1    A fast, invariant representation for human action in the visual system

2    Leyla Isik*, Andrea Tacchetti*, and Tomaso Poggio

3    Center for Brains, Minds, and Machines, MIT

4

5    Corresponding author: Leyla Isik

6    77 Massachusetts Avenue

7    Bldg 46-4141D

8    Cambridge, MA 02139

9    617.258.6933

10   lisik@mit.edu

11

12

* These authors contributed equally to this work

13 **Abstract**

14　Humans can effortlessly recognize others' actions in the presence of complex transformations,

15　such as changes in viewpoint. Several studies have located the regions in the brain involved in

16　invariant action recognition, however, the underlying neural computations remain poorly

17　understood. We use magnetoencephalography (MEG) decoding and a dataset of well-

18　controlled, naturalistic videos of five actions (run, walk, jump, eat, drink) performed by different

19　actors at different viewpoints to study the computational steps used to recognize actions across

20　complex transformations. In particular, we ask when the brain discriminates between different

21　actions, and when it does so in a manner that is invariant to changes in 3D viewpoint. We

22　measure the latency difference between invariant and non-invariant action decoding when

23　subjects view full videos as well as form-depleted and motion-depleted stimuli. We were unable

24　to detect a difference in decoding latency or temporal profile between invariant and non-

25　invariant action recognition in full videos. However, when either form or motion information is

26　removed from the stimulus set, we observe a decrease and delay in invariant action decoding.

27　Our results suggest that the brain recognizes actions and builds invariance to complex

28　transformations at the same time, and that both form and motion information are crucial for fast,

29　invariant action recognition.

30 **New and Noteworthy**

31 The human brain can quickly recognize actions despite transformations that change their visual

32 appearance. We use neural timing data to uncover the computations underlying this ability. We

33 find that within 200ms action can be read out of MEG data, and that this representation is

34 invariant to changes in viewpoint. We find form and motion are needed for this fast action

35 decoding, suggesting that the brain quickly integrates complex spatiotemporal features to form

36 invariant action representations.

37 **Keywords:** Action recognition, Magnetoencephalography, Neural decoding, Vision

**Introduction**

As a social species, humans rely on recognizing the actions of others in their everyday lives. We quickly and effortlessly extract action information from rich dynamic stimuli, despite variations in the visual appearance of action sequences, due to transformations such as changes in size, position, actor, and viewpoint (e.g., is this person running or walking towards me, regardless of which direction they are coming from). The ability to recognize actions, the middle ground between action primitives (e.g., raise the left foot and move it forward) and activities (e.g., playing basketball) (Moeslund and Granum 2001), is paramount to humans' social interactions and even survival. The computations driving this process, however, are poorly understood. This lack of computational understanding is evidenced by the fact that even state of the art computer vision algorithms, convolutional neural networks, which match human performance on object recognition tasks (He et al. 2015), still drastically underperform humans on action recognition tasks (Le et al. 2011; Karpathy et al. 2014). In particular, what makes action and other visual recognition problems challenging are transformations (such as changes in scale, position and 3D viewpoint) that alter the visual appearance of actions, but are orthogonal to the recognition task (DiCarlo and Cox 2007).

Several studies have attempted to locate the regions in the brain involved in processing actions, and in some cases, locate regions in the brain containing viewpoint-invariant representations. In humans and nonhuman primates, the extrastriate body area (EBA) has been implicated in recognizing human form and action (Downing et al. 2001; Michels et al. 2005; Lingnau and Downing 2015), and the superior temporal sulcus (STS) has been implicated in recognizing biological motion and action (Perrett et al. 1985; Oram and Perrett 1996; Grossman et al. 2000; Vaina et al. 2001; Grossman and Blake 2002; Beauchamp et al. 2003; Peelen and Downing 2005; Vangeneugden et al. 2009). The posterior portion of the STS (pSTS) represents particular types of biological motion data in a viewpoint invariant manner (Grossman et al. 2010;

63     Vangeneugden et al. 2014). Beyond visual cortex, action representations have been found in

64     human parietal and premotor cortex when people perform and view certain actions, particularly

65     hand grasping and goal-directed behavior (analogous to monkey "mirror neuron" system)

66     (Hamilton and Grafton 2006; Dinstein, Gardner, et al. 2008; Dinstein, Thomas, et al. 2008;

67     Oosterhof et al. 2010, 2012, 2013; Freeman et al. 2013). However, recent work suggests that

68     these "mirror neuron" regions do not code the abstract, invariant representations of actions,

69     which are coded in visual regions (Wurm et al. 2015, 2016).

