

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

The temporal evolution of conceptual object representations revealed through models of behavior, semantics and deep neural networks

Bankson, B. B.^{1*}, Hebart, M. N.^{1*}, Groen, I. I. A.¹, & Baker, C. I.¹

¹Section on Learning and Plasticity, Laboratory of Brain and Cognition, National Institute of Mental Health, National Institutes of Health, Bethesda, MD 20892, USA.

* equal contribution

Correspondence should be addressed to:

Brett B. Bankson
Laboratory of Cognitive Neurodynamics
UPMC Presbyterian
Suite B-400
200 Lothrop Street
Pittsburgh, PA 15213
bbb17@pitt.edu

Conflict of Interest: The authors declare no competing financial interests.

Acknowledgements: This work was supported by the Intramural Research Program of the National Institute of Mental Health (ZIA-MH-002909) - National Institute of Mental Health Clinical Study Protocol 93-M-0170, NCT00001360, a Feodor-Lynen fellowship of the Humboldt Foundation to M.N.H., and a Rubicon Fellowship from the Netherlands Organization for Scientific Research to I.I.A.G.

36 **Abstract**

37
38 Visual object representations are commonly thought to emerge rapidly, yet it has remained
39 unclear to what extent early brain responses reflect purely low-level visual features of these
40 objects and how strongly those features contribute to later categorical or conceptual
41 representations. Here, we aimed to estimate a lower temporal bound for the emergence of
42 conceptual representations by defining two criteria that characterize such representations: 1)
43 conceptual object representations should generalize across different exemplars of the same
44 object, and 2) these representations should reflect high-level behavioral judgments. To test these
45 criteria, we compared magnetoencephalography (MEG) recordings between two groups of
46 participants ($n = 16$ per group) exposed to different exemplar images of the same object
47 concepts. Further, we disentangled low-level from high-level MEG responses by estimating the
48 unique and shared contribution of models of behavioral judgments, semantics, and different
49 layers of deep neural networks of visual object processing. We find that 1) both generalization
50 across exemplars as well as generalization of object-related signals across time increase after 150
51 ms, peaking around 230 ms; 2) behavioral judgments explain the most unique variance in the
52 response after 150 ms. Collectively, these results suggest a lower bound for the emergence of
53 conceptual object representations around 150 ms following stimulus onset.

54

55 **Introduction**

56
57 There is enormous variability in the visual appearance of objects, yet we can rapidly recognize
58 them without effort, even under difficult viewing conditions (DiCarlo & Cox, 2007; Potter et al.,
59 2013). Evidence from neurophysiological studies in human suggests the emergence of visual
60 object representations within the first 150 ms of visual processing (Thorpe et al., 1996; Carlson
61 et al., 2013, Cichy et al., 2014). For example, the specific identity of objects can be decoded from
62 the magnetoencephalography (MEG) signal with high accuracy around 100 ms (Cichy et al., 2014).
63 However, knowing when discriminative information about visual objects is available does not
64 inform us about the nature of those representations, in particular whether they primarily reflect
65 (low-level) visual features or (high-level) conceptual aspects of the objects (Clarke et al., 2015).
66 To address this issue, in this study we employed multivariate MEG decoding and model-based
67 representational similarity analysis (RSA) to elucidate the nature of object representations over
68 time.

69 Previous studies have demonstrated increasing category specificity (Cichy et al., 2014),
70 tolerance for position and size (Isik et al., 2014) and semantic information (Clarke et al., 2013)
71 over the first 200ms following stimulus onset, suggesting some degree of abstraction from low-
72 level visual features. However, identifying the nature of object representations is an inherently
73 difficult problem: low-level features may be predictive of object identity, making it hard to
74 disentangle the relative contribution of low and high-level properties to measured brain signals
75 (Groen et al., 2017). In this study, we addressed this problem by combining tests for the
76 generalization of object representations with methods to separate the independent
77 contributions of low- and high-level properties. We focused on two specific criteria that would
78 need to be fulfilled for a representation to be considered conceptual. First, a conceptual

79 representation should generalize beyond the specific exemplar presented, not just variations of
80 the same exemplar. Second, a conceptual representation should also reflect high-level behavioral
81 judgments about objects (Clarke & Tyler, 2015). We consider fulfillment of these two properties
82 to provide a lower bound at which a representation could be considered conceptual.

83 We collected MEG and behavioral data from 32 participants allowing us to probe the
84 temporal dynamics of conceptual object representations according to the two criteria above. To
85 test for generalization across specific exemplars, we assessed the reliability of object
86 representations across two independent sets of objects. Further, we assessed the relation of
87 those object representations to behavior by comparing participants' behavioral judgments with
88 the MEG response patterns using RSA. Importantly, to isolate the relative contributions of low-
89 level and conceptual properties to those MEG responses, we identified the variance uniquely
90 explained by behavioral judgments, isolating low-level representations using deep neural
91 networks that have been shown to capture low- to mid-level responses in fMRI and monkey
92 ventral visual cortex (Cadieu et al., 2014; Cichy et al., 2016a; Eickenberg & Thirion, 2017; Güçlü
93 & van Gerven, 2015; Khaligh-Razavi & Kriegeskorte, 2014; Yamins et al., 2014; Wen et al., 2017).

94

95 **Methods**

96

97 *Participants*

98 32 healthy participants (18 female, mean 25.8, range 19-47) with normal or corrected-to-normal
99 vision took part in this study. As a part of a pilot experiment used for purely illustrative purposes
100 (see Figure 4a), 8 participants (5 overlap) completed the same behavioral task with a different
101 set of stimuli. All participants gave written informed consent prior to participation in the study
102 as a part of the study protocol (93-M-0170, NCT00001360). The study was approved by the
103 Institutional Review Board of the National Institutes of Health and was conducted according to
104 the Declaration of Helsinki.

105

106 *Stimuli*

107 We created two independent sets of 84 object images each that were cropped and placed on a
108 grey background. Each stimulus set contained a unique exemplar for each of the 84 object
109 concepts, as shown in Figure 1a. We selected object concepts by using a combination of two
110 word databases, one of word frequency (Corpus of Contemporary American English, Davies,
111 2008) and the other of word concreteness (Brysbaert et al., 2014). First, based on our corpus we
112 selected the 5000 most frequent nouns in American English. From this set of words, we then
113 selected nouns with concreteness ratings > 4/5. Finally, for words that would be difficult or
114 impossible to distinguish when presented as an image (e.g. 'woman', 'mother', 'wife'), we used
115 only the most frequent entry. This selection left us with a set of 112 objects.

116 To evaluate whether those categories would be labeled consistently, we generated three
117 distinct images of each object concept and asked three individuals who were not involved in the
118 study to provide a verbal label for each of the three versions of the 112 objects. Images that were
119 not labeled correctly by all raters were discarded, leaving us with 84 object concepts. From the
120 three sets of object images, we then randomly sampled two per object concept. This generated

121 two sets of unique object exemplars for 84 object concepts, divided into Image Set 1 and Image
122 Set 2. The two sets of object stimuli are shown in Supplemental Figure S1.

123

124 *Procedure*

125

126 *MEG*

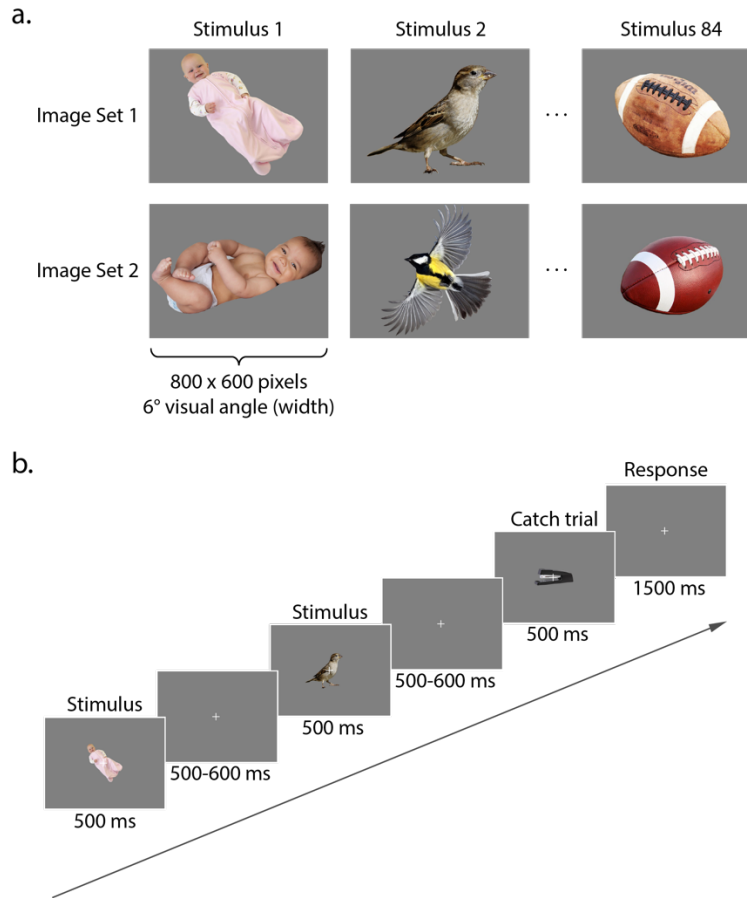
127 During MEG recordings, participants were seated upright in an electromagnetically shielded MEG
128 chamber. Stimuli were presented using the Psychophysics Toolbox (Brainard, 1997) in MATLAB
129 (version 2016a, Mathworks, Natick, MA). Visual stimulation was controlled by a Panasonic PT-
130 D3500U DLP projector with an ET-DLE400 lens, located outside of the chamber and projected
131 through a waveguide and series of mirrors onto a back-projection screen in front of the
132 participant. Participants were assigned to one of two groups and completed the experiment with
133 either Image Set 1 or Image Set 2. All stimuli were presented on a grey background with a white
134 fixation cross in the center (viewing distance: 70 cm, stimulus width: 6° of visual angle).
135 Participants completed an oddball detection task, pressing a button in response to catch trials
136 containing the oddball stimulus (desk stapler) that appeared pseudorandomly every 2-6 trials
137 (average 4, flat distribution). On each trial (Figure 1b), an object stimulus was presented at
138 fixation for 500 ms, followed by a variable fixation period (regular trials: pseudorandomly 500-
139 600 ms, catch trials: 1500 ms). In addition, participants were instructed to blink their eyes only
140 as they pressed the button of the MEG-compatible button box during catch trials, in order to
141 avoid any eye blink artifacts at other points of the experiment. Participants completed 18 runs,
142 viewing each of the 84 images 36 times over the course of the experiment.

