1    **Disentangling object category representations driven by dynamic and static**

2    **visual input**

3
4    **Abbreviated Title**:  Representation of dynamic and static object information

5
6    Sophia Robert[1,2], Leslie G. Ungerleider[1], & Maryam Vaziri-Pashkam[1]

7    [1] Lab of Brain and Cognition, National Institute of Mental Health, Bethesda, MD, USA

8    [2] Department of Psychology and Neuroscience Institute, Carnegie Mellon University, Pittsburgh, PA,

9    USA

10    Corresponding authors: Sophia Robert srobert@andrew.cmu.edu and Maryam Vaziri-Pashkam

11    maryam.vaziri-pashkam@nih.gov

12

13    **Number of Pages**: 41

14    **Number of Figures**: 8

15    **Number of Words in Abstract**: 243

16    **Number of Words in Significant Statement**: 120

17    **Number of Words in Introduction**: 635

18    **Number of Words in Discussion**: 1494

19

20    **Conflict of interest statement**: The authors declare no competing financial interests.

21

26

## Abstract

Humans can label and categorize objects in a visual scene with high accuracy and speed—a capacity well-characterized with neuroimaging studies using static images. However, motion is another cue that could be used by the visual system to classify objects. To determine how motion-defined object category information is processed in the brain, we created a novel stimulus set to isolate motion-defined signals from other sources of information. We extracted movement information from videos of 6 object categories and applied the motion to random dot patterns. Using these stimuli, we investigated whether fMRI responses elicited by motion cues could be decoded at the object category level in functionally defined regions of occipitotemporal and parietal cortex. Participants performed a one-back repetition detection task as they viewed motion-defined stimuli or static images from the original videos. Linear classifiers could decode object category for both stimulus formats in all higher order regions of interest. More posterior occipitotemporal and ventral regions showed higher accuracy in the static condition and more anterior occipitotemporal and dorsal regions showed higher accuracy in the dynamic condition. Significantly above chance classification accuracies were also observed in all regions when training and testing the SVM classifier across stimulus formats. These results demonstrate that motion-defined cues can elicit widespread robust category responses on par with those elicited by luminance cues in regions of object-selective visual cortex. The informational content of these responses overlapped with, but also demonstrated interesting distinctions from, those elicited by static cues.

## Significance Statement

Much research on visual object recognition has focused on recognizing objects in static images. However, motion cues are a rich source of information that humans might also use to categorize objects. Here, we present the first study to compare neural representations of several animate and inanimate objects when category information is presented in two formats: static cues or isolated dynamic cues. Our study shows that while higher order brain regions differentially process object categories depending on format, they also contain robust, abstract category representations that generalize across format. These results expand our previous understanding of motion-derived animate and inanimate object category processing and provide useful tools for future research on object category processing driven by multiple sources of visual information.

## Introduction

Humans can categorize objects with striking speed and accuracy. Previous research on the neural basis of visual object recognition has largely focused on the processing of static features from images along the ventral visual hierarchy of the primate brain (reviewed in Peissig & Tarr, 2007). However, real-world scenes are not static. In fact, decades of behavioral research have shown that motion cues can contain category-relevant information that humans use to make judgements about objects. Behavioral studies using point-light displays (PLDs, Johansson, 1973; Johansson, 1976) have established that, even with the impoverished motion information available in PLDs, humans can quickly perceive a moving person, identify the action being performed, and even determine the actor's age, gender, and affect (e.g., Barclay et al., 1978; Bassili, 1978; Cutting and Kozlowski, 1977; Dittrich et al., 1996).

The majority of biological motion research has focused on the perception of human motion due to the significant role that it plays in our social lives. However, our sensitivity to information in motion cues is not restricted to perceiving humans. Humans can also infer animacy and complex social relations from the movements of basic geometric shapes (Schultz & Bülthoff, 2013; Heider & Simmel, 1944; Scholl & Gao, 2013) and can recognize animal categories such as chickens, dogs, horses and cats in PLDs (Mitkin & Pavlova, 1990; Mather & West, 1993; Pinto & Shiffrar, 2009; Pinto, 1994; Pavlova et al., 2001).

Investigations of the neural underpinnings of object categorization from motion information with neuroimaging have identified the superior temporal sulcus (STS) as a key region involved in processing biological motion. The STS has been shown to track animacy signals in motion cues from simple shapes and to process dynamic movements of human faces and bodies (Schultz & Bulthoff, 2013; Hirai & Hiraki, 2006; Pitcher et al. 2011, Pavlova et al., 2004).

4

80    Neuropsychological studies have also suggested the involvement of parietal regions in the

81    integration of motion and form information during form-from-motion identification tasks (Schenk

82    & Zihl, 1997).

83          Despite extensive research into neural substrates of human motion processing (Giese,

84    2013), there have been comparatively few studies that have investigated how non-human motion

85    is processed in the brain. Previous studies suggest preferential processing of human motion over

86    that of one or two other classes, e.g., mammals or tools, in regions in lateral occipito-temporal

87    cortex (LOTC) including the posterior STS (Papeo et al., 2017), human middle temporal complex

88    (Kaiser et al., 2012), and fusiform gyrus (Grossman & Blake, 2002), as well as the inferior parietal

89    lobe, inferior frontal gyrus (Saygin et al., 2004), the posterior and anterior cingulate cortices and

90    the amygdala (Bonda et al., 1996; Ptito et al., 2003).

91          The limited neuroimaging studies that have directly compared object representations

92    driven by motion to those driven by static images have focused on human (or monkey) faces and

93    bodies (Furl et al., 2012; Hafri et al., 2017; Pitcher et al., 2011) or have only compared humans

94    with tools (Beauchamp et al., 2003). Furthermore, these studies (with the exception of Beauchamp

95    et al., 2003), have used videos containing both static and dynamic cues as their dynamic condition

96    and thus have not been able to carefully separate the contributions of motion- and image-

97    information to the responses. Thus, a systematic comparison of several object category

98    representations driven by isolated motion and static cues has yet to be undertaken.

99          Here, we devised a novel method to generate stimuli that only contained motion cues. We

100   extracted motion signals from videos of objects and simulated object movements using flow fields

101   of moving dots. We first demonstrated that humans can recognize a wide variety of animate and

102   inanimate objects in our dynamic stimuli. We then used these stimuli, along with static images, in

103 an fMRI study to compare object category representations derived from dynamic and static cues

104 in occipito-temporal and parietal regions of interest across visual cortex.

## Materials and Methods

105

**Stimuli**

106

*Stimulus creation pipeline*

107

108 Eight categories were selected to sample a wide range of animate and inanimate object

109 categories: human, non-human mammal, bird, reptile, vehicle, tool, pendulum/swing, and ball. We

110 sought videos of objects performing a wide range of movements. Video clips were downloaded

111 from various sources on the Internet or shot with in-house equipment in accordance with the

112 following criteria: 1) contained a single moving object, 2) contained the entire object in frame

113 without occlusion, 3) shot without camera movement (no zooming, panning, tracking), 4)

114 contained no movement in the background, and 5) lasted at least 3 seconds.

115 We used in-house Matlab code, the Psychtoolbox extension, and in-house python code to

116 generate moving dot patterns that followed the movement of the objects in the videos. To do this,

117 first, all videos were trimmed to 3 seconds, cropped with a 3:2 x/y aspect ratio to center the object,

118 and resized to 720 x 360 pixel resolution. Videos with 30 frames per second were then up-sampled

119 so that all videos had a frame rate of 60 fps. The local, frame-by-frame motion of the objects in

120 each video in x and y directions was then extracted using the Farneback optical flow algorithm

121 (Farneback, 2003).

122 Next, object movements extracted from the full videos were projected on moving dot

123 patterns. To create the moving dot stimuli, 2500 white dots (2 pixel diameter) were randomly

124 initialized on a grey background (360 x 720 pixels). Dots that fell within pixels with nonzero

125 motion vector values were moved in the direction and magnitude specified by the extracted motion

126    matrix in the next frame. The lifetime (number of contiguous frames of movement) of any dot was

127    randomly sampled from a uniform distribution between 1 and 17 frames. The lifetime value

128    decreased on every frame. If the lifetime of a dot reached 0 or they reached the boundaries of the

129    frame, they were reinitialized with a lifetime of 17 frames.

130    The number of dots for a given frame and their lifetime was set to mitigate the formation

131    of dot clusters that could induce perception of an edge in individual frames of the video. The

132    frames were qualitatively examined to see if they induced a perception of any kind of edge or form.

133    Videos that produced such artifacts were removed from the stimulus set. For the fMRI experiment,

134    these moving dot videos were rendered live for each trial so that the dot initializations were always

135    random.

136

137    *Stimulus Validation Experiment*

138    To ensure that the stimuli contained clear category information, we conducted an online

139    experiment. 430 participants (223 women, aged 18-65) were recruited on Amazon Mechanical

140    Turk to perform an object categorization task on the dynamic stimuli. Participants each performed

141    between 10-11 trials. For each trial, participants were asked 3 questions about the object in a looped

142    video: 1) whether the object in the video was of an animal or non-animal, 2) which of 8 listed

143    categories the object belonged to, and 3) whether they could label the object. If subjects responded

144    'yes' for the third question, they were required to type the label in a response text box. Each of the

145    three questions contained an "I don't know" option. Subjects had to answer all three questions to

146    complete each trial.

147    Overall, subjects categorized objects based on their motion in the moving dot stimuli with

148    an average accuracy of 76% (202 total videos). The three animate (human, mammal, reptile) and

149    three inanimate (tool, ball, pendulum/swing) categories with the highest accuracy were used for

150    the fMRI experiment. For each category, the 6 videos with the highest accuracy were selected

151    (mean accuracy = 96%).

152          The overall 'motion energy' of each video was calculated by averaging the motion vectors

153    across all pixels in all frames. Non-zero motion vectors were also used to calculate the average

154    non-zero 'motion energy'. The average overall and non-zero motion energy for the 6 videos in

155    each category were entered into pairwise two-sample heteroscedastic t-test comparisons to ensure

156    that there were no significant differences between categories for either metric. Neither the overall

157    nor the non-zero motion energies were significantly different across categories (all $p$s > 0.05, even

158    without correction for multiple comparisons).

