
Comparing Representations in Static and Dynamic Vision Models to the Human Brain

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We compared neural responses to naturalistic videos and representations in deep
2 network models trained with static and dynamic information. Models trained with
3 dynamic information showed greater correspondence with neural representations
4 in all brain regions, including those previously associated with the processing of
5 static information. Among the models trained with dynamic information, those
6 based on optic flow accounted for unique variance in neural responses that were not
7 captured by Masked Autoencoders. This effect was strongest in ventral and dorsal
8 brain regions, indicating that despite the Masked Autoencoders' effectiveness at a
9 variety of tasks, their representations diverge from representations in the human
10 brain in the early stages of visual processing.

11 1 Introduction

12 The human visual system is organized into distinct processing streams Ungerleider et al. [1982],
13 Pitcher and Ungerleider [2021]: a ventral stream that extends from early visual regions into the
14 inferior portions of temporal cortex, and a dorsal stream that extends into lateral occipital cortex and
15 branches into a lateral stream (along the superior temporal sulcus) and a parietal stream (reaching
16 the inferior parietal lobule). This organization likely results from the computational requirements of
17 visual perception. Therefore, understanding the representations encoded by different visual streams
18 could offer insights about the human brain and also about more general principles of vision.

19 The ventral stream has been proposed to encode static object identity Grill-Spector and Weiner [2014],
20 while dynamic information has been associated with the dorsal, lateral and parietal streams Ganel and
21 Goodale [2003], Culham et al. [2003]. Indeed, static images of objects are known to drive responses
22 in ventral temporal regions in macaques Pasupathy and Connor [2002], Logothetis et al. [1995],
23 Tanaka [1996], Hung et al. [2005] and in humans Edelman et al. [1998], Haxby et al. [2001]. Moving
24 stimuli drive stronger responses in dorsal and lateral regions Zeki et al. [1991], Tootell et al. [1995],
25 Saito et al. [1986]. In addition, disruption to lateral regions using TMS affects the processing of
26 dynamic information Beckers and Hömberg [1992], Pitcher et al. [2014] as well as motion prediction
27 Vetter et al. [2015].

28 However, other studies have challenged the hypothesis that visual streams in the human brain differ
29 based on whether they encode static or dynamic visual features. These studies suggested that both
30 static and dynamic features are represented in multiple visual streams Kourtzi et al. [2002], Freud
31 et al. [2017], Cornette et al. [1998], Sunaert et al. [1999], Robert et al. [2023]. Here, we investigated
32 the contribution of static and dynamic information to the representations encoded by different visual
33 streams, by quantifying the convergence between neural representations and representations learned
34 by deep network models.

35 Previous work compared neural responses to deep network models trained with static images Yamins
36 et al. [2013], Khaligh-Razavi and Kriegeskorte [2014], Zhuang et al. [2021], Konkle and Alvarez

37 [2022]. The present work studies the additional contribution of dynamic information during the
 38 observation of quasi-naturalistic videos, by comparing neural responses to deep networks whose
 39 inputs are static images (e.g. convolutional ResNets He et al. [2016], image masked autoencoders He
 40 et al. [2022]) and to deep networks whose inputs are videos (e.g. hidden two-stream networks, video
 41 masked autoencoders Zhu et al. [2019], Tong et al. [2022], Feichtenhofer et al. [2022]). Critically, we
 42 include in our analyses a family of self-supervised models that are widely used in Computer Vision,
 43 but that are understudied in Cognitive Neuroscience: masked autoencoders (MAEs). We investigate
 44 the correspondence of representations in MAEs and Video MAEs to neural representations in the
 45 human brain.

46 2 Methods

47 2.1 Data

48 BOLD fMRI responses ($3 \times 3 \times 3$ mm) to eight movie segments of ‘Forrest Gump’ were obtained
 49 from the publicly available *studyforrest* audiovisual dataset (<http://studyforrest.org>). Fifteen
 50 right-handed participants took part in the study (6 females; age range 21-39 years, mean 29.4 years).
 51 The data was acquired with a T2*-weighted echo-planar imaging sequence, using a whole-body 3
 52 Tesla Philips Achieva dStream MRI scanner equipped with a 32 channel head coil.

53 2.2 Preprocessing

54 Data were first preprocessed using fMRIPrep (<https://fmriprep.readthedocs.io/en/latest/index.html>): a robust pipeline for the preprocessing of diverse fMRI data. Anatomical
 55 images were skull-stripped with ANTs (<http://stnava.github.io/ANTs/>), and FSL FAST
 56 was used for tissue segmentation. Functional images were corrected for head movement with FSL
 57 MCFLIRT (<https://fsl.fmrib.ox.ac.uk/fsl/fslwiki/MCFLIRT>), and were subsequently
 58 coregistered to their anatomical scan with FSL FLIRT. Finally, the skull-stripped anatomical images
 59 were normalized to the MNI template using SPM. We denoised the data with CompCor Behzadi et al.
 60 [2007] using 5 principal components extracted from the union of cerebrospinal fluid and white matter.
 61

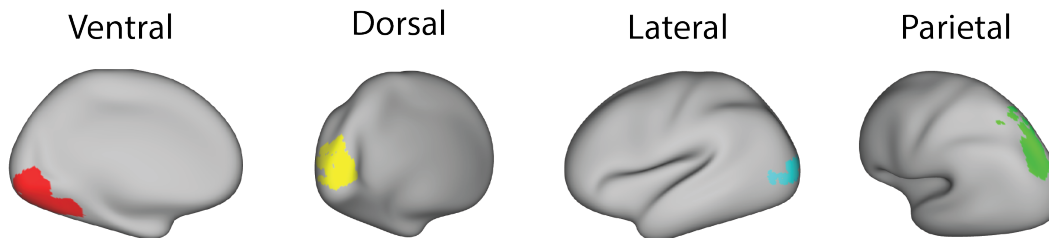


Figure 1: Masks of the visual streams in the human brain projected on an inflated cortical surface in MNI space.

62 2.3 Regions of Interest (ROI)

63 To identify the regions of interest (ROI), we used an atlas of probabilistic maps of visual topography
 64 in the human cortex from a previous study Wang et al. [2015]. The atlas contains twenty-five cortical
 65 regions and spans multiple visual streams: ventral, dorsal, parietal, and lateral (Figure 1).

66 A list of probabilities is associated with each voxel to reflect the likelihood of that voxel being part of
 67 each of the twenty-five brain regions ($R_i, i = 1, \dots, 25$). We calculated the transformation from MNI
 68 space to each participant’s native space and co-registered the probability maps with each participant’s
 69 anatomy. To prevent overlap between the regions of interest in the participants’ native space, we
 70 followed a procedure analogous to Wang et al. [2015]. Specifically, we calculated the maximum
 71 probability map for each participant, using which we exclusively classified each voxel as either
 72 belonging to a specific ROI or as being outside of all the ROIs.