143

144 *Behavior: Object arrangement task*

145 Within two days of completing the MEG session, participants took part in a follow-up behavioral
146 experiment to provide us with behavioral estimates of the representational similarity between
147 all possible object pairs. This was done using the object arrangement method (Goldstone 1994;
148 Kriegeskorte & Mur, 2012). In this method, participants arrange objects in a 2D “arena” based on
149 their subjective similarity, and the distance between the items is used to generate $(n \times n-1)/2$
150 pairwise distance estimates between object pairs. Participants were seated in front of a monitor
151 and completed the object arrangement task on the same 84 object images used in the MEG
152 experiment. All items were presented at once around the circular arena. Participants were
153 instructed to use the computer mouse and arrange the items according to their similarity at their
154 own pace, taking ~20 minutes on average to complete the task. We deliberately did not provide
155 participants with an explicit strategy or instructions on what object features to focus, so as to not
156 bias them to focus on any specific aspect of the stimuli. To facilitate the task, when a participant
157 clicked on a certain image around the arena, an enlarged version spanning 150×200 pixels was
158 displayed in the top right of the computer screen. After completion of the experiment, we
159 extracted the pixel-wise distance between each pair of items, yielding an 84×84 distance matrix
160 for each participant.

161



162
163 **Figure 1.** Stimulus format and trial progression. **a.** Two unique object exemplars were selected for each of the 84
164 object concepts used in the study. **b.** Stimuli were presented on a grey background for 500 ms, followed by fixation
165 for 500-600 ms (catch trials: 1500 ms). All 84 stimuli from both image sets are shown in Supplemental Figure S1.

166

167 *MEG acquisition and preprocessing*

168 MEG data were recorded continuously at a sampling rate of 1200 Hz with a 275-channel CTF
169 whole-head MEG system (MEG International Services, Ltd., Coquitlam, BC, Canada). All analyses
170 were conducted in MATLAB (version 2016a, The Mathworks, Natick, MA). Preprocessing was
171 carried out using Brainstorm 3.4 (version 02/2016, Tadel et al., 2011) and custom-written code,
172 using similar preprocessing steps as previously published MEG decoding work (Cichy et al., 2014;
173 Grootswagers et al., 2016, Hebart et al., 2017). Recordings were available from 272 channels
174 (dead channels: MLF25, MRF43, MRO13). The whole-head array consists of radial first-order
175 gradiometer channels equipped with synthetic third-gradient balancing to remove background
176 noise online. At the beginning of the experiment and after every third experimental run,
177 participants' head position was localized based on fiducial coil placement at the nasion, left and
178 right preauricular points. Data were bandpass filtered between 0.1 and 300 Hz, and bandstop
179 filtered at 60 Hz and harmonics. We segmented the data into single trial bins, with each trial
180 consisting of 100 ms baseline for normalization purposes and 1000 ms post-stimulus activity,
181 yielding a total of 1321 time samples for each trial. Oddball trials were discarded.

182 Three pre-analysis steps allowed us to increase SNR and reduce computational demand:
183 PCA dimensionality reduction, temporal smoothing on PCA components, and data

184 downsampling. Principal components analysis (PCA) was run to reduce the number of channels
185 into the set of most descriptive components. All data for an MEG channel across trials were
186 concatenated for PCA, and the components explaining the lowest 1 % of variance after PCA were
187 removed to speed-up further processing, with a minimum of 136 components chosen a priori as
188 a cut-off. Data across all time points were normalized according to the baseline period of -100 to
189 0 ms relative to stimulus presentation. To do so, the mean and standard deviation of the baseline
190 period for each component were computed, and the mean was subtracted from the data before
191 dividing by the standard deviation. We then used a Gaussian kernel of ± 15 ms half duration at
192 half maximum (HDHM) to temporally smooth the remaining components, and downsampled the
193 components to 120 Hz (132 samples / trial).

194

195 *Multivariate decoding and temporal generalization analysis*

196

197 *Multivariate MEG decoding*

198 Our goal was to study the representational dynamics during visual object recognition and the
199 emergence of generalizable, conceptual object representations over time. To determine the
200 amount of object information contained in the MEG signal over time, we ran time-resolved
201 multivariate decoding of MEG data using a linear support vector machine classifier (SVM; Chang
202 & Lin, 2011). The analysis steps were chosen according to general recommendations
203 (Grootswagers et al., 2016) and a recent study from our lab (Hebart et al., 2017). Multivariate
204 analyses were conducted using functions from The Decoding Toolbox (Hebart et al., 2015) and
205 custom-written code. The following analysis steps were applied to all participants, regardless of
206 experimental group.

207 First, we created supertrials by averaging 6 trials of the same object concept drawn
208 randomly without replacement (Isik et al., 2014). For each time point, preprocessed MEG data
209 within each supertrial were arranged as P dimensional measurement vectors (corresponding to
210 the number of components from PCA preprocessing), yielding K pattern vectors for each time
211 point and object concept. For each pair of object concepts and each time point, we then trained
212 the classifier on $K-1$ pattern vectors and tested it on the pair of left-out pattern vectors, yielding
213 a decoding accuracy for each pair of object categories at each time point. The assignment to
214 training and testing sets and resulting classification procedure was repeated 100 times for each
215 pair of object concepts and each time point, with a new random generation of supertrials in each
216 iteration. The resulting decoding accuracies were averaged across the 100 iterations and
217 presented as an 84×84 matrix at every time point, with rows and columns indexed according to
218 object conditions, and with the diagonal undefined. We used these matrices to evaluate average
219 decoding accuracy at each time point by computing the average of the lower triangular matrix.

220 Significance for the decoding analysis was assessed using a sign permutation test. A null
221 distribution of group means was generated by running the decoding procedure 1000 times,
222 randomly assigning a positive or negative sign value to decoding accuracies and averaging those
223 values. P -values were determined as one minus the percentile of the original group mean in this
224 null distribution. Those p -values were corrected according to the false-discovery rate (FDR) and
225 were deemed significant if the corrected p -value did not exceed 0.05.

226

227 *Temporal generalization of object representation*

228 While time-resolved multivariate decoding can reveal when specific mental representations are
229 present in patterns of neural activity, it cannot identify how said patterns at one time point relate
230 to other time points. We were interested in investigating the extent to which object-related
231 information is static or dynamic over time, which can give us an index of how rapidly neural
232 signals evolve. To investigate this, we conducted a cross-classification analysis over time, also
233 known as the temporal generalization method (King & Dehaene, 2014; Meyers et al., 2008). If a
234 classifier can successfully generalize from one time point to another, this shows that
235 representational content has not changed between these two time points. Conversely, if the
236 classifier does not generalize, this shows that patterns of neural activity have evolved to an extent
237 that representational content is no longer similar.

238 To carry out this temporal generalization analysis, we used the same classification
239 approach described above; however, instead of only testing the classifier at the same time point
240 we also tested its performance at all other time points. We repeated the analysis with all time
241 points each serving as training data once for the classifier, and generated a 132 x 132 time-time
242 decoding matrix that shows the extent to which our classifier generalizes across time.

243

244 *Representational similarity analysis (RSA)*

245

246 RSA is a method to analyze and compare data patterns, for example brain activity patterns with
247 behavioral judgments or computational models (Kriegeskorte et al., 2008). Instead of comparing
248 these patterns directly, in RSA patterns are converted to representational similarity matrices
249 (RSMs), quantifying all pairwise similarities of all patterns. These RSMs can then be compared to
250 other RSMs based on other data.