159          After the dynamic video stimulus set was finalized, the static image stimulus set was

160    generated by randomly selecting three frames of the full form video from which the moving dot

161    stimulus was created. The frame with the object in clearest view was selected and further processed

162    to extract the object from the frame. For the fMRI experiment, the isolated object was pasted onto

163    a background of 2500 randomly initialized white dots on a grey background, to mimic a frame of

164    the dynamic moving dot stimuli.

165

166    **Functional MRI experiment**

167    *Participants*

168          Fifteen healthy human subjects (six women, age range 19-42) with normal or corrected to

169    normal vision were recruited for the fMRI experiment. Participants were brought in for a 2 h fMRI

170    session that included the main experiment and three localizer tasks. Prior to entering the scanner,

171    all participants practiced the tasks for the main experiment and localizer runs and underwent a

172    short behavioral task to familiarize themselves with the stimuli. All subjects provided informed

173    consent and received compensation for their participation. The experiments were approved by the

174    NIH ethics committee.

175    *Training Session*

176         The independent norming study performed with mTurk demonstrated that people can

177    recognize the objects in these stimuli with high accuracy after minimal instruction. However, to

178    avoid introducing any random factors across subjects and differential processing during the first

179    run of the session relative to the rest, participants participated in a training session prior to entering

180    the scanner. During the training session, they familiarized themselves with the 36 dynamic stimuli

181    and were subsequently tested to ensure accurate recognition. Each video was shown on loop until

182    subjects could verbally report which of the 6 categories the object belonged to. If the subject

183    categorized the object correctly, the experimenter advanced to the next stimulus; incorrect

184    categorizations were verbally corrected by the experimenter. After all stimuli had been verbally

185    categorized, subjects underwent a testing session. In each trial, a random video was shown once

186    without looping, followed by a grey screen with 6 category labels placed in a circle around the

187    center of the screen. Subjects were instructed to categorize the object in the video by clicking on

188    the corresponding category label. No feedback was provided during the testing session. If a subject

189    performed above 90% accuracy, they continued on to the fMRI experiment. The training and

190    testing session took no longer than 15 minutes. Subjects required little to no correction during the

191    training session and performed with an average of 99% accuracy in the test session on the first

192    iteration (n = 13, data for two subjects were lost due to technical problems).

193

194

195    *MRI Methods*

196        MRI data were collected from a Siemens MAGNETOM Prisma scanner at 3 Tesla

197    equipped with a 32-channel head coil. Subjects viewed the display on a BOLDscreen 32 LCD

198    (Cambridge Research Systems, 60 Hz refresh rate, 1600 x 900 resolution, at an estimated distance

199    of 187 cm) through a mirror mounted on the head coil. The stimuli were presented using a Dell

200    laptop with MATLAB and Psychtoolbox extensions (Brainard, 1997; Kleiner, Brainard, & Pelli,

201    2007).

202        For each participant, a high resolution (1.0 x 1.0 x 1.0 mm) T1-weighted anatomical scan

203    was obtained for surface reconstruction. All functional scans were collected with a T2*-weighted

204    single-shot, multiple gradient-echo EPI sequence (Kundu et al., 2012) with a multiband

205    acceleration factor of 2 slices/pulse. 50 slices (3 mm thick, $3 \times 3$ mm$^2$ in-plane resolution) were

206    collected to cover the whole brain (TR 2 s, TE = 12 ms, 28.28 ms, 44.56 ms, flip angle = 70°, FoV

207    = 216 mm).

208    *Experimental Design*

209        **Main Experiment**: The main task of the experiment included 6 categories: human,

210    mammal, reptile, tool, pendulum/swing, and ball and 2 stimulus conditions: dynamic (moving dot

211    videos) and static (object images pasted on dot background). Both dynamic and static stimuli were

212    presented at the same size and location (subtending 9.6° x 4.8° visual angle). We used a block

213    design to present alternating blocks of dynamic and static stimuli while also alternating between

214    animate and inanimate blocks. The order of the six categories and the two formats were

215    counterbalanced within and across runs. Four different counterbalancing designs were created and

216    each subject was randomly assigned one of the designs.

10

217 Each run contained 12 condition blocks, one for each condition (2 formats x 6 categories),

218 began with an initial fixation block of 8 s, and ended with a final fixation of 12 s. Each condition

219 block began with an 8 s fixation period in which a red fixation dot (5 pixels in radius) was shown

220 on a grey background. The fixation period was then followed by the stimulus presentation period

221 in which 4 stimuli were presented from the same condition, each for 2.8 s followed by a 200 ms

222 inter-stimulus interval, resulting in 12 s of stimulus presentation. The duration of each condition

223 block was 20 s (8 s fixation and 12 s stimulus presentation). For each run, the 12 condition blocks

224 and the initial and final fixation blocks lasted 252 s (4 min 12 s). Each participant completed 12

225 runs.

226 To maintain their attention, subjects were given a one-back repetition detection task in

227 which they were instructed to press a button on an MRI-compatible button box (fORP, Cambridge

228 Research Systems) to indicate detection of a repeated stimulus within each block. There was one

229 stimulus repetition per block and the repeated stimulus of each block type was changed across

230 runs. Because there were only 3 unique trials per block but each condition had 6 unique stimuli,

231 half of the stimuli of each category were shown on odd runs and the other half were shown on the

232 even runs. These blocks were later combined during analysis. Average performance on this task

233 was 94%. To ensure proper fixation, eye movements were monitored using an ASL eye-tracker.

234 **Object Localizer task**: To localize functional ROIs in ventral and lateral occipito-temporal

235 cortex, we presented images of objects in 6 conditions: faces, scenes, head-cropped bodies, central

236 objects, peripheral objects (4 objects per image), and phase-scrambled objects in a block design

237 paradigm. Subjects were instructed to fixate while 20 images were presented in each block for

238 750ms with a 50ms fixation screen in between. Each block lasted 16 s and was repeated 4 times

239 per condition. Each run started with a 12s fixation period. Additional 8 s fixation periods were

240    presented after every 5 blocks. Total run duration was 436 s (7 min 16 s). Subjects performed a

241    motion detection task. During each block, a random image would jitter by rapidly shifting 4 pixels

242    back and forth horizontally from the center of the screen. Subjects indicated detection of motion

243    with a button press. Each participant completed 1-2 runs of this task.

244        **Motion localizer task**: To localize functional ROIs related to the perception of biological

245    and non-biological motion, we presented blocks of point light display (PLD) videos of humans

246    performing various actions in four conditions: 1) biological motion: normal PLD video (e.g.

247    walking, riding a bicycle), 2) random motion: the points in the PLD were spatially scrambled in

248    each frame, 3) translation: randomly positioned dots translated across the screen in a random

249    direction with the speed set to the average speed of the movement from the PLD videos, and 4)

250    static: a random frozen frame of the PLD was shown as an image. There were 8 exemplars per

251    condition, each presented for 1.5 s followed by a 500 ms interstimulus fixation period. Each block

252    lasted 16 s and was presented 4 times per condition. Each run began with a 6s fixation period and

253    8 s fixation periods were interspersed between each block making the total run duration 422.7 s (7

254    min 3 s). Subjects performed a one-back repetition detection task, in which they indicated detection

255    of a repeated stimulus during each block by pressing a button. Each subject completed 1-2 runs of

256    this task.

257        **Topographic mapping**: Topographic visual region V1 was mapped using 16 s blocks of

258    a vertical or horizontal polar angle wedge with an arc of 60° flashing black and white

259    checkerboards at 6 Hz. During the stimulus blocks, subjects fixated on a red fixation dot (5 pixel

260    radius) and detected a dimming on the wedge, that occurred randomly either at the inner, middle,

261    or outer ring of the wedge at 4 random times within the 16 s block. There was a 16 s fixation period

12

262     after each block and each run began with a 16 s period of fixation. Each run lasted 272 s (4 min

263     and 40 s), and subjects completed 1-2 runs of this task.

264     *Data Analysis*

265             fMRI data were analyzed using AFNI (Cox, 1996) and in-house MATLAB codes. The data

266     were pre-processed by removing the first 2 TRs of each run, motion correction, slice timing

267     correction, smoothing with 5mm FWHM, and intensity normalization. The EPI scans were

268     registered to the anatomical volume. The three echoes were combined using a weighted average

269     (Posse et al., 1999; Kundu et al., 2012). TRs with motion exceeding 0.3 mm as well as outliers

270     were excluded from further analysis. A general linear model analysis with 12 factors (2 stimulus

271     conditions x 6 categories) was used to extract t-values for each condition in each voxel. The 6

272     degrees of freedom movement parameters was used as an external regressor. To account for the

273     effect of residual autocorrelation on statistical estimates, we applied a generalized least squares

274     time series fit with restricted maximum likelihood (REML) estimation of the temporal auto-

275     correlation structure in each voxel. The t-values were calculated across all runs for the univariate

276     analysis and per-run for the multivariate analysis.

277     *ROI Definition: Group-constrained subject specific method*

278             We used a systematic, unbiased method for creating individualized regions of interest

279     constrained by group responses to our localizer experiments, basing our approach on a method of

280     region of interest definition developed by Kanwisher and Fedorenko (described in Pitcher et al.,

281     2011).

282             First, t-values were extracted from generalized linear models (GLMs) of individual

283     activation maps from the localizer experiments. All subjects' statistical activation maps (N = 15)

13

284    were converted to Talairach space. For each subject, the individual localizer contrast maps were

285    thresholded at $p < 0.0001$. Group overlap proportion maps were then created for each contrast.

