.  
820 Thus, the mask of “stimulus-driven” voxels included those voxels that achieved an inter-group  
821 correlation of 0.15 or above. We computed mean quantitative metrics over this mask in Fig-  
822 ure 3E (main text) to compare different models.

## 823 **Model architectures and implementation**

824 The base feature extraction networks and convolutional response model in Figure 1 had the archi-  
825 tecture as detailed in Figure S3. The feature extraction networks are reminiscent of the fea-  
826 ture pyramid network, which has shown significant improvements as a generic feature extractor  
827 across various applications. These networks comprise a parallel top-down pathway with lateral  
828 connections which grants them the ability to characterize both “what” and “where” in cluttered  
829 scenes, thereby enhancing object detection. We note that similar models with top-down and  
830 skip connections have been popular in vision research, since they can enrich low-level features  
831 with high-level semantics. The output of the feature extractor is fed into the convolutional re-  
832 sponse model to predict the evoked fMRI activation. This enables us to train both components  
833 of the network simultaneously in an end-to-end manner. Since the output response is differen-  
834 tiable with respect to network weights, the weights are adjusted via a first-order gradient-based  
835 optimization method to minimize the *mean squared error* between the predicted and target ac-  
836 tivation values across the entire brain.

837 For ResNet-50, we use activations of the last residual block of each stage, namely, res2, res3,  
838 res4 and res5 to construct our stimulus descriptions  $\mathbf{s}$ . From the VGG-ish network, we use  
839 the activations of each convolutional block, namely, conv2, conv3, conv4 and the penulti-  
840 mate dense layer fc2<sup>1</sup>. The first three set of activations are refined through a top-down path to  
841 enhance their semantic content, while the last activation is concatenated into  $\mathbf{s}$  directly (res4 ac-  
842 tivations are vectorized using global average pool). The top-down path comprises three feature

---

<sup>1</sup>Pre-trained tensorflow/keras models for the visual and auditory backbone were available at <https://keras.io/applications> <https://github.com/tensorflow/models/tree/master/research/audioset/vggish> respectively

843 maps at different resolutions with an up-sampling factor of 2 successively from the deepest layer  
844 of the bottom-up path. Each such feature map comprising 256/128 channels (in visual/auditory  
845 models respectively) is merged with the corresponding feature map in the bottom-up path (re-  
846 duced to 256/128 channels by 1x1 convolutions) by element-wise addition. Subsequently, the  
847 feature map at each resolution is collapsed into a 256/128 dimensional feature vector through  
848 a global average pool operation and concatenated into  $\mathbf{s}$ , leading to a 1024-D and 512-D fea-  
849 ture representation for the visual and auditory stimuli respectively. The aggregated features are  
850 then passed onto a CNN comprising the following feedforward computations: a fully connected  
851 layer to map the features into a vector space which is reshaped into a 1024-channel cuboid of  
852 size 6x7x6 followed by four 3x3x3 transposed convolutions (conv.T) with a stride of 2 and  
853 exponential linear unit activation function to up-sample the latter. Each convolution reduces  
854 the channel count by half with the exception of the last convolution which outputs the single-  
855 channel predicted fMRI response.

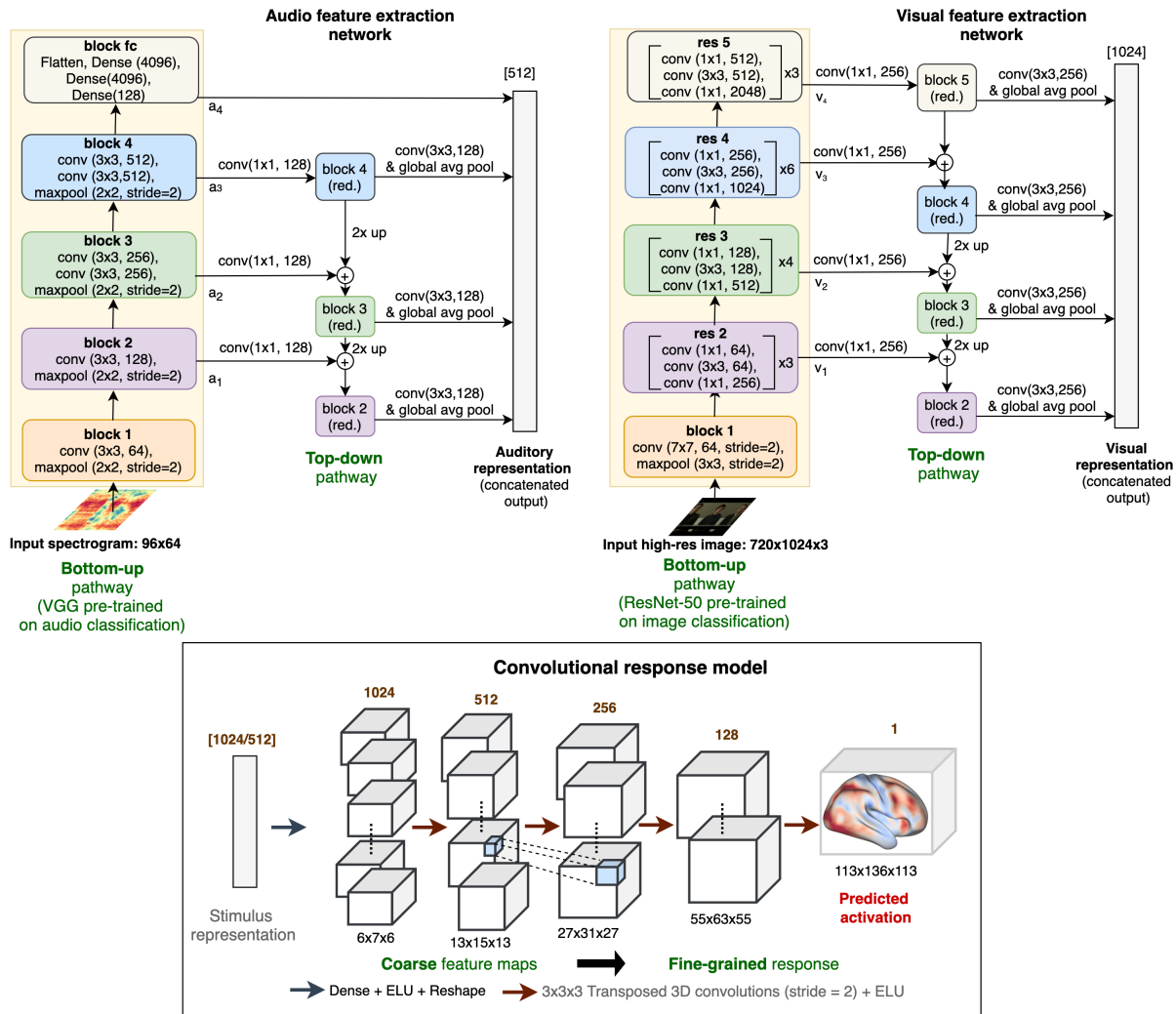
856 The 20-second models additionally comprised an LSTM layer to model the temporal propaga-  
857 tion of features across the contiguous sequence of input frames and/or spectrograms. The LSTM  
858 module has driven success across varied sequence modeling tasks due to its ability to efficiently  
859 regulate the flow of information across cells through gating. The memory cell in LSTM is modu-  
860 lated by three gates, namely, the input, forget and output gates. We note that the LSTM layer did  
861 not change the dimensionality of the input features so that equitable comparisons can be made  
862 against 1-sec models. The Audiovisual-1sec model concatenated features obtained from the  
863 base visual (1024-D) and audio (512-D) feature extraction networks, reduced their combined  
864 dimensionality to the higher value among the two (1024-D) by passing through a bottleneck  
865 dense layer followed by the same convolutional response model. The Audiovisual-20sec model  
866 additionally incorporated modal-specific LSTM networks prior to feature concatenation.