251 In this study, we used RSA for two purposes. First, across participants we directly
252 compared the time courses of MEG RSMs evoked by the *same* exemplar with MEG RSMs evoked
253 by *different* exemplars. This allows an estimate of the generalizability of representations across
254 exemplars and thus the extent to which a representation reflects high-level versus low-level
255 properties, assuming that a generalized representation indicates a more high-level, conceptual
256 representation. Second, we used RSA to study the relationship between evoked MEG activity
257 patterns and computational, semantic, and behavioral models. In particular, we wanted to
258 identify time periods at which the MEG responses reflected predominantly behavioral
259 judgments, which we take as an index of high-level conceptual processing. To do this, we
260 quantified the unique and shared variance of each model RSM with RSMs based on MEG activity
261 patterns.

262

263 *Construction of MEG similarity matrices*

264 MEG RSMs were constructed as follows. For each time point, we averaged the preprocessed MEG
265 data for all 36 trials of each object concept, yielding 84 object concept MEG patterns. Then we
266 computed the similarity between all pairs of those 84 patterns, yielding an 84 x 84 MEG RSM for
267 each time point. We then analyzed these RSMs further for the two purposes described above.

268

269 *Generalization of MEG similarity patterns across exemplars*

270 To determine time periods that generalize between representations of object exemplars, we
271 compared the time courses of similarity of RSMs *within* each image set to the similarity *between*

272 image sets. To this end, we split data between the groups for Image Set 1 and Image Set 2 and
273 conducted within- and between-group split-half correlation analyses with the RSMs for each
274 participant. We chose a repeated subsampling procedure within group to allow us to use the
275 same analysis within and between groups. The following analyses are described for one RSM at
276 one time point, but were repeated for all time points.

277 Within each group of participants ($n = 16$), we randomly assigned participants' RSMs to
278 one of two arbitrary subsets of 8 participants and averaged participants' RSMs within subsets.
279 Next, we calculated the Spearman rank correlation coefficient between the lower triangular part
280 of each 84×84 matrix, separately for every time point. We repeated this split-half analysis 1000
281 times with novel assignments of participants and averaged across repetitions, yielding a time
282 course of within-exemplar correlation. The same procedure was completed for the between-
283 group split-half analysis, but here the two subsets were each drawn from eight randomly selected
284 participants in each group, yielding a time-course of between-exemplar correlations.

285 To assess statistical significance, we conducted a randomization test. We repeated the
286 analysis above 1000 times (i.e. a total of 10^6 split-half analyses, for both within-exemplar and
287 between-exemplar comparisons). For each of those 1000 randomizations, we randomly
288 permuted the rows and columns of the matrices in one of the subgroups before calculating
289 Spearman's r . P -values were determined as one minus the percentile of the original split-half
290 analysis, and FDR-corrected to $p < 0.05$.

291
292 *Representational similarity matrices for computational models and behavior*
293 To access the representational content of the MEG data across time, we chose multiple
294 behavioral and computational models that we later compared to MEG data: a behavioral model
295 based on the group mean behavioral similarity, a semantic model to capture similarity at the
296 semantic level, and two layers of a deep neural network to capture different visual processing
297 stages. For a first comparison, we characterized the pairwise similarity of these models to assess
298 their general similarity irrespective of MEG. We calculated Spearman's r for each pair of models.
299 Significance of correlations was tested using a randomization test: The rows and columns of one
300 model RSM were randomly permuted before computing the Spearman's r between with the
301 other model RSM. This procedure was repeated 1000 times to generate a null distribution of
302 correlation coefficients, and results were deemed significant if they showed a higher correlation
303 coefficient than the distribution cut-off determined by a level of $p < 0.05$.

304
305 Behavior
306 We generated an RSM for behavioral judgments by extracting the 84×84 distance matrices from
307 each participant within a group and averaging them together. Next, we converted this
308 dissimilarity matrix to an RSM by subtracting the dissimilarities from 1. This step yielded two
309 group-level behavior RSMs corresponding to Image Set 1 and Image Set 2.

310 Semantic model: Global Vectors for Word Representation (GloVe)
311 Global Vectors for Word Representations (GloVe) is an unsupervised algorithm that is trained on
312 corpus word co-occurrence statistics to yield vector representations for words in the corpus,
313 representing semantic relationships between words (Pennington et al., 2014). As a distributional
314 measure of the semantic relatedness of words based on their shared linguistic contexts, GloVe is
315 similar to other traditional co-occurrence models of word meaning but is particularly well-suited

316 to the analysis here because of the high-dimensional similarity structure that shows semantic
317 similarity between pairs of individual words, outperforming similar models in similarity tasks. As
318 such, the structure of GloVe provides a fine-grained metric to evaluate how the representational
319 space of MEG signals reflects semantic relationships as derived from shared lexical contexts. We
320 chose 50-dimensional word vectors pre-trained on a 6-billion token Wikipedia database,
321 extracted them for each object concept in the stimulus set and calculated Spearman's r between
322 each pair of vectors, generating an 84×84 RSM.

323 Visual model: Deep neural network VGG-F

324 We used the MatConvNet toolbox (Vedaldi & Lenc, 2015) to implement a pre-trained version of
325 the Visual Geometry Group-Fast deep neural network (VGG-F DNN) (Chatfield et al., 2014) that
326 was trained to perform the ImageNet ILSVRC 2012 object classification task. This network was
327 chosen based on its high classification performance, ease of implementation, and suitability for
328 our visual object concept stimuli. DNN representations for each image in both image sets were
329 extracted from both convolutional layers (1-5) and fully-connected layers (6-8) of the network.
330 We focused on representative examples of the convolutional and fully connected layers (3 and
331 7, respectively) to reflect low-to-midlevel vision and high-level vision, respectively. Within each
332 layer, we calculated Spearman's r between each of the object conditions that yielded an 84×84
333 RSM for both layers within each participant group. This yielded four distinct RSMs: DNN Layer 3
334 and Layer 7 for Image Set 1, and DNN Layer 3 and Layer 7 for Image Set 2.

335

336 *Representational similarity analysis: Model comparisons to MEG*

337 To directly compare each model to MEG activity patterns, we calculated Spearman's r between
338 the lower diagonals of the model variables and MEG RSMs at each time point within each group.
339 These group-specific correlations were averaged together to yield a time course showing the
340 level of correlation between the model and MEG responses. Upper and lower bounds for noise
341 ceilings were determined within each of the two groups of participants according to Nili et al.
342 (2014): The upper bound was estimated by calculating the correlation between each participant's
343 RSM and the mean group RSM *including* that participant, while the lower bound was estimated
344 by calculating the correlation between each participant's RSM and the mean group RSM
345 *excluding* that participant. The upper and lower bounds from each group were averaged together
346 to yield a mean noise ceiling across all participants. The statistical significance of this suite of
347 representational similarity analyses was determined using randomization tests as described
348 above, permuting the rows and columns of a given model RSM (behavior, GloVe, DNN Layer 3,
349 DNN Layer 7) and for each randomization computing correlation time courses with the original
350 MEG RSMs. Correlations were deemed significant if they exceeded a correlation cut-off
351 determined by a level of $p < 0.05$ (FDR-corrected).

352

353 *Establishing the unique and shared contributions of individual models*

354 To determine the unique and shared variance between models and MEG signals, we conducted
355 multiple linear regression analyses using the behavior RSM, DNN Layer 3 RSM, and DNN Layer 7
356 RSM as model variables to predict MEG RSMs from these variables. Given the complexities of
357 describing the unique and shared variance partitions of more than three model variables, we
358 decided to exclude the GloVe model, which showed the weakest correlation with MEG. By
359 conducting a series of different multiple regressions with different combinations of model

360 variables, this approach allows us to determine not only the unique MEG variance explained by
361 each model individually, but also the variance shared between any combination of models.
362 Before conducting variance partitioning analyses, we averaged the group-specific RSMs of both
363 image sets for behavior and DNN models, which yielded very similar results as compared to
364 calculating them separately and averaging results afterwards. We extracted the lower diagonal
365 from the mean MEG RSM at each time point as dependent variables, and assigned each of the
366 models as independent variables. In sum, 7 variance partitioning models were tested: 1) all
367 models combined (behavior, DNN Layer 3, DNN Layer 7), (2-4) all pairwise combinations of two
368 models (behavior and DNN Layer 3, behavior and DNN Layer 7, DNN Layer 3 and DNN Layer 7),
369 and (5-7) each model alone. Comparing the explained variance (R^2) values of a single model and
370 the R^2 of the same model in conjunction with another model yields the amount of variance
371 independently explained by that model (see also Groen et al., 2012; Lescroart et al., 2015; Greene
372 et al., 2016; Hebart et al., 2017). Statistical significance was determined using a randomization
373 test as described above, randomizing columns and rows of model matrices 1000 times and
374 repeating the original analysis. For a given iteration, the same randomization was used across all
375 models to fulfill the assumptions of the randomization test. Significance cutoffs for R^2 were set to
376 $p < 0.05$ (FDR-corrected).
377
378

