

1 **A large and rich EEG dataset for modeling human visual object recognition**
2 Alessandro T. Gifford^{1,2,3}, Kshitij Dwivedi⁴, Gemma Roig⁴, Radoslaw M. Cichy^{1,2,3,5}

3
4 ¹Department of Education and Psychology, Freie Universität Berlin, Berlin, Germany

5 ²Charité - Universitätsmedizin Berlin, Einstein Center for Neurosciences Berlin,
6 Berlin, Germany

7 ³Bernstein Center for Computational Neuroscience Berlin, Berlin, Germany

8 ⁴Department of Computer Science, Goethe University, Frankfurt am Main, Germany

9 ⁵Berlin School of Mind and Brain, Humboldt-Universität zu Berlin, Berlin, Germany

10

11 **Author ORCIDs:**

12 Alessandro T. Gifford: <https://orcid.org/0000-0002-8923-9477>

13 Kshitij Dwivedi: <https://orcid.org/0000-0001-6442-7140>

14 Gemma Roig: <https://orcid.org/0000-0002-6439-8076>

15 Radoslaw M. Cichy: <https://orcid.org/0000-0003-4190-6071>

16

17 **Correspondence:**

18 alessandro.gifford@gmail.com

19

20 **Abstract**

21 The human brain achieves visual object recognition through multiple stages of
22 nonlinear transformations operating at a millisecond scale. To predict and explain
23 these rapid transformations, computational neuroscientists employ machine learning
24 modeling techniques. However, state-of-the-art models require massive amounts of
25 data to properly train, and to the present day there is a lack of vast brain datasets
26 which extensively sample the temporal dynamics of visual object recognition. Here
27 we collected a large and rich dataset of high temporal resolution EEG responses to
28 images of objects on a natural background. This dataset includes 10 participants,
29 each with 82,160 trials spanning 16,740 image conditions. Through computational
30 modeling we established the quality of this dataset in five ways. First, we trained
31 linearizing encoding models that successfully synthesized the EEG responses to
32 arbitrary images. Second, we correctly identified the recorded EEG data image
33 conditions in a zero-shot fashion, using EEG synthesized responses to hundreds of
34 thousands of candidate image conditions. Third, we show that both the high number
35 of conditions as well as the trial repetitions of the EEG dataset contribute to the
36 trained models' prediction accuracy. Fourth, we built encoding models whose
37 predictions well generalize to novel participants. Fifth, we demonstrate full end-to-
38 end training of randomly initialized DNNs that output M/EEG responses for arbitrary
39 input images. We release this dataset as a tool to foster research in visual
40 neuroscience and computer vision.

41 **Introduction**

42 Visual object recognition is a complex cognitive function that is computationally
43 solved in multiple nonlinear stages by the human brain (Marr, 1980; Goodale &
44 Milner, 1992; Van Essen et al., 1992; Riesenhuber & Poggio, 1999; Ullman, 2000;
45 Grill-Spector et al., 2001; Malach et al., 2002; Carandini et al., 2005). Through these
46 stages information is transformed from representations of simple visual features
47 such as oriented edges to representations of object categories (Tanaka, 1996;
48 Logothetis & Sheinberg, 1996). To understand the principles of these
49 transformations, computational neuroscientists build and employ mathematical
50 models that predict the brain responses to arbitrary visual stimuli and explain their
51 underlying neural mechanisms (Wu et al., 2006; Guest & Martin, 2021). The
52 performance of these models benefits from training with large datasets: as an
53 example, deep neural networks (DNNs) (Fukushima et al., 1982), the current state-
54 of-the-art computational models of the visual brain (Yamins & DiCarlo, 2016; Cichy &
55 Kaiser, 2019; Kietzmann et al., 2019a; Richards et al., 2019; Saxe et al., 2021), are
56 trained on hundreds of thousands of different data points (Russakovsky et al., 2015).
57 Yet, due to the difficulty of brain data acquisition, neuroscientific datasets usually
58 comprise no more than a few thousand trials per participant and a limited number of
59 conditions (Kay et al., 2008; Cichy et al., 2014; Horikawa & Kamitani, 2017).

60 To address the data hunger of current modeling goals, recently pioneering
61 efforts have been taken to record large datasets of functional magnetic resonance
62 imaging (fMRI) responses to images (Chang et al., 2019; Allen et al., 2021).
63 However, while providing excellent spatial resolution, fMRI data lacks the temporal
64 resolution to resolve neural dynamics at the level at which they occur. Since neurons
65 communicate at millisecond scales, high temporal resolution neural data is a crucial
66 component for building models of the visual brain (Thorpe et al., 1996; van de
67 Nieuwenhuijzen et al., 2013; Cichy et al., 2014; Harel et al., 2016; Seeliger et al.,
68 2017; Bankson et al., 2018; Dijkstra et al., 2018). Thus, in the present study we
69 collected a large millisecond resolution electroencephalography (EEG) dataset of
70 human brain responses to images of objects on a natural background. We
71 extensively sampled 10 participants, each being presented with 16,740 image
72 conditions repeated over 82,160 trials from the THINGS database (Hebart et al.,
73 2019) by using a time-efficient rapid serial visual presentation (RSVP) paradigm
74 (Intraub, 1981; Keyser et al., 2001; Grootswagers et al., 2019).

75 We then leveraged the unprecedented size and richness of our dataset to
76 train and evaluate DNN-based linearizing and end-to-end encoding models (Wu et
77 al., 2006; Kay et al., 2008; Naselaris et al., 2011; van Gerven, 2017; Seeliger et al.,
78 2017; Kriegeskorte & Douglas, 2019; Seeliger et al., 2021; Khosla et al., 2021; Allen
79 et al., 2021) that synthesize EEG responses to arbitrary images. The results
80 showcase the quality of the dataset and its potential for computational modeling in
81 five ways. First, the synthesized EEG data is strongly resemblant to its biological
82 counterpart, with robust predictions even at single participants' level. Second, we
83 built zero-shot identification algorithms (Kay et al., 2008; Seeliger et al., 2017;
84 Horikawa & Kamitani, 2017) that achieved high performance accuracies even when

85 identifying among very large candidate image conditions set sizes: 81.3% for a set
86 size of 200 candidate image conditions, 21.15% for a set size of 150,000 candidate
87 image conditions, and extrapolated accuracy $> 10\%$ for a set size of 3,650,000
88 candidate image conditions, where chance $\leq 0.5\%$. Third, we show that both the high
89 number of conditions as well as the trial repetitions of the dataset contribute to the
90 trained models' prediction accuracy. Fourth, we demonstrate that the encoding
91 models' predictions generalize to novel participants. Fifth, for the first time to our
92 knowledge we demonstrate full end-to-end training (Seeliger et al., 2021; Khosla et
93 al., 2021; Allen et al., 2021) of randomly initialized DNNs that output M/EEG
94 responses for arbitrary input images.

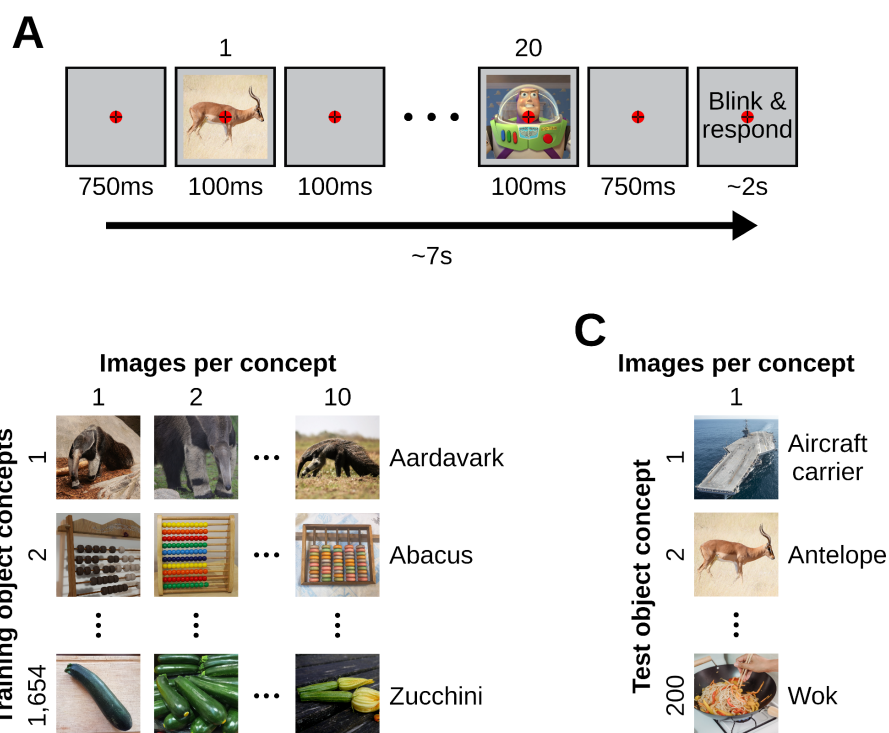
95 We release the dataset as a tool to foster research in computational
96 neuroscience and to bridge the gap between biological and artificial vision. We
97 believe this will be of great use to further understanding of visual object recognition
98 through the development of high-temporal resolution computational models of the
99 visual brain, and to optimize artificial intelligence models through biological
100 intelligence data (Sinz et al., 2019; Hassabis et al., 2017; Ullman, 2019; Toneva &
101 Wehbe, 2019; Yang et al., 2022). Also all code used to generate the presented
102 results accompanies the data release.

103 **Results**

104 **A large and rich EEG dataset of visual responses to objects on a natural**
 105 **background**

106 We used a RSVP paradigm (Intraub, 1981; Keyser et al., 2001; Grootswagers et
 107 al., 2019) to collect a large EEG dataset of visual responses to images of objects on
 108 a natural background (**Figure 1A**). This dataset contains data for 10 participants who
 109 viewed 16,540 training image conditions (**Figure 1B**) and 200 test image conditions
 110 (**Figure 1C**) coming from the THINGS database (Hebart et al., 2019). To allow for
 111 unbiased modeling the training and test images did not have any overlapping object
 112 concepts. We presented each training image condition 4 times and each test image
 113 condition 80 times, for a total of 82,160 image trials per participant over the course of
 114 four sessions. Thanks to the time-efficiency of the RSVP paradigm we collected up
 115 to 15 times more data than other typical recent M/EEG datasets used for modeling
 116 (Cichy et al., 2014; Seeliger et al., 2017). This allowed us to extensively sample
 117 single participants while drastically reducing the experimental time. During
 118 preprocessing we epoched the EEG recordings from -200ms to 800ms with respect
 119 to image onset, downsampled the resulting image epoch trials to 100 time points,
 120 and retained only the 17 occipital and parietal channels. As the basis of all further
 121 data assessment we aggregated the EEG recordings into a *biological training*
 122 (BioTrain) data matrix of shape (16,540 training image conditions × 4 condition
 123 repetitions × 17 EEG channels × 100 EEG time points) and a *biological test*
 124 (BioTest) data matrix of shape (200 test image conditions × 80 condition repetitions
 125 × 17 EEG channels × 100 EEG time points), for each participant. Providing this EEG
 126 data in its raw as well as preprocessed form is the major contribution of this
 127 resource.

128



129
130

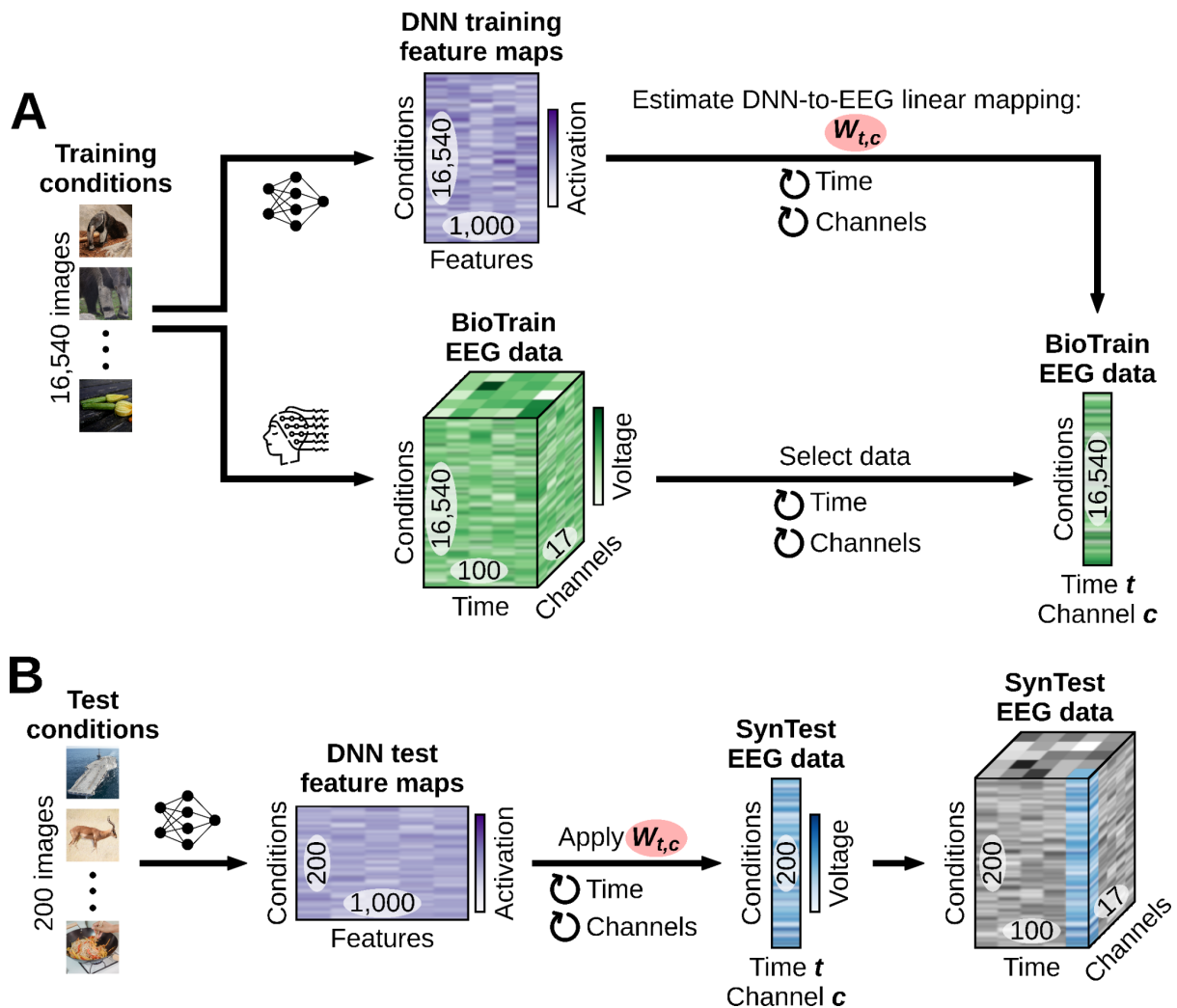
Figure 1. Experimental paradigm and stimuli images. (A) We presented participants with

131 images of objects on a natural background using a RSVP paradigm. The paradigm
132 consisted of rapid serial sequences of 20 images. Every sequence started with 750ms of
133 blank screen, then each image was presented centrally for 100ms and a stimulus onset
134 asynchrony (SOA) of 200ms, and it ended with another 750ms of blank screen. After
135 every rapid sequence there were up to 2s during which we instructed participants to first
136 blink and then report, with a keypress, whether the target image appeared in the
137 sequence. We asked participants to gaze at a central bull's eye fixation target present
138 throughout the entire experiment. **(B)** The training image partition contains 1,654 object
139 concepts of 10 images each, for a total of 16,540 image conditions. **(C)** The test image
140 partition contains 200 object concepts of 1 image each, for a total of 200 image
141 conditions.