286          Second, we thresholded the group proportion maps for each contrast separately to

287    counteract contrast- or localizer-specific differences in spatial variability or overall activation. The

288    thresholds for specific contrast maps were as follows: For the object localizer experiment, the

289    thresholds were $N \geq 0.7$ for objects vs scrambled (lateral occipital, LO; posterior fusiform sulcus,

290    pFS), $N \geq 0.5$ for bodies vs objects (extrastriate body area, EBA), and $N \geq 0.25$ for peripheral

291    objects vs scrambled (inferior intraparietal sulcus, infIPS). For the biological motion experiment,

292    the threshold for biological motion vs translation was $N \geq 0.5$ (lateral occipito-temporal biomotion

293    region, LOT-biomotion). For the retinotopy experiment, positive and negative maps were created

294    separately and thresholded at $N \geq 0.5$.
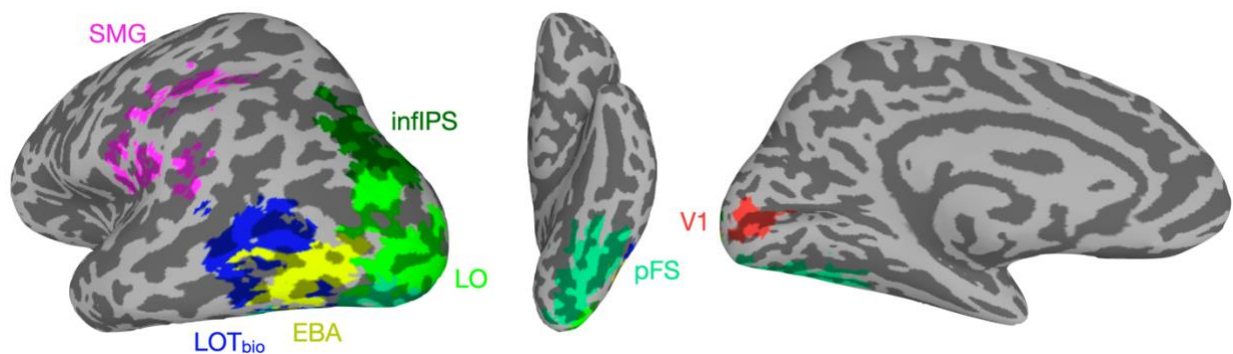
295          Third, we used a Gaussian blur of 1mm FWHM. The blurred maps were then clustered

296    using the nearest neighbors method and a minimum cluster size of 20 voxels. For V1, positive and

297    negative maps were clustered separately and then combined with a step function. Two steps were

298    required to finalize the group-constrained ROIs. Anatomical landmarks were used to separate pFS

299    from LO, and LO from infIPS. V1 was separated from V2 using a hand-drawn region based on the

300    group map. All ROIs were then selected to have no overlapping voxels.

301          The final nonoverlapping group-constrained ROIs were made subject specific by creating

302    masks based on the individual subject's activity during the localizer experiments (localizer contrast

303    threshold: $p < 0.05$). For example, for each subject's EBA, the group-constrained EBA was masked

304    by the subject's response to bodies > objects with a threshold of $p < 0.05$. If this process did not

305    yield an ROI with at least 100 voxels across the two hemispheres, the ROI was instead created

306     with a mask made from the mean response during the main experiment (task vs fix, p < 0.0001

307     uncorrected).

308         The supramarginal (SMG) region of interest was anatomically defined using a Freesurfer

309     parcellation (Desikan et al, 2006). To make the subject specific supramarginal ROIs, individual

310     masks were made from the mean response during the main experiment (task vs fixation, p < 0.0001

311     uncorrected) and intersected with the template SMG region.

312



313     **Figure 1**. Regions of interest of a single example subject generated by the group-constrained single-subject
314     method. The supramarginal area (SMG) is colored in pink, the inferior intraparietal sulcus (infIPS) is
315     colored in dark green, the lateral occipital complex (LO) is colored in light green, the extrastriate body area
316     (EBA) is colored in yellow, the biological motion related lateral occipito-temporal area (LOT-bio) is
317     colored in dark blue, the posterior fusiform sulcus (pFS) is colored in teal, and primary visual cortex (V1)
318     is colored in red.
319
320     *Univariate analysis*

321         To calculate the average fMRI response per condition for each ROI, using a general linear

322     model analysis, whole brain t-value maps were extracted for each of the 12 conditions and masked

323     with a task > fixation threshold of p < 0.0001 for each subject. The group-constrained subject-

324     specific ROIs were intersected with these maps, resulting in a t-value response per voxel in each

325     ROI for all 12 conditions in each subject. The average responses for four conditions were then

326     calculated from these ROI responses: dynamic animate, dynamic inanimate, static animate, and

327     static inanimate. The animacy preference in each ROI was calculated as the difference between

328     the animate and inanimate conditions, separately for the static and dynamic stimulus formats. One-

15

329    sample and paired t-tests were conducted to determine respectively: 1) if the animacy preference

330    in each ROI and each format was significantly different from 0, and 2) if the animacy preference

331    was significantly different across stimulus formats within each ROI. All t-tests were corrected for

332    multiple comparisons with False Discovery Rate correction (Benjamini and Hochberg, 1995)

333    across ROIs.

334    *Multivariate pattern analysis (MVPA)*

335    We performed multivariate pattern analyses to investigate whether object category

336    information was present in the fMRI responses to the dynamic and static stimuli. We extracted t-

337    values in each voxel for every condition in each run using a GLM analysis. To perform pairwise

338    object category decoding, we used a linear support vector machine classifier (SVM; Chang and

339    Lin, 2011) with feature selection. The SVM was trained using leave-one-out cross validation on

340    data that was normalized with z-scoring to avoid magnitude differences between conditions. Using

341    t-tests, we calculated the top 100 most informative voxels per ROI (Mitchell et al., 2004) to equate

342    the number of voxels analyzed per ROI and facilitate comparisons between them. This feature

343    selection was performed separately for each iteration of training. Results did not qualitatively

344    change when the analysis was performed without feature selection.

345    We trained and tested the linear SVM in two conditions: 1) within-classification, in which

346    the SVM was trained and tested on the same stimulus format, and 2) cross-classification, in which

347    SVM was trained in one stimulus format and tested on the other format. The classification was

348    performed on all unique pairs of object categories to obtain classification accuracy matrices. The

349    off-diagonal values of the matrices were averaged to produce two within-format and two cross-

350    format average object category decoding accuracies per subject. The two cross-format values were

351    then averaged to obtain one cross-classification accuracy. One-sample and paired t-tests were

16

352    conducted to determine respectively: 1) if the decoding accuracy in each ROI and each format was

353    significantly different from chance (0.5), and 2) if the decoding accuracy was significantly

354    different across stimulus formats within each ROI. All p-values listed from t-tests and ANOVAs

355    were corrected for multiple comparisons with False Discovery Rate correction across ROIs

356    (Benjamini and Hochberg, 1995). For ANOVAs, effect sizes were calculated with generalized eta

357    squared ($\eta_G^2$), for the one sample and paired t-tests, Cohen's *d* was used.

358    *Multidimensional scaling of fMRI responses*

359          To visualize how stimulus format and object category impact the responses in our regions

360    of interest, we quantified the similarities between the patterns of fMRI responses to the 12

361    conditions in each ROI by calculating all pairwise Euclidean distances. The individual subject

362    Euclidean distances per ROI were averaged across subjects to create group Euclidean distances,

363    which will be referred to as the fMRI-Euclidean matrix. We then visualized these similarities by

364    applying classical multidimensional scaling (Shepard, 1980) on the fMRI-Euclidean matrix and

365    plotting the first two dimensions for each ROI.

366          We measured the reliability of the fMRI-Euclidean matrix by performing a permutation

367    analysis wherein the individual subject matrices were split into two groups, averaged to create two

368    group matrices, and then correlated to get a measure of the split-half reliability. Correlations for

369    every possible combination of subjects in the two groups were measured and averaged to produce

370    a final reliability score. The reliabilities of the dynamic and static fMRI-Euclidean matrices were

371    evaluated separately.

372

373

374

17

**Object similarity behavioral experiment**

353 participants (32% female among the 85% who responded to the demographic survey) were recruited on Amazon Mechanical Turk to perform an object similarity task on the dynamic or static stimuli. All participants were located in the United States.

For each trial, participants were presented with three stimuli on a grey screen and were instructed to select the 'odd-one-out' stimulus (the stimulus that was most distinct among the three) by clicking on it. Dynamic and static stimuli were tested separately. Participants performed blocks of 15 trials to complete the task and were permitted to perform more than one block. To ensure data quality, trials with RTs smaller than 0.6 s and 1.2 s and larger than 10 s or 20 s were removed for the image and video tasks, respectively. These cutoffs were decided based on the distributions of RTs. If 5 or more trials in a block were eliminated, the entire block (or HIT in mTurk terminology) was removed. The eliminated blocks were resubmitted to mTurk to ensure that we had at least 2 repetitions for each unique triplet allowing for 68 trials for each pair of stimuli.

To build a dissimilarity matrix based on the odd-one-out image and video tasks, a response matrix of the pairwise dissimilarity judgments was constructed for each task by treating each triplet as three object pairs and assigning 1's to dissimilar pairs (i.e. the two pairs that included the selected odd object) and a 0 to the similar pair (i.e. the pair that did not include the selected odd object). We also constructed a count matrix to determine how many times each pair was shown together in a triplet. By dividing the response matrix by the count matrix, we obtained a dissimilarity matrix with values ranging from 0-1 with higher values denoting higher dissimilarity. To produce a category level behavioral dissimilarity matrix, we took the off-diagonal upper triangle of the 36 x 36 matrix and averaged the item distances that belonged to the same category, resulting in a 6 x 6 matrix, which will be referred to as the behavioral-dissimilarity matrix. The

18

398    diagonal was nonzero due to nonzero distances between exemplars within each category. Only the

399    off-diagonal of this matrix was used in further analyses.

400         To gauge the stability of the behavioral-dissimilarity matrix, we performed a split-half

401    reliability analysis. Because each subject only saw a small set of all possible triplets, instead of

402    splitting the data by subject, we split based on repeats of stimulus pairs (3 pairs per triplet) into

403    two groups. The binary similarity values for all pairs were correlated across the two groups to

404    produce a measure of reliability of the similarity judgments.