379 **Results**

380
381 Our aim in this study was to characterize the emergence of conceptual representations for visual
382 objects. We applied multivariate decoding and representational similarity analysis to MEG data
383 to examine (1) how object representations generalize across time and object exemplars, and (2)
384 to elucidate the unique and shared contributions of behavioral judgments to measured MEG
385 responses. The resulting temporal profiles inform us about stages of object processing from low-
386 level visual to conceptual representations.
387

388 *Time-resolved representation of object identity*

389 To characterize the time course by which neural signals in the human brain convey information
390 about object identity, we used time-resolved multivariate decoding, conducting pairwise
391 classification between MEG patterns in response to object stimuli (Figure 2a). Object identity
392 information rose rapidly in response to stimulus presentation, with decoding accuracy peaking
393 at 100 ms (mean accuracy: 91.1 %), followed by a slow decay of information that remained
394 significantly above chance after stimulus offset and for the duration of the trial time window
395 (1000 ms post stimulus onset). These results indicate that we were able to detect the temporal
396 unfolding of object-identity information encoded in MEG signals with high accuracy, establishing
397 a correspondence to previous research demonstrating that discriminable object representations
398 emerge well within 100 ms of visual recognition (Carlson et al., 2013; Cichy et al., 2014). Further,
399 these results lay an important foundation for the following analyses in which we delineate what
400 information specifically contributes to these discriminable object representations.
401
402

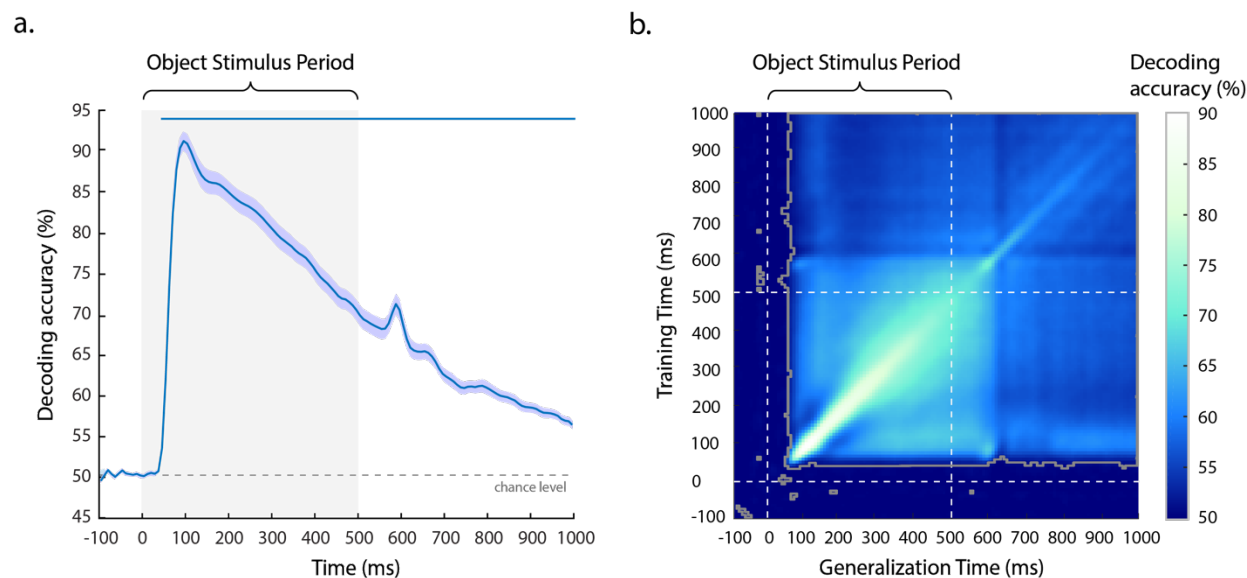
403 *Temporal generalization of object information*

404 While time-resolved multivariate decoding reveals the temporal evolution of discriminable
405 object representations, it does not inform about the dynamics and stability of those
406 representations across time. To identify the degree to which object representations generalize
407 across time, we ran a temporal generalization analysis by training a classifier on data at every
408 time point and testing it at all other time points. This yielded a temporal generalization matrix
409 (Figure 2b), with the diagonal representing training and testing at the same time points, mirroring
410 the results presented in Figure 2a. In a temporal generalization matrix, a dynamic representation
411 would be characterized by high accuracies around the diagonal and low accuracies everywhere
412 else, indicating little generalization across time. In contrast, a stable neural representation would
413 exhibit high decoding around the diagonal but also in the off-diagonal time points, demonstrating
414 a similar representation across time.

415 Our results exhibited significant generalization from ~70 ms onward, demonstrating a
416 shared representational format across the entire trial. While this result reveals a persistent
417 representation across time, the strength of generalization varies. Focusing on the first half of the
418 stimulus presentation period, the results revealed a period of increased temporal dynamics
419 between ~70-250 ms, indicated by the high decoding accuracy on the diagonal and lower
420 decoding accuracies away from the diagonal. This result suggests a relatively dynamic
421 representational format in this phase of visual processing. After ~250 ms, there was increased
422 generalization away from the diagonal, indicating a more persistent, shared representational
423 format during this later phase of visual processing. Interestingly, there was a generalization
424 period between time windows of ~70-100 ms and ~250-550 ms, suggesting an overlap of
425 representations between early visual and later conceptual processing. The markedly lower
426 information generalization between 150-250 ms and all other time points suggests the
427 information dynamics at these points are computationally dissimilar from other stages of
428 processing.

429 Taken together, these results reveal relatively weak but significant persistence of stable
430 object information throughout the entire trial. On top of this, the results reveal a general
431 broadening of information generalization after an early phase of visual processing. This
432 broadening suggests early dynamic neural activity followed by the emergence of more stable
433 object representations around 250 ms.

434



435
436

437 **Figure 2. a.** Time-resolved multivariate decoding of object identity across the trial. After onset of the object stimulus
438 (Object Stimulus Period), pairwise object decoding accuracy increased rapidly, followed by a slow decay towards
439 chance over the duration of the trial. Error bars reflect SEM across participants for each time point separately.
440 Significance is indicated by colored lines above the plot (non-parametric cluster-correction at $p < 0.05$). **b.** Temporal
441 generalization matrix for object identity. The y-axis depicts the classifier training time relative to stimulus onset, and
442 the x-axis classifier generalization time relative to stimulus onset. Dotted lines indicate stimulus onset and offset.
443 Areas bounded by a grey line contain significant temporal cross-decoding accuracy values ($p < 0.05$, FDR corrected).
444

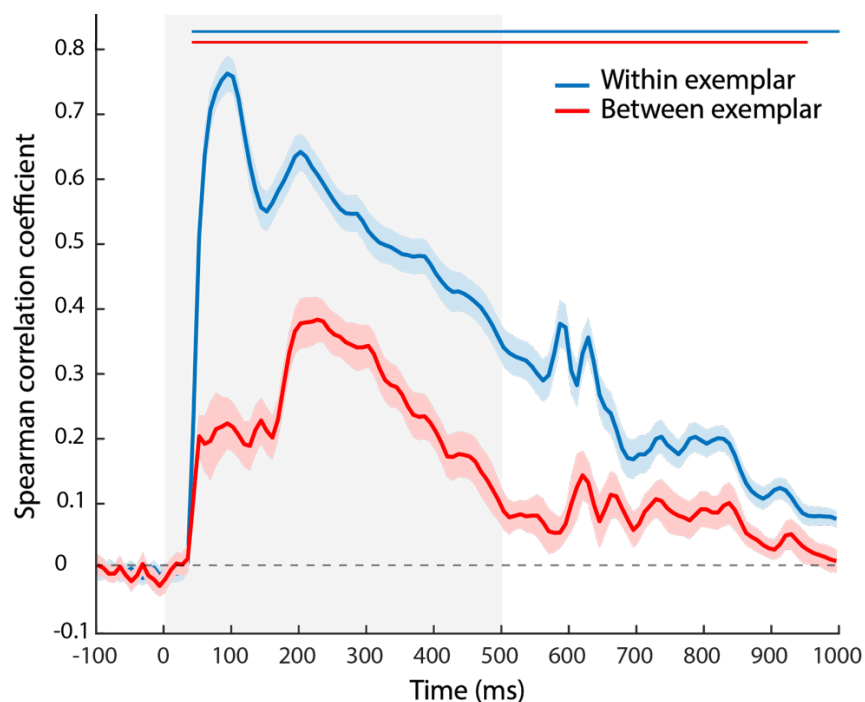
445 *Criterion 1 for conceptual object representation: Generalization between object exemplars*
446 Having established the time course of object identity-specific information, we investigated when
447 those brain responses reflect conceptual object representations. One prerequisite of a
448 conceptual object representation is a similar representational format between multiple
449 exemplars of the same object, since a conceptual representation is expected to generalize
450 beyond each individual exemplar. The data collected from Image Set 1 and 2 allow direct
451 comparison of representational similarity across exemplars for the same visual object concept
452 (Figure 3). We measured this generalization of object concept-specific information by (1)
453 calculating the correlation of within-exemplar MEG RSMs for participants who were shown the
454 same object exemplar and (2) calculating the generalization of between-exemplar MEG RSMs for
455 participants who were shown different object exemplars. Then we compared the shape of these
456 MEG correlation time courses.

