

Research Articles: Behavioral/Cognitive

Task-Dependent Warping of Semantic Representations During Search for Visual Action Categories

<https://doi.org/10.1523/JNEUROSCI.1372-21.2022>

Cite as: J. Neurosci 2022; 10.1523/JNEUROSCI.1372-21.2022

Received: 1 July 2021

Revised: 29 June 2022

Accepted: 6 July 2022

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.jneurosci.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

Task-Dependent Warping of Semantic Representations During Search for Visual Action Categories

Mo Shahdloo^{1,2,3}, Emin Çelik^{2,5}, Burcu A. Ürgen^{2,4,5}, Jack L. Gallant⁶, Tolga Çukur^{2,3,5,6*}

1. Wellcome Centre for Integrative Neuroimaging, Department of Experimental Psychology, University of Oxford, Oxford/UK
2. National Magnetic Resonance Research Centre (UMRAM), Bilkent University, Ankara/Turkey
3. Department of Electrical and Electronics Engineering, Bilkent University, Ankara/Turkey
4. Department of Psychology, Bilkent University, Ankara/Turkey
5. Neuroscience Program, Aysel Sabuncu Brain Research Centre, Bilkent University, Ankara/Turkey
6. Helen Wills Neuroscience Institute, University of California, Berkeley, CA/USA

* Correspondence to:

Tolga Çukur, cukur@ee.bilkent.edu.tr

Abbreviated Title: **Search for visual actions warps semantic representations**

Manuscript summary:

Pages: 34
Abstract: 169 words
Significance Statement: 102 words
Introduction: 650 words
Discussion: 1497 words
Figures: 15
Extended Data Figures: 11

1 **Abstract**

2 Object and action perception in cluttered dynamic natural scenes relies on efficient allocation of
3 limited brain resources to prioritize the attended targets over distractors. It has been suggested
4 that during visual search for objects, distributed semantic representation of hundreds of object
5 categories is warped to expand the representation of targets. Yet, little is known about whether
6 and where in the brain visual search for action categories modulates semantic representations. To
7 address this fundamental question, we studied brain activity recorded from five subjects (1
8 female) via functional magnetic resonance imaging while they viewed natural movies and
9 searched for either *communication* or *locomotion* actions. We find that attention directed to
10 action categories elicits tuning shifts that warp semantic representations broadly across
11 neocortex, and that these shifts interact with intrinsic selectivity of cortical voxels for target
12 actions. These results suggest that attention serves to facilitate task performance during social
13 interactions by dynamically shifting semantic selectivity towards target actions, and that tuning
14 shifts are a general feature of conceptual representations in the brain.

15
16 **Keywords:** action representation, attention, fMRI, voxelwise modelling, natural stimuli

19 **Significance Statement**

20
21 The ability to swiftly perceive the actions and intentions of others is a crucial skill for humans,
22 which relies on efficient allocation of limited brain resources to prioritise the attended targets
23 over distractors. However, little is known about the nature of high-level semantic representations
24 during natural visual search for action categories. Here we provide the first evidence showing
25 that attention significantly warps semantic representations by inducing tuning shifts in single
26 cortical voxels, broadly spread across occipitotemporal, parietal, prefrontal, and cingulate
27 cortices. This dynamic attentional mechanism can facilitate action perception by efficiently
28 allocating neural resources to accentuate the representation of task-relevant action categories.

29 Introduction

30 The ability to swiftly perceive the actions and intentions of others is a crucial skill for all social
31 animals. In the human brain this ability has been attributed to a network of occipitotemporal,
32 parietal and premotor areas collectively called the action observation network (AON) (Caspers et
33 al., 2010; Molinari et al., 2013; Oberman et al., 2007; Rozzi and Fogassi, 2017). Recent reports
34 suggest that the AON hierarchically represents diverse information pertaining to actions, ranging
35 from shape and kinematics to action-effector interactions and action categories (Grafton and de
36 C Hamilton, 2007; Handjaras et al., 2015; Oosterhof et al., 2010, 2012, 2013; Urgen et al., 2019;
37 Wurm et al., 2017; Cavina-Pratesi et al., 2018; Lingnau and Downing, 2015). Low-level shape
38 and movement kinematics are represented in occipitotemporal cortex and in the posterior bank of
39 inferior temporal cortex (Jastorff and Orban, 2009). Effector type (e.g., foot, hand) is represented
40 in ventral premotor cortex (Corbo and Orban, 2017; Jastorff et al., 2010), while parietal cortex
41 represents higher-level action categories (Abdollahi et al., 2012; Ferri et al., 2015).

42
43 Evidence suggests that selective attention alters population responses to actions across this
44 representational hierarchy. Prior electrophysiology (Muthukumaraswamy and Singh, 2008;
45 Muthukumaraswamy et al., 2004; Puglisi et al., 2017, 2018; Schuch et al., 2010) and
46 neuroimaging studies (Herrington et al., 2012; de Lange et al., 2008; Nicholson et al., 2017;
47 Rowe et al., 2002; Safford et al., 2010) have examined attention to low-level action features.
48 Schuch et al. (2010) reported increased electroencephalography (EEG) responses in AON with
49 attention to action kinematics. Safford et al. (2010) reported enhanced blood oxygen level
50 dependent (BOLD) responses in superior temporal sulcus (STS) with attention to animate actors
51 (i.e., humans) presented via simplified point-light displays (Johansson, 1973). Nicholson et al.
52 (2017) reported enhanced responses in inferior frontal gyrus (IFG), occipitotemporal cortex, and
53 middle frontal gyrus (MFG) with attention to action goals, in parietal cortex and fusiform gyrus
54 with attention to manipulated objects. Few reports have further investigated the effects of
55 attention to higher-level action features (Nastase et al., 2017, 2018). Presenting movie clips from
56 various animal taxonomies performing several actions, Nastase et al. (2017) reported that
57 attending to performed actions versus taxonomy alters multi-variate response patterns across
58 anterior intraparietal sulcus (IPS) and premotor cortex.

59
60 Current electrophysiology and neuroimaging findings on visual actions suggest that attention
61 increases AON responses to target features ranging from action kinematics and goals to actors.
62 That said, high-level semantic representations during visual search for specific action categories
63 remain understudied. Furthermore, prior studies did not question whether attending to action
64 features causes baseline and gain changes, or rather elicits dynamic tuning shifts that can alter
65 cortical representation. Recent evidence indicates that visual search for object categories shifts
66 single-voxel category tuning toward target objects (Çukur et al., 2013). Therefore, it is likely that
67 attention to action categories also causes tuning shifts to facilitate visual search. Here we
68 hypothesised that natural visual search for action categories induces semantic tuning shifts in
69 single cortical voxels toward targets. Tuning shift towards target categories elevates the local
70 sampling density near the target actions, and expands target-action representations while
71 compressing behaviourally-irrelevant-action representations by increasing the discriminability in
72 the semantic neighbourhood of the finely-sampled action categories (Fig. 1).

73

74 To test the tuning-shift hypothesis, we recorded whole-brain BOLD responses while human
75 subjects viewed 60min of natural movies and covertly searched for either 14 *communication*
76 actions or 30 *locomotion* actions among 109 action categories in the movies (Fig. 2, Fig. 2-1).
77 Using spatially-informed voxelwise modelling (Çelik et al., 2019), we measured category
78 responses for hundreds of objects and actions in the movies separately for each individual subject
79 and for each search task. We estimated a semantic space underlying action-category responses,
80 and semantic tuning for action categories were measured by projecting voxel-wise model
81 weights onto this space. Finally, semantic tuning profiles during the two search tasks were
82 compared to quantify the magnitude and direction of tuning shifts in single voxels.

85 **Materials and Methods**

86 **Subjects**

87 Five healthy adult volunteers with normal or corrected-to-normal vision participated in this
88 study: S1 (male, age 31), S2 (male, age 27), S3 (female, age 32), S4 (male, age 33), S5 (male,
89 age 27). Data were collected at the University of California, Berkeley. The experimental protocol
90 was approved by the Committee for the Protection of Human Subjects at the University of
91 California, Berkeley. All participants gave written informed consent before scanning.

93 **Stimuli and experimental design**

94 Data for the main experiment were collected in six 10min 50s runs in a single session.
95 Continuous natural movies were used as the stimulus in the main experiment. Three distinct
96 10min movie segments were compiled from short movie clips (10-20secs) without sound. Movie
97 clips were selected from a diverse set of natural movies (see Nishimoto et al. (2011) for details).
98 Movie clips were cropped into a square frame and downsampled to 512×512px. The movie
99 stimulus was displayed at 15Hz on an MRI-compatible projector screen that covered 24°x24°
100 visual angle. Subjects were instructed to covertly search for target categories in the movies while
101 maintaining fixation. A set of instructions regarding the experimental procedure and exemplars
102 of the search targets were provided to the subjects before the experiment. A colour square of
103 0.16°x0.16° at the centre with colour changing at 1Hz was used as the fixation spot. A cue word
104 was displayed before each run to indicate the attention target: *communication* or *locomotion*. The
105 *communication* target contained actions with the intent of communication, including both verbal
106 communication actions and nonverbal gestural communication actions (e.g., talking, shouting,
107 smirking). The *locomotion* target contained locomotion-related actions with the intent of moving
108 animate entities, including humans and anthropomorphized animals (e.g., moving, running,
109 driving). The same movie stimuli were used during each of the two attention tasks. The order of
110 attention conditions was interleaved across runs to minimize subject expectation bias. This
111 resulted in presentation of 1800sec of movies without repetition in each attention condition. Data
112 from the first 20secs and last 30secs of each run were discarded to minimize effects of transient
113 confounds. Following these procedures, 900 data samples for each attention condition were
114 obtained.

116 A separate set of functional data were collected while the same set of subjects passively viewed
117 120min of natural movies (i.e., passive-viewing data; this dataset was also used in Huth et al.,
118 2012 but here it was reanalysed with a focus on action categories). This dataset was used to
119 construct the semantic space and to select voxels subjected to further analyses. Data for the
120 passive-viewing experiment were collected in twelve 10min 50s runs in which 12 separate movie
121 segments were displayed. Presentation procedures were the same between the main experiment
122 and passive-viewing experiment, save for the number of runs. The passive-viewing dataset
123 contained 3600 data samples.

124 **fMRI data collection**

125 Data were collected on a 3T Siemens Tim Trio MRI scanner (Siemens Medical Solutions) via a
126 32-channel receiver coil. Functional data were collected using a T2*-weighted gradient-echo
127 echo-planar-imaging pulse sequence with the following parameters: TR=2sec, TE=33msec,
128 water-excitation pulse with flip angle=70°, voxel size=2.24mm×2.24mm×4.13mm, field of
129 view=224mm×224mm, 32 axial slices. To construct cortical surfaces, anatomical data were
130 collected using a three-dimensional T1-weighted magnetization-prepared rapid-acquisition
131 gradient-echo (MPRAGE) sequence with the following parameters: TR=2.3sec, TE=3.45 msec,
132 flip angle=10°, voxel size=1mm×1mm×1mm, field of view=256mm×212mm×256mm. Surface
133 flattening and visualisation were done via Freesurfer and PyCortex (Dale et al., 1999; Reuter
134 et al., 2012; Gao et al., 2015).

135 **fMRI data preprocessing**

136 Motion correction was performed using Statistical Parametric Mapping toolbox (SPM12; Friston
137 et al., 1995). Functional volumes were aligned to the first image from the first run in each
138 subject. Brain tissue was identified using the brain extraction tool (BET) from the FSL software
139 package (Smith, 2002). Low-frequency response components were detected using a third order
140 Savitzky-Golay low-pass filter with 240sec temporal window and were removed from voxel
141 responses. Voxel responses were then z-scored to attain zero mean and unit variance. Voxels
142 within the 2mm neighbourhood of the cortical sheet were identified as cortical voxels in each
143 subject (S1, 37791 voxels; S2, 32671 voxels; S3, 36942 voxels; S4, 42090 voxels; S5, 39254
144 voxels).

145 **Definition of regions of interest**

146 To define the anatomical regions of interest (ROIs) in each subject, the cortical surface was
147 segmented into 156 regions of the Destrieux atlas (Destrieux et al., 2010) via Freesurfer.
148 Segmentation results were projected from the anatomical space onto the functional space using
149 PyCortex, and each voxel was assigned an anatomical label based on the projections. Functional
150 ROIs were identified in each subject using visual category and retinotopic localizers (Huth et al.,
151 2012). Localizer experiments for visual category-selective areas (fusiform face area, FFA;
152 occipital face area, OFA; parahippocampal place area, PPA; retrosplenial cortex, RSC) were
153 performed in six 4.5 min runs of 16 blocks (Huth et al., 2012). Subjects passively viewed 20
154 random static images from one of the objects, scenes, body parts, faces, or spatially scrambled

155 objects groups in each block. Each image was shown for 300ms following a 500ms blank period.
156 PPA and RSC were identified as voxels with positive scene versus objects contrast (t -test, $p < 10^{-4}$,
157 uncorrected). FFA and OFA were defined using face-versus-object contrast (t -test, $p < 10^{-4}$,
158 uncorrected). The boundaries of these areas were hand drawn on the cortical surfaces along the
159 contours at which the contrast level reached half of the maximum. Localizer experiment for early
160 visual areas (RET: V1, V2, V3) contained four 9min runs. Subjects viewed clockwise and
161 counterclockwise rotating polar wedges in two runs. In the remaining two runs, subjects viewed
162 expanding and contracting rings. Visual angle and eccentricity maps were used to define visual
163 areas V1-3. Finally, ROIs were refined to voxels inside the drawn boundaries near a 2mm
164 neighbourhood of the cortical sheet.

165 **Abbreviations for regions of interest and important sulci**

166 Several regions of interest and important sulci were labelled on the flattened cortical surfaces to
167 guide the reader.