142

143 **Building linearizing encoding models of EEG visual responses**

144 We then assessed the suitability of this dataset for the development of computational
145 models of the visual brain. We employed the training and test data, respectively, to
146 build and evaluate linearizing encoding models which predict individual participant's
147 EEG visual responses to arbitrary images (Wu et al., 2006; Kay et al., 2008;
148 Naselaris et al., 2011; van Gerven, 2017; Kriegeskorte & Douglas, 2019). We based
149 our encoding algorithm on deep neural networks (DNNs), connectionist models
150 which in the last decade have excelled in predicting human and non-human primate
151 visual brain responses (Cadieu et al., 2014; Yamins et al., 2014; Güçlü & van
152 Gerven, 2015; Storrs et al., 2021). The building of encoding models involved two
153 steps. In the first step we non-linearly transformed the image pixel values using four
154 DNNs pre-trained on ILSVRC-2012 (Russakovsky et al., 2015) commonly used for
155 modeling brain responses: AlexNet (Krizhevsky, 2014), ResNet-50 (He et al, 2016),
156 CORnet-S (Kubilius et al., 2019) and MoCo (Chen et al., 2020). Separately for each
157 DNN we fed the training and test images, extracted the corresponding feature maps
158 across all layers, appended the layers' data together and downsampled it to 1,000
159 principal components using principal component analysis (PCA), resulting in the
160 training DNN feature maps matrix of shape (16,540 training image conditions \times
161 1,000 features) and the test DNN feature maps matrix of shape (200 test image
162 conditions \times 1000 features). In the second step we fitted the weights $\mathbf{W}_{t,c}$ of several
163 linear regressions that independently predicted each EEG feature's response (i.e.,
164 the EEG activity at each combination of time points (t) and channels (c)) to the
165 training images by linearly combining the training feature maps of each DNN (**Figure**
166 **2A**). We then multiplied the learned $\mathbf{W}_{t,c}$ with the test DNN feature maps, obtaining
167 the *synthetic test* (SynTest) EEG data matrix of shape (200 test image conditions \times
168 17 EEG channels \times 100 EEG time points) (**Figure 2B**). Following this procedure we
169 obtained different instances of SynTest data for each participant and DNN.

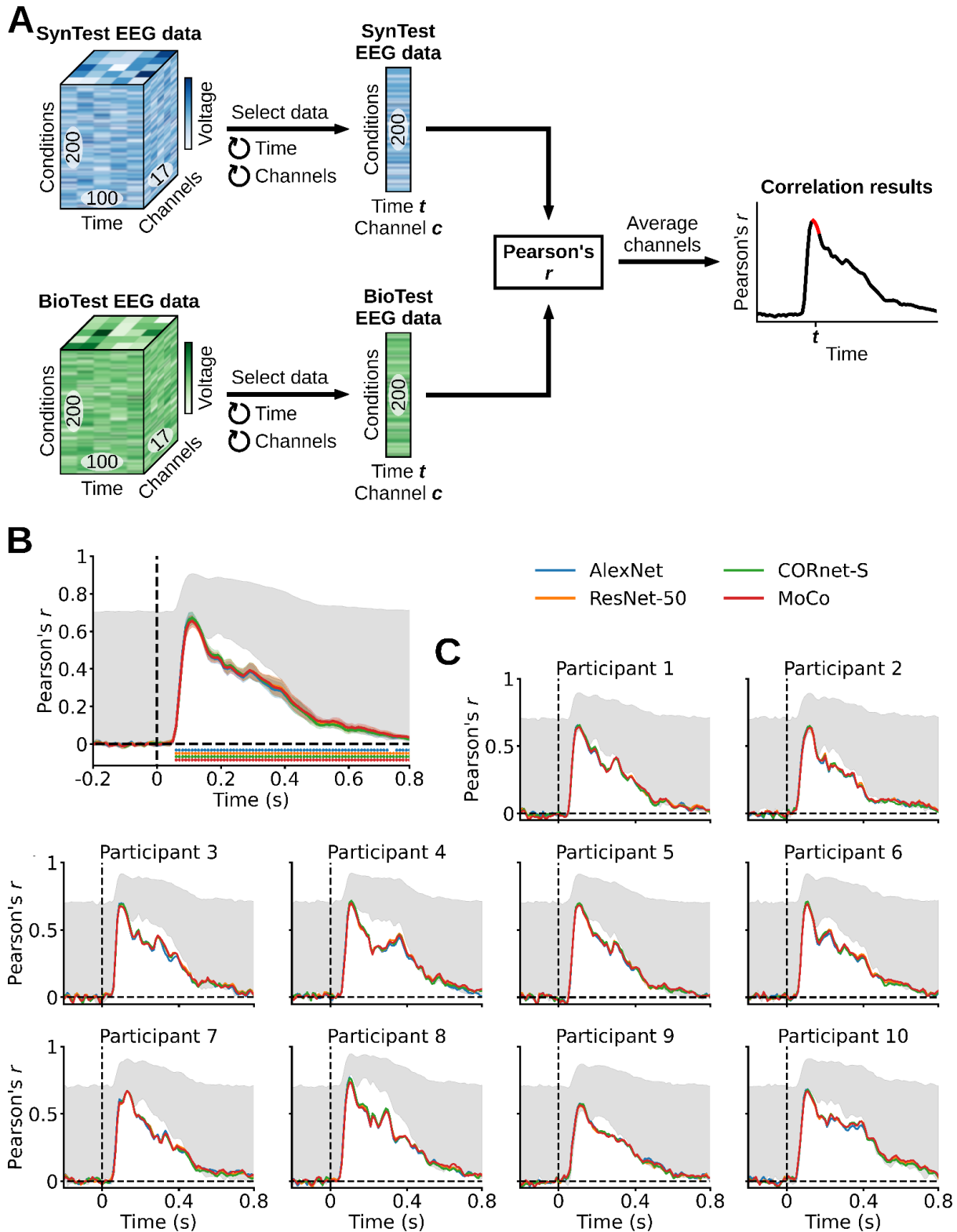


170
171
172
173
174
175
176
177
178
179
180

Figure 2. Linearizing encoding algorithm. For ease of visualization, here and in the following figures we omit the EEG condition repetitions dimension. **(A)** Through the training image conditions we obtained the training DNN feature maps and the BioTrain EEG data, and used them to build linearizing encoding models of EEG visual responses. For each combination of EEG features (time points (t) and channels (c)) we estimated the weights $W_{t,c}$ of a linear regression using the corresponding single-feature BioTrain data as criterion and the training images DNN feature maps as predictors. **(B)** To obtain the SynTest EEG data we extracted the DNN feature maps of the test images, and multiplied them with the estimated $W_{t,c}$.

181 The BioTest EEG data is well predicted by linearizing encoding models

182 To evaluate the linearizing encoding models' predictive power we quantified the
183 similarity between the SynTest data and the BioTest data through a Pearson's
184 correlation (**Figure 3A**). We correlated each SynTest data EEG feature (i.e., each
185 combination of EEG time points (t) and channels (c)) with the corresponding BioTest
186 data feature (across the 200 test image conditions), resulting in a correlation
187 coefficient matrix of shape (17 EEG channels \times 100 EEG time points). We then
188 averaged this matrix across the channels dimension, obtaining a correlation
189 coefficient result vector with 100 components, one for each EEG time point.



190
191
192
193
194
195
196
197
198

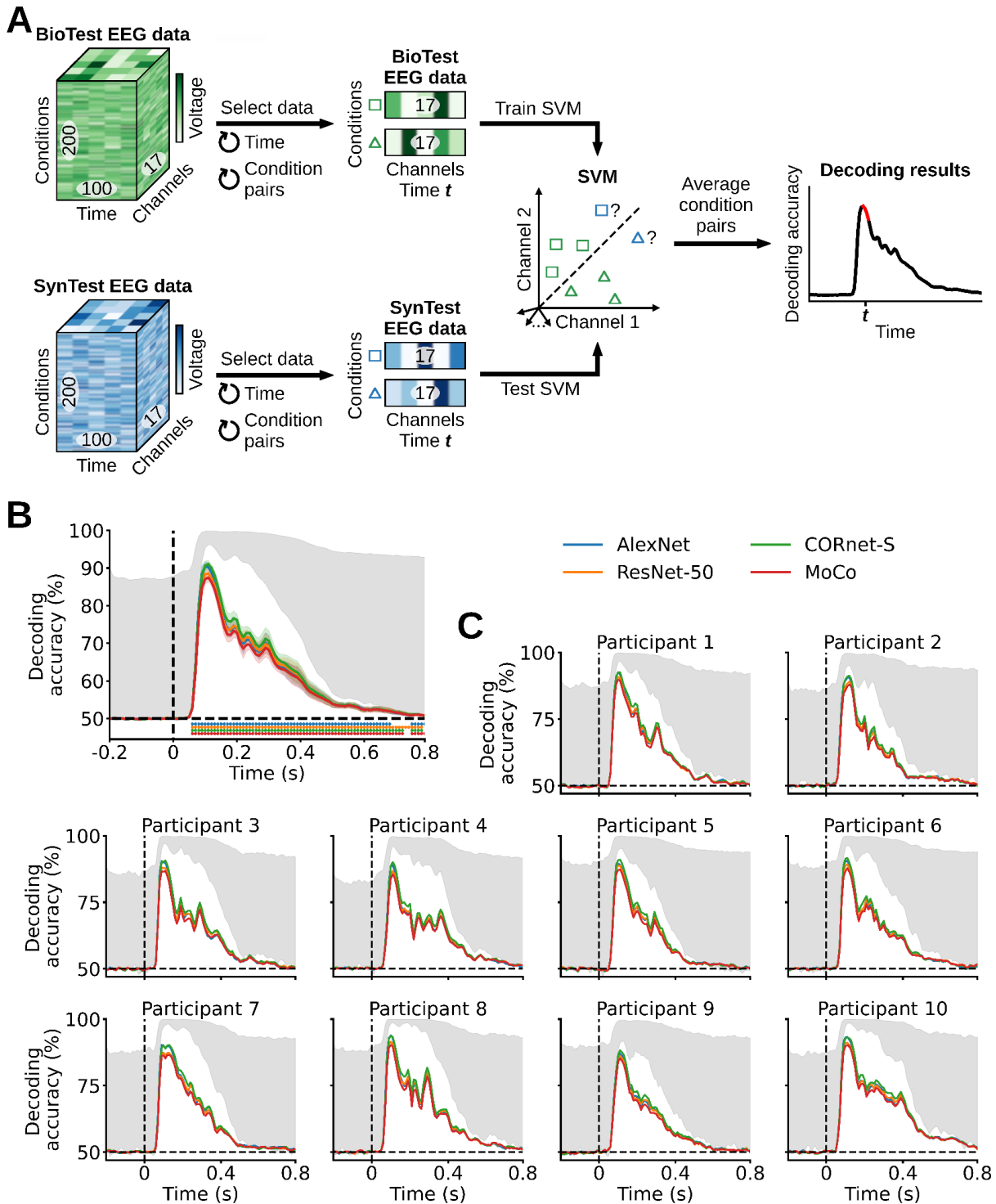
Figure 3. Evaluating the linearizing encoding models' prediction accuracy through a correlation analysis. (A) We correlated each combination of SynTest EEG data features (time points (t) and channels (c)) with the corresponding combination of BioTest EEG data features, across the 200 test image conditions, and then averaged the correlation coefficients across channels. (B) Correlation results averaged across participants. The SynTest data is significantly correlated to the BioTest data from 60ms after stimulus onset until the end of the EEG epoch ($P < 0.05$, one-sample one-sided t-test, Bonferroni-corrected), with a peak at 110ms. (C) Individual participants' results. Error margins reflect

199 95% confidence intervals. Rows of asterisks indicate significant time points ($P < 0.05$,
200 one-sample one-sided t-tests, Bonferroni-corrected). In gray is the area between the
201 noise ceiling lower and upper bounds, the black dashed vertical lines indicate onset of
202 image presentation, and the black dashed horizontal lines indicate the chance level of no
203 experimental effect.

204

205 As a complementary way to evaluate the linearizing encoding models'
206 predictive power we quantified the similarity between the SynTest data and the
207 BioTest data through decoding (**Figure 4A**). Decoding is a commonly used method
208 in computational neuroscience which exploits similar information present between
209 the trials of each experimental condition to classify neural data (Haynes & Rees,
210 2006; Mur et al., 2009). If the SynTest data and the BioTest data have similar
211 information, a decoding algorithm trained on the BioTest data would generalize its
212 performance also to the SynTest data. We tested this through pairwise decoding: we
213 trained linear support vector machines (SVMs) to perform binary classification
214 between each pair of the 200 BioTest data image conditions, and then tested them
215 on the corresponding pairs of SynTest data image conditions. We performed this
216 analysis independently for each time point (t), resulting in a decoding accuracy
217 matrix of shape (19,900 image condition pairs \times 100 EEG time points). We then
218 averaged this matrix across the image condition pairs dimension, obtaining a
219 decoding accuracy result vector with 100 components, one for each EEG time point.

220 We observe that the correlation results averaged across participants start
221 being significant at 60ms after stimulus onset, and remain significantly above chance
222 until the end of the EEG epoch at 800ms ($P < 0.05$, one-sample one-sided t-test,
223 Bonferroni-corrected). Significant correlation peaks occur for all DNNs at 110ms after
224 stimulus onset, with AlexNet, ResNet-50, CORnet-S and MoCo having correlation
225 coefficients of, respectively, 0.67, 0.66, 0.67 and 0.66 ($P < 0.05$, one-sample one-
226 sided t-test, Bonferroni-corrected), where the chance level is 0 (**Figure 3B**).
227 Similarly, the pairwise decoding results averaged across participants start being
228 significant at 60ms after stimulus onset, with significant effects present until the end
229 of the EEG epoch at 800ms ($P < 0.05$, one-sample one-sided t-test, Bonferroni-
230 corrected). Significant decoding peaks occur for all DNNs at 100-110ms after
231 stimulus onset, with AlexNet, ResNet-50, CORnet-S and MoCo having decoding
232 accuracies of, respectively, 90.37%, 88.57%, 91.06% and 87.45% ($P < 0.05$, one-
233 sample one-sided t-test, Bonferroni-corrected), where the chance level is 50%
234 (**Figure 4B**). All participants yielded qualitatively similar results (**Figure 3C**, **Figure**
235 **4C**). Taken together, these results show that the linearizing encoding models are
236 successful in predicting EEG data which robustly and significantly resembles its
237 biological counterpart. Further, they show that each participant's neural responses
238 can be consistently predicted in isolation, thus highlighting the quality of the visual
239 information contained in our EEG dataset and its potential for the development of
240 new high-temporal resolution models and theories of the visual brain.