#### 867 *Implementation:*

868 We note that all 6 models have roughly the same order of trainable parameters in the range of  
869 242M-362M. All parameters were optimized using Adam with a learning rate of  $1e-4$ . Audi-  
870 tory and visual models were trained for 50 epochs with unit batch size. The stimulus as well  
871 as subject whose fMRI response is used as the target in the loss (“mean squared error”) are  
872 randomly sampled over each step of the training but kept consistent across models. We found  
873 this method to work better than using the group-averaged response as target, presumably be-  
874 cause this sampling provides information about both the cross-subject mean and the variance  
875 of response. Given the noise characteristics at each voxel, we hypothesize that this enables the  
876 model to focus on regions that can be well predicted with the given stimulus. Validation curves  
877 were monitored for all models to ensure convergence.

#### 878 **Regularized linear regression: WordNet features**

879 Another popular approach in voxel-wise forward encoding beyond primary sensory cortices is  
880 the semantic category encoding model that is based on high-level semantic features (38). This



**Fig. S3.** Implementation details for the audio (top left) and visual (top right) feature extraction networks as well as the convolutional response model (bottom). All layers and blocks outside the yellow rectangle (bottom-up pathway) are trained from scratch. The blocks inside the yellow rectangular window are initialized with networks pre-trained on image or sound recognition. Further, ResNet-50 is frozen during the training of all encoding models, whereas VGG is fine-tuned. The sequence of operations within each block are defined from top to bottom, while the number of repetitions for each sequence within the block are indicated with the multiplicative symbol on the right.