73 The inclusion probability was computed as the probability of a voxel of being in any of the defined
 74 regions ($P(\cup_{i=1}^{25} v \in R_i)$), and The exclusion probability is the probability of a voxel not belonging to

Table 1: Models of visual cortex

Model	Input	Output	Training dataset	#Selected layers
Supervised static	image	object identity	Image-net	11
Supervised static	image	action identity	HAA-500	11
Self-supervised dynamic	video	optic flows	HAA-500	11
Supervised dynamic	optic flow	action identity	HAA-500	11
pre-trained Masked Autoencoder	(masked) image	(unmasked) image	Image-net	12
fine-tuned Masked Autoencoder	image	object identity	Image-net	12
pre-trained Masked Video Autoencoder	(masked) video	(unmasked) video	Kinetics-400	12
fine-tuned Masked Video Autoencoder	video	action identity	Kinetics-400	12
pre-trained Masked Video Distillation	(masked) video	MAE & VideoMAE high-level features	Kinetics-400	12

75 any of the ROIs ($P(\cap_{i=1}^{25} v \notin R_i)$). If the exclusion exceeded the inclusion probability, we discarded
76 the voxel. Otherwise, we classified the voxel as belonging to the region with the highest probability.
77 The resulting ROIs were grouped into four sets associated with distinct visual streams. The ventral
78 stream contains V1v, V2v, V3v, hV4, VO1, VO2, PHC1, PHC2; the dorsal stream V1d, V2d, V3d,
79 V3a, V3b; the lateral stream LO1, LO2, hMT, and finally the parietal stream IPS0, IPS1, IPS2, IPS3,
80 IPS4, IPS5, SPL1, and FEF. While often the term “dorsal stream” is used to refer to the combination
81 of the regions we labeled as “dorsal stream” and the regions we labeled as “parietal stream”, here
82 we sought to distinguish between the initial branch of the dorsal stream and its parietal and lateral
83 temporal continuations, without implying that the initial segment is disproportionately associated
84 with one or the other.

85 2.4 Models of human visual cortex

86 To study representations of quasi-naturalistic visual stimuli, we used a variety of vision models,
87 including feed-forward convolutional neural networks, as well as state-of-the-art foundation vision
88 models. The models vary in architecture, learning objective, and training data (Table 1). Here, we
89 propose an overview of the models. Training details for the HAA-trained CNNs are presented in
90 supplementary materials. The trained versions of all other models are adopted from their official
91 implementation repository. For model details, refer to the original papers.

93 **Supervised (sup) static net** is the spatial stream of the hidden two-stream convolutional neural
94 network model Zhu et al. [2019]. The sup static net has a resnet18 architecture and encodes static
95 features of visual stimulus. Two versions of the model were included in the models’ pool: one is
96 trained on Image-Net Deng et al. [2009] and predicts object identity, and the other is trained on
97 HAA-500 action dataset Chung et al. [2021] and predicts action label. Both versions take a single
98 frame as input.

99 **Self-supervised (s-sup) dynamic net** is the first part of the temporal stream (i.e., motion net) in
100 the hidden two-stream convolutional neural network model Zhu et al. [2019]. The self-supervised
101 dynamic net takes 11 consecutive frames as input and infers the optic flow between each pair of

consecutive frames. The network is trained to minimize a self-supervised learning objective obtained by combining three loss functions: 1) a pixel-wise reconstruction error, 2) a smoothness loss addressing the ambiguity problem of optic flow estimation (also known as the aperture problem), and 3) a structural dissimilarity between the original and the reconstructed image patches (see Zhu et al. [2019] for details of loss functions). The models’ pool contains one version of the self-supervised dynamic net, trained on the HAA-500 action dataset Chung et al. [2021].

Supervised (sup) dynamic net is the second part of the temporal stream in the hidden two-stream convolutional neural network model Zhu et al. [2019]. The model has resnet18 architecture and takes optic flows from the self-supervised dynamic net as input. We used the HAA-500 dataset Chung et al. [2021] and trained the supervised dynamic net to predict action labels using optic flows.

Masked Autoencoder (MAE) learn image representations, required to reconstruct original uncorrupted images from corrupted (masked) input through a series of transformer blocks He et al. [2022]. The models’ pool contains two versions of the MAE model: 1) a pre-trained version, where the model is trained to reconstruct pixel values, and 2) a fine-tuned version, where the pre-trained model is further fine-tuned to predict object identities. Both versions were trained on Image-net Deng et al. [2009].

Video Masked Autoencoder (VMAE) learns a spatiotemporal representation of videos, required to reconstruct original uncorrupted videos, from corrupted (tube masked) input through a series of transformer blocks Tong et al. [2022]. We added two versions of the VMAE to our models’ pool. The first is a pre-trained version, where the model is trained to reconstruct missing pixels of the input set of frames. The second version is the fine-tuned version obtained by fine-tuning the pre-trained version to predict action labels of input videos. Both models take a consecutive set of frames as input, and were trained on the Kinetics-400 action dataset Kay et al. [2017].

Masked Video Distillation (MVD) learns a higher-level spatial and spatiotemporal representation of the input video, required to reconstruct the representation of teacher MAE and VMAE while taking corrupted (tube-masked) videos as input Wang et al. [2023a]. Unlike VMAE and MAE, the MVD model does not learn pixel-level features. Rather, it learns high-level features of the input video using pre-trained MAE and VMAE models’ features as masked prediction targets. Using the Kinetics-400 action dataset Kay et al. [2017], a pre-trained version was obtained and added to the models’ pool.

2.5 Models’ Representational Dissimilarity Matrices (RDM)

In order to compare the models and the fMRI data, we computed representational dissimilarity matrices (RDMs) for the models’ layers with a multi-step procedure. First, since the temporal resolution of the models’ representations (25Hz) is much higher than the temporal resolution of fMRI data, we down-sampled each layer’s activation timecourses over time by selecting one data point every five time points (down to 5 Hz). Then, we convolved the layer’s activations with a standard Hemodynamic Response Function (HRF). Given that the fMRI data’s repetition time (TR) is 2 seconds, we took a layer’s activation every $25 \times 2 = 50$ time points.

Finally, for each layer we computed the dissimilarities between all pairs of timepoints, obtaining RDMs in which the entry at column j and row i contains correlation dissimilarity ($1 - r$) between the layer activations at time i and time j . We repeated this procedure for BOLD responses to all eight movie segments, resulting in eight RDMs.

2.6 Brain Representational Dissimilarity Matrices (RDM)

RDMs were constructed separately for each brain stream in the subject’s native space. The voxels for each brain stream were obtained as the union of the region voxels for individual regions within that stream. For each brain stream, we calculated the correlation dissimilarity ($1 - r$ where r is Pearson’s correlation) of fMRI response patterns for all pairs of TRs. This yielded eight RDMs, corresponding to BOLD responses in eight video segments.

2.7 Measuring models similarity with brain data

To evaluate how well each model accounts for the activity in the brain streams, we used a cross-validated linear regression to predict the left-out movie segment brain stream RDM and computed the correlation between the predicted and the true RDM in each brain stream. The correlation captures how well a model’s layers can predict a brain stream’s responses to the visual stimuli. First, we

154 used each model’s layers’ RDMs corresponding to seven (out of eight) video segments to train a
 155 linear regression model that predicts the corresponding seven RDMs in each brain stream. Then, we
 156 averaged the linear regression model’s coefficients along the seven segments and used the averaged
 157 coefficients to predict the brain stream RDM of the left-out segment, using the model layers’ RDMs
 158 of the corresponding segment. Finally, we calculated the Pearson’s correlation between the predicted
 159 and the true RDMs. We repeated the leave-one-out cross-validation process for all the segments and
 160 averaged over the obtained correlations.