168 **Regions of interest:** pMTG, posterior middle temporal gyrus; pSTS, posterior superior temporal
169 sulcus; AG, angular gyrus; SMG, supramarginal gyrus; IPS, intraparietal sulcus; aIP, anterior
170 intraparietal cortex; PrCu, precuneus; dPMC, dorsal premotor cortex; BA44/45, Brodmann area
171 44/45; MFG, middle frontal gyrus; SFG, superior frontal gyrus; ACC, anterior cingulate cortex;
172 RET, early visual areas V1-3; FFA, fusiform face area; OFA, occipital face area; PPA,
173 parahippocampal place area; RSC, retrosplenial cortex.

174 **Sulci:** TOS, temporo-occipital sulcus; STS, superior temporal sulcus; SF, Sylvian fissure; IFS,
175 inferior frontal sulcus; MFS, middle frontal sulcus; SFS, superior frontal sulcus.

176 **Head motion, eye-movement, and physiological noise**

177 To prevent head motion and physiological noise confounds, estimates of these nuisance factors
178 were regressed out of the BOLD responses. Six affine motion time courses estimated during the
179 motion-correction stage were taken as the head-motion regressors. The cardiac and respiratory
180 activity during the main experiment were recorded using a pulse oximeter and a pneumatic belt.
181 These data were then used to estimate two regressors to capture respiration and nine regressors to
182 capture cardiac activity (Verstynen and Deshpande, 2011).

183 To ensure that eye-movements did not unduly bias the results, several control analyses were
184 performed. ViewPoint EyeTracker (Arrington Research) was used to monitor subjects' eye
185 positions at 60Hz, after getting calibrated at the beginning of each experimental run. Kruskal-
186 Wallis tests were used to detect systematic differences in the distribution of eye position and
187 movement. The distribution of eye position during search for *communication* and *locomotion*
188 tasks were examined. We find that the distribution of eye position is not affected by search task
189 ($p=0.17$), or by target presence or absence ($p=0.74$), and no significant interactions are present
190 between these two factors ($p=0.60$). To test whether eye movement is affected by target or
191 distractor detection, the distribution of eye position during a 1 sec window around target onset
192 and target offset was studied. The eye position distribution is not affected by target onset
193 ($p=0.73$) or offset ($p=0.17$), and there is no significant interaction between the aforementioned

194 factors ($p=0.83$). Furthermore, the moving-average standard deviation of eye position was
195 studied in a 200ms window to determine systematic differences in rapid moment-to-moment
196 variations in eye position across the two search tasks. There are no significant effects of search
197 task ($p=0.11$), target presence or absence ($p=0.32$), target onset ($p=0.49$), or target offset
198 ($p=0.36$), and there are no significant interactions between these factors ($p=0.16$). Finally,
199 moving-average standard deviation of eye position was included in the model as a nuisance
200 regressor and was regressed out of the BOLD responses.

201 To maintain subject vigilance, the subjects were instructed to depress a button whenever they
202 detected a member of the target category in the stimulus (i.e., either a communication or a
203 locomotion action depending on the search task). The behavioural responses were initially
204 analysed to ensure that subjects performed the tasks, and that task difficulty was balanced across
205 search targets. The target detection rate was $89\pm 9\%$ for the communication and $91\pm 8\%$ for the
206 locomotion targets (mean \pm std across subjects), with no significant difference between the two
207 tasks (bootstrap test, $p>0.05$).

208

209 **Category features**

210 A category feature space was constructed to encode the information pertaining to object and
211 action categories in the movies. Each second of the movie stimulus was manually labelled using
212 the WordNet lexicon (Miller, 1995) to find the time course for presence of 922 different object
213 and action categories in the movie stimulus. This yielded an indicator matrix where each row
214 represents a one-second clip of the movie stimulus, and each column represents a category.
215 Finally, category features were obtained by downsampling the indicator matrix to 0.5Hz to
216 match the acquisition rate of fMRI.

217

218 **Motion-energy features**

219 To infer cortical selectivity for low-level scene features, local spatial frequency and orientation
220 information of each frame of the movie stimulus were quantified using a motion-energy filter
221 bank. The filter bank contained 2139 Gabor filters that were computed at eight directions (0 to
222 350° , in 45° steps), three temporal frequencies (0, 2, and 4Hz), and six spatial frequencies (0, 1.5,
223 3, 6, 12, and 24 cycles/image). Filters were placed on a square grid spanning the $24^\circ \times 24^\circ$ field of
224 view. The luminance channel was extracted from the movie frames and passed through the filter
225 bank. The outputs were then passed through a compressive nonlinearity to yield the motion-
226 energy features (Lescroart and Gallant, 2019; Nishimoto et al., 2011). Finally, the motion-energy
227 features were temporally downsampled to match the fMRI acquisition rate.

228

229 **Space-time Interest Points (STIP) features**

230 Intermediate-level kinematic information of the movies were quantified by constructing the
231 Space-Time Interest Point (STIP) features using STIP toolbox (Laptev, 2005; Laptev et al.,
232 2008). STIP features have been successfully leveraged in many computer vision applications to
233 recognize human actions. As detailed in Laptev (2005) and Laptev et al. (2008), Harris operators
234 (Harris and Stephens, 1988) were used to identify spatiotemporal interest points in the movie
235 stimulus at multiple scales $(\sigma_i^2, \tau_j^2) = (2^{1+i}, 2^j)$, $i \in \{1, \dots, 6\}$, $j \in \{1, 2\}$, where σ and τ are the

236 standard deviations of the Gaussian kernels in spatial and temporal domains respectively.
237 Histograms of oriented gradients (HoG; Dalal and Triggs, 2005), and histograms of optical flow
238 (HoF; Holte et al., 2010) were calculated in the $(\Delta_{x,i}, \Delta_{y,i}, \Delta_{t,j})$ spatiotemporal neighbourhood of
239 each interest point, where $\Delta_{x,i} = \Delta_{y,i} = 2k\sigma_i$ and $\Delta_{t,j} = 2k\tau_j$, and k is the scale factor. Scale
240 factor was set to 9 according to the default configuration of the toolbox. Finally, normalized
241 histograms were concatenated to construct the collection of 162 STIP features and were
242 downsampled to match the acquisition rate of fMRI.

243

244 **Model estimation and testing**

245 Separate linearized models were fit in each voxel to estimate model weights that map each set of
246 features (i.e., category, motion-energy, or STIP features) to the measured BOLD responses in
247 each search task in individual subjects. Banded-ridge regression (Nunez-Elizalde et al., 2019)
248 was used to fit the models. To capture the hemodynamic response, delayed feature time-courses
249 were concatenated. Delays of two, three, and four samples, corresponding to 4, 6, and 8secs were
250 used. To account for potential correlations between target detection and BOLD responses, a
251 nuisance target-presence regressor was included in the model. The target-presence regressor
252 contained the category regressor for *communication* during search for *communication* task and
253 the category regressor for *locomotion* during search for *locomotion* task. Model fitting for the
254 two search tasks was performed concurrently by concatenating the features and BOLD responses
255 across search tasks (see Fig. 2). This procedure ensured consistency between the assigned
256 regularization parameters across search tasks and enabled utilisation of the target regressor
257 (Shahdloo et al., 2020).

258 A nested cross-validation (CV) procedure was used to choose the regularization parameters and
259 estimate model weights. Data from the main experiment were segmented into 60 30-sec blocks.
260 In each of the 10 outer folds, 4 randomly chosen blocks were held-out as validation data. Then,
261 in each of the 10 inner folds, 54 randomly chosen blocks were used as training data and the 2
262 remaining blocks were used as test data. To fit models for the passive-viewing data, data were
263 segmented into 144 50-sec blocks. In each fold, 8 randomly chosen blocks were held-out as
264 validation data, 132 randomly chosen blocks were used as training data and the 4 remaining
265 blocks were used as test data. For each feature set, regularisation parameters were selected with a
266 random-search; a thousand normalized regularisation parameter candidates were sampled from a
267 Dirichlet distribution and were scaled by 30 log- spaced values ranging from 10^{-5} to 10^{20} .
268 Training data were used to fit models for each set of regularisation parameters independently.
269 Model weights were then used to predict responses in the test data and prediction scores of the fit
270 models were assessed. Prediction scores were taken as product-moment correlation coefficient
271 between measured and predicted voxel responses. The set of regularisation parameters
272 maximizing the average prediction score across inner CV folds was chosen in each voxel.
273 Finally, the optimal set of parameters were used to fit models on the union of training and test
274 data in each outer fold and model weights were averaged across the outer folds.

275 Finally, prediction performance of the fit models were evaluated. In each outer fold, after
276 discarding the nuisance regressors, responses were predicted for the validation data using the fit
277 models and prediction scores were averaged across the search tasks. Prediction scores were then
278 averaged across the outer folds.

279

280 For each voxel, separate linearized models were estimated to relate each feature representation to
281 the BOLD responses. Specifically, category models were fit to estimate category responses that
282 represented the contribution of each category to single-voxel BOLD responses separately for the
283 data in the main experiment and the passive-viewing data in individual subjects. Furthermore, a
284 motion energy model and a STIP model were fit in each voxel to represent the contribution of
285 the low- and intermediate-level stimulus features to the responses. These alternative models were
286 further used to select analysis voxels (i.e., semantic voxels).

287

288 **Variance partitioning**

289 Object-action categories can be correlated with low-level visual features of natural movies
290 (Lescroart and Gallant, 2019), and there is evidence for representation of intermediate-level
291 action features (e.g., action kinematics) across cortex (Jastorff et al., 2010). Therefore, there is a
292 possibility that the estimated category responses are confounded by selectivity for low- and
293 intermediate-level scene features. To control for potential confounds, we performed a variance
294 partitioning analysis. This analysis estimates the response variance that is uniquely explained by
295 the category model after accounting for variance that can be attributed to low- and intermediate-
296 level features captured by the motion-energy and STIP models. To do this, we separately
297 measured the variance explained when all three models (category, motion-energy, and STIP) are
298 fit simultaneously (i.e., combined model), and variance explained when only motion-energy and
299 STIP models are fit simultaneously (i.e., control model). Banded ridge regression was used to fit
300 the combined and control models to prevent bias in assigning regularisation parameters across
301 different feature sets. The explained variance (R^2) was calculated as squared prediction scores,
302 separately for the combined and control models. Note that from a model fitting perspective,
303 negative prediction scores correspond to zero explained variance. Finally, unique variance
304 explained by the category model was calculated as

$$305 \quad \widehat{R}^2_{cat} = R^2_{comb} - R^2_{cont} \quad (1)$$

306 Here \widehat{R}^2_{cat} is the variance uniquely explained by the category model after accounting for low-
307 and intermediate-level features, R^2_{comb} is the variance explained by the combined model, and
308 R^2_{cont} is the variance explained by the union of motion-energy and STIP models in each voxel.

309

310 **Action category responses**

311 The fit category responses reflect voxel tuning for each of the 922 object and action categories in
312 the movie stimulus. To infer tuning for action categories, 922-dimensional category responses
313 were masked to select only the 109 action categories. This yielded the voxelwise 109-
314 dimensional action category responses.

315

316 Semantic representation of actions

317 Passive-viewing data were used to construct a continuous semantic space for action category
 318 representation. In this space, semantically similar action categories would project to nearby
 319 points, whereas semantically dissimilar categories would project to distant points (Huth et al.,
 320 2012). Category models were fit and action category responses during passive viewing were
 321 estimated. A group semantic space was then obtained using principal component analysis (PCA)
 322 on the action category responses of cortical voxels pooled across all subjects. To maximize the
 323 quality of the semantic space, voxels in which the category model predicted unique response
 324 variance after accounting for the variance attributed to low- and intermediate-level stimulus
 325 features were selected. These voxels were further refined to include only the top 3,000 best
 326 predicted voxels within each subject. The top 12 principal components (PCs) that explained more
 327 than 95% of the variance in responses were selected. Subsequent analyses were also repeated
 328 using the top 8 PCs that explained more than 90% of the response variance but the results
 329 remained consistent. Semantic tuning profile for each voxel under each search task was then
 330 obtained by projecting the respective action category responses onto the PCs. To illustrate the
 331 semantic content of the PCs, characteristic actions of the movie stimulus were clustered in the
 332 semantic space, and cluster centres were projected onto the PCs after getting labelled (see Fig.
 333 6).

335 Consistency of the semantic space across subjects

336 To test whether the estimated semantic space is consistent across subjects, we used a leave-one-
 337 out cross-validation procedure. In each cross-validation fold, voxels from four subjects were
 338 used to derive 12 PCs to construct a semantic space. In the left-out subject, semantic tuning
 339 profile for each voxel was obtained by projecting action category responses during passive
 340 viewing onto the derived PCs. Next, product-moment correlation coefficient was calculated
 341 between the tuning profiles in the derived space and the tuning profiles in the original semantic
 342 space. Results were averaged across semantic voxels in the left-out subject. The cross-validated
 343 semantic spaces consistently correlate with the original semantic space (Fig. 7).