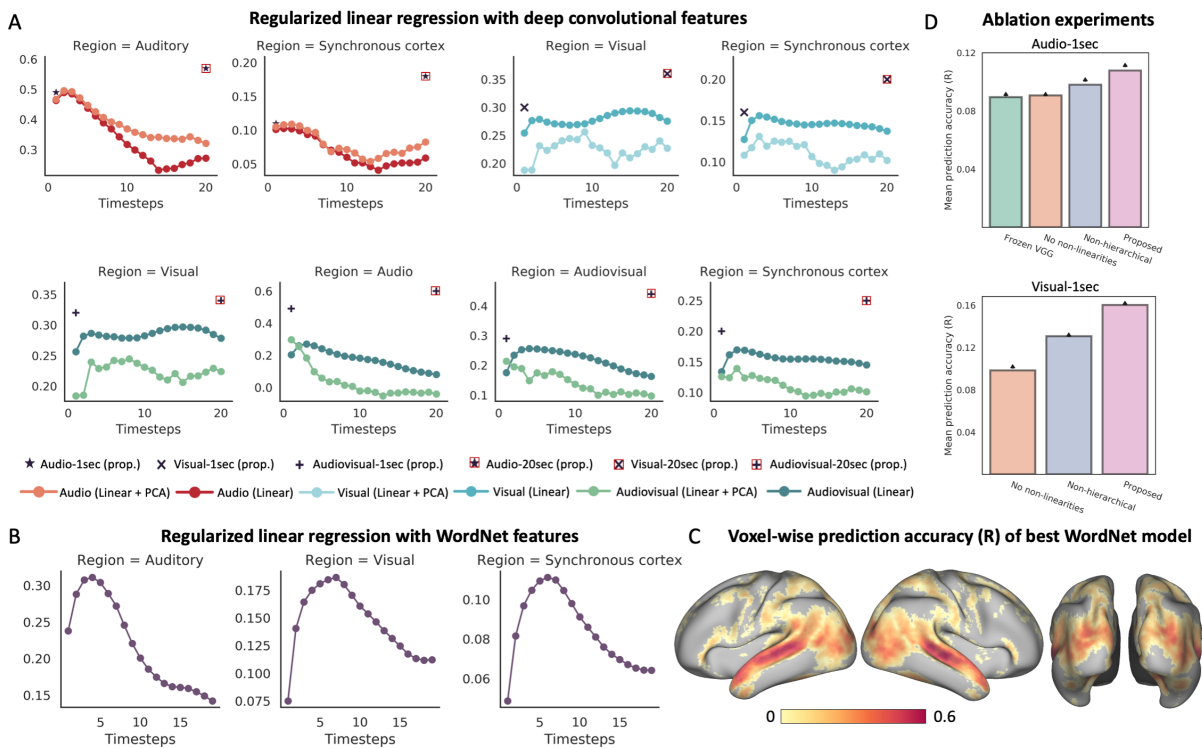


881 approach relies on labels that indicate the presence of semantic object and action categories in  
882 each movie frame. In this analysis, we employed WordNet labels that were provided as part  
883 of the HCP movie-watching data pipeline. The semantic labels were manually assigned by the  
884 Gallant lab team using the WordNet semantic taxonomy and subsequently converted to Word-  
885 Net synsets to build an 859-D semantic representational space (corresponding to 859 WordNet  
886 synset names). Following (38), we fitted  $l_2$  regularized linear regression models (known as  
887 ridge regression) to find weights corresponding to different input features for every voxel. The  
888 regularization parameter,  $\alpha$  was optimized independently for each voxel by testing among 10  
889 log-space values in [1, 1000]. The optimal alpha is obtained by averaging across 15 boot-  
890 strapped held-out sets. In addition to fitting models with WordNet features extracted 4s prior to  
891 the measured neural response, we developed longer timescale linear models by concatenating  
892 the WordNet features extracted for each second (as described above) over T-second windows  
893 with T ranging from 1 to 20 seconds and presented these aggregated features to the bootstrapped  
894 regularized regression model. Figure S4 (B) demonstrates the performance of WordNet models  
895 across different groups of regions as a function of T, and (C) depicts the voxel-level prediction  
896 accuracy (R) of the best performing WordNet model that stacks features from 4-12s (at an inter-  
897 val of 1s) prior to the encoded cortical response. While simple and interpretable, the WordNet  
898 models clearly under-perform in terms of prediction accuracy (R) in comparison to the models  
899 proposed in the present study.

## 900 **Regularized linear regression: deep convolutional features**

901 We also trained group-level encoding models using a linear response model since this consti-  
902 tutes the dominant state-of-the-art approach to neural encoding (5, 4, 7). To enable a fair com-  
903 parison against the proposed 1-sec uni-modal models, we extract hierarchical features from the  
904 same layers of the ResNet-50 and VGG-ish architectures as employed by the proposed mod-  
905 els. The only difference here is the lack of a top-down pathway (since it is not a part of the  
906 pre-trained network but is trained with random initialization on the neural response prediction  
907 task), which prevents the refinement of coarse feature maps before aggregation. Pooling the  
908 outputs of different layers channel-wise using the global average pooling operation (namely  
909  $\{v_1, v_2, v_3, v_4\}$  for the visual model and  $\{a_1, a_2, a_3, a_4\}$  for the audio model in Figure S3) leaves  
910 us with and 1024 and 3840 features to present to the auditory and visual models, respectively.  
911 Further, to compare against the longer-duration 20-sec models, we adopted two approaches: (1)  
912 we simply concatenated the stimulus features extracted for each second (as described above)  
913 over T-second windows with T ranging from 1 to 20 seconds and presented these aggregated  
914 features to the linear response model; alternatively, (2) we reduced the dimensionality of the  
915 aggregated features to a fixed length (set to 128) as in (1) using principal component analy-  
916 sis run on the training data. We added this comparison to rule out the fact that the temporal  
917 trend in performance of linear models is simply driven by a higher-dimensional feature space.  
918 We note that even after dimensionality reduction, the components retained at least 80% of the  
919 explained variance in all cases. Audio-visual encodings with linear response models were ob-

920 tained similarly by simply fusing the respective audio and visual hierarchical features through  
 921 concatenation before linear regression. We apply  $l_2$  regularization on the regression coeffi-  
 922 cients and adjust the optimal strength of this penalty through cross-validation on the training  
 923 data using log-spaced values in  $\{1e-14, 1e14\}$  for each model. We report performance of the  
 924 best models in Figure S4(A). Note that unlike the WordNet models, we found that optimizing  
 925 a single regularization penalty  $\alpha$  common across all voxels outperformed independent voxel-  
 926 wise fitting with bootstrap in this case. Thus, we only present the results for the former. We  
 927 note here that the convolutional response model in our proposed approach (instead of a fully-  
 928 connected approach) allowed us to keep the learnable parameters manageable, facilitating joint  
 929 optimization/fine-tuning of the feature extractor and response models. The consistently superior  
 930 performance of the proposed models against linear regression based approaches strongly sug-  
 931 gests that there is merit in end-to-end learning for encoding responses to dynamic, multi-sensory  
 932 stimuli.



**Fig. S4.** Performance of linear response models with (A) deep convolutional features and (B) semantically rich WordNet features. The x-axis depicts the length of the windows (in seconds) over which the stimulus features are concatenated and y-axis shows the mean Pearson's correlation coefficient between the predicted and measured responses across the stimulus-driven voxels. (C) shows the cortical map of the prediction accuracy (R) for the best WordNet model. (D) shows results of the ablation study and highlights the importance of different components of the proposed model architecture.