161 **2.8 Measuring combined models similarity with brain data**

162 We sought to study whether a combination of features from two models can improve similarity with
 163 brain data. We followed the procedure in 2.7 and used RDMs of all the layers in a pair of combined
 164 models to estimate the coefficients of a linear regression model that best predicts the RDM of a
 165 brain stream in seven (out of eight) of the video segments. Using leave-one-out cross-validation, we
 166 predicted the brain stream RDM of the left-out video segment using the average of the coefficients
 167 obtained from the seven video segments during training. Finally, we measured the correlation between
 168 the predicted RDM and the actual brain stream RDM to measure the correspondence between the
 169 combined models’ features and the brain activity.

170 **2.9 Measuring unique and shared similarity of a pair of models with brain data**

171 To evaluate how well unique and shared features among a pair of computational models correspond
 172 to the brain data, we used Pearson’s r to measure the accuracy of a ”target” model’s layers prediction
 173 of a brain stream RDM while controlling for the variation of a ”control” model layers. Using leave-
 174 one-out cross-validation, first, we estimated the coefficients of a linear regression model that predicts
 175 a brain stream’s RDM from the control model’s layers in training video segments (seven out of eight).
 176 Second, we subtracted the predicted from the actual brain stream RDM in the training and the left-out
 177 video segments to obtain training and left-out residuals. Third, we estimated the coefficients of a
 178 linear regression model that predicts training residuals of each video segment using the target model
 179 layers. Finally, we measured Pearson’s correlation between the target model’s prediction of the
 180 left-out video segment residuals and the residuals obtained from the prediction of the control model.

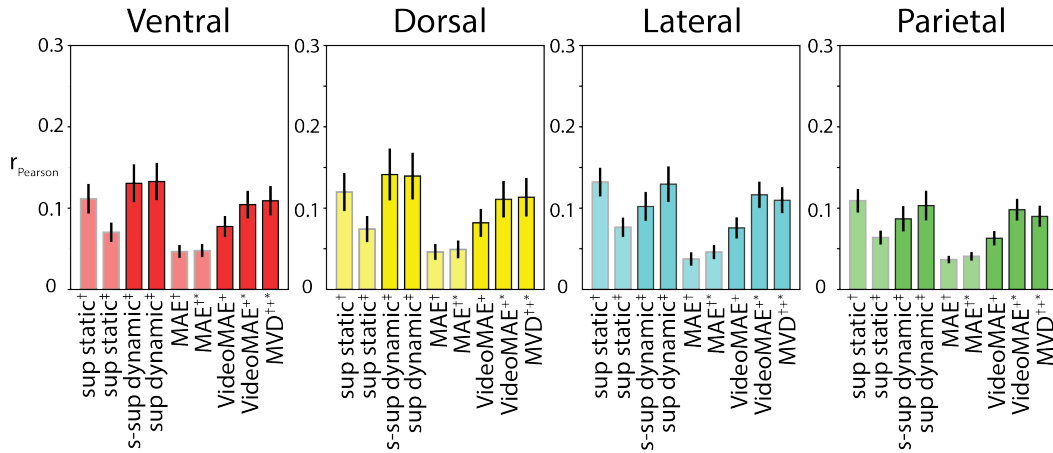


Figure 2: Pearson’s correlation between actual and predicted brain stream RDMs, averaged over participants. Predicted RDMs were obtained by training and test a leave-one-out cross-validation linear regression model using each model’s layers. Error bars show standard deviation over participants. Lighter bars correspond to models containing static, and darker ones to models containing dynamic visual information (sup: supervised, s-sup: self-supervised, †: Image-net-trained, ‡: HAA-500-trained, +: Kinetics-400-trained, *: fine-tuned; MVD was trained on pre-trained MAE (Image-net) and VideoMAE (Kinetics-400))

181 **3 Results**

182 The human visual system does not consist of a single processing stream. Instead, it is organized into
183 distinct neural pathways. To the extent that the structure of the human visual system is shaped by
184 computational optimality, understanding the visual representations encoded in these pathways can
185 offer insights into more general principles of vision. The contribution of this work is to quantify
186 the similarity between the representations in the different visual pathways in the human brain and
187 representations in models of vision that are widely used in Computer Science but understudied in
188 Cognitive Neuroscience (e.g. masked autoencoders, two-stream networks, masked video distillation).
189 In a first set of analyses 3, we leverage differences between models to reveal differences between
190 the information encoded in different visual pathways in the brain. In a second set of analyses, we
191 quantify the unique contribution of different deep network models to account for neural responses
192 3.2.

193 **3.1 DNN models similarity with human brain streams**

194 While numerous research studies have been conducted on model-to-brain correspondence using static
195 images Rose et al. [2021], Doshi and Konkle [2023], Tsao et al. [2006], the impact of dynamic
196 information on neural responses to naturalistic videos is understudied. To fill this gap, we tested
197 the correspondence between neural representations in different visual pathways (ventral, dorsal,
198 lateral, and parietal) and deep network models that can process dynamic information (two-stream
199 networks, video masked autoencoders, and masked video distillation). Comparing the correspondence
200 of neural responses with these models and their correspondence with models that only process static
201 information (standard convolutional ResNets, masked autoencoders) made it possible to study the
202 contribution of dynamic information independently of whether the learning objective is supervised
203 (as in two-stream networks) or unsupervised (as in masked autoencoders). In addition, the study of
204 the correspondence between neural representations and representations in masked autoencoders is of
205 interest in its own right: masked autoencoders are effective and widely used, but little is known about
206 their similarity to neural representations.

207 **3.1.1 Static and dynamic information in brains and feed-forward convolutional neural** 208 **networks (CNN)**

209 Functional MRI responses recorded during the observation of naturalistic videos in the ventral,
210 dorsal, lateral and parietal visual pathways were compared to the representations in feed-forward
211 convolutional neural networks. The same dataset (HAA-500) was used to train the different branches
212 of a hidden-two-stream network: the “supervised static” branch (a ResNet that takes as input
213 individual frames of a video and computes as output the action category), the “unsupervised dynamic”
214 branch (a convolutional network trained to compute optic flow by minimizing a self-supervised loss),
215 and the “supervised dynamic” branch (a ResNet that takes as input optic flow and computes as output
216 the action category). In addition, to facilitate parallels with prior work, we compared neural responses
217 to a widely studied feed-forward model: a ResNet trained with Image-net.

218 Comparing deep network models trained with the same dataset (HAA-500) showed that models
219 including dynamic information correlated with neural responses more than the Spatial model, that does
220 not use dynamic information (Figure 2). This effect was observed for all visual pathways. In addition,
221 representations in the lateral and parietal pathways correlated more with the supervised dynamic
222 model than with the unsupervised dynamic model (fisher-transformed t-values with Bonferroni-
223 corrected threshold). Lateral and parietal regions are located downstream compared to the dorsal
224 regions, thus this result is complementary to earlier work that reported a correspondence between
225 subsequent stages of processing in deep neural networks and in neural pathways in the case of static
226 visual stimuli [Khaligh-Razavi and Kriegeskorte, 2014] and in the case of auditory stimuli [Kell et al.,
227 2018].