345 Characterizing tuning shifts

346 Attentional tuning shifts toward or away from targets would be reflected in modulation of
 347 semantic selectivity for *communication* or *locomotion* action categories. Thus, the magnitude and
 348 direction of tuning shifts can be assessed by comparing the semantic selectivity for these
 349 categories between the two search tasks. Semantic selectivity for the two target categories was
 350 quantified as the similarity between semantic tuning profiles and idealized templates tuned solely
 351 for *communication* or *locomotion* action categories. First, idealized category responses were
 352 constructed as 109-dimensional vectors that contained ones for target categories (either
 353 *communication* or *locomotion* categories) and zeros elsewhere. Idealized templates were then
 354 obtained by projecting these idealized category responses onto the semantic space. Semantic
 355 selectivity for each target category was quantified as product-moment correlation coefficient
 356 between voxel semantic tuning profile and the corresponding template

$$357 \quad T_{i,C} = \text{corr}(s_i, s^t_C) \quad (1)$$

$$T_{i,L} = \text{corr}(s_i, s'_L) \quad (2)$$

where $T_{i,C}$ and $T_{i,L}$ are the tuning strength for *communication* and *locomotion* during condition $i \in \{C, L\}$ denoting attend to *communication* or attend to *locomotion*, s_i is the semantic tuning profile during condition i , and s'_C and s'_L denote the idealized semantic tuning templates for *communication* and *locomotion*, respectively. Finally, voxelwise tuning shift index (TSI_{all}) was quantified as

$$TSI_{\text{all}} = \frac{(T_{C,C} - T_{C,L}) + (T_{L,L} - T_{L,C})}{2 - \text{sign}(T_{C,C} - T_{C,L})T_{C,L} - \text{sign}(T_{L,L} - T_{L,C})T_{L,C}} \quad (3)$$

The numerator of TSI captures the difference in semantic selectivity for the attended versus unattended category, summed over the two attention tasks (i.e., search for communication and search for locomotion). Observing that the maximum possible selectivity for the attended category is 1, obtained when voxel tuning is equivalent to the idealized template, the denominator is cast to normalize the potential range of the TSI metric between 1 and -1 without affecting its sign. Tuning shifts toward the attended category would yield positive values where a TSI_{all} of 1 indicates a complete match between voxel semantic tuning and idealized templates, whereas negative values would indicate shifts away from the attended category where a TSI_{all} of -1 indicates a complete mismatch between voxel tuning and idealized templates. A TSI_{all} of 0 would indicate that the voxel tuning did not shift between the two search tasks.

The TSI metric in Eq. 3 can also be adopted to calculate tuning changes for any given set of action categories. To do this, the 922-dimensional category responses measured during attention tasks were masked to keep only the responses for the given set of actions. The masked tuning vectors and the idealised template for the given set were then projected onto the 12-dimensional semantic space. Semantic selectivity of a voxel to the given set was taken as the correlation coefficient between the projections of voxel tuning and the idealised template in the semantic space. Attentional modulation of semantic tuning for nontarget categories was examined by calculating a separate tuning shift index (TSI_{nt}). Note that this index can be calculated based on Eq. 3, but by zeroing out the category responses for communication and locomotion actions prior to projection onto the semantic space. To study the tuning shifts in an ROI, TSIs were averaged across semantic voxels within the ROI.

The change in voxelwise tuning during attending to the first target (e.g., communication) versus to the second target (e.g., locomotion) was defined as the l_1 -norm of the tuning difference between the two conditions. This calculated tuning change can be linearly decomposed into a component explained by the target features (i.e., the union of communication and locomotion features) and a component explained by the nontarget features (i.e., all features excluding the target features). The fraction of tuning change for target/nontarget features was computed by taking the ratio of the respective component to the overall tuning change.

399

400 **Characterizing target preference during visual search**

401 To investigate the interaction between tuning shifts and intrinsic selectivity for individual target
 402 action categories, we quantified a target preference index ($PI \in [-1,1]$) separately during search
 403 for communication actions (PI_{com}) and during search for locomotion actions (PI_{loc}). PI during
 404 search for each target action was taken as the difference in selectivity for the attended versus the
 405 unattended target

$$406 \quad PI_{com} = \frac{T_{C,C} - T_{C,L}}{1 - \text{sign}(T_{C,C} - T_{C,L})T_{C,L}} \quad (4)$$

$$407 \quad PI_{loc} = \frac{T_{L,L} - T_{L,C}}{1 - \text{sign}(T_{L,L} - T_{L,C})T_{L,C}} \quad (5)$$

409

410 where PI_{com} denotes the relative tuning preference for communication actions during search for
 411 communication, and PI_{loc} denotes the relative tuning preference for locomotion actions during
 412 search for locomotion. In this scheme, a PI of 1 indicates a complete match between voxel
 413 semantic tuning and the idealized template for the target, whereas a PI of -1 indicates a complete
 414 mismatch between voxel tuning and the idealized template for the target. Finally, a PI of 0
 415 indicates that the voxel semantic tuning does not shift toward any of the target actions.

416

417 **Characterizing action category preference during passive viewing**

418 To investigate the interaction between calculated preference index for individual targets and
 419 intrinsic selectivity for action categories, we quantified a selectivity index ($SI \in [-1,1]$), separately
 420 for communication actions (SI_{com}) and for locomotion actions (SI_{loc}). SI for each target in each
 421 voxel was calculated as the product-moment correlation coefficient between the voxel category
 422 response and idealised template category tuning for the given target.

423

424 **Action clustering**

425 To facilitate interpretation of stimulus information captured by individual PCs, characteristic
 426 action content of the movies was clustered and labelled. Action content (\bar{C}) for each short clip
 427 was calculated as the number of frames where each of the 109 actions were present ($\overline{N_a}$), and
 428 was normalized by the total number of clip frames (N)

$$429 \quad \bar{C} = \frac{\overline{N_a}}{N} \quad (2)$$

430 This yielded a 109-dimensional action content vector for each clip. The action content vectors
431 were then projected onto the semantic space and were grouped into 10 clusters using k-means.
432 The number of clusters was optimized using the elbow method (Thorndike, 1953). Average
433 action content of each clip (A) was calculated as the mean of the clip's action content vector

$$434 \quad A = \frac{\sum c}{109} \quad (3)$$

435 where $\bar{C} = [c_1, c_2, c_3, \dots, c_{109}]$. To label the clusters, five clips with highest average action
436 contents within each cluster were selected. Four candidate labels for each cluster were manually
437 assigned and 15 evaluators were asked to score (from 1 to 5) the correspondence of the selected
438 clips to each of the four candidate labels. Finally, the label with the highest score was selected to
439 represent each cluster.

440

441 **Statistical Analyses**

442 Bootstrap tests were used to assess statistical significance. To assess significance of the
443 prediction scores, single-voxels predicted responses were resampled 5000 times with
444 replacement. For each bootstrap sample, prediction score was computed. Significance level (p-
445 value) of the prediction scores was taken as the fraction of bootstrap samples in which the
446 prediction scores was greater than 0. Significance level of the unique response variance (Eq. 1)
447 was taken as the fraction of bootstrap samples in which the unique variance explained by the
448 category model was greater than 0. All single-voxel significance levels were corrected to account
449 for multiple comparisons using false-discovery rate correction (FDR; Benjamini and Hochberg
450 1995).

451

452 Significance of TSI_{all} , TSI_{nt} , PI_{com} , and PI_{loc} were assessed for each ROI across subjects. To do
453 this, ROI-wise metrics were resampled across subjects with replacement 10000 times.
454 Significance level was taken as the fraction of bootstrap samples where the test metric averaged
455 across resampled subjects is less than 0 (for right-sided tests) or greater than 0 (for left-sided
456 tests). This procedure was performed in a total of 21 functional ROIs separately. All ROI
457 significance levels were corrected to account for multiple comparisons using FDR.

458

459 In ROIs with a significant metric across subjects, the metric was further tested within individual
460 subjects. To do this, semantic voxels within a given ROI were resampled with replacement
461 10000 times. For each bootstrap sample, mean value of a given metric was computed across
462 resampled voxels. Significance level was taken as the fraction of bootstrap samples in which the
463 tested metric was less than 0 (for right-sided tests) or greater than 0 (for left-sided tests).

464

465 **Results**

466 **Visual search modulates category responses**

467 Little is known on whether and where in the brain natural visual search for action categories
468 warps semantic representations. To answer this question, we investigated voxel-wise tuning for

469 hundreds of object and action categories across cortex. Human subjects viewed natural movies
470 and covertly searched for *communication* or *locomotion* actions. Category regressors were
471 constructed to label presence of 922 distinct object and action categories in the movies. Separate
472 category models were then fit in each voxel for each search task. These models enabled us to
473 measure single-voxel category responses during each search task (Fig. 2a, see *Materials and*
474 *Methods*).

475
476 As natural stimuli contain correlations among various levels of features, there is a possibility that
477 estimated category responses are confounded by voxel tuning for low- and intermediate-level
478 scene features. To rule out this potential confound, we measured the response variance explained
479 by low-level motion-energy features, and intermediate-level spatiotemporal interest point (STIP)
480 features. Motion-energy features were constructed using a pyramid of spatiotemporal Gabor
481 filters (Nishimoto and Gallant, 2011). STIP features, providing an intermediate representational
482 basis for human actions, were constructed by measuring optical flow over interest points with
483 significant spatiotemporal variation (Laptev et al., 2008). We identified voxels in which the
484 category model explained unique response variance after accounting for these alternative
485 features via variance partitioning, and subsequent analyses were conducted on this set of
486 uniquely explained voxels. To prevent bias in voxel selection due to attention, variance
487 partitioning was performed on a separate dataset collected for this purpose (i.e., *passive-viewing*
488 *dataset*, see *Materials and Methods*). We find that the category model explains unique response
489 variance after accounting for low- and intermediate-level features in $25.7\pm 1.6\%$ of cortical
490 voxels (mean \pm sem across five subjects; bootstrap test, $q(FDR)<0.05$; Figs. 4 and 5), yielding
491 8,613-13,435 voxels in individual subjects (henceforth called *the semantic voxels*).

492
493 Comparison of estimated category responses across search tasks would be justified only if the fit
494 models can accurately predict BOLD responses that were held-out during model fitting. To
495 assess prediction performance of the fit category models, we measured average prediction scores
496 across the two search tasks, taken as product-moment correlation coefficient between the
497 predicted and measured held-out responses (Fig. 2b). Category models have high prediction
498 scores (greater than 1 std above the mean) in $46.9\pm 0.6\%$ of the semantic voxels. These include
499 many voxels spread across the AON comprising occipitotemporal, parietal, and premotor
500 cortices, as well as voxels in prefrontal and cingulate cortices (Fig. 3).

501
502 A recent study provided the first evidence that attention can alter single-voxel category tuning
503 profiles during search for object categories (Çukur et al., 2013). We thus hypothesised that visual
504 search for action categories can also cause changes in voxel-wise category tuning. If attentional
505 tuning changes are significant, the category models fit to individual search tasks should yield
506 higher prediction scores than a null model fit by pooling data across the two search tasks. To test
507 this prediction, we compared the prediction scores obtained from the category and null models.
508 We find that the category model significantly outperforms the null model in $46.1\pm 1.8\%$ of
509 semantic voxels (bootstrap test, $q(FDR)<0.05$). Additional control analyses further ensured that
510 these attentional changes cannot be attributed to residual eye-movements, head-motion,
511 physiological noise, or target-detection biases (see *Methods*). Taken together, these results
512 suggest that many cortical voxels in occipitotemporal, parietal, and prefrontal cortices encode
513 high-level category information, and that action-based visual search significantly modulates
514 category responses in single voxels.

515

516 **Visual search warps semantic representation of actions**

517 Previous studies suggest that the human brain represents visual categories by embedding them in
518 a continuous semantic space (Huth et al., 2012). Here, we used linear encoding models to map
519 category features of natural movies onto the recorded BOLD responses in single voxels. The
520 model features, namely actions, are fundamental semantic concepts in both language and vision.
521 The models successfully predict brain activity in cortical voxels, after controlling for lower-
522 levels of features (i.e., motion energy and STIP features). Thus, from a quantitative perspective,
523 it could be argued that there is an explicit representation of the semantic categories of actions in
524 the voxel responses (Naselaris et al., 2011). Note that a theoretical characterization of
525 relationships among semantic concepts is difficult. In computational semantics, an empirical
526 approach is adopted instead that is rooted in the distributional hypothesis. This hypothesis states
527 that concepts with similar statistical distributions have similar meaning. Accordingly, co-
528 occurrence statistics of concepts in corpora are used as a proxy metric for similarity of meaning
529 in many methods for learning semantic relationships (Jurafsky and Martin, 2021). Here, to derive
530 a semantic space underlying action category representations, we performed principal component
531 analysis (PCA) on the model weights for action categories. Visual search for actions alters
532 category model weights as reported here, so performing PCA on data from search tasks can bias
533 estimates of the semantic space. Instead, we derived the semantic space using the passive-
534 viewing dataset. Action categories that are semantically close to each other should project to
535 nearby points in this space, whereas semantically dissimilar categories should project to distant
536 points. The top twelve principal components (PCs) that explained more than 95% of the variance
537 in responses were selected, which showed a high degree of inter-subject consistency
538 ($r=0.52\pm 0.02$ mean \pm sem across subjects; Fig. 7). To visually examine the semantic information
539 captured by this space, we projected action categories onto the PCs (Fig. 6a; projections onto the
540 first three dimensions that accounted for 72.8% of the response variance is shown in Fig. 9;
541 loadings for all PCs are shown in Fig. 10). All further quantitative analyses regarding tuning
542 shifts were instead conducted in the full semantic space of 12 dimensions, including all the
543 identified PCs.

544

545 Previous evidence suggests that visual search shifts single-voxel tuning profiles to expand the
546 representation of the targets (Çukur et al., 2013). Thus, it is possible that action-based visual
547 search also shifts semantic tuning in single-voxels towards the target category. To investigate
548 this possibility, we projected action category responses onto the semantic space. The first and
549 third PCs maximally differentiated between actions belonging to the target categories (i.e.,
550 *communication* versus *locomotion* categories, Fig. 8). Therefore, we visually compared the
551 projections onto these PCs across the two search tasks. We observe that attention causes
552 semantic tuning modulations broadly across cortex (Fig. 6b; see Figs. 6-1 to 6-5 for results in
553 individual brain spaces). Specifically, voxels in inferior posterior parietal cortex (PPC), cingulate
554 cortex, and anterior inferior prefrontal cortex shift their tuning toward communication during
555 search for communication actions. Meanwhile, voxels in superior PPC and medial parietal cortex
556 shift their tuning toward locomotion during search for locomotion actions. Several reports
557 suggest involvement of superior PPC in representing locomotion actions (Corbo and Orban,
558 2017), and inferior PPC in representing communication actions (Abdollahi et al., 2012, Rizzolatti
559 and Matelli, 2003). Therefore, our findings suggest that during search for a given action

560 category, tuning shifts toward the target category are most prominent in voxels that are primarily
561 selective for the target.