405

406    *Multi-dimensional scaling and hierarchical clustering of object similarity responses*

407         We visualized the structure of the object similarity judgments from the odd-one-out tasks

408    at the category level using classical multidimensional scaling on the behavioral-dissimilarity

409    matrices of the dynamic and static stimuli separately (Shepard, 1980). The two behavioral-

410    dissimilarity matrices were also correlated to quantify their degree of similarity. To investigate the

411    structure of the object similarity judgments at the exemplar level, we used a hierarchical or

412    agglomerative clustering algorithm available in the Python package *scipy* (Virtanen et al., 2020)

413    on the dynamic and static behavioral-dissimilarity matrices separately. For visualization purposes,

414    images of the individual exemplars, which were adapted from the static stimuli used in the

415    experiment, were included under the resultant dendrograms for both static and dynamic conditions

416    (note that dynamic stimuli are not recognizable in static frames).

417

418    *Brain-behavior correlation*

419         To determine the relationship between the multivariate information for the six categories

420    in each region of interest (fMRI-Euclidean matrix) with behavioral assessments of the category

421    similarity (behavioral-dissimilarity matrix), we correlated the two measures. For each subject, the

422    off-diagonal of the fMRI-Euclidean matrix was correlated with the off-diagonal behavioral-

423    dissimilarity matrix using Pearson's linear correlation coefficient, separately for the dynamic and

424    static experiments. The correlations were then averaged across subjects. The noise ceiling of these

425    correlations was then calculated for each ROI as the square root of the product of the reliabilities

426    of the fMRI-Euclidean matrix and the behavioral-dissimilarity matrix. As the reliability of the

427    behavioral-dissimilarity matrix was calculated with only one split, the standard error of the noise

428    ceiling was calculated based on the mean and standard deviation of the reliability scores generated

429    on each permutation of the fMRI-Euclidean reliability analysis.

430

431    *Brain-optic flow correlation*

432        To ensure that optic flow information from the six object categories was not predictive of

433    the multivariate fMRI responses in any of the regions of interest, we performed a control analysis.

434    We first calculated the Euclidean distances between the dynamic stimulus information of each

435    category by vectorizing the 4-dimensional stimuli (x-coordinates, y-coordinates, x- and y-

436    magnitudes of optic flow, and time) and averaging the distances between stimuli of the same

437    category, creating the optic flow-Euclidean matrix. We then correlated the optic flow-Euclidean

438    matrix with the dynamic fMRI-Euclidean matrix of each ROI for each subject. The correlations

439    were averaged across subjects to generate group mean correlations and one-sampled t-tests were

440    used to determine whether any positive correlations were significantly above zero.

441

442

443

20

444    Results

445    **Effect of stimulus format on univariate animacy preference**

446    We first looked at the mean amplitude of responses to the two superordinate object

447    categories (animate/inanimate) in the two stimulus formats (static/dynamic). We extracted

448    individual subjects' t-values from the GLM analysis and averaged the response for the three

449    animate and the three inanimate categories within each image format to get 4 values per subject.

450    Figure 2 shows the pooled results of this analysis across subjects. A two-way ANOVA with

451    stimulus format and animacy as factors showed a significant main effect of stimulus format in all

452    ROIs ($f$s > 7.26, $p$s ≤ 0.02, $\eta_G^2$s > 0.02) with higher response amplitude in the dynamic compared

453    to the static condition. A main effect of animacy was also found in LO, pFS, EBA, LOT-biomotion,

454    and left SMG ($f$s > 7.68, $p$s < 0.03, $\eta_G^2$s > 0.02), but not in V1, infIPS, or right SMG ($f$s < 3.38, $p$s

455    > 0.12, $\eta_G^2$s < 0.009). For the four ventrotemporal cortical areas, average responses were

456    significantly higher for the animate object categories, while in left SMG the average response was

457    higher for the inanimate object categories. The pattern of responses in SMG was not solely driven

458    by the tool category as removing tools from the inanimate objects did not qualitatively change the
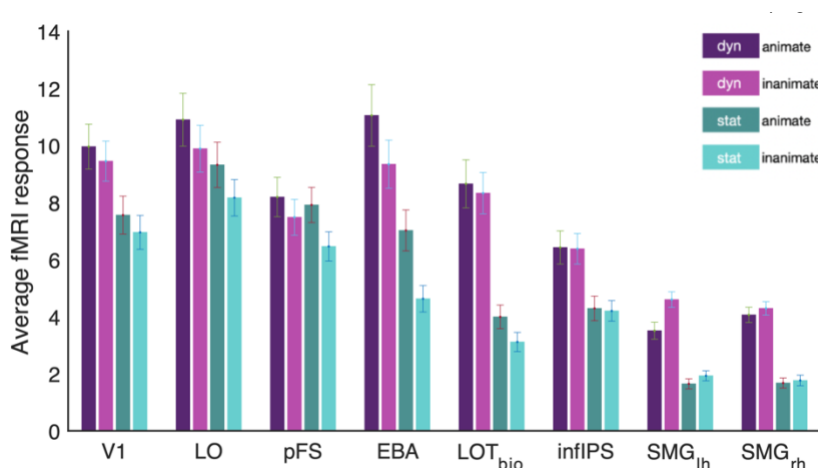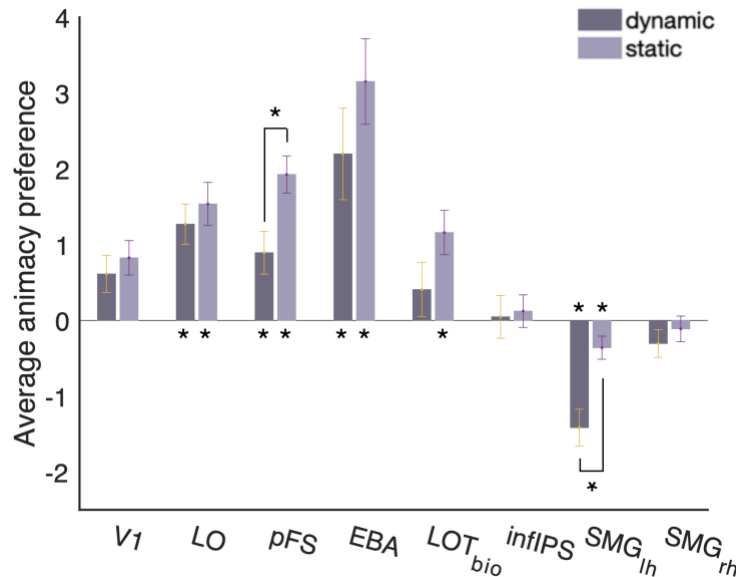
459    results (data not shown).



**Figure 2**. Univariate fMRI responses to dynamic and static stimuli averaged within animate and inanimate categories for each region of interest. Results do not qualitatively differ when removing the human and tool categories from the analysis. Error bars represent standard errors.

21

**Figure 3**. Univariate fMRI response preference for animate compared to inanimate object categories in dynamic and static stimuli for each region of interest. $*ps < 0.05$. Error bars represent standard errors.

To better visualize and investigate the interaction between stimulus format and animacy, we subtracted inanimate responses from animate responses to produce a measure of animacy preference within each stimulus format (Figure 3). Unpaired t-tests evaluating animacy preference against 0 revealed that there was no animacy preference in V1, inferior IPS, and the right SMG area in either stimulus format (dynamic: $ts < 1.56$, $ps > 0.21$, Cohen's $ds < 0.42$, static: $ts < 0.76$, $ps > 0.55$, Cohen's $ds < 0.20$). In contrast, for both stimulus formats, LO, pFS, and EBA showed a preference for animate categories (dynamic: $ts > 3.15$, $ps < 0.02$, Cohen's $ds > 0.84$, static: $ts > 5.05$, $ps < 0.0002$, Cohen's $ds > 1.35$) while left SMG preferred inanimate categories (dynamic: $t(14) = 5.59$, $p = 0.0005$, Cohen's $d = 1.49$). LOT-biomotion had significant preference for animate categories in the static ($t(14) = 3.97$, $p = 0.003$, Cohen's $d = 1.06$) but not in the dynamic condition ($t(14) = 1.14$, $p = 0.31$, Cohen's $d = 0.31$). All regions showed a preference in the same direction for dynamic and static conditions.

pFS and left SMG further showed a significant difference in the magnitude of their animacy preference across formats. pFS, a ventral region known to be involved in object recognition,

22

490    showed a stronger preference for animate object stimuli in the static compared to the dynamic

491    condition (paired t-test: $t(14) = 3.07$, $p = 0.03$, Cohen's $d = 0.79$), while left SMG, a parietal region

492    thought to be involved in tool processing and action observation had a stronger preference for

493    inanimate object stimuli in the dynamic compared to the static condition (paired t-test: $t(14) =$

494    3.73, $p = 0.02$, Cohen's $d = 0.96$). These significant interactions between stimulus format and

495    animacy preference suggest that the category preference responses in pFS and left SMG are

496    modulated by the format through which the category information is provided. The most ventral

497    region, pFS, is more sensitive to static form presentations of animate objects and the most dorsal

498    lateral region, left SMG, is more sensitive to dynamic motion information about inanimate objects.