457 A comparison of within-exemplar and between-exemplar MEG RSM correlations revealed
458 a generally higher correlation within-exemplar than between-exemplar (mean difference across
459 time: Spearman's r : 0.18, $p < 0.001$, randomization test), indicating that differences between
460 exemplars persisted throughout most of the trial. Reliable structure for within-exemplar MEG
461 RSMs emerged rapidly, peaking at 93 ms (mean Spearman's r : 0.77). This was followed by a fast
462 drop in correlation, and then another rise beginning around 160 ms and peaking at 202 ms (mean
463 Spearman's r : 0.65), after which within-exemplar correlations decreased steadily for the
464 duration of the trial while remaining significantly above chance. The correlation of between-
465 exemplar MEG RSMs also initially increased rapidly, but then reached a plateau at a comparably
466 low level of correlation between ~ 70 and ~ 160 ms (mean Spearman's r : 0.21). Importantly,

467 between-exemplar reliability then increased again after ~160 ms, peaking at 227 ms (mean
468 Spearman's r : 0.39). Between-exemplar correlation then slowly decayed back to 0, but remained
469 significantly above chance until 960 ms after stimulus presentation.

470 These results reveal an important dissociation: While within-exemplar correlations
471 reached their maximum around 100 ms, between-exemplar generalization was maximal around
472 200 ms. Thus, this analysis reveals an early processing stage during which generalization is limited
473 by the variable visual features of each individual exemplar, and a later processing stage where
474 the increased generalization likely reflects the development of a more conceptual object
475 representation that is consistent across exemplars.

476

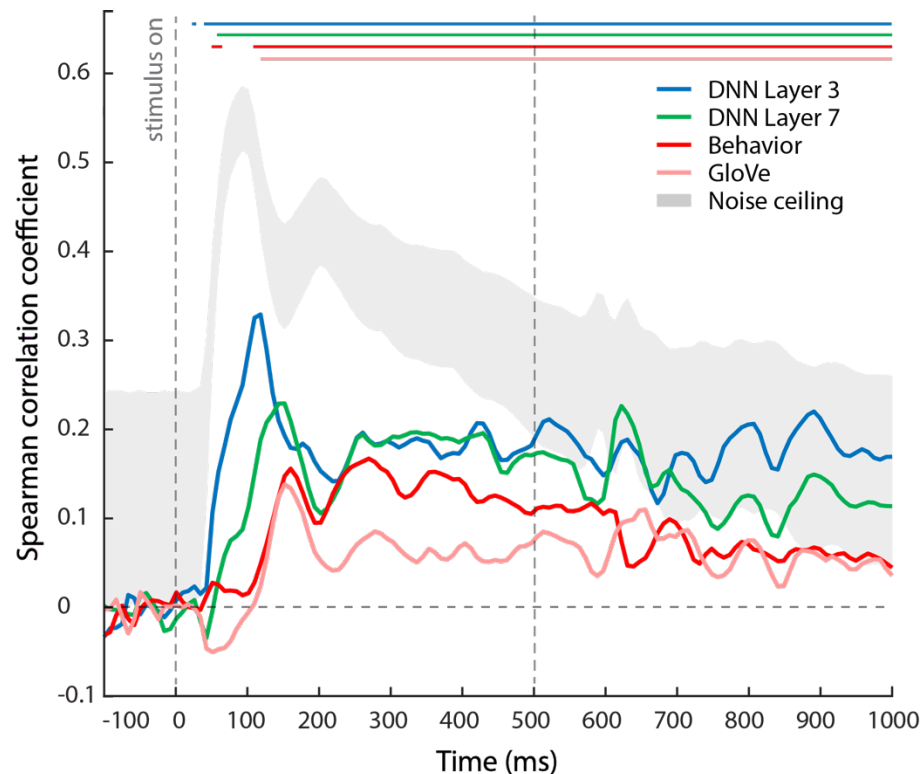


477 **Figure 3.** Within and between exemplar correlation of MEG RSMs. Within-exemplar correlation was generally higher
478 than between-exemplar correlation. Both within and between-exemplar correlations revealed an early peak (93 ms)
479 and a late peak (202 and 227 ms, respectively), with the early peak being higher than the late peak for within-
480 exemplar correlations, and the late peak being higher than the early peak for between-exemplar correlations. Error
481 bars reflect SEM. Significance is indicated by colored lines above the accuracy plot (non-parametric cluster-
482 correction at $p < 0.05$).
483

484 *Comparison of behavior and computational models of low-level and high-level processing*

485 To quantify how the RSMs derived from behavior (perceptual judgments, visualized in Figure 4b),
486 GloVe (lexical semantics), DNN Layer 3 (low/mid-level visual information), and DNN Layer 7 (high-
487 level visual information) relate to one another, we computed the correlation between each pair
488 of model RSMs (Figure 4a). For visualization purposes, we applied hierarchical clustering to
489 independent pilot data of the behavioral task to sort objects depicted in the model RSMs (Figure
490 4a). All model correlations were significant at a level of $p < 0.001$ (randomization test). An
491 estimate of the upper noise ceiling for possible model correlation values was calculated by the
492 correlation between behavior RSMs for the two groups of participants (Spearman's $r = 0.64$). The
493 greatest similarity to behavior was shown by the GloVe model. There was low similarity of
494

526 behavior with MEG signals after 150 ms raises the question whether the behavioral correlations
527 can be fully explained by the features represented in the DNN models.
528



529
530
531 **Figure 5.** Results of model-based representational similarity analysis with MEG data. Comparison includes models
532 based on DNN Layer 3, DNN Layer 7, GloVe and behavior. The results exhibit a progression of peaks from DNN Layer
533 3 to behavior, suggesting a temporal evolution of the underlying representation from more low-level to higher-
534 level/conceptual. Grey shaded area depicts the noise ceiling.

535
536 *Variance Partitioning: Shared and unique model contributions*

537 To provide a deeper understanding of the unique contributions of models to MEG variance, we
538 conducted a variance partitioning analysis in which we compared the results of different multiple
539 regression analyses applied to MEG RDMS (see Methods; Figure 6a). We first considered the total
540 percent of variance in the MEG RDMS explained when all three predictors are combined in a
541 single regression model ('full model') in comparison to the percent variance explained by each
542 model separately (Figure 6b). Since variance explained by each model separately is identical to
543 the square of the model correlation, the results of this analysis are very similar to those of the
544 previous section presented in Figure 5, with the only difference that these results were collapsed
545 across groups before conducting variance partitioning. Explained variance of DNN Layer 3 peaked
546 at 118 ms (R^2 : 11.0 %), DNN Layer 7 at 151 ms (R^2 : 7.0 %), and behavior at 160 ms (R^2 : 4.8 %).
547 Importantly, however, the dashed line indicates how these contributions relate to the total
548 variance accounted for by all three models combined. At its peak at 118 ms, the full model
549 explains 11.6 % of the variance, which is similar to the amount of variance explained by DNN
550 Layer 3 alone, suggesting that all variance captured at this time-point can be attributed uniquely
551 to DNN Layer 3, with limited additional contribution of DNN Layer 7 or behavior. At later time