241
242
243
244
245
246
247
248
249
250
251
252

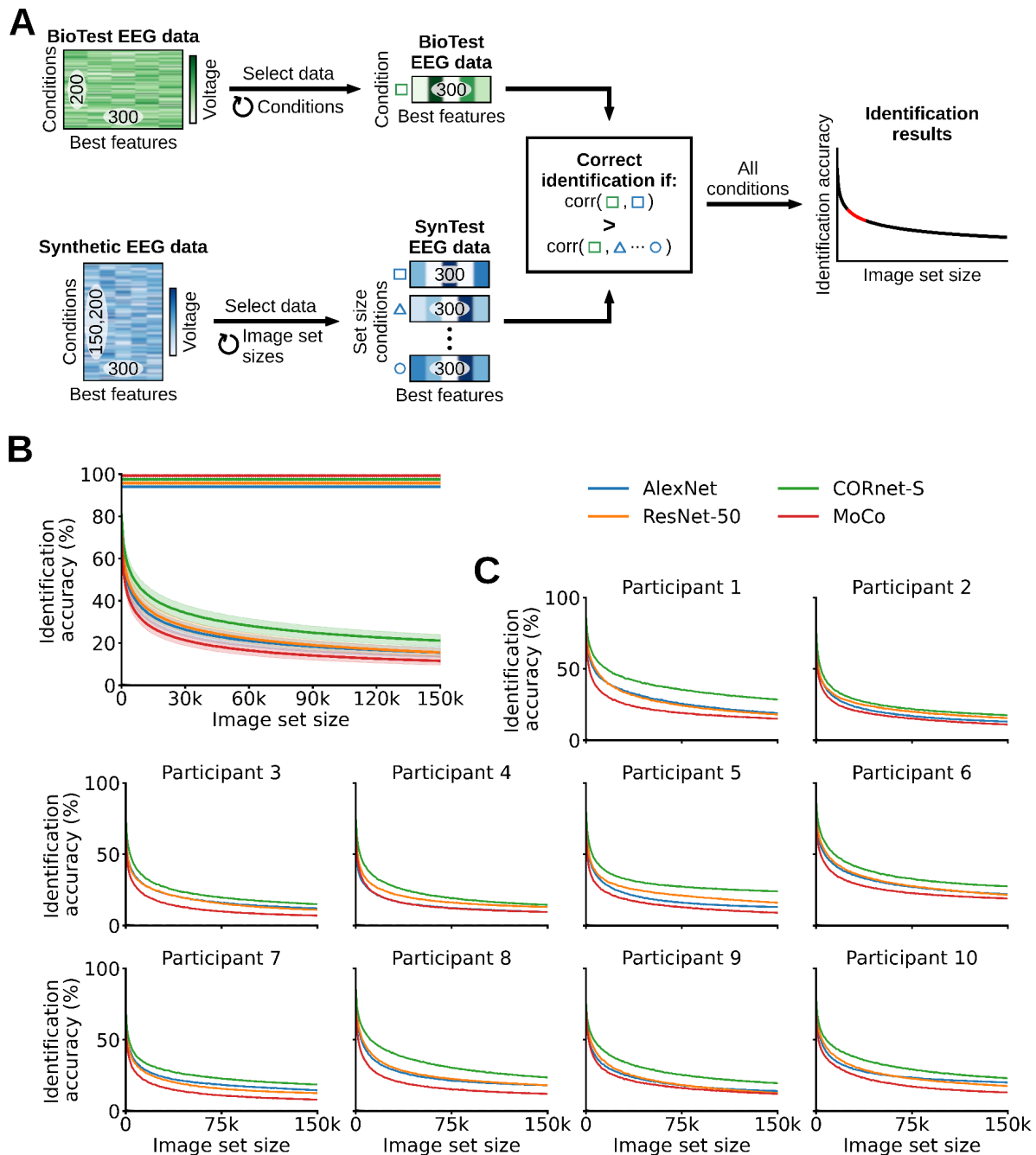
Figure 4. Evaluating the linearizing encoding models' prediction accuracy through a pairwise decoding analysis. (A) At each time point (t) we trained an SVM to classify between two BioTest EEG data image conditions (using the channels vectors) and tested it on the two corresponding SynTest EEG data image conditions. We repeated this procedure across all image condition pairs, and then averaged the decoding accuracies across pairs. (B) Pairwise decoding results averaged across participants. The linear classifiers trained on the BioTest data significantly decode the SynTest data from 60ms after stimulus onset until the end of the EEG epoch ($P < 0.05$, one-sample one-sided t-test, Bonferroni-corrected), with peaks at 100-110ms. (C) Individual participants' results. Error margins, asterisks, gray area and black dashed lines as in Figure 3.

253 **The BioTest data is significantly identified in a zero-shot fashion using**
254 **synthesized data of up to 150,200 candidate images**

255 The previous analyses showed that our linearizing encoding models synthesize EEG
256 data that significantly resembles its biological counterpart. Here we explored whether
257 we can leverage this high prediction accuracy to build algorithms that identify the
258 image conditions of the BioTest data in a zero-shot fashion, namely, that identify
259 arbitrary image conditions without prior training. If possible, this would contribute to
260 the goal of building models capable of identifying potentially infinite neural data
261 conditions on which they were never trained (Kay et al., 2008; Seeliger et al., 2017;
262 Horikawa & Kamitani, 2017) (**Figure 5A**). For the identification we used the SynTest
263 and the *synthetic Imagenet* (SynImagenet) data, where the latter consisted of the
264 synthesized EEG responses to the 150,000 validation and test images coming from
265 the ILSVRC-2012 image set (Russakovsky et al., 2015), organized in a data matrix
266 of shape (150,000 image conditions \times 17 EEG channels \times 100 EEG time points).
267 Importantly, those images did not overlap with either the image set for which EEG
268 data was recorded. The further analysis involved two steps: feature selection and
269 identification.

270 In the feature selection step we retained the 300 EEG channels and time
271 points best predicted by the encoding models, as narrowing down the EEG data to
272 these features improved the identification accuracy. In detail, we synthesized the
273 EEG responses to the 16,540 training images, obtaining the *synthetic train*
274 (SynTrain) data matrix of shape (16,540 training image conditions \times 17 EEG
275 channels \times 100 EEG time points). We then correlated each BioTrain data feature
276 (i.e., each combination of EEG channels and EEG time points) with the
277 corresponding SynTrain data feature (across the 16,540 training image conditions),
278 and only retained the 300 SynTest, BioTest and SynImagenet data EEG features
279 corresponding to the 300 highest correlation scores. This resulted in feature vectors
280 of 300 components for each image condition.

281 In the identification step we correlated the feature vectors of each BioTest
282 data image condition with the feature vectors of all the candidate image conditions,
283 where the candidate image conditions corresponded to the SynTest data image
284 conditions plus a varying amount of SynImagenet data image conditions. We
285 increased the set sizes of the SynImagenet candidate image conditions from 0 to
286 150,000 with steps of 1,000 images (for a total of 151 set sizes), and performed the
287 identification at every set size. At each set size a BioTest data image condition is
288 considered correctly identified if the correlation coefficient between its feature vector
289 and the feature vector of the corresponding SynTest data image condition is higher
290 than the correlation coefficients between its feature vector and the feature vectors of
291 all other candidate image conditions. We calculated identification accuracies through
292 the ratio of successfully decoded image conditions over all 200 BioTest image
293 conditions, obtaining a zero-shot identification result vector with 151 components,
294 one for each candidate image set size. The results of the correct SynTest data
295 image condition falling within the three or ten most correlated image conditions can
296 be seen in **Supplementary Figures 3-4** and **Supplementary Tables 1-2**.



297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312

Figure 5. Zero-shot identification of the BioTest data using the SynTest data and the synthesized EEG visual responses to the 150,000 ILSVRC-2012 validation and test image conditions (SynImagenet). **(A)** We correlated the best features of each BioTest data condition with different image set sizes of candidate synthetic image conditions (SynTest + SynImagenet data). At each image set size, a BioTest data condition is correctly identified if it is mostly correlated to its corresponding SynTest data condition, among all other synthetic data conditions. **(B)** Zero-shot identification results averaged across participants. With a SynImagenet set size of 0 the synthesized data of AlexNet, ResNet-50, CORnet-S, MoCo significantly identify the BioTest data with accuracies of, respectively, 75.05%, 75.85%, 81.3%, 70.9%. ($P < 0.05$, one-sample one-sided t-test, Bonferroni-corrected). With a SynImagenet set size of 150,000 the synthesized data of AlexNet, ResNet-50, CORnet-S, MoCo significantly identify the BioTest data with accuracies of, respectively, 15.5%, 15.55%, 21.15%, 11.55%. **(C)** Individual participants' results. Rows of asterisks indicate significant image set sizes ($P < 0.05$, one-sample one-sided t-tests, Bonferroni-corrected). Error margins and black dashed lines as in **Figure 3**.

313 The zero-shot identification results averaged across participants are
 314 significant for all SynImagenet set sizes ($P < 0.05$, one-sample one-sided t-test,
 315 Bonferroni-corrected). With a SynImagenet set size of 0 (corresponding to using only
 316 the 200 SynTest data image conditions as candidate image conditions) the BioTest
 317 data image conditions are identified by AlexNet, ResNet-50, CORnet-S and MoCo
 318 with accuracies of, respectively, 75.05%, 75.85%, 81.3%, 70.9%, where the chance
 319 level is equal to $1 / 200$ test image conditions = 0.5%. As the SynImagenet set size
 320 increases the identification accuracies monotonically decrease. With a SynImagenet
 321 set size of 150,000 (corresponding to using the 200 SynTest data plus the 150,000
 322 SynImagenet data image conditions as candidate image conditions) the BioTest data
 323 image conditions are identified by AlexNet, ResNet-50, CORnet-S and MoCo with
 324 accuracies of, respectively, 15.5%, 15.55%, 21.15%, 11.55%, where the chance
 325 level is equal to $1 / (200 \text{ test image conditions} + 150,000 \text{ ILSVRC-2012 image}$
 326 $\text{conditions}) < 10^{-5}\%$ (**Figure 5B**). To extrapolate the identification accuracies to
 327 potentially larger candidate image set sizes we fit a power-law function to the results.
 328 We averaged the extrapolations across participants, and found that the identification
 329 accuracy would remain above 10% with a candidate image set size of 914,000 for
 330 AlexNet, 588,000 for ResNet-50, 3,650,000 for CORnet-S and 348,000 for MoCo,
 331 and above 0.5% (the original chance level) with a candidate image set size of
 332 $2.18\text{E}+11$ for AlexNet, $3.43\text{E}+09$ for ResNet-50, $1.62\text{E}+13$ for CORnet-S and
 333 $1.11\text{E}+10$ for MoCo (**Table 1**). All participants yielded qualitatively similar results
 334 (**Figures 5C; Table 1**). These results demonstrate that our dataset allows building
 335 algorithms that reliably identify arbitrary neural data conditions, in a zero-shot
 336 fashion, among millions of possible alternatives.
 337

	Identification accuracy < 10%				Identification accuracy < 0.5%			
	AlexNet	ResNet-50	CORnet-S	MoCo	AlexNet	ResNet-50	CORnet-S	MoCo
Participant 1	7.35E+05	6.31E+05	5.68E+06	5.30E+05	1.10E+09	8.44E+08	1.66E+11	4.56E+09
Participant 2	2.97E+05	5.38E+05	9.02E+05	1.93E+05	7.27E+08	2.76E+09	1.24E+10	1.96E+08
Participant 3	2.35E+05	1.79E+05	4.38E+05	7.31E+04	3.53E+08	7.34E+07	1.12E+09	2.56E+07
Participant 4	1.33E+05	3.28E+05	3.67E+05	1.27E+05	4.70E+08	3.22E+09	5.38E+08	5.18E+08
Participant 5	3.12E+05	6.15E+05	1.65E+07	1.24E+05	2.27E+09	3.22E+09	1.61E+14	6.03E+07
Participant 6	2.01E+06	1.41E+06	7.06E+06	1.58E+06	3.91E+10	7.95E+09	6.80E+11	1.03E+11
Participant 7	5.27E+05	2.80E+05	1.25E+06	9.11E+04	8.31E+09	1.80E+09	3.54E+10	6.71E+07
Participant 8	1.27E+06	9.67E+05	1.63E+06	2.50E+05	8.13E+10	1.20E+10	6.64E+09	1.23E+09
Participant 9	3.50E+05	2.46E+05	9.08E+05	2.39E+05	8.76E+08	9.52E+07	2.83E+09	3.40E+08
Participant 10	3.28E+06	6.93E+05	1.84E+06	2.71E+05	2.04E+12	2.31E+09	1.46E+10	3.17E+08
Average	9.14E+05	5.88E+05	3.65E+06	3.48E+05	2.18E+11	3.43E+09	1.62E+13	1.11E+10

338 **Table 1.** Extrapolation of the zero-shot identification accuracy as a function of candidate
 339 image set sizes. The values in the table indicate the candidate image set sizes required
 340 for the identification accuracy to drop below 10% and 0.5%.
 341

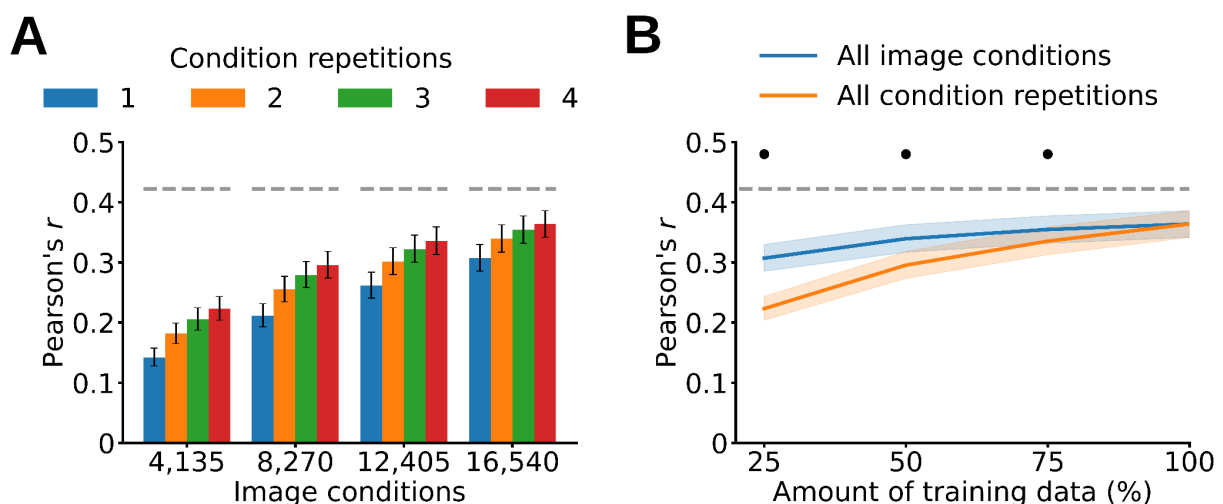
342
 343 **The amount of training image conditions and condition repetitions both**
 344 **contribute to modeling quality**

345 To understand which aspects of our EEG dataset contribute to its successful
 346 modeling we examined the linearizing encoding models' prediction accuracy as a

347 function of the amount of trials with which they are trained. The amount of training
348 trials is determined by two factors: the number of image conditions and the number
349 of EEG repetitions per each image condition. Both factors may improve the modeling
350 of neural responses in different ways, as high numbers of image conditions lead to a
351 richer training set which more comprehensively samples the representational space
352 underlying vision, and high numbers of condition repetitions increase the signal to
353 noise ratio (SNR) of the training set.