499

500    **Effect of stimulus format on multivariate object category representations**
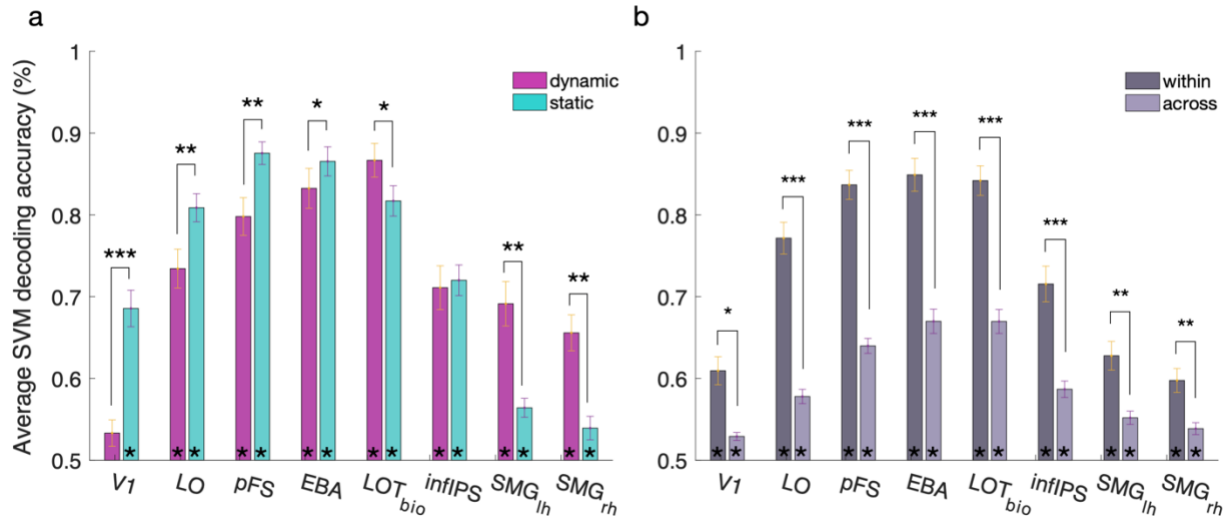
501         We next examined the multivariate patterns of each of our regions of interest to further

502    explore how object category information is represented in the brain when sourced from dynamic

503    movements and static images. We first sought to test if each of our regions contained information

504    about the 6 object categories within each stimulus format. To do this, we calculated average

505    pairwise classification accuracy for the 6 object categories for the static and dynamic conditions

506    using a linear SVM classifier (Chang and Lin, 2011). Figure 4a shows the pooled results of this

507    analysis across subjects. Unpaired t-tests revealed that the object categories were decoded

508    significantly above chance in both dynamic and static formats in all regions but V1 (dynamic: $t$s >

509    7.04, $p$s < 0.00001, Cohen's $d$s > 1.82; static: $t$s > 2.73, $p$s < 0.02, Cohen's $d$s > 0.71). In V1,

510    significant decoding was only found in the static stimulus condition (static: $t(14) = 8.31$, $p =$

511    0.00001, Cohen's $d = 2.15$; dynamic: $t(14) = 2.05$, $p = 0.06$, Cohen's $d = 0.53$). In all regions but

512    infIPS, there were significant differences between the decoding accuracies across stimulus format

23

513 (infIPS: $t(14) = 0.59$, $p = 0.57$, Cohen's $d = 0.15$). In V1, LO, pFS, and EBA decoding accuracies

514 were higher in the static condition than the dynamic ($t$s $> 2.32$, $p$s $< 0.001$, Cohen's $d$s $> 0.60$),

515 while in LOT-biomotion and bilateral SMG, decoding accuracies were higher in the dynamic

516 condition ($t$s $> 3.24$, $p$s $< 0.008$, Cohen's $d$s $> 0.84$).

517   To ensure that the significant decoding of object category from dynamic information was

518 due to differences in the responses to object categories and not contingent upon optic flow

519 information differences that were confounded with category in our stimulus set, we performed a

520 control analysis in which we correlated the dynamic stimulus information with the multivariate

521 fMRI responses (see Methods). No significant positive correlations were observed for any of the

522 regions of interest ($t$s $< 2.8$, $p$s $> 0.06$).

523   We next used a cross-classification method to determine if abstract responses to object

524 categories irrespective of stimulus format exist in our ROIs. The SVM classifier was trained in

525 one stimulus format and then tested in the other format. Decoding accuracies when training on

526 static and testing on dynamic and training on dynamic and testing on static were averaged to

527 produce the light grey bars shown in Figure 4b. We also calculated the within-classification

528 accuracy for training and testing within stimulus format (dark grey bars in Figure 4b; average of

529 the two bars in Figure 4a). Significant cross-classification was observed in all regions of interest

530 ($t$s $> 5.31$, $p$s $< 0.0001$, Cohen's $d$s $> 1.37$), and was significantly lower than within-classification

531 in all ROIs ($t$s $> 5.24$, $p$s $< 0.0001$, Cohen's $d$s $> 1.35$). This suggests that the information about

532 object categories in the multivariate pattern responses to the dynamic and static stimuli was

533 sufficiently similar to allow for significant decoding in one stimulus format after being trained on
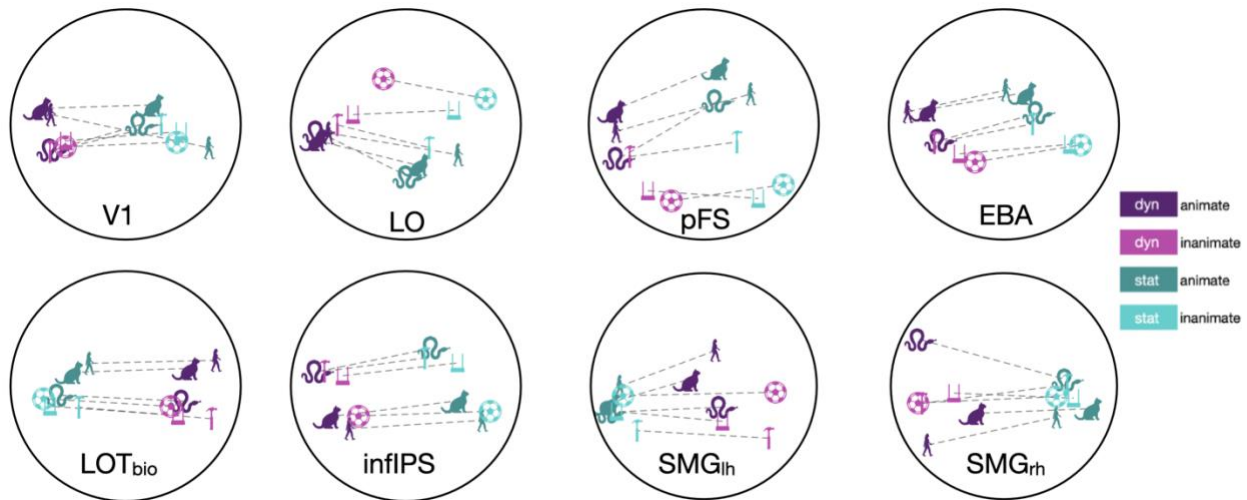
534 the other.

535

24

**Figure 4.** Object category SVM decoding accuracies in each ROI. a) Average SVM decoding accuracies when training and testing within the dynamic (pink) and static (teal) conditions. Asterisks within the bars represent significance in t-tests against chance. All average decoding accuracies were significantly above chance except for the dynamic condition in V1. Asterisks above bars represent paired t-tests across format. In all regions but infIPS, accuracies were significantly higher for one of the formats—LO, pFS, and EBA had significantly higher accuracy in the static condition while LOT-biomotion and bilateral SMG had significantly higher accuracy in the dynamic condition. b) The within stimulus format decoding accuracies, depicted in dark grey bars, were produced by averaging the dynamic and static decoding accuracies in A. The cross-format decoding accuracies are shown in light grey bars. Cross classification was significantly above chance in all regions of interest. Within classification was significantly higher than cross classification in all regions of interest. Error bars represent standard errors. Asterisk notation: * $p < 0.05$, ** $p < 0.001$, *** $p < 0.0001$.

To further visualize the similarity between the fMRI responses to the object categories in the dynamic and static conditions, we calculated the pairwise Euclidean distances between the patterns of responses to the 6 object categories and the two stimulus formats in each ROI. We then performed a multidimensional scaling analysis on the resultant dissimilarity matrix and visualized the first two dimensions in each of the regions of interest (Figure 5). In all regions, there was a clear separation between the responses to the dynamic (shown in purple and pink) and static stimuli (shown in green and teal). In addition, the ventro-temporal regions and inferior parietal cortex showed a separation amongst the individual object categories. The nearly parallel lines connecting the dynamic and static conditions of the same category indicate that categories with responses that were similar to each other in one condition were also similar to each other in

25

559    the other condition and is in line with the results of the cross-classification analysis performed

560    earlier. In bilateral supramarginal areas, this object category separation was evident for the

561    dynamic stimulus responses, but the static stimulus responses remained clustered together.  In

562    V1, while there was a separation between dynamic and static, the arrangement of categories does
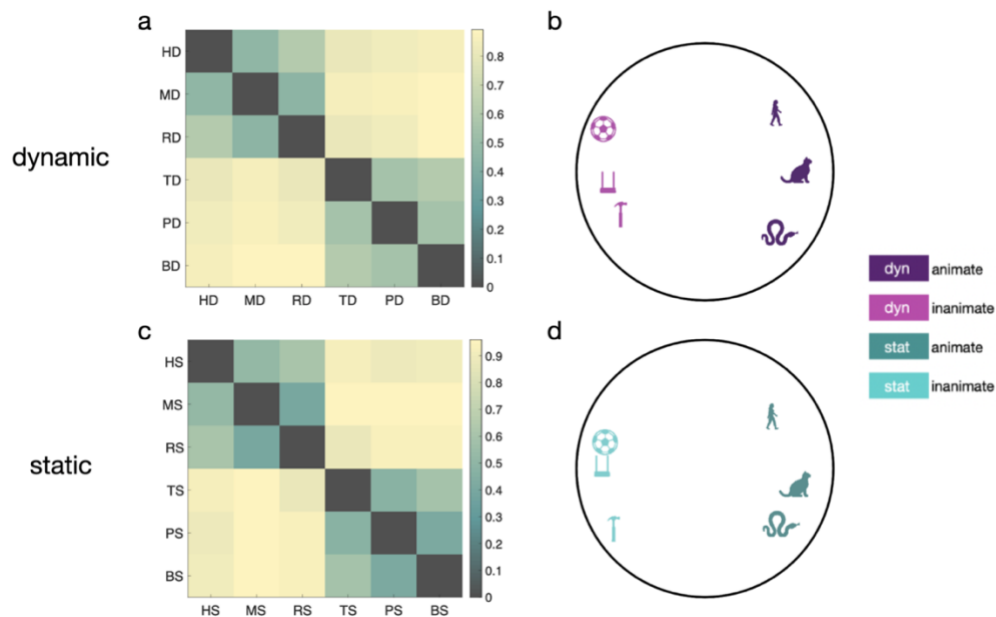
563    not appear to be consistent across conditions.