### 933 **Ablation study**

934 To determine the influence of different architectural components on prediction performance of  
935 the proposed models, we performed an ablation study to investigate the individual contributions  
936 of (i) non-linearities in the response model, (ii) hierarchical (multi-scale) feature maps and (iii)  
937 fine-tuning audio sub-network (VGG). We selectively removed each of these components from  
938 the respective 1sec models and compared the resulting performance against the proposed model  
939 that employs all (i)-(iii) components. There are several interesting observations to make from  
940 this ablation analysis (Figure S4D). (i) First, we find that encoding models with a frozen VGG  
941 network that is not updated during training incur a loss in performance compared to the pro-  
942 posed model where VGG layers are trainable during neural response prediction. This clearly  
943 demonstrates the advantages of altering these pre-trained models and suggests that fine-tuning  
944 is both feasible and beneficial in improving neural response prediction. (ii) Next, we find that  
945 prediction performance deteriorates after removing the non-linearities in both the Audio-1sec  
946 and Visual-1sec models. In the context of the Visual-1sec model with a frozen pre-trained back-  
947 bone (ResNet-50) and coupled with (i), this observation further highlights that it is possible to  
948 develop models of human sensory processing that are quantitatively more precise in matching  
949 brain activity than task-driven neural networks. (iii) Finally, we assessed the benefit of using  
950 hierarchical feature maps over selecting the single best-performing layer for each model (audio  
951 or visual) based on cross-validation. For both audio and visual models, we find that features  
952 from the last layer (i.e.,  $a_4$  and  $v_4$ , respectively) yield the highest mean prediction accuracy (R)  
953 across the synchronous cortex. However, although the convolutional response model architec-  
954 ture is common across these encoding models, it is important to note that this analysis is still  
955 plagued by confounds such as the different dimensionality of feature spaces across different  
956 layers that feed into the response model. The best performing single-layer encoding model,  
957 however, still performs worse than the hierarchical approach.

### 958 **Computing significance estimates**

959 The statistical significance of individual voxel predictions (Figure 3) was computed as the p-  
960 value of the obtained sample *correlation coefficient* for the null hypothesis of uncorrelation (i.e.,  
961 true correlation coefficient is zero) under the assumptions of a bivariate normal distribution. We  
962 employed the false-discovery procedure of Benjamini & Hochberg (1995) (51) to control for  
963 multiple comparisons under assumptions of dependence. For statistical comparison of model  
964 performance within each group of regions in Figure 2 (main text), we performed paired t-test on  
965 ROI-level average performance metrics and corrected for multiple comparisons among models  
966 (Bonferroni).

### 967 **Sensory-sensitivity index**

968 Distorting the input to the audio-visual model at test time allows us to interrogate the sensory-  
969 sensitivity of different brain regions. We developed a sensory-sensitivity index of each ROI

970 based upon predictive performance of the model with distorted inputs, as shown in Figure 5. Let  
971  $SV_r$  and  $SA_r$  denote the mean prediction accuracy of the model in region  $r$  after shuffling (tem-  
972 porally) the input order of the visual and auditory stimuli, respectively. The sensory-sensitivity  
973 index for region  $r$  is then defined as  $s_r = \frac{SA_r - SV_r}{SA_r + SV_r}$ . Note that positive values of this index indi-  
974 cate that region  $r$  incurs a greater loss in predictivity upon distortion of visual information than  
975 auditory information, suggesting a higher visual sensitivity for this voxel. Similarly, negative  
976 values signal towards a higher auditory-sensitivity.

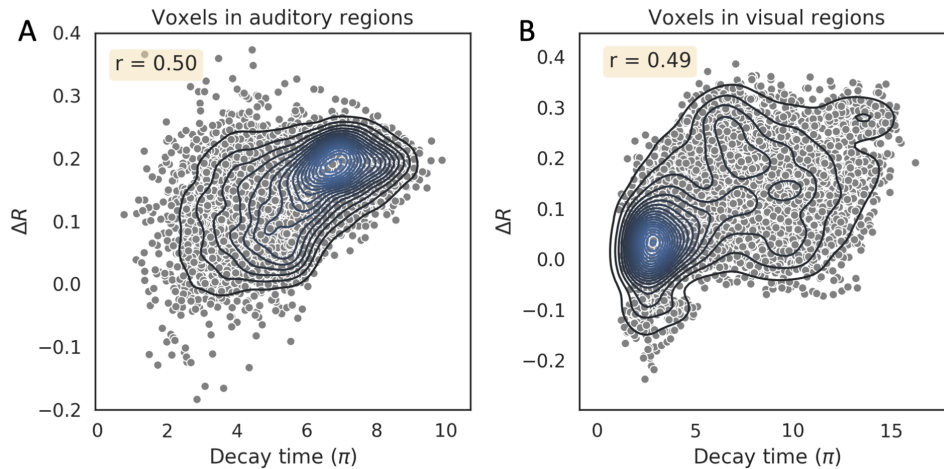
## 977 **Stimuli for synthetic contrasts**

978 Synthetic contrasts were generated to study the generalization of our models to new experi-  
979 mental paradigms (Figure 6). We focus on predicting task-based contrasts for three semantic  
980 categories, namely, *faces*, *places* and *speech*, since these are the most well-studied categories in  
981 the context of their distinct functional signatures. The stimuli for visual contrasts were derived  
982 from the HCP Working Memory paradigm, which combines category specific representation  
983 tasks (including faces and places) and working memory tasks. After excluding gray-scale im-  
984 ages, we were left with 102, 77, 97 and 103 images for the categories of faces, places, body  
985 parts and tools, respectively. Since these are static image without any dynamic content, we  
986 employed the Visual-1sec model to derive the visual contrasts (Figure 6(C),(D)).