354 To disentangle the effect of both factors we trained linearizing encoding
355 models using different quartiles of training image conditions (4,135, 8,270, 12,405,
356 16,540) and condition repetitions (1, 2, 3, 4), and tested their predictions through the
357 correlation analysis. We performed an ANOVA on the correlation results averaged
358 over participants, EEG features (all channels; time points between 60-500ms) and
359 DNN models, and observed a significant effect of both number of image conditions
360 and condition repetitions, along with a significant interaction of the two factors ($P <$
361 0.05, two-way repeated measures ANOVA) (**Figure 6A**). All participants yielded
362 qualitatively similar results (**Supplementary Figure 5**). This suggests that the
363 amount of image conditions and condition repetitions both improve the modeling of
364 neural data.

365



366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

Figure 6. Linearizing encoding models' prediction accuracy as a function of training data. **(A)** Training linearizing encoding models using different quartiles of image conditions and condition repetitions result in a significant effect of both factors ($P < 0.05$, two-way repeated measures ANOVA). **(B)** Training linearizing encoding models using all image conditions leads to higher prediction accuracies than training them using all condition repetitions ($P < 0.05$, repeated measures two-sided t-test, Bonferroni-corrected). The gray dashed line represents the noise ceiling lower bound. The asterisks indicate a significant difference between all image conditions and all condition repetitions ($P < 0.05$, repeated measures two-sided t-test, Bonferroni-corrected). Error margins and gray dashed lines as in **Figure 3**.

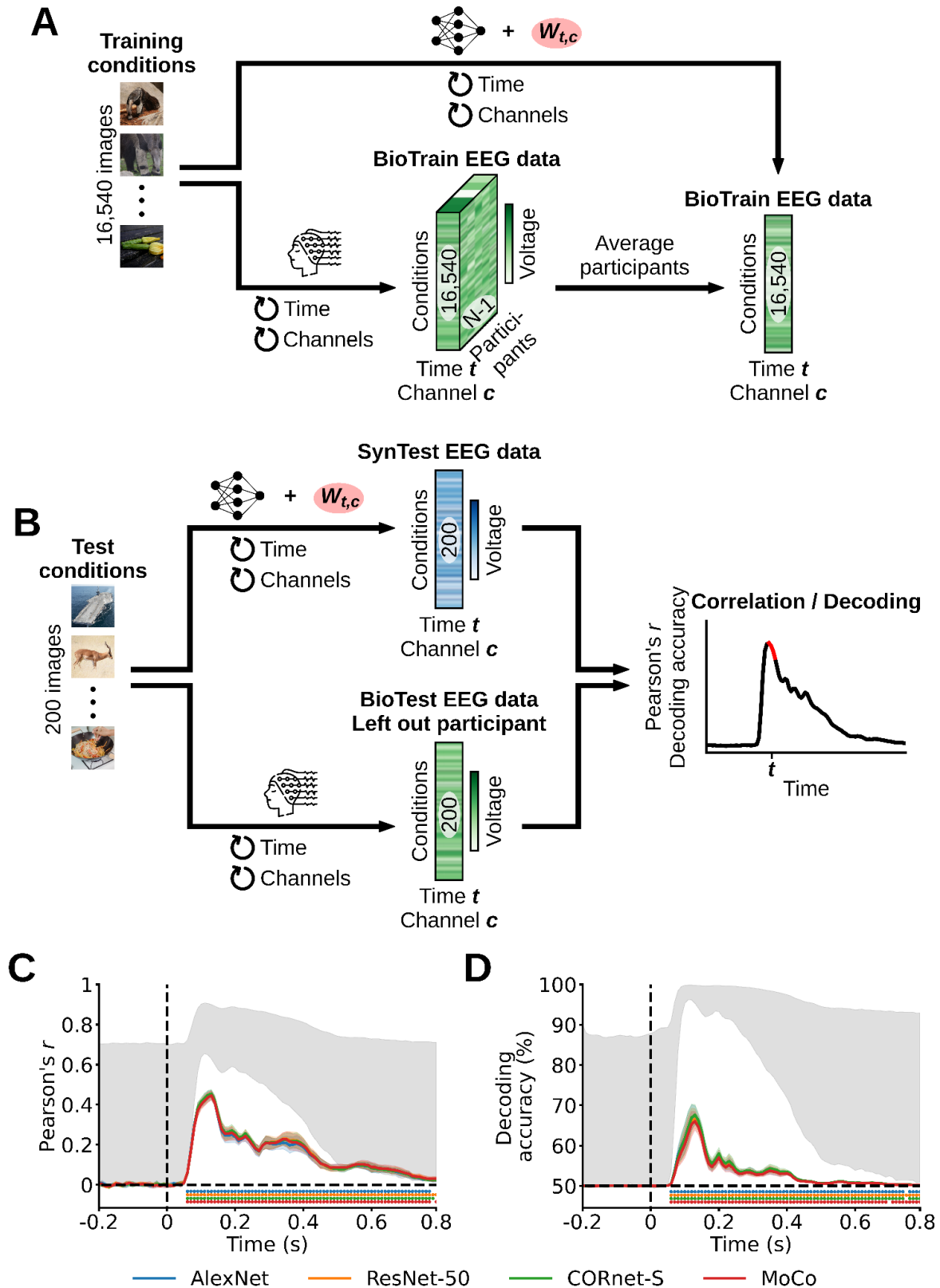
378 We then asked which of the two factors contributes more to the linearizing
379 encoding models' prediction accuracy. For this we compared model prediction
380 accuracy for cases where the number of repetitions or conditions differed, but the
381 total number of trials was the same. As we had four trial repetitions, we divided the

382 total amount of training trials into quartiles (25%, 50%, 75% and 100% of the total
383 training trials). At each quartile we trained linearizing encoding models using all
384 image conditions and the quartile's percentage of condition repetitions, and tested
385 their predictions through the correlation analysis. For example, at the first quartile we
386 trained linearizing encoding models using all image conditions and one condition
387 repetition, corresponding to 25% of the total training data. To compare, we repeated
388 the same procedure while using all condition repetitions and the quartile's
389 percentage of image conditions. The correlation results averaged across
390 participants, EEG features (all channels; time points between 60-500ms) and DNNs
391 show that using all image conditions (and quartiles of condition repetitions) leads to
392 higher prediction accuracies than using all condition repetitions (and quartiles of
393 image conditions) ($P < 0.05$, repeated measures two-sided t-test, Bonferroni-
394 corrected) (**Figure 6B**). All participants yielded qualitatively similar results
395 (**Supplementary Figure 6**). This indicates that although both factors improve the
396 modeling of neural data, the amount of image conditions does so here to a larger
397 extent.

398

399 **The linearizing encoding models' predictions generalize across participants**

400 Next we explored whether our linearizing encoding models' predictions generalize to
401 new participants. We asked: Can we accurately synthesize a participant's EEG
402 responses without using any of their data for the encoding models' training? If
403 possible, our dataset could serve as a useful benchmark for the development and
404 assessment of methods that combine EEG data across participants (Haxby et al.,
405 2020; Richard et al., 2020; Kwon et al., 2019; Zhang et al., 2021). To verify this we
406 trained linearizing encoding models on the averaged SynTrain EEG data of all minus
407 one participants (**Figure 7A**), and tested their predictions against the BioTest data of
408 the left out participant through the correlation and pairwise decoding analyses
409 (**Figure 7B**). We repeated this procedure for all participants.



410
411
412
413
414
415
416
417

Figure 7. Evaluating the prediction accuracy of linearizing encoding models which generalize to novel participants, through correlation and pairwise decoding analyses. (A) We trained linearizing encoding models on the averaged SynTrain EEG data of all minus one participants. (B) We tested the encoding models' predictions against the BioTest data of the left out participant through the correlation and pairwise decoding analyses. (C) Correlation results averaged across participants. The SynTest data is significantly correlated to the BioTest data from 60ms after stimulus onset until the end of the EEG

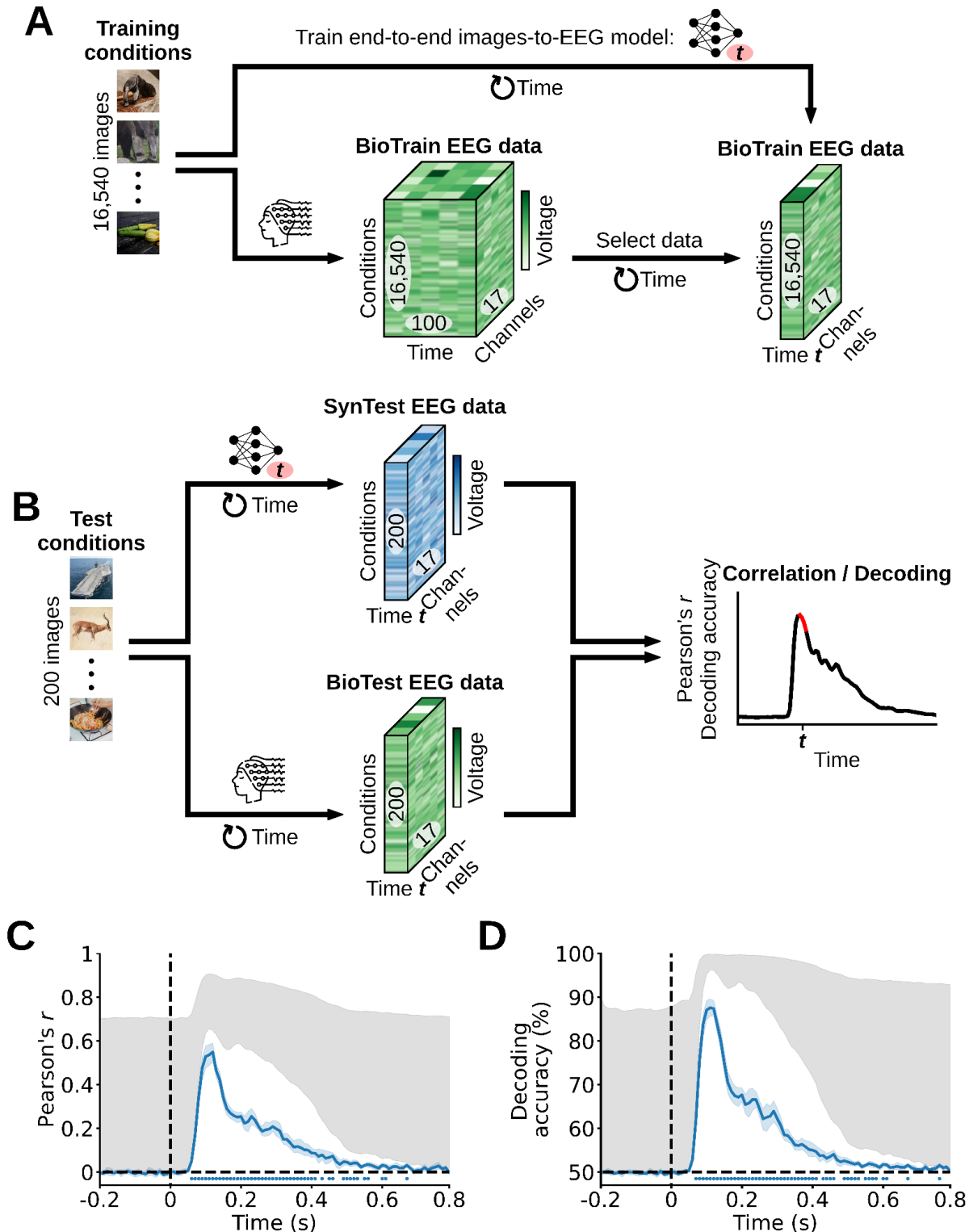
418 epoch ($P < 0.05$, one-sample one-sided t-test, Bonferroni-corrected), with a peak at
419 130ms. **(D)** Pairwise decoding results averaged across participants. The linear classifiers
420 trained on the BioTest data significantly decode the SynTest data from 60ms after
421 stimulus onset until the end of the EEG epoch ($P < 0.05$, one-sample one-sided t-test,
422 Bonferroni-corrected), with a peak at 130ms. Error margins, asterisks, gray area and
423 black dashed lines as in **Figure 3**.

424
425 When averaging the Pearson correlation coefficients across participants we
426 observe that the correlation between the SynTest data and the BioTest data starts
427 being significant at 60ms after stimulus onset, and remains significantly above
428 chance until the end of the EEG epoch at 800ms ($P < 0.05$, one-sample one-sided t-
429 test, Bonferroni-corrected). Significant correlation peaks occur for all DNNs at 130ms
430 after stimulus onset, with AlexNet, ResNet-50, CORnet-S and MoCo having
431 correlation coefficients of, respectively, 0.45, 0.46, 0.46, 0.44 ($P < 0.05$, one-sample
432 one-sided t-test, Bonferroni-corrected), where the chance level is 0 (**Figure 7C**).
433 Likewise, the decoding accuracies averaged across participants start being
434 significant at 60ms after stimulus onset, with significant effects present until the end
435 of the EEG epoch at 800ms ($P < 0.05$, one-sample one-sided t-test, Bonferroni-
436 corrected). Significant decoding peaks occur for all DNNs at 130ms after stimulus
437 onset, with AlexNet, ResNet-50, CORnet-S and MoCo having decoding accuracies
438 of, respectively, 67.44%, 66.62%, 67.63%, 66.01% ($P < 0.05$, one-sample one-sided
439 t-test, Bonferroni-corrected), where the chance level is 50% (**Figure 7D**). In both
440 analyses all participants yielded qualitatively similar results (**Supplementary**
441 **Figures 7-8**). This shows that our EEG dataset is a suitable testing ground for
442 methods which generalize and combine EEG data across participants.