562

563 **Visual search for action categories shifts single-voxel semantic tuning profiles**

564 Our inspection of semantic representations during visual search reveals that attention broadly
565 modulates high-level action representations by shifting semantic tuning profiles in single voxels.
566 To quantify the magnitude and direction of these tuning changes, we separately measured
567 semantic selectivity for *communication* and *locomotion* action categories in each search task.
568 The 922-dimensional category responses for individual voxels measured during attention tasks,
569 and idealised template vectors for the targets were projected onto the semantic space. The
570 template vector for a target is constructed as a 109-dimensional indicator vector containing ones
571 for the target category and all its subordinate categories, and zeros for remaining categories. For
572 instance, the “locomotion” template has ones for *locomotion*, and for *walk*, *run*, *crawl*, *move*,
573 *ride*, *etc.* As such, the target template vector indexes the target action as well as actions that are
574 semantically related to the target according to the WordNet hierarchy (see *Methods*). For each
575 attention task, semantic selectivity of a given voxel for a target category was then quantified as
576 the correlation coefficient between projected 12-dimensional vectors characterizing the voxel-
577 wise tuning profile and the idealised template in the semantic space. For each voxel, a tuning
578 shift index ($TSI_{all} \in [-1,1]$) was taken as the difference in semantic selectivity for targets when
579 they were attended versus unattended. A positive TSI_{all} indicates shifts towards the target, a
580 negative TSI_{all} indicates shifts away from the target, and a TSI_{all} of 0 suggests no change in
581 between tasks (see *Materials and Methods*).

582

583 We find that voxels across many cortical regions shift their tuning toward the attended category
584 (Fig. 11a; see Figs. 11-1a to 11-5a for results in individual brain spaces). Figure 15a shows
585 respective tuning shifts in relevant regions of interest (ROIs). Tuning shifts are significantly
586 greater than zero in many areas across AON including occipitotemporal cortex (posterior STS,
587 pSTS; posterior MTG, pMTG), posterior parietal cortex (intraparietal sulcus, IPS; AG, SMG),
588 and premotor cortex (Brodmann’s areas 44, 45, BA44/45; bootstrap test $q(FDR) < 0.05$; Fig. 15a).
589 This result suggests that focused attention to specific action categories shifts semantic tuning
590 toward targets in single-voxels, and that these attentional modulations are present at all levels of
591 the AON hierarchy including occipitotemporal cortex.

592

593 Prior evidence suggests that during category-based visual search, semantic tuning shifts grow
594 stronger toward later stages of semantic processing (Çukur et al., 2013). Here, we find that
595 semantic tuning shifts in AG and SMG are significantly stronger than those in occipitotemporal
596 (pSTS, pMTG) and premotor cortices (i.e., averaged over AG and SMG, compared with the
597 average over pSTS and pMTG, and with the average over dorsal premotor cortex (dPMC) and
598 BA44/45; Cohen’s $d=1.36$, $p < 0.05$). Therefore, the tuning shifts reported here could indicate that
599 AG and SMG are higher nodes in the hierarchy of semantic representation of action categories.
600 In a previous study, we reported that in medial prefrontal cortex visual search for object
601 categories causes tuning shifts toward targets while it causes tuning shifts away from targets in
602 voxels in precuneous (PrCu) and temporo-parietal junction (TPJ; Çukur et al., 2013). Similarly,
603 by qualitative inspection of the flatmaps, here we observe that visual search for action categories
604 causes negative tuning shifts in many voxels across PrCu and TPJ. These results suggest that

605 these areas might be involved in distractor detection and in error monitoring during visual search
606 for actions (Corbetta and Shulman, 2002).

607 608 **Visual search shifts semantic tuning for nontarget action categories**

609 Natural visual search for object categories was previously suggested to cause changes in
610 representations of not only targets but also nontarget categories (Çukur et al., 2013; Seidl et al.,
611 2012). Thus, it is likely that action-based visual search shifts semantic tuning for nontarget
612 categories. To address this important question, we first examined the separate contributions of
613 tuning changes for target versus nontarget categories to the overall tuning shifts. Specifically, we
614 measured the fraction of overall tuning shifts that can be attributed to the target categories versus
615 nontarget categories (i.e., all categories excluding *communication* and *locomotion* actions). We
616 find that both target and nontarget categories significantly contribute to the overall tuning shifts
617 (bootstrap test, $q(FDR) < 0.05$).

618
619 However, as would be expected, target categories account for a relatively larger fraction of the
620 overall tuning shifts compared to nontarget categories in all studied ROIs, except in early visual
621 cortex ($q(FDR) < 0.05$; Fig. 13). Next, to explicitly quantify tuning shifts for nontarget categories,
622 we calculated a separate tuning shift index exclusively on nontarget categories (TSI_{nt}). To
623 calculate TSI_{nt} , the 109-dimensional action category response vectors were masked to select
624 nontarget categories, prior to projection onto the semantic space (see *Materials and Methods*).
625 We observe that tuning shift for nontarget categories is generally smaller than the overall tuning
626 shift (Fig. 11b versus Fig. 11a and Fig. 12; see Figs. 11-1b to 11-5b for results in individual brain
627 spaces). Yet, TSI_{nt} is non-significant in all ROIs except AG, SMG, and BA45 ($q(FDR) < 0.05$;
628 Fig. 15b). Note that an insignificant TSI_{nt} does not necessarily suggest that attention has not
629 altered tuning for non-target categories, but rather the direction of tuning changes could be
630 merely not aligned towards or away from the target categories in the semantic space. Thus, these
631 results suggest that, compared to occipitotemporal areas, attention more diversely warps
632 semantic representations in parietal and premotor AON nodes by shifting tuning for both target
633 and nontarget categories.

634 635 **Tuning shifts interact with intrinsic selectivity of cortical voxels for action categories**

636 A recent study on visual attention has reported that in strongly object-selective regions voxel
637 tuning for a preferred object might be robust against attention directed to a nonpreferred object
638 (e.g., *houses* for fusiform face area, FFA, and *faces* for parahippocampal place area, PPA; Çukur
639 et al., 2013). This previous result suggests that the degree of response modulations in a brain
640 region might depend on the alignment between the search target and the intrinsically preferred
641 object. It is thus likely that tuning shifts during search for an action category also interact with
642 the intrinsic selectivity of cortical voxels for the target category. Tuning shifts as measured by
643 TSI signal an overall increase in relative selectivity for target versus nontarget categories,
644 aggregated across search tasks. Yet, interaction of tuning shifts with intrinsic selectivity for
645 action categories is task-specific by definition. Therefore, to examine potential interactions, we
646 calculated a target preference index ($PI \in [-1, 1]$) separately during search for communication
647 actions (PI_{com}) and during search for locomotion actions (PI_{loc}). PI_{com} was taken as the difference
648 in selectivity for *communication* versus *locomotion*, during search for communication actions.

649 Analogously, PI_{loc} was taken as the difference in selectivity for locomotion versus
 650 *communication*, during search for locomotion actions.

651

652 Voxel-wise PI_{com} and PI_{loc} values were projected onto cortical flat maps for visual inspection
 653 (Fig. 14; see Figs. 11-1c to 11-5c for results in individual brain spaces) and quantitatively
 654 examined in ROIs (Fig. 15c, d). We observe that semantic tuning in areas with indiscriminate
 655 selectivity for behaviourally relevant action categories (e.g., selective for low-level visual
 656 features or static object categories) show insignificant shifts regardless of the search task.
 657 Meanwhile, many voxels across anterior parietal, occipital, and cingulate cortices –with intrinsic
 658 action category preferences– show differential preference for one of the two target action
 659 categories as indicated by high PI index during either search for communication or search for
 660 locomotion actions. Lastly, semantic tuning in voxels across posterior parietal and anterior
 661 prefrontal cortices with broad selectivity for actions shift toward the attended category
 662 irrespective of the search target. These specific cases are discussed in detail below.

663

664 ***Areas where both PI_{com} and PI_{loc} are non-significant***

665 We find that PI_{com} and PI_{loc} are non-significant in retinotopic early visual areas (RET; bootstrap
 666 test, $q(FDR)>0.05$) that represent low-level stimulus features, low-level motion-selective area
 667 (hMT; $q(FDR)>0.05$), and object-selective areas (FFA; occipitotemporal face area, OFA; PPA;
 668 retrosplenial cortex, RSC; and EBA; $q(FDR)>0.05$). Furthermore, PI_{com} and PI_{loc} are non-
 669 significant in anterior intraparietal cortex (aIP; $q(FDR)>0.05$), which is not involved in
 670 representing communication or locomotion actions (non-significant SI_{com} and SI_{loc} ,
 671 $q(FDR)>0.05$) (Noppeney, 2008; Rizzolatti et al., 1997; Urgen and Orban, 2021). These results
 672 suggest that during action-based search, semantic tuning does not change substantially in cortical
 673 areas that are selective for lower-level visual features or for neutral high-level action categories
 674 irrelevant to the task.

675

676 ***Areas where either PI_{com} or PI_{loc} are significant***

677 Several previous studies suggest that lateral and medial prefrontal cortices are causally involved
 678 in representing communication actions (Van Overwalle, 2009; Wilson-Mendenhall et al., 2013).
 679 Here, we find that PI_{loc} is non-significant while PI_{com} is significantly greater than zero in anterior
 680 inferior frontal gyrus (BA44/45; $d=1.94$, $q(FDR)<0.05$; $SI_{com}=0.12$, $q(FDR)<0.05$), in superior
 681 frontal gyrus (SFG; $d=1.94$, $q(FDR)<0.05$; $SI_{com}=0.18$, $q(FDR)<0.05$), and in anterior cingulate
 682 cortex (ACC; $d=.34$, $q(FDR)<0.05$; $SI_{com}=0.18$, $q(FDR)<0.05$). On the other hand, previous
 683 reports provide evidence for representation of animate locomotion actions in PPC, including IPS
 684 ($SI_{loc}=0.15$, $q(FDR)<0.05$) (Abdollahi et al., 2012; Battelli et al., 2003; Bremmer et al., 2001; Ilg
 685 et al., 2004). In accord, we find that PI_{com} is non-significant while PI_{loc} is significantly greater
 686 than zero in IPS ($d=3.95$, $q(FDR)<0.05$). Taken together, our findings suggest that in areas that
 687 are strongly selective for specific action categories, visual search for the preferred action shifts
 688 tuning more vigorously towards the preferred target category. It is also worth noting that these
 689 attentional effects are not limited to the AON, but rather extend to higher-order cortical areas
 690 involved in social cognition. Lastly, we find that PI_{loc} is significantly less than zero while PI_{com} is
 691 non-significant ($d=0.73$, $q(FDR)>0.05$; $SI_{loc}=-0.23$, $q(FDR)<0.05$) in dPMC. This result supports

692 the view that dPMC enhances the representation of distractors during search for locomotion
693 actions (Anticevic et al., 2010; Toepper et al., 2010; Zhou et al., 2012).

694

695 *Areas where both PI_{com} and PI_{loc} are significant*

696 Posterior STS (pSTS), posterior middle temporal gyrus (pMTG), and SMG are considered as
697 AON nodes that maintain representation of actions regardless of their semantic category
698 (Caspers et al., 2010; Jastorff et al., 2016; Lui et al., 2008). We find that both PI_{com} and PI_{loc} are
699 significantly greater than zero in pSTS, pMTG, and SMG, consistent with their generic action
700 selectivity. In addition, several previous studies suggest that MFG –as a node in dorsal attention
701 network– facilitates visual search by maintaining the representation of targets (Corbetta and
702 Shulman, 2002; Mars and Grol, 2007; Paneri and Gregoriou, 2017; Ptak et al., 2017).
703 Accordingly, here we find that PI_{com} and PI_{loc} are significantly greater than zero in MFG
704 ($q(FDR)<0.05$). Overall, these results indicate that in areas with generic action selectivity and in
705 high-level cortical areas, attention facilitates action-based search by shifting representations
706 toward targets irrespective of their semantic category.

707

708 The results presented here can be explored online via an interactive brain viewer at
709 http://www.icon.bilkent.edu.tr/brainviewer/shahdloo_etal/.

710

711

712 **Discussion**

713 Several previous studies have reported response modulations during action-based attention in
714 parietal and prefrontal cortices, but not in occipitotemporal areas (Nastase et al., 2017, 2018;
715 Nicholson et al., 2017). Yet we observe significant attentional tuning shifts in occipitotemporal
716 cortex. Unlike previous studies, our analysis approach enables us to measure single-voxel tuning.
717 Our movie stimulus contains a large set of action categories in natural context in contrast to
718 controlled stimuli with a handful of actions on a homogeneous background. Lastly, we
719 investigate actions that are performed by animate actors, known to elicit robust responses across
720 the occipitotemporal cortex (Isik et al., 2017; Thompson and Parasuraman, 2012; Walbrin and
721 Koldewyn, 2019; Walbrin et al., 2018). These design factors might have enabled us to detect
722 tuning shifts in early stages of AON comprising occipitotemporal areas.