987 Stimuli for the speech and non-speech contrast were extracted from large popular datasets for  
988 these categories. Speech stimuli were extracted from a human speech-utterance dataset com-  
989 prising short audio clips of interviews recorded on YouTube (52). Non-speech stimuli were  
990 extracted from another large dataset comprising short clips of environmental sounds (53). We  
991 randomly extracted  $\sim 100$  minutes of audio waveforms from these datasets for both categories.  
992 The stimuli were processed for mel-spectrogram extraction in the same manner as the HCP  
993 audio-visual movies. Since the non-speech stimuli only comprised contiguous clips of roughly  
994 3 – 5 second duration, we employed the Audio-1sec model to obtain the speech contrast (Fig-  
995 ure 6(B)).

## 996 **Performance improvement and autocorrelation decay**

997 In the past, processing timescales in the brain have been probed using several different means (46).  
998 In one of the proposed approaches, the decay time of temporal autocorrelation is used as a proxy  
999 measure to understand the variation in processing timescales across different brain regions.  
1000 With this approach, it was shown that decay times increased progressively along the temporal  
1001 hierarchy. Following this line of work, we estimated the autocorrelation decay time constant  
1002 ( $\pi$ ) for each voxel by fitting an exponential,  $A \exp\{-t/\pi\}$ , to the autocorrelation function (au-  
1003 tocorrelation computed at different lags). The exponential model was first independently fit for  
1004 each movie run and each voxel and the estimated  $\pi$  were subsequently averaged across runs to  
1005 obtain one decay time constant per voxel. Here, we were primarily interested in understanding  
1006 whether there is any relationship between the performance improvement of the 20-sec model



**Fig. S5.** Performance boost of the 20-sec model over 1-sec model is higher in voxels with longer autocorrelation decay times. (A) & (B) depict the performance improvement ( $\Delta R$ ) against decay time constants for voxels associated with auditory and visual regions, respectively (Table S2). The  $r$  value indicates the Pearson's correlation coefficient between the two quantities. Each dot in the scatterplot represents an individual voxel. Bivariate kernel density estimates are overlaid on top of the scatterplot as contours to depict the probability distribution of observations.

1007 over 1-sec model,  $\Delta R$ , computed as the difference between the prediction accuracies of the  
1008 Audiovisual-20sec and Audiovisual-1sec at every voxel, and the temporal autocorrelation prop-  
1009 erties of that voxel. We hypothesized that in voxels with longer processing timescales, the au-  
1010 tocorrelation would persist for longer durations (resulting in larger  $\pi$ ) and the longer timescale  
1011 model (20-sec) would yield more substantive improvement over the 1-sec model. As shown in  
1012 Figure S5, we observed a significantly positive correlation between performance improvement  
1013 and the autocorrelation decay time constant ( $r = 0.49$  and  $0.50$  across voxels in auditory and  
1014 visual regions as defined in Table S2), in line with our hypothesis. This suggests that the benefit  
1015 of employing the 20-sec model, as quantified in terms of performance improvement, is indeed  
1016 more remarkable in regions with longer processing timescales.

## 1017 Surface visualization

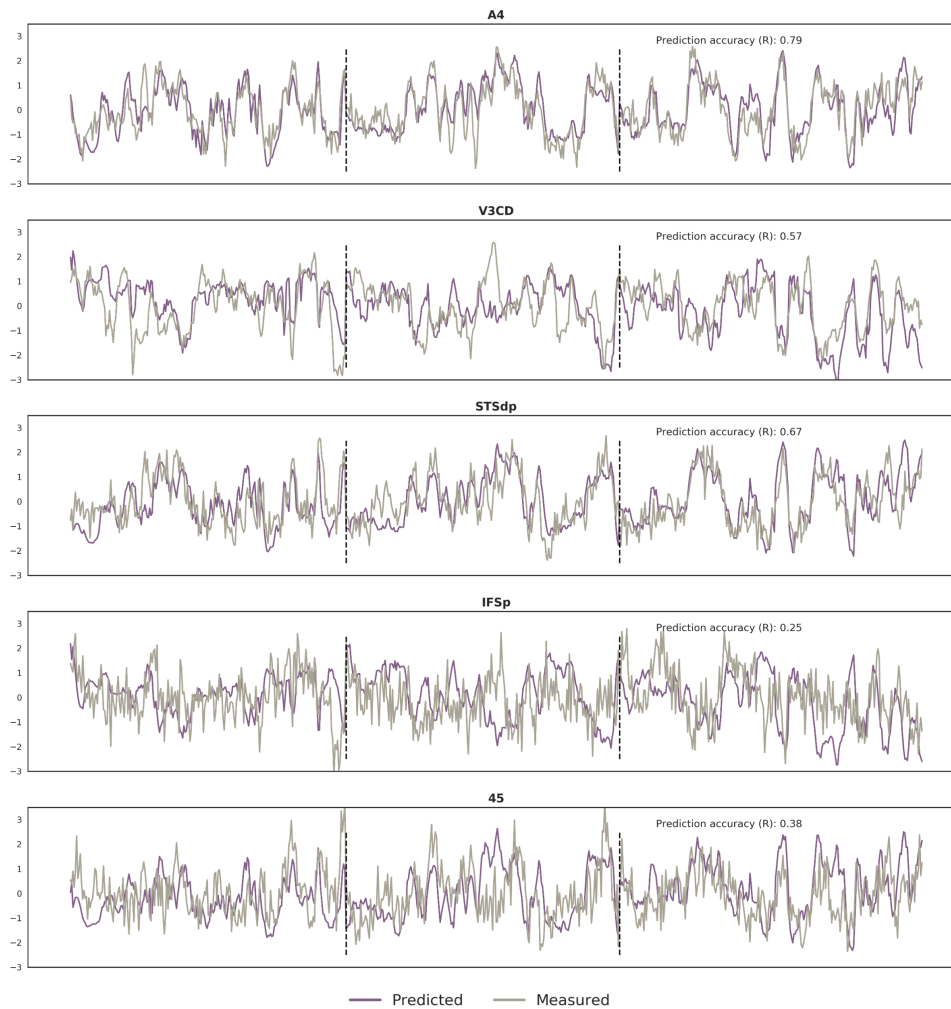
1018 All input fMRI data, as well as response predictions in this study are volume based. In order  
1019 to be consistent with prior research on encoding models that employ surface visualizations, we  
1020 created surface versions of volumetric predictability and synthetic contrast maps, as shown in  
1021 Figures 3, 5 and 6. We employed the 3D trilinear mapping method from connectome workbench  
1022 that computes the result on each vertex based on linear interpolation from voxels on each side  
1023 of the vertex<sup>2</sup>. However, since volume to surface mappings are an approximation, we only  
1024 employ this conversion for visualizations. All reported metrics are computed on volumes only