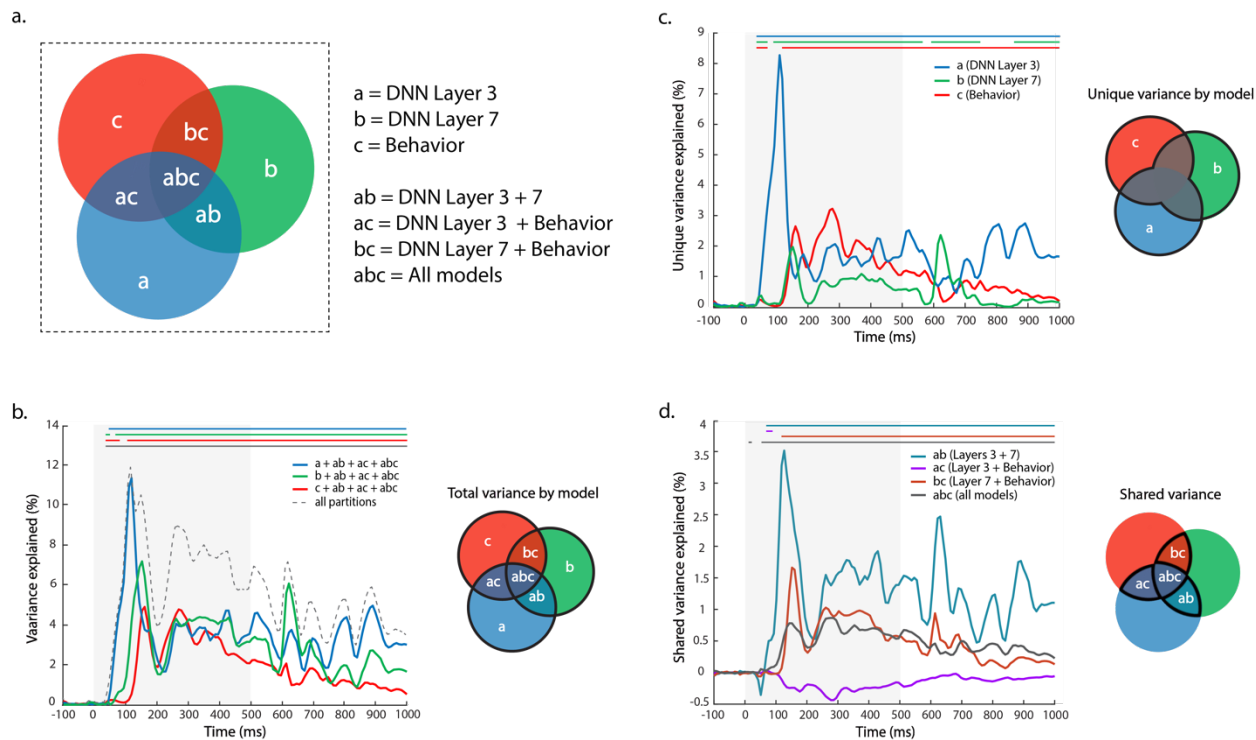
552 points, however, the full model always substantially explains more variance than the individual
553 predictors, providing a first clue that some or all of these predictors contribute unique (i.e.,
554 additive) variance.

555 To directly quantify the unique and shared variance of each model, we compared the
556 regression outcomes with different model variables included (Figures 6c, 6d). The unique MEG
557 variance explained by DNN Layer 3 peaked very early in time, at 109 ms (R^2 : 8.3%). DNN Layer 7
558 peaked next at 151 ms (R^2 : 2.0 %), followed closely by behavior at 160 ms (R^2 : 2.6 %), with a
559 second peak at 277 ms (R^2 : 3.2 %). Importantly, DNN Layer 3 explained the most unique variance
560 until 143 ms, after which behavior predicted the most unique variance until ~400 ms. Thus, while
561 all three models (DNN Layer 3, DNN Layer 7 and behavior) captured some unique variance in
562 MEG activity throughout the trial, behavior dominated after around 150 ms.

563 Finally, to complete the picture, we partitioned the variance into shared contributions
564 from combinations of the different models. Both DNN Layers contributed the most shared
565 variance across all time points after stimulus onset, which is perhaps not surprising considering
566 that both layers are derived from the same computational model. This shared variance between
567 DNN Layer 3 and Layer 7 peaked at 126 ms (R^2 : 3.5 %). Interestingly, the shared variance between
568 behavior and DNN layer 7 demonstrated a clear peak at 151 ms (R^2 : 1.7 %), suggesting that it is
569 around this time-point that DNN Layer 7 best captures neural information that is also reflected
570 in behavior. The shared variance between DNN Layer 3 and behavior was slightly negative, a
571 result that is not untypical for variance partitioning, indicative of small suppression effects
572 (Pedhazur, 1997) and suggesting that DNN layer 3 does not capture information that is relevant
573 for behavioral judgments.

574 Collectively, the variance partitioning results indicate that behavioral judgments are
575 reflected in the MEG response above and beyond what is captured by the DNN, with behavioral
576 judgments explaining the most unique variance between 200 and 400 ms after stimulus onset.
577 Further, before 150 ms, DNN layer 3 explains the most variance, suggesting that representations
578 prior to this point are unlikely to be conceptual in nature.

579



580
581
582
583
584
585
586
587

Figure 6. Time-resolved variance partitioning: Total, shared, and unique MEG variance explained by models: DNN Layer 3, DNN Layer 7, and behavior. **a.** Schematic of unique and shared variance components using a Venn diagram. **b.** Percent MEG variance explained by each model independently (colored lines), and total MEG variance explained at all time points (dotted line). **c.** Unique variance explained by each model. **d.** Shared variance between different model combinations.

588 Discussion

589

590 In this study, we investigated the temporal evolution of visual object representations. In
 591 particular we focused on determining a lower bound for the emergence of conceptual
 592 representations of objects. We proposed two criteria that would reflect conceptual
 593 representations: 1) generalization of representations between different exemplars of the same
 594 object, and 2) relationship to high-level behavioral judgments. We find qualitatively different
 595 processing of objects over time: Early responses (< 150 ms) were characterized by exemplar-level
 596 representations and similarity with computational visual models, whereas later responses (> 150
 597 ms) showed increasing generalization across exemplars and similarity with behavioral judgments,
 598 with greater stability of representations over time.

599 To evaluate generalization of representations reflecting conceptual processing, we
 600 compared the representational structure of MEG responses, both within exemplar and between
 601 sets of exemplars. This analysis revealed two interesting features. First, between-exemplar
 602 generalization was found to be consistently lower than within-exemplar generalization,
 603 demonstrating the persistence of exemplar-specific responses. This reduced between-exemplar
 604 generalization likely reflects the impact of low-level features varying between different
 605 exemplars. The fact that this advantage is maintained throughout the trial, suggests some

606 persistence of low-level feature representation. This interpretation is supported by the temporal
607 generalization even for very early time points and the variance explained by DNN Layer 3 (which
608 likely corresponds to early to mid-level visual processing, Cichy et al., 2016a; Güçlü & van Gerven,
609 2015; Wen et al., 2017), throughout the trial. Second, both within and between-exemplar
610 generalization showed two distinct peaks, one early around 100 ms, and another late around 200
611 ms. However, their relative amplitude was reversed: While the early peak was stronger than the
612 second within-exemplar, this pattern was reversed between-exemplar. This striking increase in
613 generalization between exemplars that occurs for the later peak suggests the emergence of a
614 common representation across exemplars, a key marker for conceptual representations.
615 Together these results suggest that the earliest time point for the emergence of conceptual
616 representations is around 150 ms, but also suggest a prolonged representation of low-level visual
617 features.

618 To evaluate the relationship to high-level behavioral judgments, we compared models
619 derived from behavior, semantics (GloVe), and computational vision (DNN) with the MEG
620 response to objects. We found that all models show significant correlation with the MEG
621 response throughout most of the trial. The early DNN layer showed the strongest and earliest
622 correlation, while the GloVe model showed the weakest correlation. This result highlights the
623 importance of testing multiple models rather than relying on a significant effect for a single
624 model. Since the models themselves are correlated (Figure 4), this demonstrates that testing
625 multiple models is *also* not sufficient; it is important to determine the unique and shared variance
626 explained by the different models (Lescroart et al., 2015; Groen et al., 2012; Greene et al., 2016;
627 Hebart et al., 2017), motivating our variance partitioning analysis. Given the complexities of
628 describing the unique and shared variance partitions of more than three model variables, we
629 decided to exclude one of the four. Since the GloVe model showed the weakest correlation with
630 MEG we focused on the DNN and behavioral model variables.

631 The variance partitioning revealed several important features. Focusing on the unique
632 contribution of each model variable, it becomes clear that DNN Layer 3 dominates early MEG
633 responses peaking at 100 ms, whereas behavior explains the most variance after 150 ms, peaking
634 at 270 ms. This result fulfills our second criterion – relationship with high-level behavioral
635 judgments – converging with the results of both the temporal generalization analysis and the
636 representational generalization across exemplars in identifying the time period after around 150
637 ms as reflecting a lower bound for the emergence of conceptual representations. Focusing on
638 the shared contribution of model variables, the results largely reflect the correlations between
639 model variables (Figure 4), e.g. no shared variance between DNN Layer 3 and behavior, high
640 shared variance between Layers 3 and 7 of the DNN model. However, they provide important
641 information about the timing of the shared variances. In particular, the shared variance between
642 DNN Layers 3 and 7 persisted even late in time, again suggesting a sustained representation of
643 low-level visual information.

644 Our results are generally consistent with prior work investigating how visual processing
645 of objects evolves over time, showing the gradual emergence of high-level representations
646 (Contini et al., 2017). While early signals reflect low-level visual features (e.g. Groen et al., 2013;
647 Cichy et al., 2014), later signals reflect perceptual similarity (Wardle et al., 2016), some tolerance
648 for changes in size and position (Isik et al., 2014), categorical processing (Carlson et al., 2013;
649 Cichy et al., 2014), and correlate with task performance and reaction times (Van Rullen and

650 Thorpe, 2001; Philiastides and Sajda, 2006; Martinovic et al., 2008; Ritchie et al., 2015). Further,
651 comparisons of deep neural networks with MEG have revealed a correspondence of early layers
652 with earlier MEG responses, likely reflecting initial stages of processing in early visual cortex,
653 while higher layers reflect later stages of processing in occipitotemporal cortex (Cichy et al.,
654 2016b; Seeliger et al., 2017). Our results significantly extend these results by establishing a lower
655 bound for the development of conceptual representations.