564 **Figure 5**. Multidimensional scaling visualization of fMRI response similarity between the object categories
565 presented in the dynamic and static formats. MDS was performed on the similarity matrix obtained from
566 the Euclidean distances of response patterns for the 12 conditions in each ROI. Dotted lines connect
567 dynamic and static presentations of the same object category. The dynamic condition is signified by purple
568 and the static condition is signified by green. Within each condition, the darker hues represent the animate
569 categories while the lighter hues represent the inanimate categories. The 6 object categories are symbolized
570 as with the following icons: human (person from side profile), mammal (cat), reptile (snake), tool (hammer
571

572 **Odd-one-out behavioral experiment**

573    To investigate how the responses of each ROI to the 6 object categories in each format

574    relates to the behavioral measure of similarity we performed two behavioral experiments on

575    Amazon Mechanical Turk in which we showed participants three objects (either in static condition

576    or in dynamic condition) and asked them to judge the similarity between the three objects and pick

577    the odd-one-out. We calculated two dissimilarity matrices based on the responses, one for the static

578    stimuli and one for the dynamic stimuli (see Methods). We then averaged the individual object

26

579    distances from each category to obtain dissimilarity scores between the 6 object categories for the

580    two stimulus formats (Figure 6a). The reliability of these similarity judgments was evaluated for

581    each stimulus format separately (see Methods). Participants rated both stimulus formats with

582    highly stable similarity judgments ($r = 0.98$ for both dynamic and static stimuli). We used

583    multidimensional scaling on the pairwise dissimilarities of each stimulus format to visualize the

584    distance between object categories in the first two dimensions (Figure 6b).

585         The dynamic and static similarity judgments had highly similar structure, showing a clear

586    separation between animate and inanimate categories in the first dimension. The animate (human,

587    mammal, and reptile) and inanimate (tool, pendulum/swing, and ball) categories were also

588    separated from each other along the second dimension in both tasks. Overall, the dissimilarities

589    from the dynamic and static tasks were highly correlated ($r = 0.98$, $p = 2.80\text{e-}10$), however, there

590    also appeared to be slight qualitative differences in the arrangement of the inanimate object

591    categories along the second dimension.



592    **Figure 6**. Odd-one-out similarity judgements of dynamic and static stimuli at the category level. The
593    matrices depict pairwise dissimilarity scores between object categories in dynamic (a) and static (c)
594    stimulus formats. The circle plots represent the object categories project into the first two dimensions from
595    multidimensional scaling on their dissimilarities in the dynamic (b) and static (d) stimuli.

27

596    To further explore the similarity structure of the dynamic and static stimuli at the exemplar

597    level, a hierarchical clustering algorithm was used on the odd-one-out similarity judgments (Figure

598    7). Similar to the MDS of odd-one-out judgements at the category level, a gross distinction

599    between animate and inanimate objects was observed for both the static and dynamic conditions.

600    Moreover, as in the MDS, the three object categories within the animate and inanimate

601    superordinate categories are largely distinguished in both formats. However, the clustering

602    algorithm also revealed several interesting differences in the similarity judgments of the same

603    objects when presented in either static image or dynamic optic flow format. For example, the

604    dynamic baboon stimulus, a clip of a baboon sitting and feeding, was grouped with the human

605    stimuli, while the static baboon stimulus was grouped with the mammal stimuli. Similarly, the

606    dynamic presentation of the two pendulum stimuli were grouped with the swings, presumably due

607    to their shared movement patterns, while their static presentations were grouped with the balls,

608    likely due to their shared global form. These deviations of specific exemplars from their category

609    clusters illustrate important differences in the category information provided by dynamic and static

610    visual cues and shed light on some of the heuristics that are used to guide similarity judgments in

611    the absence of either form or motion information. When luminance-defined edges are not

612    available, robust category information can be derived from dynamic motion-isolated inputs.
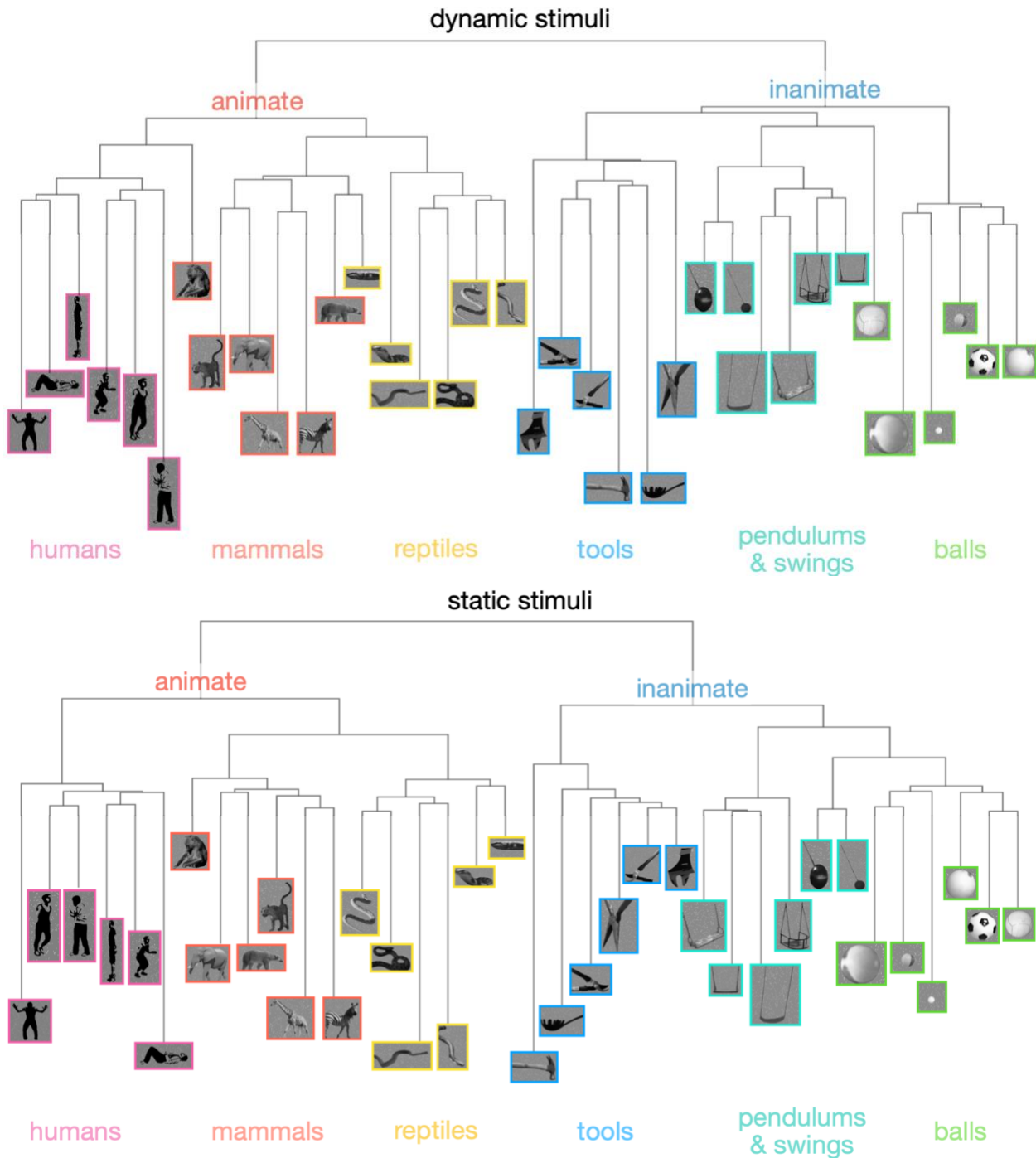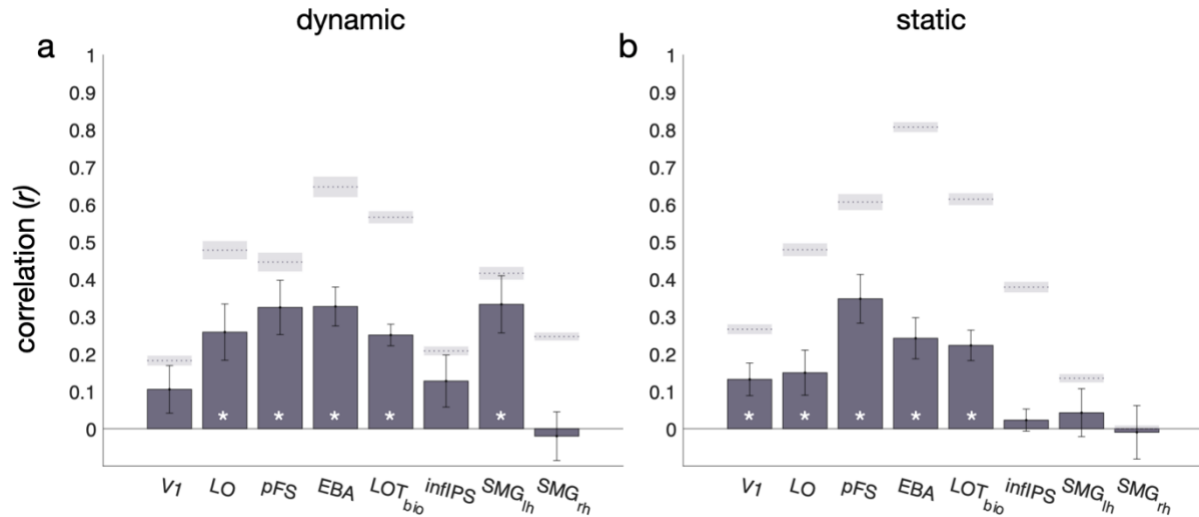
613

614

615

616

617

618

619

620



621

**Figure 7**. Hierarchical clustering of odd-one-out similarity judgments of the dynamic and static stimuli at the exemplar level. Edited versions of the static stimuli were used to visualize the similarity structure of both the dynamic (top) and static (bottom) stimuli as category of the dynamic stimuli cannot be gleaned from individual frames. The scale and position of the objects are not representative of the stimuli during presentation. Stimulus borders were colored to distinguish the six object categories. The human stimulus examples were modified into two-tone images for this figure to deidentify the individuals in the stimuli.