443 444 **The BioTest EEG data is successfully predicted by end-to-end encoding** 445 **models based on the AlexNet architecture**

446 So far we predicted the synthetic data through the linearizing encoding framework,
447 which relied on DNNs pre-trained on an image classification task. An alternative
448 encoding approach, named end-to-end encoding (Seeliger et al., 2021; Khosla et al.,
449 2021; Allen et al., 2021), is based on DNNs trained from scratch to predict the neural
450 responses to arbitrary images. This direct infusion of brain data during the model's
451 learning could lead to DNNs having internal representations that more closely match
452 the properties of the visual brain (Sinz et al., 2019; Allen et al., 2021). However, with
453 a few exceptions (Seeliger et al., 2021; Khosla et al., 2021; Allen et al., 2021), the
454 development of end-to-end encoding models has been so far prohibitive due to the
455 large amount of data needed to train a DNN in combination with the small size of
456 existing brain datasets. Thus, in this final analysis we exploited the largeness and
457 richness of our EEG dataset to train randomly initialized AlexNet architectures to
458 synthesize the EEG responses to images, independently for each participant. We
459 started by replacing AlexNet's 1000-neurons output layer with a 17-neurons layer,
460 where each neuron corresponded to one of the 17 EEG channels. Then, for each
461 EEG time point (t) we trained one such model to predict the multi-channel EEG
462 responses to visual stimuli using the 16,540 training images as inputs and the

463 corresponding BioTrain data as output targets (**Figure 8A**). Finally, we deployed the
 464 trained networks to synthesize the EEG responses to the 200 test images and
 465 evaluated their prediction accuracy through the same correlation and pairwise
 466 decoding analyses described above (**Figure 8B**).
 467



468
 469
 470
 471

Figure 8. Evaluating the end-to-end encoding models' prediction accuracy through correlation and pairwise decoding analyses. (A) At each EEG time point (t) we trained one encoding model end-to-end to predict the SynTrain data channel responses using

472 the corresponding training images as input. **(B)** We used the trained encoding models to
473 predict the SynTest data using the test images as input, and evaluated their prediction
474 accuracies by comparing the SynTest and BioTest data through correlation and pairwise
475 decoding analyses. **(C)** Correlation results averaged across participants. The SynTest
476 data is significantly correlated to the BioTest data from 60ms after stimulus onset until
477 670ms ($P < 0.05$, one-sample one-sided t-test, Bonferroni-corrected), with a peak at
478 120ms. **(D)** Pairwise decoding results averaged across participants. The linear classifiers
479 trained on the BioTest data significantly decode the SynTest data from 70ms after
480 stimulus onset until 760ms ($P < 0.05$, one-sample one-sided t-test, Bonferroni-corrected),
481 with peaks at 110ms. Error margins, asterisks, gray area and black dashed lines as in
482 **Figure 3**.

483
484 We observe that the correlation results averaged across participants start
485 being significant at 60ms after stimulus onset, with a correlation coefficient peak at
486 120ms of 0.55, and have significant effects until 670ms ($P < 0.05$, one-sample one-
487 sided t-test, Bonferroni-corrected) (**Figure 8C**). Similarly, the pairwise decoding
488 results averaged across participants start being significant at 70ms after stimulus
489 onset, with a decoding accuracy peak at 110ms of 87.59%, and have significant
490 effects until 760ms ($P < 0.05$, one-sample one-sided t-test, Bonferroni-corrected)
491 (**Figure 8D**). All participants yielded qualitatively similar results (**Supplementary**
492 **Figures 9-10**). This proves that our EEG dataset allows for the successful training of
493 DNNs in an end-to-end fashion, paving the way for a stronger symbiosis between
494 brain data and deep learning models benefitting both neuroscientists interested in
495 building better models of the brain (Seeliger et al., 2021; Khosla et al., 2021; Allen et
496 al., 2021) and computer scientists interested in creating better performing and more
497 brain-like artificial intelligence algorithms through inductive biases from biological
498 intelligence (Sinz et al., 2019; Hassabis et al., 2017; Ullman, 2019; Toneva &
499 Wehbe, 2019; Yang et al., 2022).

500 **Discussion**

501 **Summary**

502 We used a RSVP paradigm (Intraub, 1981; Keyzers et al., 2001; Grootswagers et
503 al., 2019) to collect a large and rich EEG dataset of neural responses to images of
504 real-world objects on a natural background, which we release as a tool to foster
505 research in vision neuroscience and computer vision. Through computational
506 modeling we established the quality of this dataset in five ways. First, we trained
507 linearizing encoding models (Wu et al., 2006; Kay et al., 2008; Naselaris et al., 2011;
508 van Gerven, 2017; Kriegeskorte & Douglas, 2019) that successfully synthesized the
509 EEG responses to arbitrary images. Second, we correctly identified the recorded
510 EEG data image conditions in a zero-shot fashion (Kay et al., 2008; Seeliger et al.,
511 2017; Horikawa & Kamitani, 2017), using EEG synthesized responses to hundreds
512 of thousands of candidate image conditions. Third, we show that both the high
513 number of conditions as well as the trial repetitions of the EEG dataset contribute to
514 the trained model's prediction accuracy. Fourth, we built encoding models whose
515 predictions well generalize to novel participants. Fifth, we demonstrate full end-to-
516 end training (Seeliger et al., 2021; Khosla et al., 2021; Allen et al., 2021) of randomly
517 initialized DNNs that output M/EEG responses for arbitrary input images.

518

519 **Size matters**

520 In the last years cognitive neuroscientists have drastically increased the scope of
521 their recordings from datasets with dozens of stimuli to datasets comprising several
522 thousands of stimuli per participant (Chang et al., 2019; Naselaris et al., 2021; Allen
523 et al., 2021). Compared to their predecessors, these large datasets more
524 comprehensively sample the visual space and interact better with modern data-
525 hungry machine learning algorithms. In this spirit we extensively sampled 10
526 participants with 82,160 trials spanning 16,740 image conditions, and showed how
527 this unprecedented size contributes to high modeling performances. We released the
528 data in both its raw and preprocessed format ready for modeling to allow researchers
529 of different analytical perspectives to use the dataset in their preferred way
530 immediately. We believe the largeness of this dataset holds great promise for
531 neuroscientists interested in further improving theories and models of the visual
532 brain, as well as computer scientists interested in improving machine vision models
533 through biological vision constraints (Haxby et al., 2020; Richard et al., 2020; Kwon
534 et al., 2019; Zhang et al., 2021).

535

536 **Linearizing encoding modeling**

537 We showcased the potential of the dataset for modeling visual responses by building
538 linearizing encoding algorithms (Wu et al., 2006; Kay et al., 2008; Naselaris et al.,
539 2011; van Gerven, 2017; Kriegeskorte & Douglas, 2019) that predicted EEG visual
540 responses to arbitrary images. Through correlation and decoding analyses we
541 showed that the encoding models synthesized data which significantly resembles its
542 biological counterpart robustly and consistently across all participants. These results
543 highlight the signal quality of the visual information present in the EEG dataset,

544 making it a promising candidate for the development of new high-temporal resolution
545 models and theories of the neural dynamics of vision capable of predicting, decoding
546 and even explaining visual object recognition.

547

548 **Zero-shot identification**

549 Decoding models in neuroscience typically classify between only a few data
550 conditions, while relying on data exemplars from these same conditions to train
551 (Haynes & Rees, 2006; Mur et al., 2009). As a result, their performance fails to
552 generalize to the unlimited space of different brain states. Here we exploited the
553 prediction accuracy of the synthesized EEG responses to build zero-shot
554 identification algorithms that identify potentially infinite neural data image conditions,
555 without the need of prior training (Kay et al., 2008; Seeliger et al., 2017; Horikawa &
556 Kamitani, 2017). Through this framework we identified the BioTest EEG image
557 conditions between hundreds of thousands of candidate image conditions. Even
558 when the identification algorithm failed to assign the correct image condition to the
559 biological EEG responses, we show that it nevertheless selected a considerable
560 amount (up to 45%) of the correct image conditions as the first three or ten choices
561 (**Supplementary Figures 3-4**). These results suggest that our dataset is a good
562 starting ground for the future creation of zero-shot identification algorithms to be
563 deployed not only in research, but also in cutting-edge brain computer interface
564 (BCI) technology (Abiri et al., 2019; Petit et al., 2021).

565

566 **Both number of image conditions and condition repetitions determine dataset 567 quality**

568 Building linearizing encoding algorithms with different amounts of training data
569 revealed that the encoding models' prediction accuracies are significantly affected by
570 both the amount of EEG image conditions (to a higher extent) and repetitions of
571 measurements (to a lower extent). Given that the noise ceiling lower bound estimate
572 is not reached, these findings suggest that the prediction accuracy of the linearizing
573 encoding models would have benefited from either more training data trials, or from
574 a training dataset with the same amount of trials but having more image conditions
575 and less repetitions of measurements. Based on these observations, for future
576 dataset collections we recommend prioritizing the amount of stimuli conditions over
577 the amount of repetitions of measurements.

578

579 **Inter-participant predictions**

580 Typically, computational models in neuroscience are trained and evaluated on the
581 data of single participants (Kay et al., 2008; Yamins et al., 2014; Güçlü & van
582 Gerven, 2015; Seeliger et al., 2017; Horikawa & Kamitani, 2017). While this
583 approach is well motivated by the neural idiosyncrasies of every individual (Charest
584 et al., 2014), it fails to produce models that leverage the shared information across
585 multiple brains. Here we show that our encoding models well predict out-of-set
586 participants, indicating that our dataset is a suitable testing ground for methods
587 which generalize and combine neural data across participants, as well as for BCI

588 technology that can be readily used on novel participants without the need of
589 calibration (Haxby et al., 2020; Richard et al., 2020; Kwon et al., 2019; Zhang et al.,
590 2021).

591

592 **End-to-end encoding**

593 So far limitations in neural dataset sizes led computational neuroscientists to model
594 brain data mostly using pre-trained DNNs (Cadieu et al., 2014; Yamins et al., 2014;
595 Güçlü & van Gerven, 2015; Naselaris et al., 2015; Seeliger et al., 2017). Here, we
596 leveraged the largeness and richness of our dataset to demonstrate, for the first time
597 to our knowledge with EEG data, the feasibility of training a randomly initialized
598 AlexNet architecture to predict the neural responses to arbitrary images in an end-to-
599 end fashion (Seeliger et al., 2021; Khosla et al., 2021; Allen et al., 2021). The end-to-
600 end approach opens the doors to training complex computational algorithms directly
601 with brain data, potentially leading to models which more closely mimic the internal
602 representation of the visual system (Sinz et al., 2019; Allen et al., 2021). This will in
603 turn make it possible for computer scientists to use the neural representations of
604 biological systems as inductive biases to improve artificial systems under the
605 assumption that increasing the brain-likeness of computer models could increase
606 their performance in tasks at which humans excel (Sinz et al., 2019; Hassabis et al.,
607 2017; Ullman, 2019; Toneva & Wehbe, 2019; Yang et al., 2022).

608

609 **Dataset limitations**

610 A major limitation of our dataset is the backward and forward noise introduced by the
611 very short (200ms) stimulus onset asynchronies (SOAs) of the RSVP paradigm
612 (Intraub, 1981; Keyser et al., 2001; Grootswagers et al., 2019). The forward noise
613 at a given EEG image trial comes from the ongoing neural activity of the previous
614 trial, whereas the backward noise coming from the following trial starts from around
615 260ms after image onset, which corresponds to the SOA length plus the amount of
616 time required for the visual information to travel from the retina to the visual cortex.
617 Despite these noise sources, we showed that the visual responses are successfully
618 predicted during the entire EEG epoch. We believe that averaging the EEG image
619 conditions across several repetitions of measurements reduced the noise, and that
620 the backward noise was further mitigated given that the neural processing required
621 to detect and recognize object categories can be achieved in the first 150ms of
622 vision (Thorpe et al., 1996; Rousselet et al., 2002). A second limitation concerns the
623 ecological validity of the dataset. The stimuli images used consisted of objects
624 presented at foveal vision with natural backgrounds that have little clutter.
625 Furthermore, participants were asked to constantly gaze at a central fixation target.
626 This does not truthfully represent human vision, in which objects are perceived and
627 recognized also when at the periphery of the visual field, within cluttered scenes, and
628 while making eye movements. However, our results pave the way towards studies
629 aiming to provide large amounts of EEG responses recorded during more natural
630 viewing conditions.

631

632 **Contribution to the THINGS initiative**

633 The visual brain is an ensemble of billions of neurons communicating with high
634 spatial and temporal precision. However, neither current neural recording modalities,
635 nor single lab efforts can capture this complexity. This motivates the need to
636 integrate data across different imaging modalities and labs. To address this
637 challenge, the so-called THINGS initiative promotes using the THINGS database to
638 collect and share behavioral and neuroscientific datasets for the same set of images
639 - also used here - among vision researchers (<https://things-initiative.org/>). We
640 contribute to the initiative by providing rich high temporal resolution EEG data, that
641 complements other datasets in both a within- and between-modality fashion. As an
642 example of the within-modality fashion, Grootswagers and collaborators recently
643 published an EEG dataset of visual responses to images coming from the THINGS
644 database (Grootswager et al., 2022). While their dataset comprises more
645 participants and image conditions, our dataset provides more repetitions of
646 measurements, longer image presentation latencies, and an extensive assessment
647 of the dataset's potential based on the resulting high signal-to-noise ratio.
648 Researchers can choose between one or the other based on the nature,
649 requirements and constraints of their own experiments. As an example of the
650 between-modality fashion, our data can be used to make bridges from the EEG
651 temporal domain to, for example, the fMRI spatial domain through modeling
652 frameworks such as representational similarity analysis (Kriegeskorte et al., 2008;
653 Cichy et al., 2014; Cichy et al., 2016; Khaligh-Razavi et al., 2017), thus promoting a
654 more integrated understanding of the neural basis of visual object recognition.

655

656 **Comparing the modeling results of the four DNNs evaluated**

657 The size and quality of our dataset make it a good candidate for the comparison of
658 predictive and explanatory models of the visual brain (Schrimpf et al., 2020; Cichy et
659 al., 2019). Here, we built encoding models using four DNNs: despite the prediction
660 accuracies of these DNNs being overall qualitatively similar (Storrs et al., 2021), the
661 results of our analyses suggest that the EEG data is best predicted by the linearizing
662 encoding models based on the recurrent CORnet-S architecture (Kubilius et al.,
663 2019). This supports a growing amount of literature which asserts that recurrent
664 computations are critical for object recognition along the ventral stream, and
665 therefore any model of visual object processing must also take recurrency into
666 account (Kriegeskorte, 2015; Spoerer et al., 2017; Mohsenzadeh et al., 2018; Kar et
667 al., 2019; Kietzmann et al., 2019b; Kubilius et al., 2019; Rajaei et al., 2019; van
668 Bergen & Kriegeskorte, 2020). However, this interpretation should be taken with a
669 grain of salt as we compared DNNs differing not only in the hypotheses of visual
670 processing they embedded (e.g., recurrent vs. pure feedforward visual processing),
671 but also in potential confounding factors such as architecture and complexity.

672

673 **The modeling accuracy is not homogeneous across time**

674 As expected, the prediction accuracies of our encoding algorithms did not reach the
675 noise ceiling level (**Supplementary Figures 1-2**), indicating that our dataset is well

676 suited for further model improvements. Interestingly, we found that the modeling
677 accuracy is not homogeneous across time: the differences between the prediction
678 accuracy and the noise ceiling are smaller in the first 100ms after image onset, and
679 peak at 200-220ms, suggesting that the four DNNs used are more similar to the
680 brain at earlier stages of visual processing. This calls for future improvements in
681 model building (e.g., by including high-level visual semantics or improving biological
682 realism of the models) to more closely match the internal representations of the brain
683 at all time points.

684

685 **Conclusion**

686 We view our EEG dataset as a valuable tool for computational neuroscientists and
687 computer scientists. We believe that its largeness, richness and quality will facilitate
688 steps towards a deeper understanding of the neural mechanisms underlying visual
689 processing and towards more human-like artificial intelligence models.