228 Supervised CNNs trained with Image-Net performed well, achieving correspondence with neural
229 responses that was close to that of HAA-trained models that included dynamic information. This
230 could indicate that some of the variance in neural responses that correlates with dynamic models
231 might also be accounted for by models trained exclusively with static information, as long as a
232 suitable training dataset is used (in this case, Image-Net). However, an alternative possibility is that
233 the supervised static model trained with Image-Net and the dynamic models trained with HAA might

234 account for different portions of the variance in neural responses. We investigate these alternative
235 possibilities in section 3.2.

236 3.1.2 Static and dynamic information in brains and masked autoencoders

237 Masked Autoencoders (MAE, He et al. [2022]) and Video Masked Autoencoders (VideoMAE,
238 Tong et al. [2022], Feichtenhofer et al. [2022]) models are trained to reconstruct masked pixels of
239 input (image or video) during pre-training and are further fine-tuned to predict object/action labels.
240 MAE and VideoMAE models are very effective in learning visual representations and have been
241 shown to outperform competing models in several visual tasks He et al. [2022], Tong et al. [2022],
242 Feichtenhofer et al. [2022], Wang et al. [2023b], Venkatesh et al.. However, it is still unknown
243 whether the representations learned by models based on masked autoencoding are similar to visual
244 representations in the human brain. Here we investigated this question, quantifying the correlation
245 between neural responses measured with fMRI while participants watched naturalistic videos, and
246 representations learned by models trained with masked autoencoding.

247 We compared the correspondence between neural responses and MAEs trained with images (which
248 learn spatial relationships between component of an image, Wang et al. [2023a]) as well as Video-
249 MAEs (which learn temporal relationships in videos, Wang et al. [2023a]). Finally, we also compared
250 neural responses to masked video distillation (MVD, Wang et al. [2023a]), which combines image
251 MAEs and videoMAEs to better capture both spatial and temporal relationships. Unlike MAE and
252 VideoMAE, the MVD model does not aim to reconstruct missing patches at the level of pixel values.
253 Instead, MVD adopts a knowledge-distillation approach, reconstructing missing information at the
254 level of features extracted from pre-trained MAE and VideoMAE teachers.

255 As in the case of supervised models trained with the HAA dataset, models that included dynamic
256 information (VideoMAEs) outperformed models using only static information (Image MAEs). This
257 pattern was observed across all visual pathways. Image MAEs did not correlate well with neural
258 responses, even compared to supervised models trained with static inputs. Overall, the representations
259 learned by Image MAEs were very different from neural representations. By contrast, VideoMAEs
260 showed greater correspondence with neural responses. In particular, fine-tuning with an action
261 recognition task (Figure 2, VideoMAE fine-tuned) improved the correspondence between Video-
262 MAE representations and neural representations across all streams (fisher-transformed t-values with
263 Bonferroni-corrected threshold). Across all the pre-trained models, pre-trained MVD showed the
264 highest similarity to neural representations in all brain streams. Further, MVD showed comparable
265 similarity with brain streams to that of fine-tuned VideoMAE.

266 3.2 Vision models capture shared and unique neural activity variation in human brain 267 streams

268 The results described in 3 show that representations from models trained with dynamic information
269 are more correlated with neural representations compared to representations from models trained
270 with static information. This overall pattern is broken by the exception of ResNets trained with
271 ImageNet, which performed on par with models trained with dynamic information. This raises
272 the question of whether ResNets trained with ImageNet and dynamic models explain overlapping
273 variance in neural responses or whether, instead, they are complementary, capturing non-overlapping
274 portions of the variance. This question can be posed more generally for any pair of models studied in
275 section 3. We investigated this first by combining layers from two models and measuring whether a
276 combination of models can better predict the pattern of neural activity in visual pathways. Second,
277 we measured the correspondence between a “target” model’s representations and the representations
278 in each brain stream while controlling for the representations encoded in a “control” model. To this
279 end, we predicted neural representations using the representations of the control model and obtained
280 the residuals. Then, we predicted the residuals using the representations in the target model (see
281 Methods for details).

282 Each matrix in Figure 3.a shows how well a combination of models’ layers can predict the pattern
283 of neural activity in a brain stream. Each column of each row demonstrates the correlation between
284 the neural response pattern of a brain stream and the combined models’ layers’ prediction of that
285 brain stream’s neural activity pattern. Model-to-brain-stream similarity increased in all brain streams
286 when combining features from static models with features from dynamic models. Notably, the
287 correspondence between combined models’ features with both dorsal and ventral streams improved

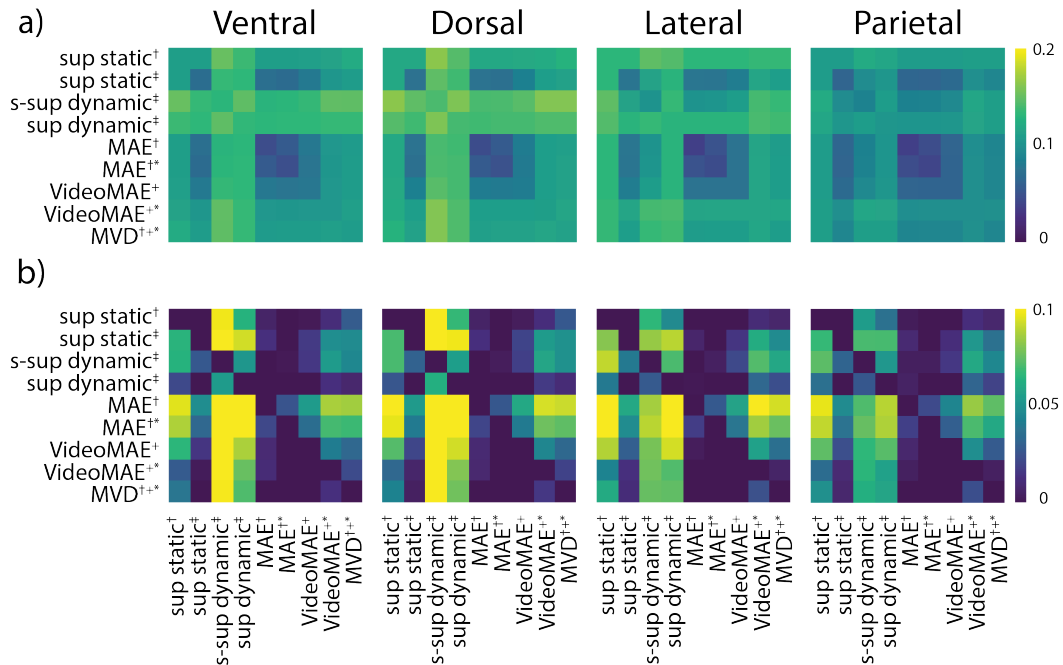


Figure 3: **a)** Model combination similarity with brain streams. The similarity was calculated using Pearson’s correlation between a brain stream’s actual and predicted RDMs. These predictions were obtained by combining layers from two models (corresponding to the row and column names), and averaged across participants. A linear regression model was trained and tested using leave-one-out cross-validation to generate the predictions.. **b)** Models unique similarity with brain streams. The similarity was calculated using Pearson’s correlation between the actual RDM of a brain stream and the RDM predicted by a target model while controlling for the variation explained by a control model in the brain stream. These correlations were averaged across participants. Each row corresponds to a different control model, and each column corresponds to a different target model used for prediction. (sup: supervised, s-sup: self-supervised, †: Image-net-trained, ‡: HAA-500-trained, +: Kinetics-400-trained, *: fine-tuned; MVD was trained on pre-trained MAE (Image-net) and VideoMAE (Kinetics-400))