628

629    To investigate how the object category fMRI responses to each format relate to behavioral

630    judgements of similarity, we correlated the dissimilarity scores from the dynamic and static

631    behavioral experiments (dynamic and static reliability: 0.985) to those obtained from the Euclidean

632    distances between the multivariate response patterns in each region of interest ($r$s: dynamic > 0.03;

633    static > 0.02, apart from right SMG, see below). As shown in Figure 8, most ventral and lateral

634    temporal regions—LO, pFS, EBA, LOT-biomotion—showed significant correlations with the

635    object similarity judgments for both the dynamic and static stimuli (dynamic: $p$s < 0.01; static: $p$s

636    < 0.05). The responses in infIPS were not correlated to object similarity judgments for either the

637    dynamic or static stimuli (dynamic: $p = 0.12$, static: $p = 0.59$). The activity in left SMG was

638    significantly correlated with the similarity judgments for the dynamic stimuli ($p = 0.001$), but not

639    for the static stimuli ($p = 0.59$). Similarly, the activity in V1 was significantly correlated with

640    similarity judgments for the static stimuli ($p = 0.02$), but not for the dynamic stimuli ($p = 0.14$).

641    The only significant difference between the correlations of the behavioral similarity judgments and

642    the fMRI responses to the two conditions was found in the left SMG area, in which the correlation

643    was significantly higher with similarity judgments of the dynamic stimuli compared to the static

644    stimuli ($t(14) = 3.32$, $p = 0.04$, Cohen's $d = 0.86$). In the right SMG area, the $r$ value was -0.0083

645    for the static condition, signifying a reliability of zero. As this suggests that the responses to the

646    static stimuli in this region were unreliable, the correlation between the multivariate fMRI

647    responses in the right SMG to the static stimuli with behavioral assessments of their similarity will

648    not be interpreted.

649

650

651

**Figure 8**. Correlation of Euclidean distance between multivariate fMRI responses and behavioral dissimilarity matrices for a) dynamic and b) static stimuli. * $p$s $< 0.05$. Error bars represent standard errors. Shaded regions represent the average noise ceiling (dotted line) and the standard error of noise (shaded region) for each ROI.

## Discussion

Motion is an important visual cue that can provide category-relevant information in the absence of luminance-defined edges and form. Here, we introduce a novel approach to systematically separate form and motion signals and study the contribution of the motion signal to object category processing in isolation. To our knowledge, our study is the first to use this approach to compare the neural processing of form and motion signals from several animate and inanimate object categories. We sought to determine whether category-relevant information from the two sources is shared across the visual system by comparing dynamic and static category processing in regions of interest across visual occipito-temporal and parietal cortices. The two highly dissimilar information sources produced distinct but overlapping representations of animate and inanimate object categories, with a shift in processing primarily static information in more ventral regions to primarily dynamic information in more dorsal regions of cortex.

31

**Categorizing Objects with Motion Information**

670

671    An object identification task was used to determine whether our method for simulating the

672    extracted motion information in dynamic flow fields could produce stimuli in which objects were

673    recognizable. Our findings illustrate that, not only do people categorize motion-defined *animate*

674    objects with high accuracy (Pinto, 2006; Pinto, 1994; Pavlova et al., 2001), this high performance

675    also holds for three *inanimate* object categories: tools, swinging objects, and balls. These results

676    extend previous research by showing that a wide range of objects spanning animate and inanimate

677    categories can be recognized from just motion information. Our odd-one-out judgment task further

678    demonstrated that the similarity judgments for the dynamic and static stimuli were highly

679    correlated. This consistency suggests that people infer the similarity of objects from the two

680    sources of information in a similar way.

681    When discussing the perception of objects from motion, it is important to distinguish

682    between two types of information that can be gleaned from motion cues: 1) structure from motion,

683    a percept of a form arising from the global integration of coherent local motion vectors, and 2)

684    types of actions that are diagnostic of a particular object category such as walking, swinging, tool

685    use, bouncing, etc. Though it was not within the scope of this study to systematically distinguish

686    these two sources, the exemplar level clustering of our odd-one-out data qualitatively suggests that

687    both factors may play an important role in subjects' judgements of object similarity. For example,

688    images of pendulums and bouncing balls maybe judged to be similar since they both contain a

689    round shape, but distinct in dynamic form because they move differently.

690

691

692

**Format-dependent processing of object categories**

Comparison of the object category information across the two stimulus formats revealed differences in many of our regions of interest. Our findings suggest that stimulus format matters for: 1) processing of animate and inanimate objects—indicated by the regions of interest with significant interactions between stimulus format and univariate animacy preference (i.e., pFS and left SMG)—and 2) discriminating object categories within format—indicated by regions with significant differences in the multivariate classification accuracy of the responses to dynamic and static stimuli (i.e., all regions but infIPS). Broadly speaking, we found that the most ventral and posterior regions we examined (LO, EBA, and pFS) showed higher classification in the static condition, while most dorsal and anterior regions (LOT-biomotion and bilateral SMG) had stronger classification in the dynamic condition. Interestingly, infIPS used both sources of information without dominance of one source over the other. Importantly, all regions of interest but V1 showed robust responses to, and significant decoding accuracies of, all categories presented in both static image and dynamic motion formats. Thus, differential multivariate processing of object category based on stimulus format in these regions is a matter of degree. These results align with predictions from the model presented by Giese and Poggio (2003), in which form and motion signals are processed by distinct neural populations that largely overlap in topographic regions across ventral and dorsal cortex.

**Animate and Inanimate Category Processing**

Relative to static images, investigation of topographic organization of object category processing driven by motion information has been largely neglected. However, an important exception can be found in the work of Beauchamp and colleagues (2003), in which they compared

33

716 univariate fMRI responses between 1) full form videos and static images of humans and tools and

717 2) full form videos and point-light displays of humans and tools. Beauchamp et al. (2003) argued

718 for two processing pathways—form and motion. Lateral temporal regions (STS and MTG),

719 respond to their preferred category, humans and tools, respectively, in both PLDs and videos,

720 suggesting category preference from motion without requiring form. Meanwhile, ventral temporal

721 cortex (lateral and medial fusiform), needed form information for category preference responses.

722 Our results are in agreement with these findings and demonstrate that the topography of animacy

723 preference is not dependent on or exclusive to the human and tool categories—it also expands to

724 other animate objects such as mammals and reptiles, and other inanimate objects such as

725 pendulums/swings, and balls. These results suggest that large-scale animacy preference maps

726 (Konkle & Caramazza, 2013, Sha et al., 2015) found with static objects in the brain might also be

727 present for motion defined stimuli. Future studies with a larger stimulus set and sufficient power

728 to perform whole-brain analyses will be crucial for expanding our findings beyond functionally

729 defined regions of interest in VOTC and parietal cortex.

730

731 **Distinct but Overlapping Representations of Object Category for Dynamic and Static**

732 **Stimuli**

733 Using linear SVM classifiers, we decoded object category with high accuracy in all regions

734 tested. In all regions but V1 and the right supramarginal area, both information sources drove

735 object representations that were sufficiently distinguishable from each other to allow for high

736 classification performance. Extracting form and motion information from the same objects and

737 presenting them separately also allowed us to investigate the extent to which the representations

738 are overlapping across stimulus formats. We used a cross-classification approach to identify

739    regions that have format independent responses. A similar analysis has been used previously to

740    study fMRI responses to human actions in full form videos and images (Hafri et al., 2017). Our

741    results are largely in qualitative agreement with those of Hafri and colleagues, with the exception

742    that we found significantly more widespread cross-classification, possibly because our static

743    stimuli were source matched to our dynamic stimuli. Cross-decoding in all regions (apart from V1)

744    suggests that the object category representations driven by static and dynamic information were

745    sufficiently distinct to allow for significant within format classification, but also sufficiently

746    overlapping that their shared information could lead to significant cross-classification. These

747    results suggest the existence of abstract object category responses that pool information about

748    object category across various cues in the visual input.

749

750    **Relationship between brain and behavior**

751          Multivariate responses to both the dynamic and static conditions in LO, pFS, EBA, and

752    LOT-biomotion—the ventral and lateral regions—were correlated with the object similarity

753    judgments of the dynamic and static stimuli, respectively, with no differences across condition.

754    This implies that the fMRI responses in these regions follow the structure of the stimulus similarity

755    characterized by our odd-one-out experiment. The only region to show a difference in correlation

756    across the stimulus conditions was the left supramarginal area, which showed higher correlations

757    for the fMRI responses to the dynamic relative to the static stimuli. By contrast, the right

758    supramarginal area showed no significant correlation to behavioral judgments of either condition,

759    which indicates a lateralization of inanimate category processing to the left supramarginal area.

760    This left lateralization has been shown previously in research on tool processing (Beauchamp et

761    al., 2003). Importantly, not all regions that showed significant animacy preference or object

762  category decoding had responses that were significantly correlated with the similarity structure of

763  the behavioral judgments. In V1 and infIPS, the fMRI responses to both conditions were unrelated

764  to the similarity judgments of both stimulus types, suggesting that these regions were extracting

765  features irrelevant to similarity judgments on the objects.