690 **Materials and methods**

691 **Participants**

692 Ten healthy adults (mean age 28.5 years, SD=4; 8 female, 2 male) participated, all
693 having normal or corrected-to-normal vision. They all provided informed written
694 consent and received monetary reimbursement. Procedures were approved by the
695 ethical committee of the Department of Education and Psychology at Freie
696 Universität Berlin and were in accordance with the Declaration of Helsinki.

697

698 **Stimuli**

699 All images came from THINGS (Hebart et al., 2019), a database of 12 or more
700 images of objects on a natural background for each of 1,854 object concepts, where
701 each concept (e.g., antelope, strawberry, t-shirt) belongs to one of 27 higher-level
702 categories (e.g., animal, food, clothing). The building of encoding models involves
703 two stages: model training and model evaluation. Since each of these stages
704 requires an independent data partition, we pseudo-randomly divided the 1,854 object
705 concepts into non-overlapping 1,654 training and 200 test concepts under the
706 constraint that the same proportion of the 27 higher-level categories had to be kept
707 in both partitions. We then selected ten images for each training partition concept
708 and one image for each test partition concept, resulting in a training image partition
709 of 16,540 image conditions (1,654 training object concepts \times 10 images per concept
710 = 16,540 training image conditions) and a test image partition of 200 image
711 conditions (200 test object concepts \times 1 image per concept = 200 test image
712 conditions). We used the training and test data partitions for the encoding model
713 training and testing, respectively. The experiment had an orthogonal target detection
714 task (see “experimental paradigm” section below), and as task-relevant target stimuli
715 we used 10 different images of the “Toy Story” character Buzz Lightyear. All images
716 were of square size. We reshaped them to 500 \times 500 pixels for the EEG data
717 collection paradigm. For the modeling with DNNs we reshaped the images to 224 \times
718 224 pixels, and normalized them.

719

720 **Experimental Paradigm**

721 The experiment consisted in a RSVP paradigm (Intraub, 1981; Keyser et al., 2001;
722 Grootswagers et al., 2019) with an orthogonal target detection task to ensure
723 participants paid attention to the visual stimuli. All 10 participants completed four
724 equivalent experimental sessions, resulting in 10 datasets of 16,540 training images
725 conditions repeated 4 times and 200 test image conditions repeated 80 times, for a
726 total of (16,540 training image conditions \times 4 training image repetitions) + (200 test
727 image conditions \times 80 test image repetitions) = 82,160 image trials per dataset.

728 One session comprised 19 runs, all lasting around 5m. In each of the first 4
729 runs we showed participants the 200 test image conditions through 51 rapid serial
730 sequences of 20 images, for a total of 4 test runs \times 51 sequences per run \times 20
731 images per sequence = 4,080 image trials. In each of the following 15 runs we
732 showed 8,270 training image conditions (half of all the training image conditions, as
733 different halves were shown on different sessions) through 56 rapid serial sequences

734 of 20 images, for a total of 15 training runs \times 56 sequences per run \times 20 images per
735 sequence = 16,800 image trials.

736 Every rapid serial sequence started with 750ms of blank screen, then each of
737 the 20 images was presented centrally with a visual angle of 7 degrees for 100ms
738 and a SOA of 200ms, and it ended with another 750ms of blank screen. After every
739 rapid sequence there were up to 2s during which we instructed participants to first
740 blink (or make any other movement) and then report, with a keypress, whether the
741 target image of Buzz Lightyear appeared in the sequence. The images were
742 presented in a pseudo-randomized order, and a target image appeared in 6
743 sequences per run. A central bull's eye fixation target (Thaler et al., 2013) was
744 present on the screen throughout the entire experiment, and we asked participants
745 to constantly gaze at it. We controlled stimulus presentation using the Psychtoolbox
746 (Brainard, 1997), and recorded EEG data during the experimental sessions.

747

748 **EEG recording and preprocessing**

749 We recorded the EEG data using a 64-channel EASYCAP with electrodes arranged
750 in accordance with the standard 10-10 system (Nuwer et al., 2998), and a
751 Brainvision actiCHamp amplifier. We recorded the data at a sampling rate of
752 1000Hz, while performing online filtering (between 0.1Hz and 100Hz) and
753 referencing (to the Fz electrode). We performed offline preprocessing in Python,
754 using the MNE package (Gramfort et al., 2013). We epoched the continuous EEG
755 data into trials ranging from 200ms before stimulus onset to 800ms after stimulus
756 onset, and applied baseline correction by subtracting the mean of the pre-stimulus
757 interval for each trial and channel separately. We then down-sampled the epoched
758 data to 100Hz, and we selected 17 channels overlying occipital and parietal cortex
759 for further analysis (O1, Oz, O2, PO7, PO3, POz, PO4, PO8, P7, P5, P3, P1, Pz, P2,
760 P4, P6, P8). All trials containing target stimuli were not analyzed further, and we
761 randomly selected and retained 4 measurement repetitions for each training image
762 condition and 80 measurement repetitions for each test image condition. Next, we
763 applied multivariate noise normalization (Guggenmos et al., 2018) independently to
764 the data of each recording session. For each participant, the preprocessing resulted
765 in the EEG *biological training* (BioTrain) data matrix of shape (16,540 training image
766 conditions \times 4 condition repetitions \times 17 EEG channels \times 100 EEG time points) and
767 *biological test* (BioTest) data matrix of shape (200 test image conditions \times 80
768 condition repetitions \times 17 EEG channels \times 100 EEG time points). We used the
769 BioTrain and BioTest EEG data for the encoding models training and testing,
770 respectively.

771

772 **DNN models used**

773 We built linearizing encoding models (Wu et al., 2006; Kay et al., 2008; Naselaris et
774 al., 2011; van Gerven, 2017; Kriegeskorte & Douglas, 2019) of EEG visual
775 responses using four different DNNs: AlexNet (Krizhevsky, 2014), a supervised
776 feedforward neural network of 5 convolutional layers followed by 3 fully-connected
777 layers that won the Imagenet large-scale visual recognition challenge in 2012;

778 ResNet-50 (He et al, 2016), a supervised feedforward 50 layer neural network with
779 shortcut connections between layers at different depths; CORnet-S (Kubilius et al.,
780 2019), a supervised deep recurrent neural network of four convolutional layers and
781 one fully-connected layer; MoCo (Chen et al., 2020), a feedforward ResNet-50
782 architecture trained in a self-supervised fashion. All of them had been pre-trained on
783 object categorization on the ILSVRC-2012 training image partition (Russakovsky et
784 al., 2015).

785

786 **Linearizing encoding models of EEG visual responses**

787 The first step in building linearizing encoding models is to use DNNs to non-linearly
788 transform the image input space onto a feature space. A DNNs feature space is
789 given by its feature maps, layerwise representations (non-linear transformations) of
790 the input images. To get the training and test feature maps we fed the training and
791 test images separately to each DNN and appended the vectorized image
792 representations of its layers onto each other. We extracted AlexNet's feature maps
793 from layers maxpool1, maxpool2, ReLU3, ReLU4, maxpool5, ReLU6, ReLU7, and
794 fc8; ResNet-50's and MoCo's feature maps from the last layer of each of their four
795 blocks, and from the decoder layer; CORnet-S' feature maps from the last layers of
796 areas V1, V2 (at both time points), V4 (at all four time points), IT (at both time
797 points), and from the decoder layer. We then standardized the appended feature
798 maps of the training and test data to zero mean and unit variance for each feature
799 across the sample (images) dimension, using the mean and standard deviation of
800 the training feature maps. Finally, we used the Scikit-learn (Pedregosa et al., 2011)
801 implementation of non-linear principal component analysis (computed on the training
802 feature maps using a polynomial kernel of degree 4) to reduce the feature maps of
803 both the training and test images to 1,000 components. For each DNN model, this
804 resulted in the training feature maps matrix of shape (16,540 training image
805 conditions \times 1,000 features) and test feature maps matrix of shape (200 test image
806 conditions \times 1,000 features).

807 The second step in building linearizing encoding models is to linearly map the
808 DNNs' feature space onto the EEG neural space, effectively predicting the EEG
809 responses to images. We performed this linear mapping independently for each
810 participant, DNN model and EEG feature (i.e., for each of the 17 EEG channels (\mathbf{c}) \times
811 100 EEG time points (\mathbf{t}) = 1,700 EEG features). We fitted the weights $\mathbf{W}_{t,c}$ of a linear
812 regression using the DNNs' training feature maps as the predictors and the
813 corresponding BioTrain data (averaged across the image conditions repetitions) as
814 the criterion: during training the regression weights learned the existing linear
815 relationship between the DNN feature maps of a given image and the EEG
816 responses of that same image. No regularization techniques were used. We
817 multiplied $\mathbf{W}_{t,c}$ with the DNNs' test feature maps. For each participant and DNN, this
818 resulted in the *synthetic test* (SynTest) EEG data matrix of shape (200 test image
819 conditions \times 17 EEG channels \times 100 EEG time points).

820

821

822 **Correlation**

823 We used a Pearson correlation to assess how similar the SynTest EEG data of each
824 participant and DNN is to the corresponding BioTest data, thus quantifying the
825 encoding models' predicted power. We started the analysis by averaging the BioTest
826 data across 40 image conditions repetitions (see "noise ceiling" section below),
827 resulting in a BioTest data matrix equivalent in shape to the SynTest data matrix
828 (200 test image conditions \times 17 EEG channels \times 100 EEG time points). Next, we
829 implemented a nested loop over the EEG channels and time points. At each loop
830 iteration we indexed the 200-dimensional BioTest data vector containing the 200 test
831 image conditions of the EEG channel (**c**) and time point (**t**) in question, and
832 correlated it with the corresponding 200-dimensional SynTest data vector. This
833 procedure yielded a Pearson correlation coefficient matrix of shape (17 EEG
834 channels \times 100 EEG time points). Finally, we averaged the Pearson correlation
835 coefficient matrix over the EEG channels, obtaining a correlation results vector of
836 length (100 EEG time points) for each participant and DNN.

837

838 **Pairwise decoding**

839 The rationale of this analysis was to see if a classifier trained on the BioTest data is
840 capable of generalizing its performance to the SynTest data. This is a
841 complementary way (to the correlation analysis) to assess the similarity between the
842 SynTest data and the BioTest data, hence the encoding models' predictive power.
843 We started the analysis by averaging 40 BioTest data image conditions repetitions
844 (see "noise ceiling" section below) into 10 pseudo-trials of 4 repeats each, yielding a
845 matrix of shape (200 test image conditions \times 10 image condition pseudo-trials \times 17
846 EEG channels \times 100 EEG time points). Next, we used the pseudo-trials for training
847 linear SVMs to perform binary classification between each pair of the 200 BioTest
848 data image conditions (for a total of 19,900 image condition pairs) using their EEG
849 channels vectors (of 17 components). We then tested the trained classifiers on the
850 corresponding pairs of SynTest data image conditions. We performed the pairwise
851 decoding analysis independently for each EEG time point (**t**), which resulted in a
852 matrix of decoding accuracy scores of shape (19,900 image condition pairs \times 100
853 EEG time points). We then averaged the decoding accuracy scores matrix across
854 the image condition pairs, obtaining a pairwise decoding results vector of length (100
855 EEG time points) for each participant and DNN.

856

857 **Zero-shot identification**

858 In this analysis we exploited the linearizing encoding models' predictive power to
859 identify the BioTest data image conditions in a zero-shot fashion, that is, to identify
860 arbitrary image conditions without prior training (Kay et al., 2008; Seeliger et al.,
861 2017; Horikawa & Kamitani, 2017). We identified each BioTest data image condition
862 using the SynTest data and an additional synthesized EEG dataset of up to 150,000
863 candidate image conditions. These 150,000 image conditions came from the
864 ILSVRC-2012 (Russakovsky et al., 2015) validation (50,000) plus test (100,000)
865 sets. We synthesized them into their corresponding EEG responses following the

866 same procedure described above, resulting in the *synthetic Imagenet* (SynImagenet)
867 data matrix of shape (150,000 image conditions \times 17 EEG channels \times 100 EEG time
868 points). The zero-shot identification analysis involved two steps: feature selection
869 and identification.

870 In the feature selection step we used the training data to pick only the most
871 relevant EEG features (out of all 17 EEG channels \times 100 EEG time points = 1,700
872 EEG features). We synthesized the EEG responses to the 16,540 training images,
873 obtaining the *synthetic train* (SynTrain) data matrix of shape (16,540 training image
874 conditions \times 17 EEG channels \times 100 EEG time points). Next, we correlated each
875 SynTrain data feature (across the 16,540 training image conditions, with a Pearson
876 correlation), with the corresponding BioTrain data feature (averaged across the
877 image conditions repetitions). We then selected only the 300 BioTest data, SynTest
878 data and SynImagenet data EEG features corresponding to the 300 highest
879 correlation scores, thus obtaining a BioTest data matrix of shape (200 test image
880 conditions \times 80 condition repetitions \times 300 EEG features), a SynTest data matrix of
881 shape (200 test image conditions \times 300 EEG features), and a SynImagenet data
882 matrix of shape (150,000 image conditions \times 300 EEG features).

883 In the identification step we started by averaging the BioTest data across all
884 the 80 image conditions repetitions: this yielded feature vectors of 300 components
885 for each of the 200 image conditions. Next, we correlated (through a Pearson
886 correlation) the feature vectors of each BioTest data image condition with the feature
887 vectors of all the candidate image conditions: the SynTest data image conditions
888 plus a varying amount of SynImagenet data image conditions. We increased the set
889 sizes of the SynImagenet candidate image conditions from 0 to 150,000 with steps of
890 1,000 images (for a total of 151 set sizes), where 0 corresponded to using only the
891 SynTest data candidate image conditions, and performed the zero-shot identification
892 at every set size. At each SynImagenet set size a BioTest data image condition is
893 considered correctly identified if the correlation coefficient between its channel vector
894 and the channel vector of the corresponding SynTest data image condition is higher
895 than the correlation coefficients between its channel vector and the channel vectors
896 of all other candidate SynTest data and SynImagenet data image conditions. Thus,
897 we calculated the zero-shot identification accuracies through the ratio of correctly
898 classified images over all 200 BioTest images, obtaining a zero-shot identification
899 results vector of length (151 candidate image set sizes). We iterated the
900 identification step 100 times, while always randomly selecting different SynImagenet
901 data image conditions at each set size, and then averaged the results across the 100
902 iterations.

903 To extrapolate the drop in identification accuracy with larger candidate image
904 set sizes we fit the power-law function to the results of each participant. The power
905 law function is defined as:

$$906 \quad f(x) = ax^b$$

907 where x is the image set size, a and b are constants learned during function fitting,
908 and $f(x)$ is the predicted zero-shot identification accuracy. We fit the function using
909 the 100 SynImagenet set sizes ranging from 50,200 to 150,200 images (along with

910 their corresponding identification accuracies), and then used it to extrapolate the
911 image set size required for the identification accuracy to drop below 10% and 0.5%.