<sup>2</sup><https://www.humanconnectome.org/software/workbench-command>

1025 on a per-voxel basis.

## 1026 Qualitative analysis

1027 To gain qualitative insights into the predictions of the most accurate model (Audiovisual-20sec)  
1028 on the held-out movie, we plot the predicted as well as measured response time-series of the  
1029 voxel with ‘median’ prediction accuracy (R) in the best performing ROI of each group (Fig-  
1030 ure S6). The latter corresponds to A4, V3CD, STSdp, IFSp and Area 45 for the auditory, visual,  
multi-sensory, frontal and language groups respectively.



**Fig. S6.** Predicted and measured response time-series of the ‘median’ predictive accuracy (R) voxel across ROIs of different functional groups. Vertical dashed lines mark the boundary of clip segments in the held-out movie.

1031

### 1032 **Group-level prediction accuracy: held-out set**

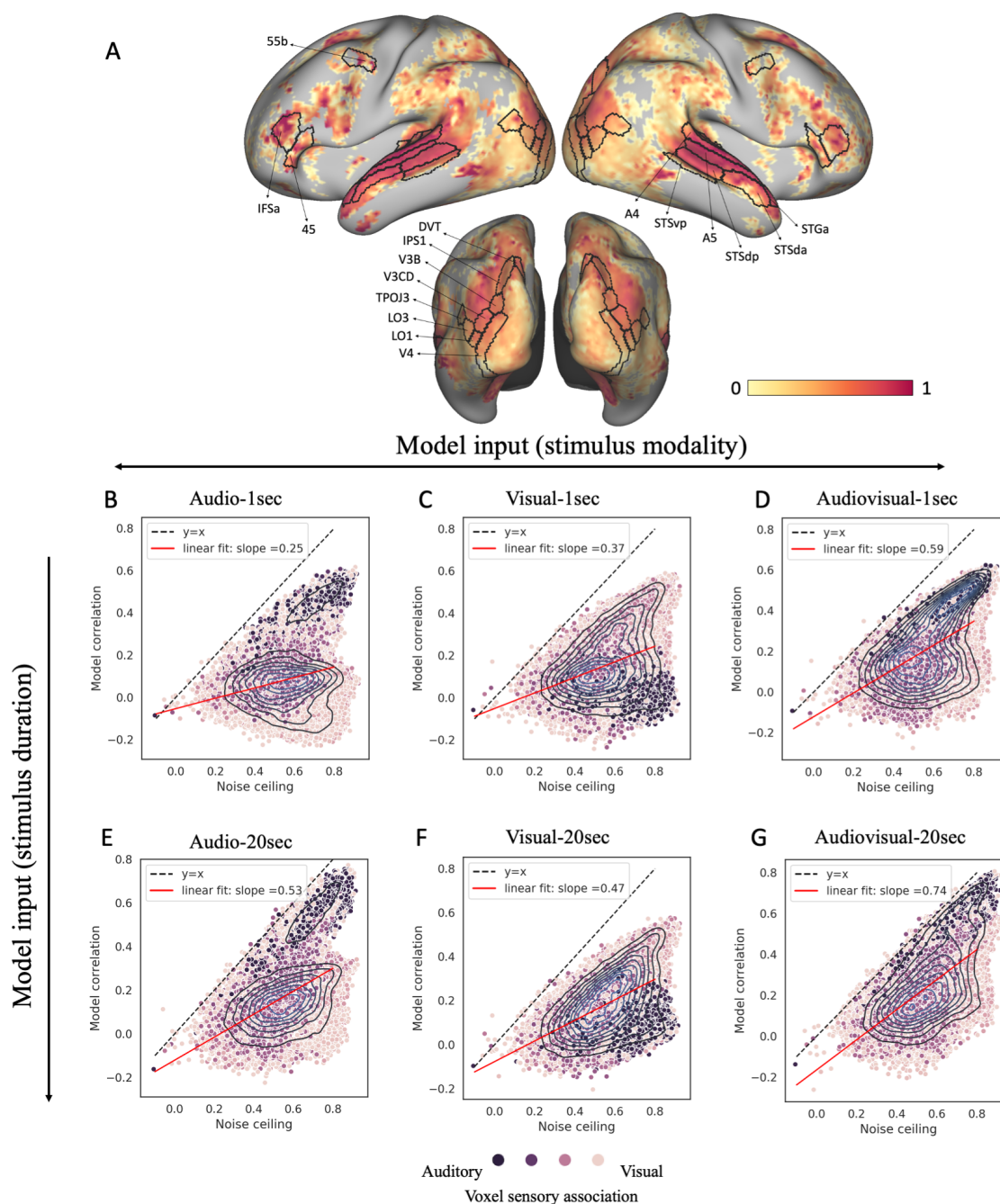
1033 To test the generality of the models, we further compared model predictions against the group-  
1034 averaged response of a held-out group within HCP comprising 20 novel subjects distinct from  
1035 the 158 individuals used in the training set, on the same independent held-out movie.