723

724 Recent studies emphasize the role of AG and SMG in multi-modal semantic representation while
725 observing actions, hearing action sounds, or reading action words (Bedny and Caramazza, 2011;
726 van Dam et al., 2010; Liljeström et al., 2008; Pizzamiglio et al., 2005). Evidence also suggests
727 that during semantic processing these areas act as central connectivity hubs, passing information
728 from low-level perceptual areas onto higher-level areas in prefrontal cortex (Farahibozorg et al.,
729 2019; Hoeren et al., 2013). We find significant tuning shifts toward targets in AG and SMG,
730 higher than that in occipitotemporal and premotor AON nodes, irrespective of the search target.
731 Our results can be taken to suggest a higher place for AG and SMG in the hierarchy of
732 semantic representations compared to remaining AON nodes. Another potential account could be
733 that areas with stronger action selectivity might undergo stronger tuning shifts and future studies
734 are warranted to investigate this issue more directly.

735

736 Cortical areas selective for an object category are suggested to retain their preferred tuning even
737 when a non-preferred category is the search target (Çukur et al., 2013; Reddy and Kanwisher,
738 2007; Shahdloo et al., 2020). We find that semantic tuning of voxels in locomotion-action-
739 selective superior parietal cortex are shifted toward locomotion actions only during search for
740 this target. Likewise, semantic tuning of voxels in communication-selective anterior prefrontal
741 cortex are shifted toward communication actions only during search for communication. These
742 results suggest that semantic tuning shifts interact with the intrinsic selectivity for target
743 categories.

744
745 We used WordNet to label action categories in the stimulus and create a one-hot-encoded
746 stimulus feature matrix. Thus, it is possible to conduct part of the reported analyses by directly
747 examining modulations of category responses. However, assessments in the 922-dimensional
748 category space would treat each category independently ignoring semantic similarities, and they
749 would be inherently noisier reducing our sensitivity for detecting tuning shifts. To assess TSI for
750 non-target categories, category responses were masked to zero out responses for communication
751 and locomotion actions. If selectivity measurements had been performed based on one-hot
752 category features, this masking would eliminate all information related to target categories. It
753 would then be impossible to quantify whether tuning for non-target categories shifts
754 towards/away from the attended category. Therefore, we performed our analyses in a dense-
755 encoded semantic space obtained via PCA. An alternative is voxel-wise modelling with a dense-
756 encoded stimulus feature matrix derived using embedding models (Mikolov et al., 2013; Devlin
757 et al., 2019). During preliminary experiments in the current study and prior studies from our lab
758 (Huth et al, 2018; Çelik et al., 2021), we compared the category model against dense embedding
759 models and estimates of attentional modulations did not vary significantly by choice of model.
760 As such, we do not expect a profound difference between results from these various models,
761 although there could be practical differences in terms of interpretation and feature similarity
762 assessments.

763
764 We employed communication and locomotion as target categories to maximize our chances for
765 detecting semantic tuning shifts, since previous studies suggest that these action categories have
766 broadly distributed and distinctive representations (Urgen and Orban, 2021). Attentional
767 modulations in multi-voxel response patterns were recently reported during search for several
768 other categories related to animal taxonomy or actions (Nastase et al., 2017). We have observed
769 in preliminary experiments that search for many salient categories in natural movies elicits
770 tuning shifts (data not shown here). Thus, it is likely that tuning shifts are a ubiquitous
771 mechanism for response modulation during natural visual search for action categories. However,
772 there may be differences in the strength and cortical distribution of tuning shifts depending on
773 the target action, and future studies are warranted to systematically examine whether and how
774 tuning shifts generalise across action categories. Evidence suggests that attending to an object
775 can modulate responses to features correlated with the target (O'Craven et al., 1999). We have
776 previously reported that attending to a target object (e.g., vehicles) enhances the representation of
777 semantically similar actions (e.g., driving) (Çukur et al., 2013). It is thus possible that attending
778 to a target action could induce tuning shifts for correlated features such as the object categories
779 pertaining to the actor. Since we restricted the target actions to be performed by the same
780 animate actors, we did not examine tuning changes for objects in this study.

781

782 The tuning profile of a voxel refers to its response levels to the examined range of features.
783 Attention can induce different modulations on this profile including baseline changes, gain
784 changes and tuning shifts. Baseline changes imply an additive offset, gain changes imply a
785 multiplicative offset to responses uniformly across features, neither changing the shape of the
786 profile. Instead, tuning shifts alter shape by shifting selectivity towards the target, changing
787 responses to both attended and unattended features. Here, we find that the overall tuning shift is
788 attributed to significant tuning changes for both target and non-target categories. Such broadly
789 distributed changes imply alteration in the shape of the tuning profile. Since our measurements
790 are naturally limited by the spatiotemporal resolution of BOLD responses, we cannot make
791 definitive inferences about the neural mechanisms underlying voxel tuning shifts, which could be
792 attributed to baseline, gain or selectivity changes in single neurons (Connor et al., 1997; David et
793 al., 2008; Reynolds et al., 2000). Further electrophysiological work would be needed to
794 characterize neural tuning shifts during action-based search.

795
796 A common practice in fMRI is to collect a relatively limited dataset from a greater number of
797 subjects to increase reliability of across-subject assessments at the expense of individual-subject
798 results. Diverting away from this practice, here we collect a larger amount of data per subject to
799 give greater focus to reliability in single subjects. This procedure substantially increased the
800 amount and diversity of fMRI data collected per subject, which enhanced the quality of resulting
801 models and thereby reliability of individual-subject results. However, we acknowledge that
802 future studies are warranted to assess to what degree the results reported in the current study
803 generalise to a broader population of subjects.

804
805 The natural movie stimuli used here have greater ecological validity compared to simplified or
806 controlled movie clips used in many action-perception studies. That said, action categories in
807 natural movies might be correlated with low-level features such as global motion-energy
808 (Nishimoto et al., 2011; Weiss et al., 2006) and intermediate-level features such as scene
809 dynamics (Grossman and Blake, 2002). Substantial correlations can confound the estimated
810 category responses and tuning shifts. We employed several procedures to control for potential
811 biases. First, to minimize correlations between category responses and global motion-energy, we
812 used a nuisance motion-energy regressor (Nishimoto et al., 2011). Second, we restricted analyses
813 to voxels uniquely predicted by the category model after accounting for motion-energy and STIP
814 features. Voxels in areas such as LOC might encode multiple levels of features ranging from
815 motion-energy and kinematics to semantics. Thus, controlling for motion-energy and STIP
816 features might reduce sensitivity for attentional modulation of perceptual selectivity in these
817 areas. Our analyses do not consider attentional tuning shifts that might be evident for motion-
818 energy and STIP features, or other features such as expected action goals (Hudson et al., 2016a,
819 2016b), and actors' perceived attitude (Bach and Schenke, 2017). Some level of ambiguity will
820 be naturally evident about what specific aspect of the correlated stimulus features are most
821 relevant for measured cortical representations. Addressing this ambiguity requires complete
822 decorrelation of all possible feature sets, yet conclusions derived using decorrelated stimuli
823 deprived from their natural context might no longer be ecologically relevant. It remains
824 important work to assess the effects of category-based search on multiple levels of feature
825 representations.

826

827 In conclusion, we showed that natural visual search for a specific action category modulates
828 semantic representations, causing tuning shifts toward the target in single voxels within and
829 beyond the AON. Attentional modulations further interact with intrinsic selectivity of neural
830 populations for search targets. This dynamic attentional mechanism can facilitate action
831 perception by efficiently allocating neural resources to accentuate the representation of task-
832 relevant action categories. Overall, these findings offer new insights into the effects of category-
833 based visual search on brain responses (Çukur et al., 2013; Erez and Duncan, 2015; Harel et al.,
834 2014; Peelen et al., 2009), as our results help explain humans' astounding ability to perceive
835 others' actions in dynamic, cluttered daily-life experiences.

836

837

838 **Data and software availability**

839 Data supporting the findings of this study are available from the corresponding authors upon
840 request. Results can be explored online via an interactive brain viewer at

841 http://www.icon.bilkent.edu.tr/brainviewer/shahdloo_etal/.

842 The codes used to estimate spatially informed voxelwise model weights is freely available on
843 GitHub at <https://github.com/icon-lab/SPIN-VM>.

844 **Acknowledgements**

845 **Funding**

846 This work was supported in part by a Marie Curie Actions Career Integration Grant (PCIG13-
847 GA-2013-618101), by a European Molecular Biology Organisation Installation Grant (IG 3028),
848 by a TUBA GEBIP 2015 fellowship, and by a BAGEP 2017 fellowship.

849 **Author contributions**

850 Conceptualization, J.L.G and T.C.; Methodology, M.S., and T.C.; Software, M.S., and E.C.;
851 Investigation, M.S., B.A.U. and T.C.; Writing-Original Draft, M.S.; Writing-Review and
852 Editing, M.S., T.C., B.A.U., J.L.G, and E.C.; Funding Acquisition, T.C.; Resources, M.S., T.C.,
853 and J.L.G; Supervision, T.C.

854 The authors declare no competing financial interests.

References

- 855
856 Abdollahi, R.O., Jastorff, J., and Orban, G.A. (2012). Common and Segregated Processing of
857 Observed Actions in Human SPL. *Cerebral Cortex* *23*, 2734–2753.
- 858 Anticevic, A., Repovs, G., and Barch, D.M. (2010). Resisting emotional interference: Brain
859 regions facilitating working memory performance during negative distraction. *Cognitive,*
860 *Affective, & Behavioral Neuroscience* *10*, 159–173.
- 861 Bach, P., and Schenke, K.C. (2017). Predictive social perception: Towards a unifying framework
862 from action observation to person knowledge. *Social and Personality Psychology Compass* *11*,
863 *e12312*.
- 864 Battelli, L., Cavanagh, P., and Thornton, I.M. (2003). Perception of biological motion in parietal
865 patients. *Neuropsychologia* *41*, 1808–1816.
- 866 Bedny, M., and Caramazza, A. (2011). Perception, action, and word meanings in the human
867 brain: The case from action verbs. *Annals of the New York Academy of Sciences* *1224*, 81–95.
- 868 Bremner, F., Schlack, A., Duhamel, J.-R., Graf, W., and Fink, G.R. (2001). Space Coding in
869 Primate Posterior Parietal Cortex. *NeuroImage* *14*, S46–S51.
- 870 Buccino, G., Binkofski, F., Fink, G.R., Fadiga, L., Fogassi, L., Gallese, V., Seitz, R.J., Zilles, K.,
871 Rizzolatti, G., and Freund, H.J. (2001). Action observation activates premotor and parietal areas
872 in a somatotopic manner: An fMRI study. *European Journal of Neuroscience* *13*, 400–404.
- 873 Caspers, S., Zilles, K., Laird, A.R., and Eickhoff, S.B. (2010). ALE meta-analysis of action
874 observation and imitation in the human brain. *NeuroImage* *50*, 1148–1167.
- 875 Cavina-Pratesi, C., Connolly, J.D., Monaco, S., Figley, T.D., Milner, D., Schenk, T., Culham,
876 J.C. (2018). Human neuroimaging reveals the subcomponents of grasping, reaching and pointing
877 actions. *Cortex* *98*, 128–148.
- 878 Çelik, E., Dar, S.U.H., Yılmaz, Ö., Keleş, Ü., and Çukur, T. (2019). Spatially informed
879 voxelwise modeling for naturalistic fMRI experiments. *NeuroImage* *186*, 741–757.
- 880 Connor, C.E., Preddie, D.C., Gallant, J.L., and Van Essen, D.C. (1997). Spatial Attention Effects
881 in Macaque Area V4. *Journal of Neuroscience* *17*, 3201–3214.
- 882 Corbetta, M., and Shulman, G.L. (2002). Control of goal-directed and stimulus-driven attention
883 in the brain. *Nature Reviews Neuroscience* *3*, 215–229.
- 884 Corbo, D., and Orban, G.A. (2017). Observing others speak or sing activates SPT and
885 neighboring parietal cortex. *Journal of Cognitive Neuroscience* *29*, 1002–1021.
- 886 Çukur, T., Nishimoto, S., Huth, A.G., and Gallant, J.L. (2013). Attention during natural vision
887 warps semantic representation across the human brain. *Nature Neuroscience* *16*, 763–770.

- 888 Dalal, N., and Triggs, B. (2005). Histograms of oriented gradients for human detection. In IEEE
889 Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05), IEEE,
890 886–893.
- 891 van Dam, W.O., Rueschemeyer, S.-A., and Bekkering, H. (2010). How specifically are action
892 verbs represented in the neural motor system: An fMRI study. *NeuroImage* 53, 1318–1325.
- 893 David, S.V., Hayden, B.Y., Mazer, J.A., and Gallant, J.L. (2008). Attention to stimulus features
894 shifts spectral tuning of V4 neurons during natural vision. *Neuron* 59, 509–521.
- 895 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of Deep
896 Bidirectional Transformers for Language Understanding. ArXiv:1810.04805.
- 897 Erez, Y., and Duncan, J. (2015). Discrimination of Visual Categories Based on Behavioral
898 Relevance in Widespread Regions of Frontoparietal Cortex. *Journal of Neuroscience* 35, 12383–
899 12393.
- 900 Farahibozorg, S.-R., Henson, R.N., Woollams, A.M., and Hauk, O. (2019). Distinct roles for the
901 Anterior Temporal Lobe and Angular Gyrus in the spatio-temporal cortical semantic network.
902 *BioRxiv*, 544114.
- 903 Ferri, S., Rizzolatti, G., and Orban, G.A. (2015). The organisation of the posterior parietal cortex
904 devoted to upper limb actions: An fMRI study. *Human Brain Mapping* 36, 3845–3866.
- 905 Grafton, S.T., and de C Hamilton, A.F. (2007). Evidence for a distributed hierarchy of action
906 representation in the brain. *Human Movement Science* 26, 590–616.
- 907 Griffiths, T. L., Steyvers, M., and Tenenbaum, J. B. (2007). Topics in semantic representation.
908 *Psychological Review*, 114(2), 211–244.
- 909 Grosbras, M.H., Beaton, S., and Eickhoff, S.B. (2012). Brain regions involved in human
910 movement perception: A quantitative voxel-based meta-analysis. *Human Brain Mapping* 33,
911 431–454.
- 912 Grossman, E.D., and Blake, R. (2002). Brain areas active during visual perception of biological
913 motion. *Neuron* 35, 1167–1175.
- 914 Hamilton, A.F. de C., and Grafton, S.T. (2006). Goal representation in human anterior
915 intraparietal sulcus. *Journal of Neuroscience* 26, 1133–1137.
- 916 Handjaras, G., Bernardi, G., Benuzzi, F., Nichelli, P.F., Pietrini, P., and Ricciardi, E. (2015). A
917 topographical organisation for action representation in the human brain. *Human Brain Mapping*
918 36, 3832–3844.
- 919 Harel, A., Kravitz, D.J., and Baker, C.I. (2014). Task context impacts visual object processing
920 differentially across the cortex. *PNAS* 111, E962-71.