288 in two cases: 1) combined features from Image-net-trained static supervised models with dynamic
 289 features from self-supervised model and 2) features from the combination of the self-supervised
 290 dynamic model with either VideoMAE or MVD. These cases shows that, first, ventral and dorsal brain
 291 streams both represent static and dynamic visual features, and second, different types of dynamic
 292 information are represented in both of these human streams.

293 Figure 3.b demonstrates the correspondence between a target model’s features and each brain stream
 294 when we controlled for the features of a control model in the brain stream’s neural responses. The
 295 results are visualized as a matrix in which each row corresponds to a control model and each column
 296 to a target model. The first row displays the correlations between models and neural responses after
 297 controlling for the Image-net-trained static model. The high values for the columns corresponding to
 298 the self-supervised dynamic and the supervised dynamic models indicate that these models and the
 299 Image-net-trained static model capture non-overlapping variance in neural responses. Representations
 300 learned by these models also capture non-overlapping variance with those learned by the unsupervised
 301 dynamic models: the VideoMAEs. This finding shows that despite VideoMAEs exhibit relatively
 302 high correlations with neural responses (outperforming Image MAEs), they nonetheless fail to capture
 303 some variance in human visual representations that is accounted for by self-supervised and supervised
 304 dynamic models.

305 VideoMAEs and MVD accounted for additional variance in neural responses compared to MAEs
 306 (as expected given the results in Figure 2) but also compared to the HAA-trained static and self-
 307 supervised dynamic models. However, they accounted for a minimal amount (if any) of additional

308 variance compared to the supervised dynamic model, suggesting some degree of convergence on
309 common representations across models trained with different learning objectives.

310 **4 Limitations**

311 This study focused on a set of models selected to enable comparing the contribution of static and
312 dynamic information and the impact of supervised and unsupervised learning objectives. The selection
313 of models in this study includes only a subset of the existing models, future work will be needed to
314 expand the set of models tested. In addition, the present work centered on the comparison between
315 models and entire visual streams. A finer-grained analysis comparing models to individual regions
316 within each stream will require further work.

317 **5 Discussion**

318 Three main findings emerged. First, models including dynamic information outperformed models
319 using exclusively static information, not only in the dorsal, lateral and parietal streams, but also in
320 the ventral stream. This is in line with recent evidence of responses to dynamic features in ventral
321 brain regions [Robert et al., 2023]. Patients with deficits for motion perception typically present
322 with lesions affecting dorsal regions (such as area V5, [McLeod, 1996, Vaina et al., 1990, Zihl et al.,
323 1983]) or parietal regions [Battelli et al., 2003]. By contrast, patients with damage to ventral regions
324 typically do not present with deficits for motion perception [Gilaie-Dotan et al., 2015]. This raises
325 the question of what might be the use of dynamic information represented in the ventral stream. We
326 hypothesize that this information might be used to support object segmentation, as proposed in recent
327 computational models [Chen et al., 2022] inspired by classic work in Developmental Psychology
328 [Spelke, 1990].

329 Second, Image MAEs showed little correspondence with neural representations, even compared
330 to other models trained exclusively with static information. These results indicate that despite the
331 effectiveness of Image MAEs for learning visual representations that can transfer to a variety of
332 visual tasks [He et al., 2022], these models do not converge on representations that are similar to
333 those observed in the human brain, suggesting that human vision and image MAEs rely on different
334 computational mechanisms.

335 Third, models based on optic flow representations accounted for unique variance in all streams, even
336 compared to video masked autoencoders that can make use of dynamic information. Fine-tuning video
337 MAEs with an action classification task increased their correspondence with neural representations,
338 but did not fully bridge the gap with neural responses compared to optic flow models, which still
339 explained additional unique variance compared to the fine-tuned video MAEs. The additional
340 contribution of optic flow models was particularly strong in ventral and dorsal regions, suggesting
341 that representations based on optic flow exhibit greater correspondence with representations in early
342 stages of visual processing in the human brain compared to both image and video MAEs.

343

344 **References**

- 345 Leslie G Ungerleider, Mortimer Mishkin, et al. Two cortical visual systems. analysis of visual
346 behavior. *Ingle DJ, Goodale MA, Mansfield RJW*, 1982.
- 347 David Pitcher and Leslie G Ungerleider. Evidence for a third visual pathway specialized for social
348 perception. *Trends in Cognitive Sciences*, 25(2):100–110, 2021.
- 349 Kalanit Grill-Spector and Kevin S Weiner. The functional architecture of the ventral temporal cortex
350 and its role in categorization. *Nature Reviews Neuroscience*, 15(8):536–548, 2014.
- 351 Tzvi Ganel and Melvyn A Goodale. Visual control of action but not perception requires analytical
352 processing of object shape. *Nature*, 426(6967):664–667, 2003.
- 353 Jody C Culham, Stacey L Danckert, Joseph FX De Souza, Joseph S Gati, Ravi S Menon, and
354 Melvyn A Goodale. Visually guided grasping produces fmri activation in dorsal but not ventral
355 stream brain areas. *Experimental brain research*, 153:180–189, 2003.