912

913 **End-to-end encoding models of EEG visual responses**

914 We based our end-to-end encoding models (Seeliger et al., 2021; Khosla et al.,
915 2021; Allen et al., 2021) on the AlexNet architecture which, once trained, predicted
916 the EEG responses to the test images. To match AlexNet's output with the channel
917 responses of our EEG data, we replaced AlexNet's 1000-neurons output layer with a
918 17-neurons layer, where each neuron represented one of the 17 EEG channels.
919 Next, we randomly initialized independent AlexNet instances for each participant and
920 EEG time point (t). We used Pytorch (Paszke et al., 2019) to train the AlexNets on a
921 regression task: given the input training images and the corresponding target
922 BioTrain EEG data channel activity (averaged across the image condition
923 repetitions), the models had to optimize their weights so to minimize the summed
924 squared error between their predictions and the BioTrain data. We trained the
925 models using batch sizes of 256 images and an Adam optimizer with a learning rate
926 of 0.0001, a weight decay term of 0.001, and the default value for the remaining
927 parameters. We implemented a cross-validation loop over the 200 test image
928 conditions to identify the optimal amount of training epochs for the synthesis of each
929 image's EEG responses. At every loop iteration we selected the image condition of
930 interest, synthesized the EEG responses to the remaining 199 test images for each
931 of 30 training epochs, and correlated the synthetic data with the corresponding 199
932 biological test EEG data conditions, resulting in one correlation score per epoch. We
933 then synthesized the EEG responses to the image condition of interest using the
934 model weights of the epoch leading to the highest correlation score. For each
935 participant, this resulted in the SynTest data matrix of shape (200 test image
936 conditions \times 17 EEG channels \times 100 EEG time points).

937

938 **Noise ceiling calculation**

939 We calculated the noise ceilings of the correlation and pairwise decoding analyses to
940 estimate the theoretical maximum results given the level of noise in the BioTest data.
941 If the results of the SynTest data reach this theoretical maximum the encoding
942 models are successful in explaining all the BioTest data variance which can be
943 explained. If not, further model improvements could lead to more accurate
944 predictions of neural data.

945 For the noise ceiling estimation we randomly divided the BioTest data into two
946 non-overlapping partitions of 40 image condition repetitions each, where the first
947 partition corresponded to the 40 repeats of BioTest data image conditions used in
948 the correlation and pairwise decoding analyses described above. We then performed
949 the two analyses while substituting the SynTest data with the second BioTest data
950 partition (averaged across image condition repetitions). This resulted in the noise
951 ceiling lower bound estimates. To calculate the upper bound estimates we
952 substituted the SynTest data with the average of the BioTest data over all 80 image
953 condition repetitions and reiterated the two analyses. We assume that the true noise

954 ceiling is somewhere in between the lower and the upper bound estimates. To avoid
955 the results being biased by one specific configuration of the BioTest data repeats we
956 iterated the correlation and pairwise decoding analyses 100 times, while always
957 selecting different repeats for the two BioTest data partitions, and then averaged the
958 results across the 100 iterations.

959

960 **Statistical testing**

961 To assess the statistical significance of the correlation, pairwise decoding and zero-
962 shot identification analyses we tested all results against chance using one-sample
963 one-sided t-tests. Here, the rationale was to reject the null hypothesis H_0 of the
964 analyses results being at chance level with a confidence of 95% or higher (i.e., with a
965 P -value of $P < 0.05$), thus supporting the experimental hypothesis H_1 of the results
966 being significantly higher than chance. The chance level differed across analyses: 0
967 in the correlation; 50% in the pairwise decoding; $(1 / (200 \text{ test image conditions} + N$
968 $\text{ILSVRC-2012 image conditions}))$ in the zero-shot identification (where N varied from
969 0 to 150,000). When analyzing the linearizing encoding models' prediction accuracy
970 using different amounts of training data we used a two-way repeated measures
971 ANOVA to reject the null hypothesis H_0 of no significant effects of number of image
972 conditions and/or condition repetitions on the prediction accuracy, and a repeated
973 measures two-sided t-test to reject the null hypothesis H_0 of no significant
974 differences between the effects of training image conditions and condition
975 repetitions.

976 We controlled familywise error rate by applying a conservative Bonferroni-
977 correction to the resulting P -values to correct for the number of EEG time points ($N =$
978 100) in the correlation and pairwise decoding analyses, for the amount of training
979 data quartiles ($N = 4$) in the analysis of the linearizing encoding models' prediction
980 accuracy as a function of training image conditions and condition repetitions, and for
981 the number of candidate images set sizes ($N = 151$) in the zero-shot identification
982 analysis.

983 To calculate the confidence intervals of each statistic, we created 10,000
984 bootstrapped samples by sampling the participant-specific results with replacement.
985 This yielded empirical distributions of the results, from which we took the 95%
986 confidence intervals.

987 **Acknowledgments**

988 A.T.G. is supported by a PhD fellowship of the Einstein Center for Neurosciences.
989 G.R. is supported by the Alfons and Gertrud Kassel Foundation. R.M.C. is supported
990 by German Research Council (DFG) grants (CI 241/1-1, CI 241/3-1, CI 241/1-7) and
991 the European Research Council (ERC) starting grant (ERC-StG-2018-803370). We
992 thank Martin Hebart for support with the THINGS database. We thank Daniel Kaiser
993 and Kendrick Kay for helpful comments on the manuscript. We thank the HPC
994 Service of ZEDAT, Freie Universität Berlin, for computing time.

995

996 **Competing interests**

997 The authors declare no competing interests.

998

999 **Author contributions**

1000 A.T.G., K.D. and R.M.C. designed research, A.T.G. acquired data, A.T.G. analyzed
1001 data, A.T.G., K.D., G.R. and R.M.C. interpreted results, A.T.G. prepared figures,
1002 A.T.G. drafted manuscript, A.T.G., K.D., G.R. and R.M.C. edited and revised
1003 manuscript. All authors approved the final version of the manuscript.

1004

1005 **Data availability**

1006 The raw and preprocessed EEG dataset, the stimuli image set and the extracted
1007 DNN feature maps are available on [OSF](#).

1008

1009 **Code availability**

1010 The code to reproduce all the results is available on [GitHub](#).

1011 **References**

- 1012 Abiri R, Borhani S, Sellers EW, Jiang Y, Zhao X. 2019. A comprehensive review of EEG-based brain–
1013 computer interface paradigms. *Journal of Neural Engineering*, 16(1):011001. DOI:
1014 <https://doi.org/10.1088/1741-2552/aaf12e>
- 1015 Allen EJ, St-Yves G, Wu Y, Breedlove JL, Prince JS, Dowdle LT, Nau M, Caron B, Pestilli F, Charest
1016 I, Hutchinson JB, Naselaris T, Kay K. 2022. A massive 7T fMRI dataset to bridge cognitive
1017 neuroscience and computational intelligence. *Nature Neuroscience*, 25(1):116–126. DOI:
1018 <https://doi.org/10.1038/s41593-021-00962-x>
- 1019 Bankson BB, Hebart MN, Groen II, Baker CI. 2018. The temporal evolution of conceptual object
1020 representations revealed through models of behavior, semantics and deep neural networks.
1021 *NeuroImage*, 178:172–182. DOI: <https://doi.org/10.1016/j.neuroimage.2018.05.037>
- 1022 Brainard DH. 1997. The psychophysics toolbox. *Spatial Vision*, 10:433–436. DOI:
1023 <https://doi.org/10.1163/156856897X00357>
- 1024 Cadieu CF, Hong H, Yamins DLK, Pinto N, Ardila D, Solomon EA, Majaj NJ, DiCarlo JJ. 2014. Deep
1025 neural networks rival the representation of primate IT cortex for core visual object recognition.
1026 *PLoS Computational Biology*, 10(12):e1003963. DOI:
1027 <https://doi.org/10.1371/journal.pcbi.1003963>
- 1028 Carandini M, Demb JB, Mante V, Tolhurst DJ, Dan Y, Olshausen BA, Gallant JL, Rust NC. 2005. Do
1029 we know what the early visual system does?. *Journal of Neuroscience*, 25(46):10577–10597.
1030 DOI: <https://doi.org/10.1523/JNEUROSCI.3726-05.2005>
- 1031 Chang N, Pyles JA, Marcus A, Gupta A, Tarr M, Aminoff EM. 2019. BOLD5000, a public fMRI dataset
1032 while viewing 5000 visual images. *Scientific Data*, 6(1):1–18. DOI:
1033 <https://doi.org/10.1038/s41597-019-0052-3>
- 1034 Charest I, Kievit RA, Schmitz TW, Deca D, Kriegeskorte N. 2014. Unique semantic space in the brain
1035 of each beholder predicts perceived similarity. *Proceedings of the National Academy of
1036 Sciences*, 111(40): 14565–14570. DOI: <https://doi.org/10.1073/pnas.1402594111>
- 1037 Chen X, Fan H, Girshick R, He K. 2020. Improved baselines with momentum contrastive learning.
1038 *arXiv preprint*, arXiv:2003.04297. DOI: <https://doi.org/10.48550/arXiv.2003.04297>
- 1039 Cichy RM, Kaiser D. 2019. Deep neural networks as scientific models. *Trends in Cognitive Sciences*,
1040 23(4):305–317. DOI: <https://doi.org/10.1016/j.tics.2019.01.009>
- 1041 Cichy RM, Khosla A, Pantazis D, Torralba A, Oliva A. 2016. Comparison of deep neural networks to
1042 spatio-temporal cortical dynamics of human visual object recognition reveals hierarchical
1043 correspondence. *Scientific Reports*, 6(1):1-13. DOI: <https://doi.org/10.1038/srep27755>
- 1044 Cichy RM, Pantazis D, Oliva A. 2014. Resolving human object recognition in space and time. *Nature
1045 Neuroscience*, 17(3):455–462. DOI: <https://doi.org/10.1038/nn.3635>
- 1046 Cichy RM, Roig G, Oliva A. 2019. The Algonauts Project. *Nature Machine Intelligence*, 1(12):613–
1047 613. DOI: <https://doi.org/10.1038/s42256-019-0127-z>
- 1048 Dijkstra N, Mostert P, de Lange FP, Bosch S, Gerven MA. 2018. Differential temporal dynamics
1049 during visual imagery and perception. *Elife*, 7:e33904. DOI: <https://doi.org/10.7554/eLife.33904>
- 1050 Fukushima K, Miyake S. 1982. Neocognitron: A self-organizing neural network model for a
1051 mechanism of visual pattern recognition. In *Competition and Cooperation in Neural Nets*: 267–
1052 285. Springer, Berlin, Heidelberg. DOI: https://doi.org/10.1007/978-3-642-46466-9_18
- 1053 Goodale MA, Milner AD. 1992. Separate visual pathways for perception and action. *Trends in
1054 Neurosciences*, 15(1): 20–25. DOI: [https://doi.org/10.1016/0166-2236\(92\)90344-8](https://doi.org/10.1016/0166-2236(92)90344-8)
- 1055 Gramfort A, Luessi M, Larson E, Engemann DA, Strohmeier D, Brodbeck C, Goj R, Jas M, Brooks T,
1056 Parkkonen L, Hämäläinen MS. 2013. MEG and EEG data analysis with MNE-Python. *Frontiers
1057 in Neuroscience*, 7(267):1–13. DOI: <https://doi.org/10.3389/fnins.2013.00267>
- 1058 Grill-Spector K, Kourtzi Z, Kanwisher N. 2001. The lateral occipital complex and its role in object
1059 recognition. *Vision Research*, 41(10-11):1409–1422. DOI: [https://doi.org/10.1016/S0042-
1060 6989\(01\)00073-6](https://doi.org/10.1016/S0042-6989(01)00073-6)
- 1061 Grootswagers T, Robinson AK, Carlson TA. 2019. The representational dynamics of visual objects in
1062 rapid serial visual processing streams. *NeuroImage*, 188:668-679. DOI:

- 1063 <https://doi.org/10.1016/j.neuroimage.2018.12.046>
- 1064 Grootswagers T, Zhou I, Robinson AK, Hebart MN, Carlson TA. 2022. Human EEG recordings for
1065 1,854 concepts presented in rapid serial visual presentation streams. *Scientific Data*, 9(1):1–7.
1066 DOI: <https://doi.org/10.1038/s41597-021-01102-7>
- 1067 Güçlü U, van Gerven MAJ. 2015. Deep neural networks reveal a gradient in the complexity of neural
1068 representations across the ventral stream. *Journal of Neuroscience*, 35(27):10005–10014. DOI:
1069 <https://doi.org/10.1523/JNEUROSCI.5023-14.2015>
- 1070 Guest O, Martin AE. 2021. On logical inference over brains, behaviour, and artificial neural networks.
1071 Guggenmos M, Sterzer P, Cichy RM. 2018. Multivariate pattern analysis for MEG: A comparison of
1072 dissimilarity measures. *NeuroImage*, 173:434–447. DOI:
1073 <https://doi.org/10.1016/j.neuroimage.2018.02.044>
- 1074 Harel A, Groen II, Kravitz DJ, Deouell LY, Baker CI. 2016. The temporal dynamics of scene
1075 processing: A multifaceted EEG investigation. *Eneuro*, 3(5). DOI:
1076 <https://doi.org/10.1523/ENEURO.0139-16.2016>
- 1077 Hassabis D, Kumaran D, Summerfield C, Botvinick M. 2017. Neuroscience-inspired artificial
1078 intelligence. *Neuron*, 95(2):245–258. DOI: <https://doi.org/10.1016/j.neuron.2017.06.011>
- 1079 Haxby JV, Guntupalli JS, Nastase SA, Feilong M. 2020. Hyperalignment: Modeling shared information
1080 encoded in idiosyncratic cortical topographies. *Elife*, 9:e56601. DOI:
1081 <https://doi.org/10.7554/eLife.56601>
- 1082 Haynes JD, Rees G. 2006. Decoding mental states from brain activity in humans. *Nature Reviews*
1083 *Neuroscience*, 7(7):523–534. DOI: <https://doi.org/10.1038/nrn1931>
- 1084 He K, Zhang X, Ren S, Sun J. 2016. Deep residual learning for image recognition. *Proceedings of the*
1085 *IEEE Conference on Computer Vision and Pattern Recognition*, 770–778. DOI:
1086 <https://doi.org/10.1109/CVPR.2016.90>
- 1087 Hebart MN, Dickter AH, Kidder A, Kwok WY, Corriveau A, Van Wicklin C, Baker CI. 2019. THINGS: A
1088 database of 1,854 object concepts and more than 26,000 naturalistic object images. *PLoS*
1089 *ONE*, 14(10): e0223792. DOI: <https://doi.org/10.1371/journal.pone.0223792>
- 1090 Horikawa T, Kamitani Y. 2017. Generic decoding of seen and imagined objects using hierarchical
1091 visual features. *Nature Communications*, 8(1):1–15. DOI: <https://doi.org/10.1038/ncomms15037>
- 1092 Intraub H. 1981. Rapid conceptual identification of sequentially presented pictures. *Journal of*
1093 *Experimental Psychology: Human Perception and Performance*, 7(3):604. DOI:
1094 <https://doi.org/10.1037/0096-1523.7.3.604>
- 1095 Kar K, Kubilius J, Schmidt K, Issa EB, DiCarlo JJ. 2019. Evidence that recurrent circuits are critical to
1096 the ventral stream’s execution of core object recognition behavior. *Nature Neuroscience*,
1097 22(6):974–983. DOI: <https://doi.org/10.1038/s41593-019-0392-5>
- 1098 Kay KN, Naselaris T, Prenger RJ, Gallant JL. 2008. Identifying natural images from human brain
1099 activity. *Nature*, 452(7185):352–355. DOI: <https://doi.org/10.1038/nature06713>
- 1100 Keyser C, Xiao DK, Földiák P, Perrett DI. 2001. The speed of sight. *Journal of cognitive*
1101 *neuroscience*, 13(1):90–101. DOI: <https://doi.org/10.1162/089892901564199>
- 1102 Khaligh-Razavi SM, Henriksson L, Kay K, Kriegeskorte N. 2017. Fixed versus mixed RSA: Explaining
1103 visual representations by fixed and mixed feature sets from shallow and deep computational
1104 models. *Journal of Mathematical Psychology*, 76:184–197. DOI:
1105 <https://doi.org/10.1016/j.jmp.2016.10.007>
- 1106 Khosla M, Ngo GH, Jamison K, Kuceyeski A, Sabuncu MR. 2021. Cortical response to naturalistic
1107 stimuli is largely predictable with deep neural networks. *Science Advances*, 7(22):eabe7547.
1108 DOI: <https://doi.org/10.1126/sciadv.abe7547>
- 1109 Kietzmann TC, McClure P, Kriegeskorte N. 2019a. Deep neural networks in computational
1110 neuroscience. In *Oxford Research Encyclopedia of Neuroscience*. DOI:
1111 <https://doi.org/10.1093/acrefore/9780190264086.013.46>
- 1112 Kietzmann TC, Spoerer CJ, Sörensen LK, Cichy RM, Hauk O, Kriegeskorte N. 2019b. Recurrence is
1113 required to capture the representational dynamics of the human visual system. *Proceedings of*
1114 *the National Academy of Sciences*, 116(43):21854–21863. DOI:
1115 <https://doi.org/10.1073/pnas.1905544116>