- 921 Harris, C., and Stephens, M. (1988). A combined corner and edge detector. In In Proc. Fourth
922 Alvey Vision Conference, 147–152.
- 923 Herrington, J., Nymberg, C., Faja, S., Price, E., and Schultz, R. (2012). The responsiveness of
924 biological motion processing areas to selective attention towards goals. *NeuroImage* 63, 581–
925 590.
- 926 Hoeren, M., Kaller, C.P., Glauche, V., Vry, M.-S., Rijntjes, M., Hamzei, F., and Weiller, C.
927 (2013). Action semantics and movement characteristics engage distinct processing streams
928 during the observation of tool use. *Experimental Brain Research* 229, 243–260.
- 929 Holte, M.B., Moeslund, T.B., and Fihl, P. (2010). View-invariant gesture recognition using 3D
930 optical flow and harmonic motion context. *Computer Vision and Image Understanding* 114,
931 1353–1361.
- 932 Hudson, M., Nicholson, T., Ellis, R., and Bach, P. (2016a). I see what you say: Prior knowledge
933 of other's goals automatically biases the perception of their actions. *Cognition* 146, 245–250.
- 934 Hudson, M., Nicholson, T., Simpson, W.A., Ellis, R., and Bach, P. (2016b). One step ahead: The
935 perceived kinematics of others' actions are biased toward expected goals. *Journal of*
936 *Experimental Psychology: General* 145, 1–7.
- 937 Huth, A.G., Nishimoto, S., Vu, A.T., and Gallant, J.L. (2012). A Continuous Semantic Space
938 Describes the Representation of Thousands of Object and Action Categories across the Human
939 Brain. *Neuron* 76, 1210–1224.
- 940 Ilg, U.J., Schumann, S., and Thier, P. (2004). Posterior parietal cortex neurons encode target
941 motion in world-centered coordinates. *Neuron* 43, 145–151.
- 942 Isik, L., Koldewyn, K., Beeler, D., and Kanwisher, N.G. (2017). Perceiving social interactions in
943 the posterior superior temporal sulcus. *PNAS* 114, E9145–E9152.
- 944 Jastorff, J., and Orban, G.A. (2009). Human Functional Magnetic Resonance Imaging Reveals
945 Separation and Integration of Shape and Motion Cues in Biological Motion Processing. *Journal*
946 *of Neuroscience* 29, 7315–7329.
- 947 Jastorff, J., Begliomini, C., Fabbri-Destro, M., Rizzolatti, G., and Orban, G.A. (2010). Coding
948 observed motor acts: Different organisational principles in the parietal and premotor cortex of
949 humans. *Journal of Neurophysiology* 104, 128–140.
- 950 Jastorff, J., Abdollahi, R.O., Fasano, F., and Orban, G.A. (2016). Seeing biological actions in
951 3D: An fMRI study. *Human Brain Mapping* 37, 203–219.
- 952 Johansson, G. (1973). Visual perception of biological motion and a model for its analysis.
953 *Perception & Psychophysics* 14, 201–211.
- 954 Jurafsky, D. and Martin, J. (2021). *Speech and language processing*. 3rd ed., Prentice Hall.

- 955 Kilintari, M., Raos, V., and Savaki, H.E. (2014). Involvement of the Superior Temporal Cortex
956 in Action Execution and Action Observation. *Journal of Neuroscience* *34*, 8999–9011.
- 957 Kiremitçi, I., Yılmaz, Ö., Çelik, E., Shahdloo, M., Huth, A.G., and Çukur, T. (2021). Attentional
958 Modulation of Hierarchical Speech Representations in a Multi-Talker Environment. *Cerebral*
959 *Cortex*, bhab136
- 960 de Lange, F.P., Spronk, M., Willems, R.M., Toni, I., and Bekkering, H. (2008). Complementary
961 Systems for Understanding Action Intentions. *Current Biology* *18*, 454–457.
- 962 Laptev, I. (2005). On space-time interest points. *International Journal of Computer Vision*, 107–
963 123.
- 964 Laptev, I., Marszałek, M., Schmid, C., and Rozenfeld, B. (2008). Learning realistic human
965 actions from movies. 26th IEEE Conference on Computer Vision and Pattern Recognition
966 (CVPR), 1-8.
- 967 Lescroart, M.D., and Gallant, J.L. (2019). Human Scene-Selective Areas Represent 3D
968 Configurations of Surfaces. *Neuron* *101*, 178-192.e7.
- 969 Liljeström, M., Tarkiainen, A., Parviainen, T., Kujala, J., Numminen, J., Hiltunen, J., Laine, M.,
970 and Salmelin, R. (2008). Perceiving and naming actions and objects. *NeuroImage* *41*, 1132–
971 1141.
- 972 Lingnau, A., and Downing, P. E. (2015). The lateral occipitotemporal cortex in action. *Trends in*
973 *cognitive sciences*, *19*(5), 268-277.
- 974 Lui, F., Buccino, G., Duzzi, D., Benuzzi, F., Crisi, G., Baraldi, P., Nichelli, P., Porro, C.A., and
975 Rizzolatti, G. (2008). Neural substrates for observing and imagining non-object-directed actions.
976 *Social Neuroscience* *3*, 261–275.
- 977 Mars, R.B., and Grol, M.J. (2007). Dorsolateral prefrontal cortex, working memory, and
978 prospective coding for action. *Journal of Neuroscience* *27*, 1801–1802.
- 979 Miller, G.A. (1995). WordNet: a lexical database for English. *Communications of ACM* *38*, 39–
980 41.
- 981 Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient Estimation of Word
982 Representations in Vector Space. ArXiv:1301.3781
- 983 Molinari, E., Baraldi, P., Campanella, M., Duzzi, D., Nocetti, L., Pagnoni, G., and Porro, C.A.
984 (2013). Human parietofrontal networks related to action observation detected at rest. *Cerebral*
985 *Cortex* *23*, 178–186.
- 986 Muthukumaraswamy, S.D., and Singh, K.D. (2008). Modulation of the human mirror neuron
987 system during cognitive activity. *Psychophysiology* *45*, 896–905.

- 988 Muthukumaraswamy, S.D., Johnson, B.W., and McNair, N.A. (2004). Mu rhythm modulation
989 during observation of an object-directed grasp. *Cognitive Brain Research* *19*, 195–201.
- 990 Nastase, S.A., Connolly, A.C., Oosterhof, N.N., Halchenko, Y.O., Guntupalli, J.S., Di Oleggio
991 Castello, M.V., Gors, J., Gobbini, M.I., and Haxby, J.V. (2017). Attention selectively reshapes
992 the geometry of distributed semantic representation. *Cerebral Cortex* *27*, 4277–4291.
- 993 Nastase, S.A., Halchenko, Y.O., Connolly, A.C., Gobbini, M.I., and Haxby, J.V. (2018). Neural
994 responses to naturalistic clips of behaving animals in two different task contexts. *Frontiers in*
995 *Neuroscience* *12*, 316.
- 996 Nelissen, K., Vanduffel, W., and Orban, G.A. (2006). Charting the Lower Superior Temporal
997 Region, a New Motion-Sensitive Region in Monkey Superior Temporal Sulcus. *Journal of*
998 *Neuroscience* *26*, 5929–5947.
- 999 Nelissen, K., Borra, E., Gerbella, M., Rozzi, S., Luppino, G., Vanduffel, W., Rizzolatti, G., and
1000 Orban, G.A. (2011). Action observation circuits in the macaque monkey cortex. *Journal of*
1001 *Neuroscience* *31*, 3743–3756.
- 1002 Newman-Norlund, R., van Schie, H.T., van Hoek, M.E.C., Cuijpers, R.H., and Bekkering, H.
1003 (2010). The role of inferior frontal and parietal areas in differentiating meaningful and
1004 meaningless object-directed actions. *Brain Research* *1315*, 63–74.
- 1005 Nicholson, T., Roser, M., and Bach, P. (2017). Understanding the Goals of Everyday
1006 Instrumental Actions Is Primarily Linked to Object, Not Motor-Kinematic, Information:
1007 Evidence from fMRI. *PLoS ONE* *12*, e0169700.
- 1008 Nishimoto, S., Vu, A.T., Naselaris, T., Benjamini, Y., Yu, B., and Gallant, J.L. (2011).
1009 Reconstructing Visual Experiences from Brain Activity Evoked by Natural Movies. *Current*
1010 *Biology* *21*, 1641–1646.
- 1011 Noppeney, U. (2008). The neural systems of tool and action semantics: A perspective from
1012 functional imaging. *Journal of Physiology* *102*, 40–49.
- 1013 Nunez-Elizalde, A. O., Huth, A. G., and Gallant, J. L. (2019). Voxelwise encoding models with
1014 non-spherical multivariate normal priors. *NeuroImage*, *197*:482–492.
- 1015
- 1016 Oberman, L.M., Pineda, J.A., and Ramachandran, V.S. (2007). The human mirror neuron
1017 system: A link between action observation and social skills. *Social Cognitive and Affective*
1018 *Neuroscience* *2*, 62–66.
- 1019 Oosterhof, N.N., Wiggett, A.J., Diedrichsen, J., Tipper, S.P., and Downing, P.E. (2010). Surface-
1020 based information mapping reveals crossmodal vision-action representations in human parietal
1021 and occipitotemporal cortex. *Journal of Neurophysiology* *104*, 1077–1089.

- 1022 Oosterhof, N.N., Tipper, S.P., and Downing, P.E. (2012). Viewpoint (in)dependence of action
1023 representations: an MVPA study. *Journal of Cognitive Neuroscience* *24*, 975–989.
- 1024 Oosterhof, N.N., Tipper, S.P., and Downing, P.E. (2013). Crossmodal and action-specific:
1025 Neuroimaging the human mirror neuron system. *Trends in Cognitive Sciences* *17*, 311–318.
- 1026 Paneri, S., and Gregoriou, G.G. (2017). Top-down control of visual attention by the prefrontal
1027 cortex. *Functional specialization and long-range interactions. Frontiers in Neuroscience* *11*, 545.
- 1028 Peelen, M.V., Fei-Fei, L., and Kastner, S. (2009). Neural mechanisms of rapid natural scene
1029 categorization in human visual cortex. *Nature* *460*, 94–97.
- 1030 Pizzamiglio, L., Aprile, T., Spitoni, G., Pitzalis, S., Bates, E., D’Amico, S., and Di Russo, F.
1031 (2005). Separate neural systems for processing action- or non-action-related sounds. *NeuroImage*
1032 *24*, 852–861.
- 1033 Popham, S. F., Huth, A. G., Bilenko, N. Y., Deniz, F., Gao, J. S., Nunez-Elizalde, A. O., and
1034 Gallant, J. L. (2021). Visual and linguistic semantic representations are aligned at the border of
1035 human visual cortex. *Nature Neuroscience*, *24*(11), 1628–1636.
- 1036 Ptak, R., Schnider, A., and Fellrath, J. (2017). The Dorsal Frontoparietal Network: A Core
1037 System for Emulated Action. *Trends in Cognitive Sciences* *21*, 589–599.
- 1038 Puglisi, G., Leonetti, A., Landau, A., Fornia, L., Cerri, G., and Borroni, P. (2017). The role of
1039 attention in human motor resonance. *PLoS ONE* *12*, e0177457.
- 1040 Puglisi, G., Leonetti, A., Cerri, G., and Borroni, P. (2018). Attention and cognitive load modulate
1041 motor resonance during action observation. *Brain & Cognition* *128*, 7–16.
- 1042 Ramsey, R., and Hamilton, A.F.D.C. (2010). Understanding actors and object-goals in the
1043 human brain. *NeuroImage* *50*, 1142–1147.
- 1044 Reddy, L., and Kanwisher, N.G. (2007). Category Selectivity in the Ventral Visual Pathway
1045 Confers Robustness to Clutter and Diverted Attention. *Current Biology* *17*, 2067–2072.
- 1046 Reynolds, J.H., Pasternak, T., and Desimone, R. (2000). Attention Increases Sensitivity of V4
1047 Neurons. *Neuron* *26*, 703–714.
- 1048 Rizzolatti, G., and Matelli, M. (2003). Two different streams form the dorsal visual system:
1049 Anatomy and functions. *Experimental Brain Research* *153*, 146–157.
- 1050 Rizzolatti, G., Fogassi, L., and Gallese, V. (1997). Parietal cortex: from sight to action. *Current*
1051 *Opinion in Neurobiology* *7*, 562–567.
- 1052 Rowe, J., Friston, K., Frackowiak, R., and Passingham, R. (2002). Attention to action: Specific
1053 modulation of corticocortical interactions in humans. *NeuroImage* *17*, 988–998.