- 356 Anitha Pasupathy and Charles E Connor. Population coding of shape in area v4. *Nature neuroscience*,
357 5(12):1332–1338, 2002.
- 358 Nikos K Logothetis, Jon Pauls, and Tomaso Poggio. Shape representation in the inferior temporal
359 cortex of monkeys. *Current biology*, 5(5):552–563, 1995.
- 360 Keiji Tanaka. Inferotemporal cortex and object vision. *Annual review of neuroscience*, 19(1):109–139,
361 1996.
- 362 Chou P Hung, Gabriel Kreiman, Tomaso Poggio, and James J DiCarlo. Fast readout of object identity
363 from macaque inferior temporal cortex. *Science*, 310(5749):863–866, 2005.
- 364 Shimon Edelman, Kalanit Grill-Spector, Tammar Kushnir, and Rafael Malach. Toward direct
365 visualization of the internal shape representation space by fmri. *Psychobiology*, 26:309–321, 1998.
- 366 James V Haxby, M Ida Gobbini, Maura L Furey, Alunit Ishai, Jennifer L Schouten, and Pietro
367 Pietrini. Distributed and overlapping representations of faces and objects in ventral temporal cortex.
368 *Science*, 293(5539):2425–2430, 2001.
- 369 Semir Zeki, JD Watson, CJ Lueck, Karl J Friston, C Kennard, and RS Frackowiak. A direct
370 demonstration of functional specialization in human visual cortex. *Journal of neuroscience*, 11(3):
371 641–649, 1991.
- 372 Roger B Tootell, John B Reppas, Kenneth K Kwong, Rafael Malach, Richard T Born, Thomas J
373 Brady, Bruce R Rosen, and John W Belliveau. Functional analysis of human mt and related visual
374 cortical areas using magnetic resonance imaging. *Journal of Neuroscience*, 15(4):3215–3230,
375 1995.
- 376 Hide-aki Saito, Masao Yukie, Keiji Tanaka, Kazuo Hikosaka, Yoshiro Fukada, and Eiichi Iwai.
377 Integration of direction signals of image motion in the superior temporal sulcus of the macaque
378 monkey. *Journal of Neuroscience*, 6(1):145–157, 1986.
- 379 G Beckers and V Hömberg. Cerebral visual motion blindness: transitory akinetopsia induced by
380 transcranial magnetic stimulation of human area v5. *Proceedings of the Royal Society of London.
381 Series B: Biological Sciences*, 249(1325):173–178, 1992.
- 382 David Pitcher, Bradley Duchaine, and Vincent Walsh. Combined tms and fmri reveal dissociable
383 cortical pathways for dynamic and static face perception. *Current Biology*, 24(17):2066–2070,
384 2014.
- 385 Petra Vetter, Marie-Helene Grosbras, and Lars Muckli. Tms over v5 disrupts motion prediction.
386 *Cerebral cortex*, 25(4):1052–1059, 2015.
- 387 Zoe Kourtzi, Heinrich H Bühlhoff, Michael Erb, and Wolfgang Grodd. Object-selective responses in
388 the human motion area mt/mst. *Nature neuroscience*, 5(1):17–18, 2002.
- 389 Erez Freud, Jody C Culham, David C Plaut, and Marlene Behrmann. The large-scale organization of
390 shape processing in the ventral and dorsal pathways. *elife*, 6:e27576, 2017.
- 391 L Cornette, Patrick Dupont, A Rosier, Stefan Sunaert, P Van Hecke, J Michiels, Luc Mortelmans, and
392 GA Orban. Human brain regions involved in direction discrimination. *Journal of Neurophysiology*,
393 79(5):2749–2765, 1998.
- 394 Stefan Sunaert, Paul Van Hecke, Guy Marchal, and Guy A Orban. Motion-responsive regions of the
395 human brain. *Experimental brain research*, 127:355–370, 1999.
- 396 Sophia Robert, Leslie G Ungerleider, and Maryam Vaziri-Pashkam. Disentangling object category
397 representations driven by dynamic and static visual input. *Journal of Neuroscience*, 43(4):621–634,
398 2023.
- 399 Daniel L Yamins, Ha Hong, Charles Cadieu, and James J DiCarlo. Hierarchical modular optimization
400 of convolutional networks achieves representations similar to macaque it and human ventral stream.
401 *Advances in neural information processing systems*, 26, 2013.

- 402 Seyed-Mahdi Khaligh-Razavi and Nikolaus Kriegeskorte. Deep supervised, but not unsupervised,
403 models may explain it cortical representation. *PLoS computational biology*, 10(11):e1003915,
404 2014.
- 405 Chengxu Zhuang, Siming Yan, Aran Nayebi, Martin Schrimpf, Michael C Frank, James J DiCarlo, and
406 Daniel LK Yamins. Unsupervised neural network models of the ventral visual stream. *Proceedings
407 of the National Academy of Sciences*, 118(3):e2014196118, 2021.
- 408 Talia Konkle and George A Alvarez. A self-supervised domain-general learning framework for
409 human ventral stream representation. *Nature communications*, 13(1):491, 2022.
- 410 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
411 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
412 pages 770–778, 2016.
- 413 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
414 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer
415 vision and pattern recognition*, pages 16000–16009, 2022.
- 416 Yi Zhu, Zhenzhong Lan, Shawn Newsam, and Alexander Hauptmann. Hidden two-stream convolu-
417 tional networks for action recognition. In *Computer Vision–ACCV 2018: 14th Asian Conference
418 on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part III 14*,
419 pages 363–378. Springer, 2019.
- 420 Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-
421 efficient learners for self-supervised video pre-training. *Advances in neural information processing
422 systems*, 35:10078–10093, 2022.
- 423 Christoph Feichtenhofer, Yanghao Li, Kaiming He, et al. Masked autoencoders as spatiotemporal
424 learners. *Advances in neural information processing systems*, 35:35946–35958, 2022.
- 425 Yashar Behzadi, Khaled Restom, Joy Liau, and Thomas T Liu. A component based noise correction
426 method (compcor) for bold and perfusion based fmri. *Neuroimage*, 37(1):90–101, 2007.
- 427 Liang Wang, Ryan EB Mrcuzek, Michael J Arcaro, and Sabine Kastner. Probabilistic maps of visual
428 topography in human cortex. *Cerebral cortex*, 25(10):3911–3931, 2015.
- 429 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
430 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
431 pages 248–255. Ieee, 2009.
- 432 Jihoon Chung, Cheng-hsin Wu, Hsuan-ru Yang, Yu-Wing Tai, and Chi-Keung Tang. Haa500:
433 Human-centric atomic action dataset with curated videos. In *Proceedings of the IEEE/CVF
434 International Conference on Computer Vision*, pages 13465–13474, 2021.
- 435 Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan,
436 Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset.
437 *arXiv preprint arXiv:1705.06950*, 2017.
- 438 Rui Wang, Dongdong Chen, Zuxuan Wu, Yinpeng Chen, Xiyang Dai, Mengchen Liu, Lu Yuan, and
439 Yu-Gang Jiang. Masked video distillation: Rethinking masked feature modeling for self-supervised
440 video representation learning. In *Proceedings of the IEEE/CVF conference on computer vision
441 and pattern recognition*, pages 6312–6322, 2023a.
- 442 Olivia Rose, James Johnson, Binxu Wang, and Carlos R Ponce. Visual prototypes in the ventral
443 stream are attuned to complexity and gaze behavior. *Nature communications*, 12(1):6723, 2021.
- 444 Fenil R Doshi and Talia Konkle. Cortical topographic motifs emerge in a self-organized map of
445 object space. *Science Advances*, 9(25):eade8187, 2023.
- 446 Doris Y Tsao, Winrich A Freiwald, Roger BH Tootell, and Margaret S Livingstone. A cortical region
447 consisting entirely of face-selective cells. *Science*, 311(5761):670–674, 2006.