- 1116 Kriegeskorte N. 2015. Deep neural networks: a new framework for modeling biological vision and
1117 brain information processing. *Annual Review of Vision Science*, 1:417–446. DOI:
1118 <https://doi.org/10.1146/annurev-vision-082114-035447>
- 1119 Kriegeskorte N, Douglas PK. 2019. Interpreting encoding and decoding models. *Current opinion in*
1120 *neurobiology*, 55, 167–179. DOI: <https://doi.org/10.1016/j.conb.2019.04.002>
- 1121 Kriegeskorte N, Mur M, Bandettini PA. 2008. Representational similarity analysis-connecting the
1122 branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2:4.8 DOI:
1123 <https://doi.org/10.3389/neuro.06.004.2008>
- 1124 Krizhevsky A. 2014. One weird trick for parallelizing convolutional neural networks. *arXiv preprint*,
1125 arXiv:1404.5997
- 1126 Kubilius J, Schrimpf M, Kar K, Rajalingham R, Hong H, Majaj N, Issa E, Bashivan P, Prescott-Roy J,
1127 Schmidt K, Nayeibi A, Bear D, Yamins DL, DiCarlo JJ. 2019. Brain-like object recognition with
1128 high-performing shallow recurrent ANNs. *Advances in neural information processing systems*,
1129 32.
- 1130 Kwon OY, Lee MH, Guan C, Lee SW. 2019. Subject-independent brain–computer interfaces based on
1131 deep convolutional neural networks. *IEEE transactions on neural networks and learning*
1132 *systems*, 31(10):3839–3852. DOI: <https://doi.org/10.1109/TNNLS.2019.2946869>
- 1133 Logothetis NK, Sheinberg DL. 1996. Visual object recognition. *Annual review of neuroscience*,
1134 19(1):577–621. DOI: <https://doi.org/10.1146/annurev.ne.19.030196.003045>
- 1135 Malach R, Levy I, Hasson U. 2002. The topography of high-order human object areas. *Trends in*
1136 *Cognitive Sciences*, 6(4):176–184. DOI: [https://doi.org/10.1016/S1364-6613\(02\)01870-3](https://doi.org/10.1016/S1364-6613(02)01870-3)
- 1137 Marr D. 1980. Visual information processing: The structure and creation of visual representations.
1138 *Philosophical Transactions of the Royal Society of London. B, Biological Sciences*,
1139 290(1038):199–218. DOI: <https://doi.org/10.1098/rstb.1980.0091>
- 1140 Mohsenzadeh Y, Qin S, Cichy RM, Pantazis D. 2018. Ultra-Rapid serial visual presentation reveals
1141 dynamics of feedforward and feedback processes in the ventral visual pathway. *Elife*,
1142 7:e36329. DOI: <https://doi.org/10.7554/eLife.36329>
- 1143 Mur M, Bandettini PA, Kriegeskorte N. 2009. Revealing representational content with pattern-
1144 information fMRI—an introductory guide. *Social Cognitive and Affective Neuroscience*,
1145 4(1):101–109. DOI: <https://doi.org/10.1093/scan/nsn044>
- 1146 Naselaris T, Allen E, Kay K. 2021. Extensive sampling for complete models of individual brains.
1147 *Current Opinion in Behavioral Sciences*, 40:45–51. DOI:
1148 <https://doi.org/10.1016/j.cobeha.2020.12.008>
- 1149 Naselaris T, Kay KN, Nishimoto S, Gallant JL. 2011. Encoding and decoding in fMRI. *NeuroImage*,
1150 56(2):400–410. DOI: <https://doi.org/10.1016/j.neuroimage.2010.07.073>
- 1151 Naselaris T, Olman CA, Stansbury DE, Ugurbil K, Gallant JL. 2015. A voxel-wise encoding model for
1152 early visual areas decodes mental images of remembered scenes. *NeuroImage*, 105:215–228.
1153 DOI: <https://doi.org/10.1016/j.neuroimage.2014.10.018>
- 1154 Nuwer MR, Comi G, Emerson R, Fuglsang-Frederiksen A, Guérit JM, Hinrichs H, Ikeda A, Luccas
1155 FJC, Rappelsburger P. 1998. IFCN standards for digital recording of clinical EEG.
1156 *Electroencephalography and Clinical Neurophysiology*, 106(3):259–261. DOI:
1157 [https://doi.org/10.1016/s0013-4694\(97\)00106-5](https://doi.org/10.1016/s0013-4694(97)00106-5)
- 1158 Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga
1159 L, Desmaison A. 2019. Pytorch: An imperative style, high-performance deep learning library.
1160 *Advances in Neural Information Processing Systems*, 32:8026–8037.
- 1161 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P,
1162 Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M,
1163 Duchesnay É. 2011. Scikit-learn: Machine learning in Python. *the Journal of Machine Learning*
1164 *Research*, 12:2825–2830.
- 1165 Petit J, Rouillard J, Cabestaing F. 2021. EEG-based brain–computer interfaces exploiting steady-state
1166 somatosensory-evoked potentials: a literature review. *Journal of Neural Engineering*,
1167 18(5):051003. DOI: <https://doi.org/10.1088/1741-2552/ac2fc4>
- 1168 Rajaei K, Mohsenzadeh Y, Ebrahimpour R, Khaligh-Razavi SM. 2019. Beyond core object

- 1169 recognition: Recurrent processes account for object recognition under occlusion. *PLoS*
1170 *computational biology*, 15(5):e1007001. DOI: <https://doi.org/10.1371/journal.pcbi.1007001>
- 1171 Richard H, Gresle L, Hyvarinen A, Thirion B, Gramfort A, Ablin P. 2020. Modeling shared responses
1172 in neuroimaging studies through multiview ica. *Advances in Neural Information Processing*
1173 *Systems*, 33:19149–19162.
- 1174 Richards BA, Lillicrap TP, Beaudoin P, Bengio Y, Bogacz R, Christensen A, Clopath C, Costa RP, de
1175 Berker A, Ganguli S, Gillon CJ. 2019. A deep learning framework for neuroscience. *Nature*
1176 *Neuroscience*, 22(11):1761–1770. DOI: <https://doi.org/10.1038/s41593-019-0520-2>
- 1177 Riesenhuber M, Poggio T. 1999. Hierarchical models of object recognition in cortex. *Nature*
1178 *neuroscience*, 2(11):1019–1025. DOI: <https://doi.org/10.1038/14819>
- 1179 Rousselet GA, Fabre-Thorpe M, Thorpe SJ. 2002. Parallel processing in high-level categorization of
1180 natural images. *Nature Neuroscience*, 5(7):629–630. DOI: <https://doi.org/10.1038/nn866>
- 1181 Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A,
1182 Bernstein M, Berg AC, Fei-Fei L. 2015. ImageNet Large Scale Visual Recognition Challenge.
1183 *International Journal of Computer Vision*, 115(3):211–252. DOI: [https://doi.org/10.1007/s11263-](https://doi.org/10.1007/s11263-015-0816-y)
1184 [015-0816-y](https://doi.org/10.1007/s11263-015-0816-y)
- 1185 Saxe A, Nelli S, Summerfield C. 2021. If deep learning is the answer, what is the question?. *Nature*
1186 *Reviews Neuroscience*, 22(1):55–67. DOI: <https://doi.org/10.1038/s41583-020-00395-8>
- 1187 Schrimpf M, Kubilius J, Lee MJ, Murty NAR, Ajemian R, DiCarlo JJ. 2020. Integrative benchmarking
1188 to advance neurally mechanistic models of human intelligence. *Neuron*, 108(3):413–423. DOI:
1189 <https://doi.org/10.1016/j.neuron.2020.07.040>
- 1190 Seeliger K, Ambrogioni L, Güçlütürk Y, van den Bulk LM, Güçlü U, van Gerven MAJ. 2021. End-to-
1191 end neural system identification with neural information flow. *PLOS Computational Biology*,
1192 17(2):e1008558. DOI: <https://doi.org/10.1371/journal.pcbi.1008558>
- 1193 Seeliger K, Fritsche M, Güçlü U, Schoenmakers S, Schoffelen J-M, Bosch S, van Gerven, MAJ. 2017.
1194 Convolutional neural network-based encoding and decoding of visual object recognition in
1195 space and time. *NeuroImage*, 180:253–266. DOI:
1196 <https://doi.org/10.1016/j.neuroimage.2017.07.018>
- 1197 Sinz FH, Pitkow X, Reimer J, Bethge M, Tolias AS. 2019. Engineering a less artificial intelligence.
1198 *Neuron*, 103(6):967–979. DOI: <https://doi.org/10.1016/j.neuron.2019.08.034>
- 1199 Spoerer CJ, McClure P, Kriegeskorte N. 2017. Recurrent convolutional neural networks: a better
1200 model of biological object recognition. *Frontiers in psychology*, 8:1551. DOI:
1201 <https://doi.org/10.3389/fpsyg.2017.01551>
- 1202 Storrs KR, Kietzmann TC, Walther A, Mehrer J, Kriegeskorte N. 2021. Diverse Deep Neural Networks
1203 All Predict Human Inferior Temporal Cortex Well, After Training and Fitting. *Journal of Cognitive*
1204 *Neuroscience*, 33(10):2044–2064. DOI: https://doi.org/10.1162/jocn_a_01755
- 1205 Tanaka K. 1996. Inferotemporal cortex and object vision. *Annual review of neuroscience*, 19:109–139.
1206 DOI: <https://doi.org/10.1146/annurev.ne.19.030196.000545>
- 1207 Thaler L, Schütz AC, Goodale MA, Gegenfurtner KR. 2013. What is the best fixation target? The
1208 effect of target shape on stability of fixational eye movements. *Vision Research*, 76:31–42. DOI:
1209 <https://doi.org/10.1016/j.visres.2012.10.012>
- 1210 Thorpe S, Fize D, Marlot C. 1996. Speed of processing in the human visual system. *Nature*,
1211 381(6582):520–522. DOI: <https://doi.org/10.1038/381520a0>
- 1212 Toneva M, Wehbe L. 2019. Interpreting and improving natural-language processing (in machines)
1213 with natural language-processing (in the brain). *Advances in Neural Information Processing*
1214 *Systems*, 32.
- 1215 Ullman S. 2000. High-level vision: Object recognition and visual cognition. *MIT press*.
- 1216 Ullman S. 2019. Using neuroscience to develop artificial intelligence. *Science*, 363(6428):692–693.
1217 DOI: <https://doi.org/10.1126/science.aau6595>
- 1218 Van Essen, D.C., Anderson, C.H. and Felleman, D.J., 1992. Information processing in the primate
1219 visual system: an integrated systems perspective. *Science*, 255(5043):419–423. DOI:
1220 <https://doi.org/10.1126/science.1734518>
- 1221 van Bergen RS, Kriegeskorte N. 2020. Going in circles is the way forward: the role of recurrence in

- 1222 visual inference. *Current Opinion in Neurobiology*, 65:176–193. DOI:
1223 <https://doi.org/10.1016/j.conb.2020.11.009>
- 1224 van de Nieuwenhuijzen ME, Backus AR, Bahramisharif A, Doeller CF, Jensen O, van Gerven MA.
1225 2013. MEG-based decoding of the spatiotemporal dynamics of visual category perception.
1226 *Neuroimage*, 83:1063–1073. DOI: <https://doi.org/10.1016/j.neuroimage.2013.07.075>
- 1227 van Gerven MA. 2017. A primer on encoding models in sensory neuroscience. *Journal of*
1228 *Mathematical Psychology*, 76:172–183. DOI: <https://doi.org/10.1016/j.jmp.2016.06.009>
- 1229 Wu MC-K, David SV, Gallant JL. 2006. Complete functional characterization of sensory neurons by
1230 system identification. *Annual Review of Neuroscience*, 29(1):477–505. DOI:
1231 <https://doi.org/10.1146/annurev.neuro.29.051605.113024>
- 1232 Yamins DLK, DiCarlo JJ. 2016. Using goal-driven deep learning models to understand sensory cortex.
1233 *Nature Neuroscience*, 19(3):356–365. DOI: <https://doi.org/10.1038/nn.4244>
- 1234 Yamins DLK, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ. 2014. Performance-optimized
1235 hierarchical models predict neural responses in higher visual cortex. *Proceedings of the*
1236 *National Academy of Sciences*, 111(23):8619–8624. DOI:
1237 <https://doi.org/10.1073/pnas.1403112111>
- 1238 Yang X, Yan J, Wang W, Li S, Hu B, Lin J. 2022. Brain-inspired models for visual object recognition:
1239 an overview. *Artificial Intelligence Review*, 1–49. DOI: [https://doi.org/10.1007/s10462-021-](https://doi.org/10.1007/s10462-021-10130-z)
1240 [10130-z](https://doi.org/10.1007/s10462-021-10130-z)
- 1241 Zhang K, Robinson N, Lee SW, Guan C. 2021. Adaptive transfer learning for EEG motor imagery
1242 classification with deep Convolutional Neural Network. *Neural Networks*, 136:1–10. DOI:
1243 <https://doi.org/10.1016/j.neunet.2020.12.013>