- 1054 Rozzi, S., and Fogassi, L. (2017). Neural coding for action execution and action observation in
1055 the prefrontal cortex and its role in the organisation of socially driven behavior. *Frontiers in*
1056 *Neuroscience 11*, 276.
- 1057 Safford, A.S., Hussey, E.A., Parasuraman, R., and Thompson, J.C. (2010). Object-based
1058 attentional modulation of biological motion processing: spatiotemporal dynamics using
1059 functional magnetic resonance imaging and electroencephalography. *Journal of Neuroscience 30*,
1060 9064–9073.
- 1061 Schuch, S., Bayliss, A.P., Klein, C., and Tipper, S.P. (2010). Attention modulates motor system
1062 activation during action observation: Evidence for inhibitory rebound. *Experimental Brain*
1063 *Research 205*, 235–249.
- 1064 Seidl, K.N., Peelen, M.V., and Kastner, S. (2012). Neural evidence for distracter suppression
1065 during visual search in real-world scenes. *Journal of Neuroscience 32*, 11812–11819.
- 1066 Shahdloo, M., Çelik, E., and Çukur, T. (2020). Biased competition in semantic representation
1067 during natural visual search. *NeuroImage 216*, 116383.
- 1068 Tarhan, L., and Konkle, T. (2020). Sociality and interaction envelope organise visual action
1069 representations. *Nature Communications 11*, 3002.
- 1070 Thompson, J., and Parasuraman, R. (2012). Attention, biological motion, and action recognition.
1071 *NeuroImage 59*, 4–13.
- 1072 Toepper, M., Gebhardt, H., Beblo, T., Thomas, C., Driessen, M., Bischoff, M., Blecker, C.R.,
1073 Vaitl, D., and Sammer, G. (2010). Functional correlates of distracter suppression during spatial
1074 working memory encoding. *Neuroscience 165*, 1244–1253.
- 1075 Urgen, B.A., and Orban, G.A. (2021). The unique role of parietal cortex in action observation:
1076 Functional organisation for communicative and manipulative actions. *NeuroImage 273*, 118220.
- 1077 Urgen, B.A., Pehlivan, S., and Saygin, A.P. (2019). Distinct representations in occipito-temporal,
1078 parietal, and premotor cortex during action perception revealed by fMRI and computational
1079 modeling. *Neuropsychologia 127*, 35–47.
- 1080 Van Overwalle, F. (2009). Social cognition and the brain: A meta-analysis. *Human Brain*
1081 *Mapping 30*, 829–858.
- 1082 Walbrin, J., and Koldewyn, K. (2019). Dyadic interaction processing in the posterior temporal
1083 cortex. *NeuroImage 198*, 296–302.
- 1084 Walbrin, J., Downing, P., and Koldewyn, K. (2018). Neural responses to visually observed social
1085 interactions. *Neuropsychologia 112*, 31–39.
- 1086 Weiss, P.H., Rahbari, N.N., Lux, S., Pietrzyk, U., Noth, J., and Fink, G.R. (2006). Processing the
1087 spatial configuration of complex actions involves right posterior parietal cortex: An fMRI study
1088 with clinical implications. *Human Brain Mapping 27*, 1004–1014.

- 1089 Wilson-Mendenhall, C.D., Simmons, W.K., Martin, A., and Barsalou, L.W. (2013). Contextual
 1090 processing of abstract concepts reveals neural representations of nonlinguistic semantic content.
 1091 *Journal of Cognitive Neuroscience* 25, 920–935.
- 1092 Wurm, M.F., and Caramazza, A. (2019). Lateral occipitotemporal cortex encodes perceptual
 1093 components of social actions rather than abstract representations of sociality. *NeuroImage* 202,
 1094 116153.
- 1095 Wurm, M.F., Caramazza, A., and Lingnau, A. (2017). Action Categories in Lateral
 1096 Occipitotemporal Cortex Are Organised Along Sociality and Transitivity. *Journal of*
 1097 *Neuroscience* 37, 562–575.
- 1098 Zhou, X., Katsuki, F., Qi, X.-L., and Constantinidis, C. (2012). Neurons with inverted tuning
 1099 during the delay periods of working memory tasks in the dorsal prefrontal and posterior parietal
 1100 cortex. *J Neurophysiology* 108, 31–38.

1101
 1102
 1103

1104 Figure Captions

1105

1106 **Figure 1. Hypothesised changes in semantic representation of action categories.** Recent evidence suggests that
 1107 the human brain organises hundreds of object and action categories in a semantic space that is distributed
 1108 systematically across the cerebral cortex (Huth et al., 2012). **a.** Semantic representation for a single subject from
 1109 Çukur et al. (2013) is shown on flattened cortical surface and on inflated hemispheres. Colours indicate tuning for
 1110 different object or action categories (see colour legend). Regions of interest identified using conventional functional
 1111 localizers are denoted by white borders. Abbreviations for regions of interest are listed in *Materials and Methods*. **b.**
 1112 In the semantic space, action categories that are semantically similar to each other are mapped to nearby points and
 1113 semantically dissimilar actions are mapped to distant points. There is evidence that visual search for object
 1114 categories warps semantic representation in favour of the targets by shifting single-voxel tuning for object categories
 1115 toward target objects (Çukur et al., 2013). Thus, we hypothesised that visual search for a given action category
 1116 should similarly expand the semantic representation of the target and semantically similar categories.

1117

1118 **Figure 2. Model fitting and validation procedure.** Undergoing fMRI, human subjects viewed 60mins of natural
 1119 movies and covertly searched for *communication* or *locomotion* action categories while fixating on a central dot. **a.**
 1120 An indicator matrix was constructed that identified the presence of each of the 922 object and action categories in
 1121 each 1-sec clip of the movies (see Fig. 2-1). Nuisance regressors were included to account for head-motion,
 1122 physiological noise, and eye-movement confounds. An additional nuisance regressor was included to account for
 1123 target detection confounds. In a nested cross-validation (CV) procedure, regularized linear regression was used to
 1124 estimate separate category model weights (i.e., category responses) for each search task that mapped each category
 1125 feature to the recorded BOLD responses in single voxels. **b.** Accuracy of the fit models was cross-validated by
 1126 measuring prediction performance on the held-out data in each CV fold, after discarding the nuisance regressors and
 1127 the target regressor. Prediction score of the fit models was taken as product-moment correlation coefficient between
 1128 estimated and measured BOLD responses, averaged across the two search tasks.

1129

1130 **Figure 3. Prediction performance of the category model.** To test the performance of fit category models,
 1131 prediction score was calculated on held-out data as the product-moment correlation coefficient between the
 1132 predicted category responses and measured BOLD responses, and it was averaged across the two search tasks. **a.**

1133 Prediction scores of the category model are plotted on flattened cortical surfaces of individual subjects. A variance
 1134 partitioning analysis was used to quantify the response variance that was uniquely predicted by the category model
 1135 after accounting for low- and intermediate-level stimulus features (see *Materials and Methods*, Fig. 4). Voxels
 1136 where the category model did not explain unique response variance after accounting for these features were masked
 1137 (bootstrap test, $q(FDR) < 0.05$; see Fig. 11). **b.** To visualise single-subject results in a common space, prediction score
 1138 values are shown following projection onto the standard brain template from Freesurfer and averaging across
 1139 subjects, after getting thresholded in single subjects. Only voxels that were identified as semantic in all individual
 1140 subjects were averaged and displayed in the template. Regions of interest are illustrated by white borders. Several
 1141 important sulci are illustrated by dashed grey lines. Abbreviations for regions of interest and sulci are listed in
 1142 *Materials and Methods*. The category model predicts responses across ventral-temporal, parietal, and frontal cortices
 1143 well, suggesting that visual categories are broadly represented across visual and nonvisual cortex. Results can be
 1144 explored via an interactive brain viewer at http://www.icon.bilkent.edu.tr/brainviewer/shahdloo_etal/.

1145
 1146 **Figure 4. Comparison of category and control models.** The prediction scores (raw product-moment correlation
 1147 coefficient) of the category and control (the collection of motion-energy and STIP regressors) models were
 1148 measured for all cortical voxels. Voxels across all subjects are displayed. Each voxel is represented with a dot. Red
 1149 versus blue dots indicate whether the category model or the control model yields higher prediction scores. Black
 1150 dots indicate voxels where none of the models has high prediction scores. The category model outperforms the
 1151 control model in $53.75 \pm 3.29\%$ of cortical voxels (mean \pm sem; average over five subjects).

1152
 1153 **Figure 5. Fraction of uniquely predicted voxels in regions of interest (ROIs).** We identified voxels in which the
 1154 category model explained unique response variance after accounting for low-level motion-energy, and intermediate-
 1155 level STIP stimulus features by performing a variance partitioning analysis (see *Materials and Methods*). Fraction of
 1156 these *semantic voxels* is shown across ROIs, in individual subjects. Asterisk shows across subject significance
 1157 (bootstrap test, $q(FDR) < 0.05$).

1158 **Figure 6. Attention warps semantic representation of action categories.** To assess attentional changes, we
 1159 projected voxel-wise tuning profiles onto a continuous semantic space. **a.** The semantic space was derived from
 1160 principal components analysis (PCA) of tuning vectors measured during a separate passive-viewing task, and was
 1161 tested to be consistent across subjects (Fig. 7). To illustrate the semantic information embedded within this space,
 1162 action categories were projected onto PC1 and PC3 that best delineate the target actions (Fig. 8; words in regular
 1163 font show projections of individual categories, see Fig. 9). To illustrate the semantic content of the PCs,
 1164 characteristic actions of the movie stimulus were clustered in the semantic space, and cluster centres were projected
 1165 onto the PCs after getting labelled (bold-italic words; see *Materials and Methods*, Fig. 10). Average location of the
 1166 *communication* and *locomotion* actions are specified with red and green dots. **b.** Action category responses during
 1167 passive viewing and during the two search tasks were projected onto the semantic space, and a two-dimensional
 1168 colourmap was used to colour each voxel based on the projection values along PC1 and PC3 (see legend).
 1169 Projections in individual subjects were mapped onto the standard brain template from Freesurfer, and average
 1170 projections across subjects are displayed (see Figs. 6-1 to 6-5 for data in individual subjects). Figure formatting is
 1171 identical to Fig. 3. Many voxels across occipitotemporal, parietal, and prefrontal cortices shift their tuning toward
 1172 targets, suggesting that attention warps semantic representations of actions. Specifically, voxels in inferior posterior
 1173 parietal cortex, cingulate cortex, and anterior inferior prefrontal cortex shift their tuning toward *communication*
 1174 during search for *communication* actions. Meanwhile, voxels in superior posterior and medial parietal cortex shift
 1175 their tuning toward *locomotion* during search for *locomotion* actions. Results can be explored via an interactive brain
 1176 viewer at http://www.icon.bilkent.edu.tr/brainviewer/shahdloo_etal/.

1177
 1178 **Figure 7. Consistency of the semantic space across subjects.** To test whether the estimated semantic space is
 1179 consistent across subjects, leave-one-out cross-validation was performed. In each cross-validation fold, best-
 1180 predicted voxels from four subjects were used to derive 12 PCs to construct a semantic space. In the left-out subject,
 1181 semantic tuning profile for each voxel was obtained by projecting action category responses during passive viewing
 1182 onto the derived PCs. Next, product-moment correlation coefficient was calculated between the tuning profiles in
 1183 the derived space and the tuning profiles in the original semantic space. Results were averaged across semantic
 1184 voxels in the left-out subject. Correlation coefficients are shown for each PC and each subject. The cross-validated
 1185 semantic spaces consistently correlate with the original semantic space.

1186

1187 **Figure 8. The distance between target actions in subspaces spanned by different pairs of PCs.** To visualise
 1188 attentional modulation of semantic representation in Fig. 6, we compared projections of action category responses
 1189 onto a pair of PCs across the search tasks. To maximize our sensitivity in visualising the attentional modulations, we
 1190 chose the pair of dimensions that maximally separates the actions belonging to the two target categories (i.e.,
 1191 *communication* and *locomotion* categories). The Mahalanobis distance between communication actions and
 1192 locomotion actions (mean±sem across communication and locomotion actions) in the subspace spanned by each pair
 1193 of PCs is shown. Target actions are maximally separated across the subspace spanned by the first and third PCs.

1194

1195 **Figure 9. Distribution of action categories across PCs.** To illustrate the distribution of action categories embedded
 1196 within the semantic space, action categories were projected onto the PCs. Projections onto the first three PCs are
 1197 shown (words in regular font show projections of individual categories). To facilitate illustration, categories were
 1198 collapsed into 10 clusters and cluster centres were also projected onto the PCs (bold-italic words; see *Materials and*
 1199 *Methods*). Average location of the *communication* and *locomotion* actions are specified with red and green dots. The
 1200 estimated semantic space captures reasonable semantic variance across action categories in natural movies.

1201

1202 **Figure 10. Projections of action category clusters onto PCs.** Each of the 109 action categories were projected
 1203 onto the twelve semantic principal components (PCs). The projections were then clustered into 10 groups using k-
 1204 means and labelled for interpretation (see *Materials and Methods*). The projections of the cluster centres onto 12
 1205 PCs are shown. The first three dimensions were used to visualise the semantic space. The first dimension
 1206 distinguishes between self-movements (e.g., swirl, consume) and actions that are targeted toward other humans or
 1207 objects (e.g., reach, talk). The second dimension distinguishes between dynamic (e.g., drive, chase) versus static
 1208 actions (e.g., consume, struggle). The third dimension distinguishes between actions that involve humans (e.g., talk,
 1209 reach) and dynamic actions (e.g., fly, swirl).