- 448 Alexander JE Kell, Daniel LK Yamins, Erica N Shook, Sam V Norman-Haignere, and Josh H
449 McDermott. A task-optimized neural network replicates human auditory behavior, predicts brain
450 responses, and reveals a cortical processing hierarchy. *Neuron*, 98(3):630–644, 2018.
- 451 Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao.
452 Videomae v2: Scaling video masked autoencoders with dual masking. In *Proceedings of the*
453 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14549–14560, 2023b.
- 454 Rahul Venkatesh, Honglin Chen, Klemen Kotar, Kevin Feigelis Wanhee Lee Daniel Bear, and Daniel
455 Yamins. Climbing the ladder of causation with counterfactual world modeling.
- 456 P McLeod. Preserved and impaired detection of structure from motion by a "motion-blind" patient.
457 *Visual Cognition*, 3(4):363–392, 1996.
- 458 Lucia M Vaina, Marjorie Lemay, Don C Bienfang, Albert Y Choi, and Ken Nakayama. Intact
459 "biological motion" and "structure from motion" perception in a patient with impaired motion
460 mechanisms: A case study. *Visual neuroscience*, 5(4):353–369, 1990.
- 461 Josef Zihl, D Von Cramon, and Norbert Mai. Selective disturbance of movement vision after bilateral
462 brain damage. *Brain*, 106(2):313–340, 1983.
- 463 Lorella Battelli, Patrick Cavanagh, and Ian M Thornton. Perception of biological motion in parietal
464 patients. *Neuropsychologia*, 41(13):1808–1816, 2003.
- 465 Sharon Gilaie-Dotan, Ayse Pinar Saygin, Lauren J Lorenzi, Geraint Rees, and Marlene Behrmann.
466 Ventral aspect of the visual form pathway is not critical for the perception of biological motion.
467 *Proceedings of the National Academy of Sciences*, 112(4):E361–E370, 2015.
- 468 Honglin Chen, Rahul Venkatesh, Yoni Friedman, Jiajun Wu, Joshua B Tenenbaum, Daniel LK
469 Yamins, and Daniel M Bear. Unsupervised segmentation in real-world images via spelke object
470 inference. In *European Conference on Computer Vision*, pages 719–735. Springer, 2022.
- 471 Elizabeth S Spelke. Principles of object perception. *Cognitive science*, 14(1):29–56, 1990.
- 472 Michael Hanke, Nico Adelhöfer, Daniel Kottke, Vittorio Iacovella, Ayan Sengupta, Falko R Kaule,
473 Roland Nigbur, Alexander Q Waite, Florian Baumgartner, and Jörg Stadler. A studyforrest extension,
474 simultaneous fmri and eye gaze recordings during prolonged natural stimulation. *Scientific*
475 *data*, 3(1):1–15, 2016.

476 **5.1 Supplementary materials**

477 **5.1.1 Training and testing the Two-stream CNN for action recognition**

478 We adopted the models in Zhu et al. [2019] and trained on the HAA500 dataset Chung et al. [2021].
479 The dataset contains over 591k labeled frames with 500 action classes. 85% of the data points were
480 used for training, 5% for validation, and 10% for testing 5.1.1. The training dataset was converted to
481 the Webdataset format, i.e., shards of tar files. We used 4 V100 GPUs and 8 workers to load the
482 dataset and train the models. All the analyses were performed on the same version of the movie that
483 was used to acquire fMRI responses in the StudyForrest dataset Hanke et al. [2016].
484 The *supervised static model* have a ResNet18 architecture He et al. [2016], and were trained for 47
485 epochs with a batch size of 128. The training was done with the stochastic gradient descent algorithm
486 with a 0.001 initial learning rate and a 0.0001 weight decay. During training, the gradients were
487 accumulated and backpropagated for every two batches. Each frame in an input batch is a 224×224
488 frame and was randomly flipped horizontally.
489 The *unsupervised dynamic model* was trained for 12 epochs with a batch size of 32 and an initial
490 learning rate of 0.01. No weight decay was used during training. Input to this model consists of a set
491 of 11 frames each with dimensions of 224×224 .
492 The *supervised dynamic model* was trained for 50 epochs with a batch size of 128 and an initial
493 learning rate of 0.001. A weight decay of 0.0005 was used to train the models, and the gradients
494 were accumulated and backpropagated every 5 batches.
495

Table 2: Test performance of models on the HAA500 dataset

Model	epochs	Performance	
		Top-1	Top-3
sup static	47	30.80%	49.38%
unsup dynamic + sup dynamic	12, 50	22.72%	37.90%

496 **NeurIPS Paper Checklist**

497 **1. Claims**

498 Question: Do the main claims made in the abstract and introduction accurately reflect the
499 paper’s contributions and scope?

500 Answer: [Yes]

501 Justification: [TODO]The abstract contains a summary of conclusions, and the introduction
502 clearly states the motivation and the contributions made in the paper.

503 Guidelines:

- 504 • The answer NA means that the abstract and introduction do not include the claims
505 made in the paper.
- 506 • The abstract and/or introduction should clearly state the claims made, including the
507 contributions made in the paper and important assumptions and limitations. A No or
508 NA answer to this question will not be perceived well by the reviewers.
- 509 • The claims made should match theoretical and experimental results, and reflect how
510 much the results can be expected to generalize to other settings.
- 511 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
512 are not attained by the paper.

513 **2. Limitations**

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

515 Answer: [Yes]

516 Justification: [TODO]A separate part is dedicated to describing the limitations.

517 Guidelines:

- 518 • The answer NA means that the paper has no limitation while the answer No means that
519 the paper has limitations, but those are not discussed in the paper.
- 520 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 521 • The paper should point out any strong assumptions and how robust the results are to
522 violations of these assumptions (e.g., independence assumptions, noiseless settings,
523 model well-specification, asymptotic approximations only holding locally). The authors
524 should reflect on how these assumptions might be violated in practice and what the
525 implications would be.
- 526 • The authors should reflect on the scope of the claims made, e.g., if the approach was
527 only tested on a few datasets or with a few runs. In general, empirical results often
528 depend on implicit assumptions, which should be articulated.
- 529 • The authors should reflect on the factors that influence the performance of the approach.
530 For example, a facial recognition algorithm may perform poorly when image resolution
531 is low or images are taken in low lighting. Or a speech-to-text system might not be
532 used reliably to provide closed captions for online lectures because it fails to handle
533 technical jargon.
- 534 • The authors should discuss the computational efficiency of the proposed algorithms
535 and how they scale with dataset size.
- 536 • If applicable, the authors should discuss possible limitations of their approach to
537 address problems of privacy and fairness.

538 • While the authors might fear that complete honesty about limitations might be used by
539 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
540 limitations that aren't acknowledged in the paper. The authors should use their best
541 judgment and recognize that individual actions in favor of transparency play an impor-
542 tant role in developing norms that preserve the integrity of the community. Reviewers
543 will be specifically instructed to not penalize honesty concerning limitations.

544 3. Theory Assumptions and Proofs

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

547 Answer: [NA]

548 Justification: [TODO]

549 Guidelines:

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

560 4. Experimental Result Reproducibility

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

564 Answer: [Yes]

565 Justification: [TODO] All the datasets used in the paper, including image and video datasets
566 as well as the fMRI dataset, are publicly available, and the proper citation is included in
567 the manuscript. All the steps required to reproduce the results are clearly explained in the
568 Methods section.