656 Other studies have also investigated high-level conceptual processing over time using
657 explicit semantic feature models (Clarke & Tyler, 2015) or behavioral judgments (Cichy et al.,
658 2017). For example, Clarke and colleagues showed semantic feature effects before 120 ms,
659 although including basic visual features based on the HMAX model revealed unique semantic
660 contributions to MEG signals only after ~200 ms (Clarke et al., 2013; Clarke et al., 2014). In
661 contrast to these studies, we used more recent deep convolutional neural networks which have
662 been shown to be more closely tied to neural and behavioral data (Khaligh-Razavi et al., 2016;
663 Jozwik et al., 2017; Cichy et al., 2016a). Further, we operationalized high-level conceptual
664 processing, using both a computational semantic model based on semantic co-occurrence
665 statistics (GloVe model), as well as behavioral judgments of object similarity that we take to more
666 broadly reflect conceptual processing. Indeed, our results suggest that MEG variance explained
667 by the GloVe model was comparably low and mostly covaried with behavioral judgments,
668 suggesting that conceptual representations extend beyond those relationships captured by the
669 GloVe model. Despite these differences, our results are generally consistent with the results of
670 Clarke and colleagues, but suggest a lower bound for conceptual processing around ~150 ms (see
671 also Cichy et al., 2017). Further, we show that the computational visual model and behavioral
672 judgments explain shared variance even prior to 150 ms. This shared variance indicates that the
673 neural activity captured by computational models is behaviorally relevant and argues against a
674 strong distinction between (low-level) visual features on the one hand, and high-level conceptual
675 processing on the other. At the same time, the presence of significant unique variance explained
676 by behavior after 150 ms suggests that not *all* aspects of conceptual object representations
677 reflected in MEG activity are explained by current generations of computational visual models.

678 While our study provides insight into the development of conceptual representations,
679 there are some important considerations. First, we used behavioral similarity judgments using
680 the multi-arrangement task (Kriegeskorte & Mur, 2012) to index conceptual processing.
681 However, this choice of method might constrain the ability to capture conceptual
682 representations. While the behavioral judgments explain more variance in the MEG signal than
683 the semantic GloVe model we tested, we do not know what aspects of conceptual processing are
684 reflected in those judgments. Further, it is unclear how sensitive those behavioral judgments are
685 to the context imposed by the stimuli and instructions. Second, we only employed two exemplars
686 per object concept to test generalization of representations and this may not have contained
687 sufficient variability to fully disentangle low-level and high-level processing. Future studies
688 should consider broader sets of stimuli, different behavioral tasks, and alternative computational
689 models that may better match the MEG signal.

690 In conclusion, by focusing on two criteria for conceptual object representations we
691 provide an estimate for a lower bound for the emergence of conceptual object representations
692 of around 150 ms. Prior to this time, our results demonstrate limited generalization across object
693 exemplars and time, and importantly little unique contributions of behavioral judgments to the

694 MEG response. The multifaceted nature of our findings here show that the combination of neural
695 data, behavior, and models are a viable method to probe the temporal dynamics of object
696 recognition and allow us to establish a novel profile of emergent conceptual representations in
697 time.

698 **References**

699

700 Brainard, D. H., 1997. The Psychophysics Toolbox. *Spat. Vis.* 10, 433-436.

701 Brysbaert, M., Wariner, A. B., & Kuperman, V., 2014. Concreteness ratings for 40 thousand
702 generally known English word lemmas. *Behav. Res. Methods* 46, 904-911.

703 doi:10.3758/s13428-013-0403-5.

704 Bullier, J., 2001. Integrated model of visual processing. *Brain Res. Rev.* 36, 96–107.

705 doi:10.1016/S0165-0173(01)00085-6.

706 Cadieu, C. F., Hong, H., Yamins, D. L. K., Pinto, N., Ardila, D., Solomon, E. A., Majaj, N. J., DiCarlo,
707 J. J., 2014. Deep neural networks rival the representation of primate IT cortex for core
708 visual object recognition. *PLoS Comp. Bio.* 10, e1003963.

709 doi:10.1371/journal.pcbi.1003963.

710 Carlson, T. A., Tovar, D., Alink, A., Kriegeskorte, N., 2013. Representational dynamics of object
711 vision: The first 1000 ms. *J. Vis.* 13, 1–19. doi:10.1167/13.10.1.doi.

712 Carlson, T. A., Hogendoorn, H., Kanai, R., Mesik, J., Turret, J., 2011. High temporal resolution
713 decoding of object position and category. *J. Vis.* 11, 1-17. doi:10.1167/11.10.9.

714 Chang, C. C., Lin, C. J., 2011. LIBSVM: a library for support vector machines. *ACM Trans. Intell.*
715 *Syst. Technol.*, 2-27. doi:10.1145/1961189.1961199.

716 Chatfield, K., Simonyan, K., Vedaldi, A., Zisserman, A., 2014. Return of the devil in the details:
717 Delving deep into convolutional nets. *Brain Mach. Vis. Conf.*

718 Cichy, R. M., Khosla, A., Pantazis, D., Oliva, A., 2016a. Comparison of deep neural networks to
719 spatio-temporal cortical dynaics of human visual object recognition reveals hierarchical
720 correspondence. *Sci. Rep.* 6, 27755. doi:10.1038/srep27755.

721 Cichy, R. M., Kriegeskorte, N., Jozwik, K. M., van den Bosch, J. J. F., Charest, I., 2017. Neural
722 dynamics of real-world object vision that guide behavior. *bioRxiv*. doi:

723 <http://dx.doi.org/10.1101/147298>.

724 Cichy, R. M., Pantazis, D., Oliva, A., 2014. Resolving human object recognition in space and time.
725 *Nat. Neurosci.* 17, 455-464. doi:10.1038/nn.3635.

726 Cichy, R.M., Pantazis, D., Oliva, A., 2016b. Similarity-based fusion of MEG and fMRI reveals
727 spatio-temporal dynamics in human cortex during visual object recognition. *Cereb. Cortex*

728 26, 3563-3579. doi:10.1093/cercor/bhw135.

729 Contini, E. W., Wardle, S. G., Carlson, T. A., 2017. Decoding the time-course of object
730 recognition in the human brain: From visual features to categorical decisions.

731 *Neuropsychologia* 105, 165-176. doi: 10.1016/j.neuropsychologia.2017.02.013.

732 Clarke, A., Devereux, B. J., Randall, B., Tyler, L. K., 2014. Predicting the time course of individual
733 objects with MEG. *Cereb. Cortex* 10, 3602-3612. doi:10.1093/cercor/bhu203.

734 Clarke, A., Taylor, K. I., Devereux, B., Randall, B., Tyler, L. K., 2013. From perception to
735 conception: how meaningful objects are processed over time. *Cereb. Cortex* 23, 187-197.

736 doi:10.1093/cercor/bhs002.