766

767  **Conclusion**

768      In sum, our study demonstrates that in regions across occipito-temporal and parietal

769  cortices, category responses driven by isolated motion signals parallel category responses to static

770  form signals in a number of interesting ways. Regions that are traditionally considered part of the

771  visual object recognition pathway that processes static information, such as the pFS, LO, and EBA,

772  also extract robust object category information from isolated motion signals relevant to behavioral

773  judgments of object similarity. Furthermore, cross-classification of object categories in all regions

774  suggests that object-category information from static and dynamic signals overlap. Lastly,

775  preferential processing of certain kinds of objects, such as animate or inanimate objects, is

776  sensitive in some regions, i.e., the pFS and left SMG, to the format of visual information. Using

777  the stimulus generation approach we have introduced, future studies can expand beyond the six

778  object categories tested here and introduce parametric manipulations of dimensions that are likely

779  to play an important role in differential processing of motion-derived object categories. Candidate

780  dimensions include the type of action or movements that the objects are performing as well as the

781  orientation from which the movements are viewed. Such studies will be important for furthering

782  our understanding of how various visual cues to object-category are processed and integrated

783  together to form rich and robust object representations in the human brain.

## References

1. Barclay, C. D., Cutting, J. E., & Kozlowski, L. T. (1978). Temporal and spatial factors in gait perception that influence gender recognition. *Perception & psychophysics, 23*(2), 145-152.

2. Bassili, J. N. (1978). Facial motion in the perception of faces and of emotional expression. *Journal of experimental psychology: human perception and performance, 4*(3), 373.

3. Beauchamp, M. S., Lee, K. E., Haxby, J. V., & Martin, A. (2003). FMRI responses to video and point-light displays of moving humans and manipulable objects. *Journal of cognitive neuroscience, 15*(7), 991-1001.

4. Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological), 57*(1), 289-300.

5. Bonda, E., Petrides, M., Ostry, D., & Evans, A. (1996). Specific involvement of human parietal systems and the amygdala in the perception of biological motion. *Journal of Neuroscience, 16*(11), 3737-3744.

6. Brainard, D. H. (1997) The psychophysics toolbox. *Spatial Vision.* 10:433–436.

7. Chang, C. C., & Lin, C. J. (2011). LIBSVM: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST), 2*(3), 1-27.

8. Cox, R. W. (1996). AFNI: software for analysis and visualization of functional magnetic resonance neuroimages. *Computers and Biomedical Research, 29*(3):162-173. doi:10.1006/cbmr.1996.0014

9. Cutting, J. E., Kozlowski, L. (1977) "Recognition of friends by their walk." *Bulletin of the Psychonomic Society, 9*, 353–356.

10. Desikan, R. S., Ségonne, F., Fischl, B., Quinn, B. T., Dickerson, B. C., Blacker, D., Buckner, R.L., … & Killiany, R. J. (2006). An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *Neuroimage, 31*(3), 968-980.

11. Dittrich, W. H., Troscianko, T., Lea, S. E., & Morgan, D. (1996). Perception of emotion from dynamic point-light displays represented in dance. *Perception, 25*(6), 727-738.

810    12. Farnebäck, G. (2003, June). Two-frame motion estimation based on polynomial expansion. In

811         Scandinavian conference on Image analysis (pp. 363-370). Springer, Berlin, Heidelberg.

812    13. Furl, N., Hadj-Bouziane, F., Liu, N., Averbeck, B. B., & Ungerleider, L. G. (2012). Dynamic and

813         static facial expressions decoded from motion-sensitive areas in the macaque monkey. J*ournal of*

814         *Neuroscience, 32*(45), 15952-15962.

815    14. Giese, M. A., & Poggio, T. (2003). Neural mechanisms for the recognition of biological

816         movements. *Nature Reviews Neuroscience, 4*(3), 179-192.

817    15. Giese, M. A. (2013). Biological and body motion perception. The Oxford handbook of perceptual

818         organization, 575-596.

819    16. Grossman, E. D., & Blake, R. (2002). Brain areas active during visual perception of biological

820         motion. *Neuron, 35*(6), 1167-1175.

821    17. Hafri, A., Trueswell, J., & Epstein, R. (2017) Neural Representations of Observed Actions

822         Generalize across Static and Dynamic Visual Input. *Journal of Neuroscience 37*(11): 3056-3071.

823    18. Hirai, M., & Hiraki, K. (2006). The relative importance of spatial versus temporal structure in the

824         perception of biological motion: an event-related potential study. *Cognition, 99*(1), B15-B29.

825    19. Ishai, A., Ungerleider, L. G., Martin, A., Schouten, J. L., & Haxby, J. V. (1999). Distributed

826         representation of objects in the human ventral visual pathway. *Proceedings of the National*

827         *Academy of Sciences, 96*(16), 9379-9384.

828    20. Johansson, G. (1976). Spatio-temporal differentiation and integration in visual motion perception.

829         *Psychological research, 38*(4), 379-393.

830    21. Johansson, G. (1973) "Visual perception of biological motion and a model of its analysis"

831         *Perception & Psychophysics, 14*, 201–211.

832    22. Kaiser, M. D., Shiffrar, M., & Pelphrey, K. A. (2012). Socially tuned: Brain responses

833         differentiating human and animal motion. *Social neuroscience, 7*(3), 301-310.

834    23. Kleiner, M., Brainard, D., Pelli, D. (2007) "What's new in Psychtoolbox-3?" *Perception, 36,* ECVP

835         Abstract Supplement.

836    24. Konkle, T., & Caramazza, A. (2013). Tripartite organization of the ventral stream by animacy and
837        object size. *Journal of Neuroscience, 33*(25), 10235-10242.

838    25. Kundu, P., Inati, S.J., Evans, J.W., Luh, W.M. & Bandettini, P.A. (2012). Differentiating BOLD
839        and non-BOLD signals in fMRI time series using multi-echo EPI. *NeuroImage, 60*, 1759-1770.

840    26. Mitchell TM, Hutchinson R, Niculescu RS, Pereira F, Wang XR, Just M, Newman S (2004)
841        Learning to decode cognitive states from brain images. *Machine Learning 57*:145–175.

842    27. Mather, G., & West, S. (1993). Recognition of animal locomotion from dynamic point-light
843        displays. *Perception, 22*(7), 759-766.

844    28. Papeo, L., Wurm, M. F., Oosterhof, N. N., & Caramazza, A. (2017). The neural representation of
845        human versus nonhuman bipeds and quadrupeds. *Scientific reports, 7*(1), 1-8.

846    29. Pavlova, M., Krägeloh-Mann, I., Sokolov, A., & Birbaumer, N. (2001). Recognition of point-light
847        biological motion displays by young children. *Perception, 30*(8), 925-933.

848    30. Pavlova, M., Lutzenberger, W., Sokolov, A., & Birbaumer, N. (2004). Dissociable cortical
849        processing of recognizable and non-recognizable biological movement: analysing gamma MEG
850        activity. *Cerebral Cortex, 14*(2), 181-188.

851    31. Peissig, J. J., & Tarr, M. J. (2007). Visual object recognition: Do we know more now than we did
852        20 years ago? *Annual Review of Psychology*, *58*, 75-96.

853    32. Pinto, J. (1994). Human infants' sensitivity to biological motion in pointlight cats. *Infant Behavior*
854        *and Development, 17*, 871.

855    33. Pinto, J. (2006). "Developing body representations: A review of infants' responses to biological-
856        motion displays". In *Perception of the human body from the inside out*, Edited by: Knoblich, G.,
857        Grosjean, M., Thornton, I. and Shiffrar, M. 305–322.

858    34. Pinto, J., & Shiffrar, M. (2009). The visual perception of human and animal motion in point-light
859        displays. *Social Neuroscience, 4*(4), 332-346.

860   35. Pitcher, D., Dilks, D. D., Saxe, R. R., Triantafyllou, C., & Kanwisher, N. (2011). Differential

861        selectivity for dynamic versus static information in face-selective cortical

862        regions. *Neuroimage*, *56*(4), 2356-2363.

863   36. Posse, S., Wiese, S., Gembris, D., Mathiak, K., Kessler, C., Grosse-Ruyken, M. L., Elghahwagi,

864        B., … & Kiselev, V. G. (1999). Enhancement of BOLD-contrast sensitivity by single-shot multi-

865        echo functional MR imaging. *Magnetic Resonance in Medicine: An Official Journal of the*

866        *International Society for Magnetic Resonance in Medicine, 42*(1), 87-97.

867   37. Ptito, M., Faubert, J., Gjedde, A., & Kupers, R. (2003). Separate neural pathways for contour and

868        biological-motion cues in motion-defined animal shapes. *Neuroimage, 19*(2), 246-252.

869   38. Saygin, A. P., Wilson, S. M., Hagler, D. J., Bates, E., & Sereno, M. I. (2004). Point-light biological

870        motion perception activates human premotor cortex. *Journal of Neuroscience, 24*(27), 6181-6188.

871   39. Schenk, T., & Zihl, J. (1997). Visual motion perception after brain damage: II. Deficits in form-

872        from-motion perception. *Neuropsychologia, 35*(9), 1299-1310.

873   40. Scholl, B. J., & Gao, T. (2013). Perceiving animacy and intentionality: Visual processing or higher-

874        level judgment. *Social perception: Detection and interpretation of animacy, agency, and intention,*

875        *4629*, 197-229.

876   41. Schultz, J., & Bülthoff, H. H. (2013). Parametric animacy percept evoked by a single moving dot

877        mimicking natural stimuli. *Journal of vision, 13*(4), 15-15.

878   42. Sha, L., Haxby, J. V., Abdi, H., Guntupalli, J. S., Oosterhof, N. N., Halchenko, Y. O., & Connolly,

879        A. C. (2015). The animacy continuum in the human ventral vision pathway. *Journal of cognitive*

880        *neuroscience, 27*(4), 665-678.

881   43. Shepard, R. N. (1980) Multidimensional scaling, tree-fitting, and clustering. *Science 210*:390 –398.

882   44. Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., ... & Van

883        Mulbregt, P. (2020). SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature*

884        *methods, 17*(3), 261-272.