569 Guidelines:

- 570 • The answer NA means that the paper does not include experiments.
- 571 • If the paper includes experiments, a No answer to this question will not be perceived
572 well by the reviewers: Making the paper reproducible is important, regardless of
573 whether the code and data are provided or not.
- 574 • If the contribution is a dataset and/or model, the authors should describe the steps taken
575 to make their results reproducible or verifiable.
- 576 • Depending on the contribution, reproducibility can be accomplished in various ways.
577 For example, if the contribution is a novel architecture, describing the architecture fully
578 might suffice, or if the contribution is a specific model and empirical evaluation, it may
579 be necessary to either make it possible for others to replicate the model with the same
580 dataset, or provide access to the model. In general, releasing code and data is often
581 one good way to accomplish this, but reproducibility can also be provided via detailed
582 instructions for how to replicate the results, access to a hosted model (e.g., in the case
583 of a large language model), releasing of a model checkpoint, or other means that are
584 appropriate to the research performed.
- 585 • While NeurIPS does not require releasing code, the conference does require all submis-
586 sions to provide some reasonable avenue for reproducibility, which may depend on the
587 nature of the contribution. For example
 - 588 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
589 to reproduce that algorithm.
 - 590 (b) If the contribution is primarily a new model architecture, the paper should describe
591 the architecture clearly and fully.

- 592 (c) If the contribution is a new model (e.g., a large language model), then there should
593 either be a way to access this model for reproducing the results or a way to reproduce
594 the model (e.g., with an open-source dataset or instructions for how to construct
595 the dataset).
- 596 (d) We recognize that reproducibility may be tricky in some cases, in which case
597 authors are welcome to describe the particular way they provide for reproducibility.
598 In the case of closed-source models, it may be that access to the model is limited in
599 some way (e.g., to registered users), but it should be possible for other researchers
600 to have some path to reproducing or verifying the results.

601 5. Open access to data and code

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

605 Answer: [Yes]

606 Justification: [TODO]All the datasets used in the paper, including image and video datasets
607 as well as the fMRI dataset, are publicly available, and the proper citation is included in the
608 manuscript. The code will be released shortly after the submission.

609 Guidelines:

- 610 • The answer NA means that paper does not include experiments requiring code.
- 611 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
612 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 613 • While we encourage the release of code and data, we understand that this might not be
614 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
615 including code, unless this is central to the contribution (e.g., for a new open-source
616 benchmark).
- 617 • The instructions should contain the exact command and environment needed to run to
618 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
619 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 620 • The authors should provide instructions on data access and preparation, including how
621 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 622 • The authors should provide scripts to reproduce all experimental results for the new
623 proposed method and baselines. If only a subset of experiments are reproducible, they
624 should state which ones are omitted from the script and why.
- 625 • At submission time, to preserve anonymity, the authors should release anonymized
626 versions (if applicable).
- 627 • Providing as much information as possible in supplemental material (appended to the
628 paper) is recommended, but including URLs to data and code is permitted.

629 6. Experimental Setting/Details

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

633 Answer: [Yes]

634 Justification: [TODO]full details of the model that were trained from scratch are explained
635 in supplementary materials. Other models were adopted from the official repository, without
636 any changes to the original setting.

637 Guidelines:

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

643 7. Experiment Statistical Significance

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

646 Answer: [Yes]

647 Justification: [TODO]error bars were demonstrated in result figures. significant statistical
648 tests were reported throughout the manuscript

649 Guidelines:

- 650 • The answer NA means that the paper does not include experiments.
- 651 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
652 dence intervals, or statistical significance tests, at least for the experiments that support
653 the main claims of the paper.
- 654 • The factors of variability that the error bars are capturing should be clearly stated (for
655 example, train/test split, initialization, random drawing of some parameter, or overall
656 run with given experimental conditions).
- 657 • The method for calculating the error bars should be explained (closed form formula,
658 call to a library function, bootstrap, etc.)
- 659 • The assumptions made should be given (e.g., Normally distributed errors).
- 660 • It should be clear whether the error bar is the standard deviation or the standard error
661 of the mean.
- 662 • It is OK to report 1-sigma error bars, but one should state it. The authors should
663 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
664 of Normality of errors is not verified.
- 665 • For asymmetric distributions, the authors should be careful not to show in tables or
666 figures symmetric error bars that would yield results that are out of range (e.g. negative
667 error rates).
- 668 • If error bars are reported in tables or plots, The authors should explain in the text how
669 they were calculated and reference the corresponding figures or tables in the text.

670 8. Experiments Compute Resources

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

674 Answer: [Yes]

675 Justification: [TODO]Yes. The details are explained in the supplementary materials

676 Guidelines:

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

685 9. Code Of Ethics

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

688 Answer: [Yes]

689 Justification: [TODO]

690 Guidelines:

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

696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747

10. Broader Impacts

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

Answer: [NA]

Justification: [TODO]

Guidelines:

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

11. Safeguards

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

Answer: [NA]

Justification: [TODO]

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: [TODO] Proper citation and URL are included in the manuscript.

Guidelines:

- 748 • The answer NA means that the paper does not use existing assets.
- 749 • The authors should cite the original paper that produced the code package or dataset.
- 750 • The authors should state which version of the asset is used and, if possible, include a
- 751 URL.
- 752 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 753 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 754 service of that source should be provided.
- 755 • If assets are released, the license, copyright information, and terms of use in the
- 756 package should be provided. For popular datasets, paperswithcode.com/datasets
- 757 has curated licenses for some datasets. Their licensing guide can help determine the
- 758 license of a dataset.
- 759 • For existing datasets that are re-packaged, both the original license and the license of
- 760 the derived asset (if it has changed) should be provided.
- 761 • If this information is not available online, the authors are encouraged to reach out to
- 762 the asset's creators.

763 13. New Assets

764 Question: Are new assets introduced in the paper well documented and is the documentation

765 provided alongside the assets?

766 Answer: [NA]

767 Justification: **[TODO]**

768 Guidelines:

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

777 14. Crowdsourcing and Research with Human Subjects

778 Question: For crowdsourcing experiments and research with human subjects, does the paper

779 include the full text of instructions given to participants and screenshots, if applicable, as

780 well as details about compensation (if any)?

781 Answer: [Yes] The human data is adopted from a publicly available dataset. For provided

782 essential instruction for referring to the study and the accompanying dataset.

783 Justification: **[TODO]**

784 Guidelines:

- 785 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 786 human subjects.
- 787 • Including this information in the supplemental material is fine, but if the main contribu-
- 788 tion of the paper involves human subjects, then as much detail as possible should be
- 789 included in the main paper.
- 790 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 791 or other labor should be paid at least the minimum wage in the country of the data
- 792 collector.

793 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human

794 Subjects

795 Question: Does the paper describe potential risks incurred by study participants, whether

796 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

797 approvals (or an equivalent approval/review based on the requirements of your country or

798 institution) were obtained?

799 Answer: [NA]

800 Justification: **[TODO]**We did not collect any data and only used publicly available repositories.
801

802 Guidelines:

- 803 • The answer NA means that the paper does not involve crowdsourcing nor research with
804 human subjects.
- 805 • Depending on the country in which research is conducted, IRB approval (or equivalent)
806 may be required for any human subjects research. If you obtained IRB approval, you
807 should clearly state this in the paper.
- 808 • We recognize that the procedures for this may vary significantly between institutions
809 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
810 guidelines for their institution.
- 811 • For initial submissions, do not include any information that would break anonymity (if
812 applicable), such as the institution conducting the review.