1210

1211

1212 **Figure 11. Cortical distribution of tuning shifts. a.** To quantify the tuning shifts for the attended versus
 1213 unattended categories, a tuning shift index ($TSI_{all} \in [-1,1]$) was calculated for each voxel. Tuning shifts toward the
 1214 attended category would yield positive TSI (red colour), whereas negative TSI would indicate shifts away from the
 1215 attended category (blue colour). TSI_{all} values from individual subjects were projected onto the standard brain
 1216 template and averaged across subjects (see Figs. 11-1a to 11-5a for data in individual subjects). Figure formatting is
 1217 identical to Fig. 3. AON is outlined by green dashed lines. Voxels across many cortical regions shifted their tuning
 1218 toward the attended category. These include regions across AON (occipitotemporal cortex, posterior parietal cortex,
 1219 and premotor cortex), lateral prefrontal cortex, and anterior cingulate cortex. **b.** To examine how representation of
 1220 nontarget action categories changes during visual search, we measured a separate tuning shift index specifically for
 1221 these categories (TSI_{nt}). TSI_{nt} values from individual subjects were projected onto the standard brain template and
 1222 averaged across subjects (see Figs. 11-1b to 11-5b for data in individual subjects). TSI_{nt} shows a similar distribution
 1223 to TSI_{all} shown in **a**, albeit with lower magnitude (Fig. 12). Tuning shift for nontarget categories is positive across
 1224 many voxels within posterior parietal cortex and anterior prefrontal cortex, suggesting a more flexible semantic
 1225 representation of actions in these cortices, compared to occipitotemporal AON nodes. Results can be explored via an
 1226 interactive brain viewer at http://www.icon.bilkent.edu.tr/brainviewer/shahdloo_etal/.

1227

1228 **Figure 12. Difference in tuning shift for target, versus non-target categories.** The difference between absolute
 1229 values of TSI_{all} and TSI_{nt} were calculated in individual ROIs. TSI_{all} is significantly larger than TSI_{nt} in all areas with
 1230 significant tuning shift.

1231

1232 **Figure 13. Fraction of the overall tuning shifts.** Fraction of the overall tuning shifts explained by shifts in tuning
 1233 for target categories (mean±sem across subjects) and nontarget categories (i.e., excluding the union of
 1234 communication and locomotion categories) is shown. Target categories explain a greater portion of the overall

1235 tuning shifts broadly across ROIs, except for early retinotopic areas. At the same time, nontarget categories
 1236 significantly contribute to the overall tuning shifts.

1237
 1238

1239 **Figure 14. Interaction of tuning shifts with intrinsic selectivity for individual targets.** To examine the
 1240 interaction between tuning shifts and the intrinsic selectivity for individual targets, separate target preference indices
 1241 (PI) were calculated during search for *communication* (PI_{com}), and *locomotion* (PI_{loc}) categories. PI during search for
 1242 a specific target action was taken as the difference in selectivity for the target versus distractor during search for that
 1243 target. PI_{com} and PI_{loc} values are shown following projection onto the standard brain template (see Figs. 11-1c to 11-
 1244 5c for data in individual subjects). A two-dimensional colourmap was used to annotate each voxel based on PI_{com}
 1245 and PI_{loc} values (see legend). Figure format is identical to Fig. 3. AON is outlined by green dashed lines. Semantic
 1246 tuning in voxels across posterior parietal and anterior prefrontal cortices shift toward the attended category
 1247 irrespective of the search target. However, tuning in many voxels in anterior parietal, occipital, and cingulate
 1248 cortices shift toward the attended category only during search for communication or only during search for
 1249 locomotion actions.

1250

1251 **Figure 15. Attentional tuning changes in regions of interest.** Average (a) TSI_{all} , (b) TSI_{nt} , (c) PI_{com} , and (d) PI_{loc}
 1252 values were examined in cortical areas (mean±sem across five subjects). Significant values are denoted by green
 1253 bars and grey bars denote non-significant values (bootstrap test, $q(FDR)>0.05$). Values for individual subjects are
 1254 indicated by dots. Grey dots show values in areas with non-significant mean, green dots show non-significant values
 1255 in areas with significant mean, and green crosses show significant values in areas with significant mean. Tuning
 1256 shift is significantly greater than zero in many regions across all levels of the AON including occipitotemporal
 1257 cortex (pSTS, pMTG), posterior parietal cortex (IPS, AG, SMG), and premotor cortex (BA44, BA45), and in
 1258 regions across prefrontal and cingulate cortices (SFG, ACC). Compared to occipitotemporal areas, attention more
 1259 diversely modulates semantic representations in parietal and premotor AON nodes, manifested as significantly
 1260 positive tuning shift for nontarget categories in posterior parietal cortex (AG, SMG) and anterior inferior frontal
 1261 cortex (BA45). PI_{com} is significantly greater than zero in BA44/45, SFG, and ACC. In contrast, PI_{loc} is significantly
 1262 greater than zero in IPS and AG and is significantly less than zero in dPMC. Both PI_{com} and PI_{loc} are significantly
 1263 greater than zero in pSTS, pMTG, SMG, and MFG. Tuning shifts interact with the attention task, and with intrinsic
 1264 selectivity of cortical areas for target action categories.

1265

1266 **Figure 2-1. Stimulus action categories.** Actions belonging to the two target categories, and non-target categories
 1267 are specified. Occurrence frequency of each action, calculated as the number of movie frames where a given action
 1268 was present, is indicated in parentheses.

1269

1270 **Figure 6-1. Cortical flat maps of semantic representation for subject S1.** Action category responses during a.
 1271 passive viewing, b. search for *communication*, and c. search for *locomotion* categories were projected onto the
 1272 semantic space in subject S1. A two-dimensional colourmap was used to colour each voxel based on the projection
 1273 values along the first and third semantic dimensions (see colour legend). Voxels where the category model does not
 1274 explain unique response variance after accounting for low- and intermediate-level stimulus features are masked
 1275 (bootstrap test, $q(FDR)<0.05$). Regions of interest are illustrated by white borders. Several important sulci are
 1276 illustrated by dashed grey lines. Abbreviations for regions of interest and sulci are listed in *Materials and Methods*.
 1277 Many voxels across occipitotemporal, parietal, and prefrontal cortices shift their tuning toward targets.

1278

1279 **Figure 6-2. Cortical flat maps of semantic representation for subject S2.** Action category responses during a.
 1280 passive viewing, b. search for *communication*, and c. search for *locomotion* categories were projected onto the
 1281 semantic space in subject S2. Formatting is identical to Fig. 6-1.

1282

1283 **Figure 6-3. Cortical flat maps of semantic representation for subject S3.** Action category responses during a.
 1284 passive viewing, b. search for *communication*, and c. search for *locomotion* categories were projected onto the
 1285 semantic space in subject S3. Formatting is identical to Fig. 6-1.

1286

1287 **Figure 6-4. Cortical flat maps of semantic representation for subject S4.** Action category responses during **a.**
1288 passive viewing, **b.** search for *communication*, and **c.** search for *locomotion* categories were projected onto the
1289 semantic space in subject S4. Formatting is identical to Fig. 6-1.

1290
1291 **Figure 6-5. Cortical flat maps of semantic representation for subject S5.** Action category responses during **a.**
1292 passive viewing, **b.** search for *communication*, and **c.** search for *locomotion* categories were projected onto the
1293 semantic space in subject S5. Formatting is identical to Fig. 6-1.

1294
1295 **Figure 11-1. Cortical flat maps of TSI and PI for subject S1.** **a.** Tuning shift index for all action categories
1296 (TSI_{all}), **b.** tuning shift for nontarget categories (TSI_{nt}), and **c.** preference index values (PI_{com} , PI_{loc}) were projected
1297 onto the semantic space in subject S1 (see legends in Fig. 11 and 14 for the colour map). Only significant voxels are
1298 shown (bootstrap test, $q(FDR) < 0.05$). Formatting is identical to Fig. 6-1.

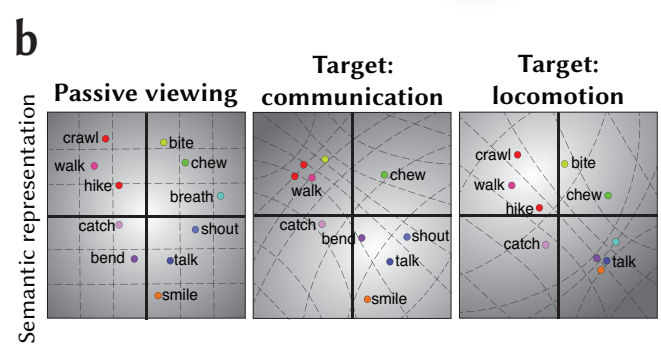
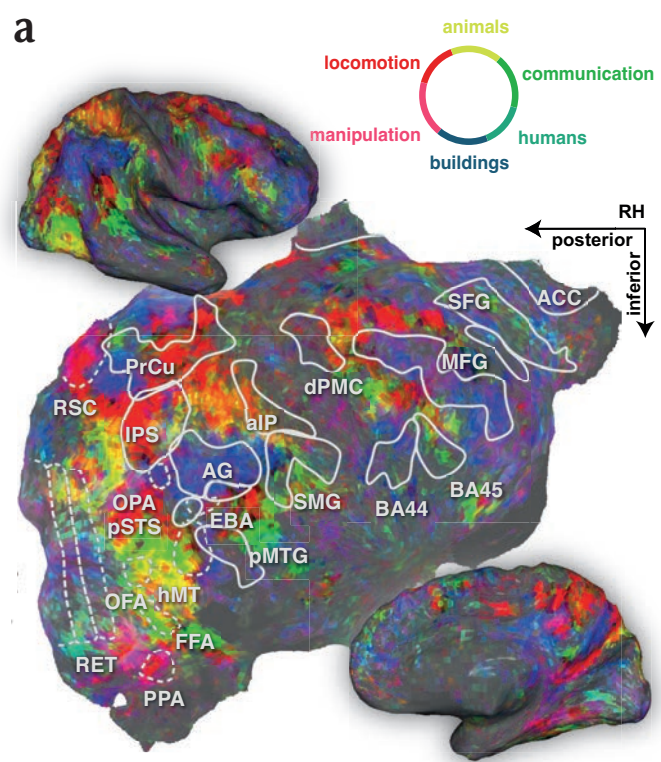
1299
1300 **Figure 11-2. Cortical flat maps of TSI and PI for subject S2.** **a.** Tuning shift index for all action categories
1301 (TSI_{all}), **b.** tuning shift for nontarget categories (TSI_{nt}), and **c.** preference index values (PI_{com} , PI_{loc}) were projected
1302 onto the semantic space in subject S2 (see legends in Fig. 11 and 14 for the colour map). Only significant voxels are
1303 shown (bootstrap test, $q(FDR) < 0.05$). Formatting is identical to Fig. 6-1.

1304
1305 **Figure 11-3. Cortical flat maps of TSI and PI for subject S3.** **a.** Tuning shift index for all action categories
1306 (TSI_{all}), **b.** tuning shift for nontarget categories (TSI_{nt}), and **c.** preference index values (PI_{com} , PI_{loc}) were projected
1307 onto the semantic space in subject S3 (see legends in Fig. 11 and 14 for the colour map). Only significant voxels are
1308 shown (bootstrap test, $q(FDR) < 0.05$). Formatting is identical to Fig. 6-1.

1309
1310 **Figure 11-4. Cortical flat maps of TSI and PI for subject S4.** **a.** Tuning shift index for all action categories
1311 (TSI_{all}), **b.** tuning shift for nontarget categories (TSI_{nt}), and **c.** preference index values (PI_{com} , PI_{loc}) were projected
1312 onto the semantic space in subject S4 (see legends in Fig. 11 and 14 for the colour map). Only significant voxels are
1313 shown (bootstrap test, $q(FDR) < 0.05$). Formatting is identical to Fig. 6-1.

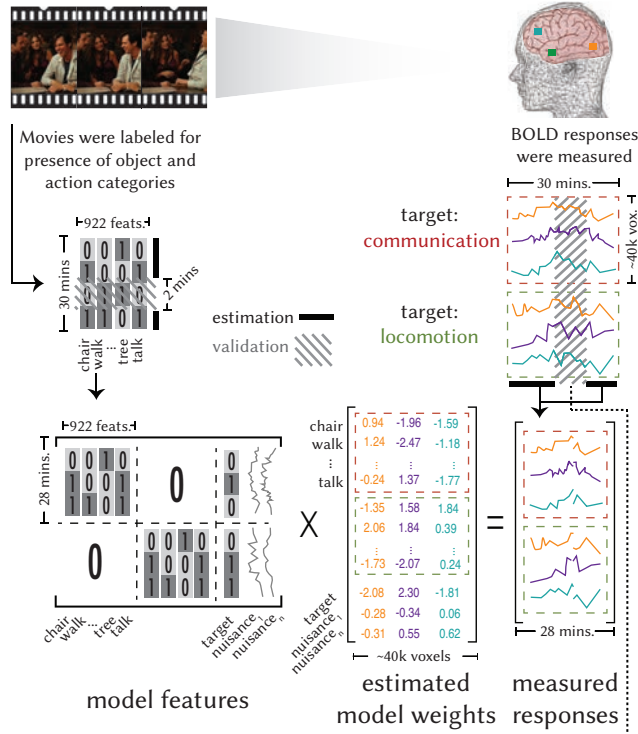
1314
1315 **Figure 11-5. Cortical flat maps of TSI and PI for subject S5.** **a.** Tuning shift index for all action categories
1316 (TSI_{all}), **b.** tuning shift for nontarget categories (TSI_{nt}), and **c.** preference index values (PI_{com} , PI_{loc}) were projected
1317 onto the semantic space in subject S5 (see legends in Fig. 11 and 14 the colour map). Only significant voxels are
1318 shown (bootstrap test, $q(FDR) < 0.05$). Formatting is identical to Fig. 6-1.

1319



a Estimating voxelwise models

Subjects viewed natural movies
and attended to *communication* or *locomotion* actions



b Validating the fit models

Estimated models were cross-validated,
after discarding the nuisance regressors
and the target regressor