1036 *Noise ceiling estimation:* For the held-out group, we obtain the noise ceiling by considering  
1037 variability across subjects. Here, the noise ceiling was computed as the correlation coefficient  
1038 between the mean measured response for the *independent* test movie across all 158 subjects  
1039 in the training set and the group-averaged response computed over the 20 new subjects. This  
1040 metric captures the response component shared across independent groups of subjects and thus  
1041 reflects the upper bound achievable by a group-level encoding model. We employ this noise  
1042 ceiling for comparison against the prediction accuracy of the model on the held-out group of  
1043 subjects (Figure S7).

1044 The models accurately predicted cortical responses evoked by the *independent* test movie as  
1045 measured in the *independent* subject population (Figure S7, S8), with the best performing  
1046 model (Audiovisual-20sec) even achieving close to perfect predictivity relative to the “noise  
1047 ceiling” in certain multi-sensory sites such as the posterior STS (Figure S7(A), (G)). Here, the  
1048 noise ceiling was computed as the correlation coefficient between the mean neural response in  
1049 the *independent* test movie, across all 158 subjects in the training set and the group-averaged  
1050 response computed over the 20 new subjects. This metric captures the response component  
1051 shared across independent subject populations and thus reflects the upper bound achievable by  
1052 a group-level encoding model. These results clearly indicate that inclusion of temporal history  
1053 and multi-sensory information pushes the prediction accuracies closer to their upper bound, as  
1054 also evidenced by a higher slope of the linear model fit on their corresponding data points. Fur-  
1055 ther, voxels that truly approach the noise ceiling are predominantly associated with the auditory  
1056 group of regions as broadly characterized within the HCP MMP parcellation. Interestingly, we  
1057 find that this regional distribution of predictivity against noise ceiling holds even for subject-  
1058 specific responses and not just the group-averaged responses, as described in the next section  
1059 and shown in Figure S9.

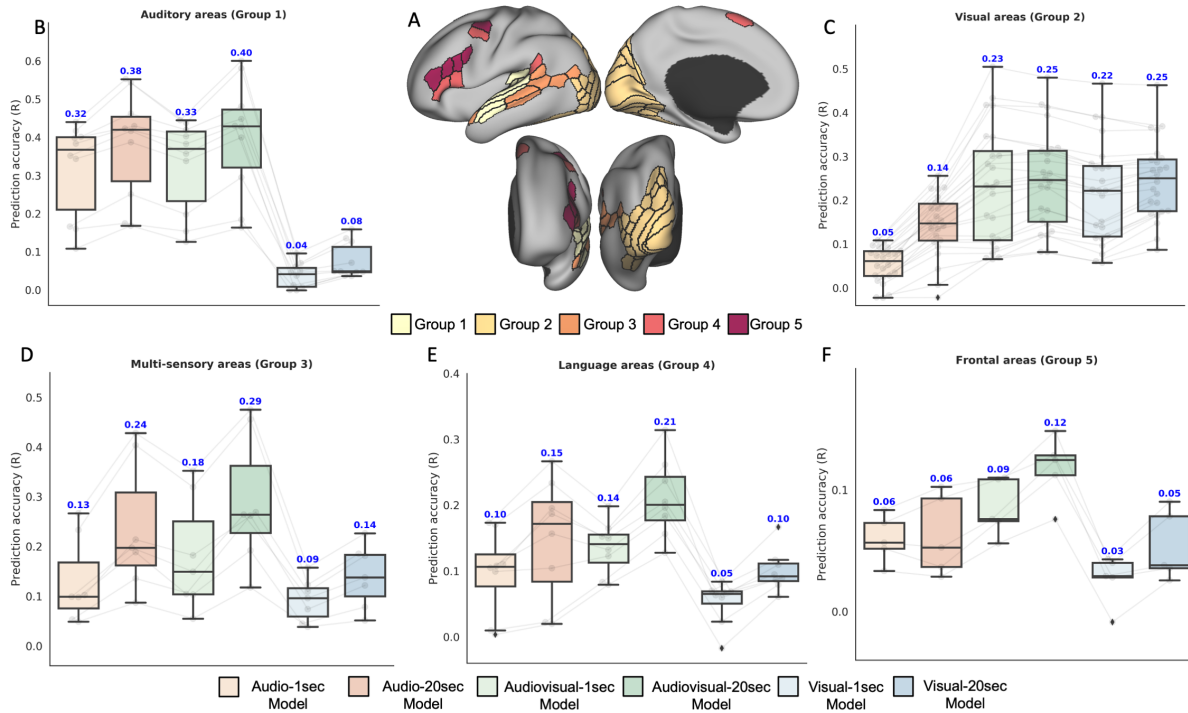
### 1060 **Subject-level prediction accuracy: held-out set**