737 Clarke, A., Tyler, L. K., 2015. Understanding what we see: How we derive meaning from vision.

738 *Trends Cog. Sci.* 19, 677-687. doi:10.1016/j.tics.2015.08.008

739 Davies, M., 2008. The Corpus of Contemporary American English (COCA): 520 million words,
740 1990-present.

- 741 Eickenberg, M., Gramfort, A., Varoquaux, G., Thirion, B., 2017. Seeing it all: Convolutional
742 network layers map the function of the human visual system. *NeuroImage* 152, 184-194.
743 doi:<https://doi.org/10.1016/j.neuroimage.2016.10.001>.
- 744 Farah, M. J., McClelland, J. L., 1991. A computational model of semantic memory impairment:
745 Modality specificity and emergent category specificity. *J. Exp. Psychol. Gen.* 120, 339–357.
746 doi:10.1037/0096-3445.120.4.339.
- 747 Goldstone, R., 1994. An efficient method for obtaining similarity data. *Behav. Res. Methods*
748 *Instrum. Comput.* 26, 381-386. doi:10.3758/BF03204653.
- 749 Güçlü, U., van Gerven, M. A. J., 2015. Deep neural networks reveal a gradient in the complexity
750 of representations across the ventral stream. *J. Neurosci.* 35, 10005-10014.
751 doi:<https://doi.org/10.1523/JNEUROSCI.5023-14.2015>.
- 752 Greene, M. R., Baldassano, C., Esteva, A., Beck, D. M., Li, F. F., 2016. Visual scenes are
753 categorized by function. *J. Exp. Psych.* 145, 82-94. doi:10.1037/xge0000129.
- 754 Groen, I. I. A., Ghebreab, S., Lamme, V. A. F., Scholte, H. S., 2012. Spatially pooled contrast
755 responses predict neural and perceptual similarity of naturalistic image categories. *PLoS*
756 *Comput. Biol.* 8, e1002726. doi:10.1371/journal.pcbi.1002726.
- 757 Groen, I. I. A., Ghebreab, S., Prins, H., Lamme, V. A. F., Scholte, H. S., 2013. From image statistics
758 to scene gist: Evoked neural activity reveals transition from natural image structure to
759 scene category. *J. Neurosci.* 33, 18814-18824. doi:
760 <https://doi.org/10.1523/JNEUROSCI.3128-13.2013>.
- 761 Groen, I. I. A., Silson, E. H., Baker, C. I., 2017. Contributions of low- and high-level properties to
762 neural processing of visual scenes in the human brain. *Phil. Trans. B.* 372.
763 doi:<https://doi.org/10.1098/rstb.2016.0102>.
- 764 Grootswagers, T., Wardle, S. G., Carlson, T. A., 2016. Decoding dynamic brain patterns from
765 evoked responses: A tutorial on multivariate pattern analysis applied to time series
766 neuroimaging data. *J. Cogn. Neurosci.* 29, 677-697. doi:10.1162/jocn_a_01068.
- 767 Hebart, M. N., Bankson, B. B., Harel, A., Baker, C. I., Cichy, R., M., 2017. Representational
768 dynamics of task context and its influence on visual object processing. *bioRxiv*. doi:
769 10.1101/153684.
- 770 Hebart M. N., Görden K., Haynes J-D., 2015. The Decoding Toolbox (TDT): a versatile software
771 package for multivariate analyses of functional imaging data. *Front. Neuroinform.* 8:88.
- 772 Isik, L., Meyers, E. M., Leibo, J. Z., Poggio, T., 2014. The dynamics of invariant object recognition
773 in the human visual system. *J. Neurophysiol.* 111, 91-102. doi: 10.1152/jn.00394.2013.
- 774 Johnson, J. S., Olshausen, B. A., 2003. Timecourse of neural signatures of object recognition. *J.*
775 *Vis.* 3, 499–512. doi: 10:1167/3.7.4.
- 776 Jozwik, K. M., Kriegeskorte, N., Storrs, K. R., Mur, M., 2017. Deep convolutional neural networks
777 outperform feature-based but not categorical models in explaining object similarity
778 judgments. *Front. Psychol.* 8, 1726. doi:<https://doi.org/10.3389/fpsyg.2017.01726>
- 779 Khaligh-Razavi S-M., Henriksson, L., Kay, K., Kriegeskorte, N., 2016. Fixed versus mixed RSA:
780 Explaining visual representations by fixed and mixed sets from shallow and deep
781 computational models. *J. Math. Psychol.* 76, 184-197.
782 doi:<https://doi.org/10.1016/j.jmp.2016.10.007>.

- 783 Khaligh-Razavi, S-M., Kriegeskorte, N., 2014. Deep supervised, but not unsupervised,
784 models may explain IT cortical representation. *PLoS Comp. Bio.* 10, e1003915. doi:
785 <https://doi.org/10.1371/journal.pcbi.1003915>.
- 786 King, J. R., Dehaene, S., 2014. Characterizing the dynamics of mental representations: the
787 temporal generalization method. *Trends Cog. Sci.* 18, 203-210.
788 doi:10.1016/j.tics.2014.01.002.
- 789 Kriegeskorte, N., Mur, M., 2012. Inverse MDS: Inferring dissimilarity structure from multiple
790 item arrangements. *Front. Psychol.* 3, 245. doi: 10.3389/fpsyg.2012.00245.
- 791 Lescroart, M. D., Stansbury, D. E., Gallant, J. L., 2015. Fourier power, subjective distance, and
792 object categories all provide plausible models of BOLD responses in scene-selective visual
793 areas. *Front. Comput. Neurosci.* 9. doi: 10.3389/fncom.2015.00135
- 794 Martinovic, J., Gruber, T., Müller, M. M., 2008. Coding of visual object features and feature
795 conjunctions in the human brain. *PLoS One* 3, e3781. doi:10.1371/journal.pone.0003781.
- 796 McRae, K., Seidenberg, M. S., De Sa, V. R., 1997. On the nature and scope of featural
797 representations of word meaning. *J. Exp. Psychol. Gen.* 126, 99–130.
798 doi:10.1080/10430719008404646.
- 799 Meyers, E. M., Freedman, D. J., Kreiman, G., Miller, E. K., Poggio, T., 2008. Dynamic population
800 coding of category information in inferior temporal and prefrontal cortex. *J. Neurophysiol.*
801 100, 1407-1419. doi:10.1152/jn.90248.2008.
- 802 Nili, H., Wingfield, C., Walther, A., Su, L., Marslen-Wilson, W., Kriegeskorte, N., 2014. A Toolbox
803 for Representational Similarity Analysis. *PLoS Comput. Biol.* 10, e1003553.
804 doi:10.1371/journal.pcbi.1003553.
- 805 Pedhazur, E. J., 1997. Multiple regression in behavioral research: Explanation and prediction, 3rd
806 Edition. Orlando, FL: Harcourt Brace.
- 807 Pennington, J., Socher, R., Manning, C. D., 2014. GloVe: Global vectors for word representation.
808 *Proc. of 2014 Conf. Empir. Methods Nat. Lang. Process.*
- 809 Philastides, M. G., Sajda, P., 2006. Temporal characterization of the neural correlates of
810 perceptual decision making in the human brain. *Cereb. Cortex* 16, 509-518.
811 doi:10.1093/cercor/bhi130.
- 812 Potter, M. C., Wyble, B., Haggmann, C. E., McCourt, E. S., 2014. Detecting meaning in RSVP at 13
813 ms per picture. *Attent. Percept. Psychophys.* 76, 270-279. doi:[10.3758/s13414-013-0605-](https://doi.org/10.3758/s13414-013-0605-z)
814 [z](https://doi.org/10.3758/s13414-013-0605-z).
- 815 Ritchie, J. B., Kaplan, D., Klein, C., 2017. Decoding the brain: Neural representation and the
816 limits of multivariate pattern analysis in cognitive neuroscience. *BioRxiv*. doi:
817 10.1101/127233.
- 818 Ritchie, J. B., Tovar, D. A., Carlson, T. A., 2015. Emerging object representations in the visual
819 system predict reaction times for categorization. *PLoS Comp. Biol.* 11, e1004316.
820 doi:10.1371/journal.pcbi.1004316.
- 821 Rosch, E., Mervis, C. B., Gray, W. D., Johnson, D. M., Boyes-Braem, P., 1976. Basic objects in
822 natural categories. *Cogn. Psychol.* 8, 382–439. doi:10.1016/0010-0285(76)90013-X.
- 823 Schendan, H. E., Maher, S. M., 2009. Object knowledge during entry-level categorization is
824 activated and modified by implicit memory after 200 ms. *NeuroImage* 44, 1423–1438.
825 doi:10.1016/j.neuroimage.2008.09.061.

- 826 Seeliger, K., Fritsche, M., Güçlü, U., Schoenmakers, S., Schoffelen, J., Bosch, S. E., Gerven, M. A.
827 J., 2017. Convolutional neural network-based encoding and decoding of visual object
828 recognition in space and time. *NeuroImage* 155. doi:
829 <https://doi.org/10.1016/j.neuroimage.2017.07.018>.
- 830 Tadel F., Baillet S., Mosher J. C., Pantazis D., Leahy R. M., 2011. Brainstorm: a user-friendly
831 application for MEG/EEG analysis. *Comput. Intell. Neurosci*, 8. doi:10.1155/2011/879716.
- 832 Thorpe, S., Fize, D., Marlot, C., 1996. Speed of processing in the human visual system. *Nat.* 381,
833 520-522. doi:10.1038/381520a0.
- 834 Tyler, L. K., Moss, H. E., 2001. Towards a distributed account of conceptual knowledge. *Trends*
835 *Cogn. Sci.* 5, 244-262. doi:10.1016/S1364-6613(00)01651-X.
- 836 VanRullen, R., 2007. The power of the feed-forward sweep. *Adv. Cogn. Psychol.* 3, 167–176.
837 doi:10.2478/v10053-008-0022-3.
- 838 VanRullen, R., Thorpe, S. J., 2001. The time course of visual processing: From early perception
839 to decision-making. *J. Cogn. Neurosci.* 13, 454-461. doi:10.1162/08989290152001880.
- 840 Vedaldi, A., Lenc, K., 2015. MatConvNet – convolutional neural networks for MATLAB. *ACM Int.*
841 *Conf. on Multimedia*.
- 842 Wen, H., Shi, J., Zhang, Y., Lu, K.-H., Cao, J., Liu, Z., 2017. Neural encoding and decoding with
843 deep learning for dynamic natural vision. *Cerebral Cortex*.
844 <https://doi.org/10.1093/cercor/bhx268>
- 845 Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., DiCarlo, J. J., 2014.
846 Performance-optimized hierarchical models predict neural responses in higher visual
847 cortex. *Proc. Natl. Acad. Sci. USA* 111, 8619:8624.