1061 For each participant in our independent subject group ( $N = 20$ ), we computed the correlation  
1062 coefficient ( $R$ ) between the predictions of the best performing model (Audiovisual-20sec) and  
1063 the subject-specific fMRI response corresponding to the independent movie. We further contrast  
1064 this cortical map of prediction performance against another map computed as the voxel-wise  
1065 correlation coefficient between the mean neural response across all 158 training subjects and the  
1066 respective subject-specific response on the independent movie. The latter places an upper bound  
1067 on the predictivity of each voxel as achievable by any group-level model. Here, we present the  
1068 results for 5 subjects with mean prediction accuracy (un-normalized) within the stimulus-driven  
1069 cortex in the  $i$ th percentile with  $i \in \{0.01, 25, 50, 75, 99.9\}$ . The results (Figure S9) suggest that

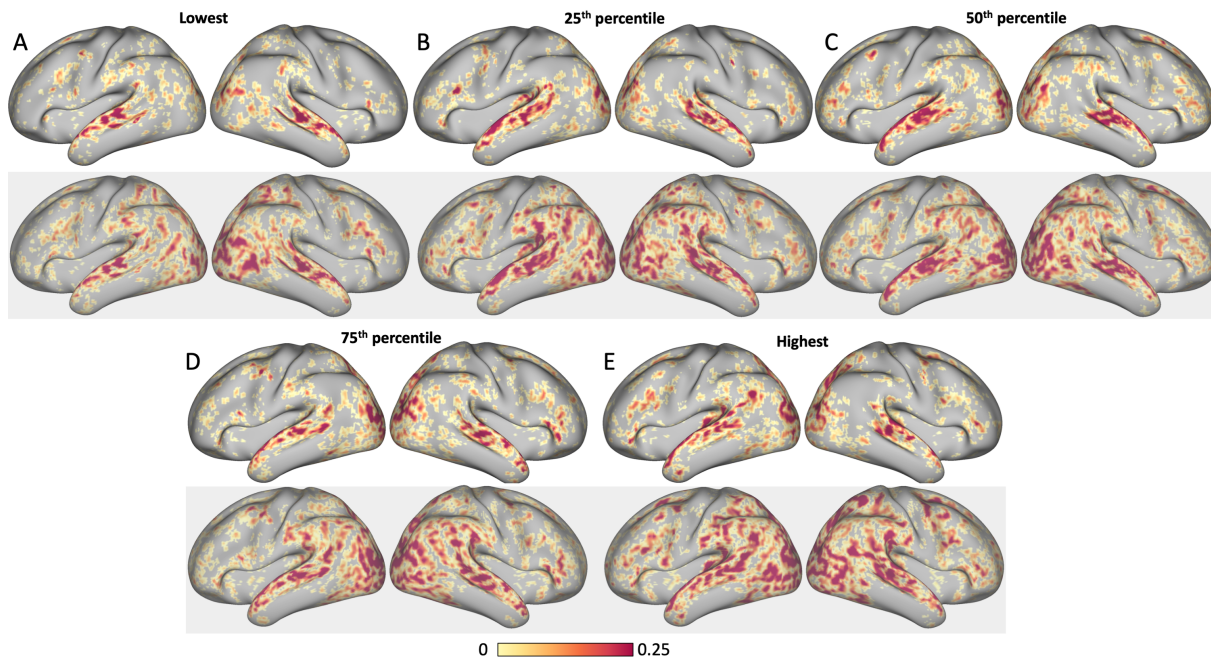


**Fig. S7.** Model performance on held-out group of subjects. (A) Pearson’s correlation coefficient (R) between the model predictions and group-averaged response of an independent subject group comprising 20 subjects, on the held-out test movie, normalized by the voxel-specific noise ceiling. (B) Predictivity against the noise ceiling for all voxels with high “synchrony” across training movies ( $>0.5$ ) (see Supplementary Information for details). This gives a total of 52,954 highly “synchronous” voxels that are colored based on their association with auditory and visual groups. This hue assignment of each voxel was derived from the coloration of the corresponding ROI in the multi-modal HCP parcellation. Each dot in the scatterplot represents an individual voxel. Bivariate kernel density estimates are overlaid on top of the scatterplot as contours to depict the probability distribution of observations (prediction accuracy/noise ceiling pair at every voxel). 40





**Fig. S8.** Quantitative evaluation metrics for all the proposed models on the independent *held-out* population comprising 20 novel subjects. (B)-(F) depict prediction accuracy (R) for all the proposed models across major groups of regions as identified in the HCP MMP parcellation (A). Predictive accuracy of all models is summarized across (B) auditory, (C) visual, (D) multi-sensory, (E) language and (F) frontal areas. Box plots depict quartiles and swarmplots depict mean prediction accuracy of every ROI in the group. For language areas (Group 4), left and right hemisphere ROIs are shown as separate points in the swarmplot because of marked differences in the prediction accuracy. Statistical significance tests (results indicated with horizontal bars) are performed to compare 1-sec and 20-sec models of the same modality (3 comparisons) or uni-modal against multi-modal models of the same duration (4 comparisons) using paired t-test ( $p$ -value  $< 0.05$ , Bonferroni corrected) on mean prediction accuracy within ROIs of each group.



**Fig. S9.** Comparison of voxel-level prediction accuracies ( $R$ ) against subject-specific noise ceiling for 5 representative subjects from the held-out set. The subjects were chosen such that their mean prediction accuracy (un-normalized) within the stimulus-driven cortex lied in the  $i$ th percentile with  $i \in \{0.01, 25, 50, 75, 99.9\}$ . Surface maps with white background in (A)-(E) depict raw correlation coefficients between model (Audiovisual-20sec) predictions and subject-specific response on the held-out movie whereas maps on gray background indicate the respective subject-specific noise ceiling. Only significantly correlated voxels ( $p < 0.05$ , FDR corrected) are colored on the surface.

1070 the model can successfully capture the response component that individual subjects share with  
1071 the population.