70         Here we investigate the neural dynamics of action processing, rather than the particular

71     brain regions involved, in order to elucidate the underlying computations. We use

72     magnetoencephalography (MEG) decoding to understand when action information is present

73     and how the brain computes representations that are invariant to complex, non-affine

74     transformations such as changes in viewpoint. Timing information can constrain the

75     computations underlying visual recognition by informing when different visual representations

76     are computed. For example, recent successes in MEG decoding have revealed interesting

77     properties about invariant object recognition in humans, mainly that it is fast and highly dynamic,

78     and that varying levels of abstract categorization and invariance increase over the first 200ms

79     following image onset (Carlson et al. 2011, 2013; Cichy et al. 2014; Isik et al. 2014).

80         Prior work has shown that biological motion can be distinguished from spatially

81     scrambled dots (Hirai et al. 2003; Hirai and Hiraki 2006; Pavlova et al. 2007) and inverted

82     figures (Jokisch et al. 2005) within 200 ms. However, it remains unknown when neural signals

83     can not only detect, but discriminate between different types of biological motion. We use timing

84     data to ask first, when the brain can discriminate between different actions, and second, when it

85     computes invariance to complex, non-affine transformations. Previous studies of invariant

86     recognition of static faces and objects suggest that 3D-viewpoint invariance develops at later

87 stages in the visual processing hierarchy (Logothetis and Sheinberg 1996; Freiwald and Tsao

88 2010; Leibo et al. 2017). Does this hold for invariant action recognition?

89      Our results show that we can read out actions as early as 200 ms after a video begins.

90 We further find that the MEG signals are already invariant to changes in viewpoint, suggesting

91 that the brain performs both action recognition and invariance at the same processing stage.

92 We further show that two types of action information, form (as tested with static images) and

93 motion (as tested with point light figures), both contribute to these immediately view-invariant

94 representations. When either form or motion information is removed, view-invariant decoding is

95 lower accuracy and delayed. These results suggest that features that are rich in form and

96 motion content drive the fast, invariant representation of the actions in the human brain.

97

98 **Materials and Methods**

99 *Action recognition dataset*

100 To study the effect of changes in view on action recognition, we used a dataset of five actors

101 performing five different actions (drink, eat, jump, run and walk) on a treadmill from two different

102 views (0 and 90 degrees from the front of the actor/treadmill; the treadmill rather than the

103 camera was rotated in place to film from different viewpoints) [Figure 1] (Tacchetti et al. 2016).

104 These actions were selected to be highly familiar, and thus something subjects would have

105 experienced under many viewing conditions, to include both reaching-oriented (eat and drink)

106 and leg-oriented (jump, run, walk) actions, as well as to span both coarse (eat and drink versus

107 run and walk) and fine (eat versus drink and run versus walk) action distinctions. Every video

108 was filmed on the same background, and the same objects were present in each video,

109 regardless of action (e.g., to avoid confounds such as "run" being detected based on the

110 presence of a treadmill and "drink" being detected based on the presence of a water bottle).

111 Each action-actor-view combination was filmed for at least 52-seconds. The videos were then

112 cut into two-second clips that each included at least one cycle of each action, and started at

113 random points in the cycle (for example, a jump may start midair or on the ground). This dataset

114 allows testing of actor and view invariant action recognition, with few low-level confounds.

115     To explore the roles of form and motion in invariant action representations, we extended

116 this video dataset with two additional components: a form only dataset, consisting of

117 representative single frames for each action, and a motion-only dataset, consisting of point light

118 figures performing the same actions. For the form dataset, the authors selected one frame per

119 video making sure that the selected frames were unambiguous for action identity (special

120 attention was paid to the actions eat and drink to ensure the food or drink was near the mouth,

121 and occluded views to ensure there was some visual information about action). For the motion

122 point light dataset, the videos were put on Amazon Mechanical Turk and workers were asked to

123 label 15 landmarks in every single frame: center of head, shoulders, elbows, hands, torso, hips,

124 knees, and ankles. Three workers labeled each video frame. We used the spatial median of the

125 three independent labels for each frame and landmark to increase the signal to noise ratio, and

126 independently low-pass filtered the time series (Gaussian Filter with a 30 frames aperture and

127 normalized convolution) for each of the 15 points to reduce the high frequency artifacts

128 introduced by single-frame labeling.

129

130 *Participants*

131 Three separate MEG experiments were conducted (see below). Ten subjects (5 female, 8 right-

132 handed, age: mean±SD = 28.6±6.1) participated in experiment one, ten subjects (7 female, 10

133 right-handed, age mean±SD = 25.2±5.0) participated experiment two, and ten subjects (7

134 female, 9 right-handed, age: mean±SD = 28.3±5.7) participated in experiment three. All subjects

135 had normal or corrected to normal vision. The MIT Committee on the Use of Humans as

136   Experimental Subjects approved the experimental protocol. Subjects provided informed written

137   consent before the experiment.

138

139   *Experimental procedure*

140   In the first experiment, we assessed if we could read out different actions both within viewpoint

141   (training and testing on videos at 0 degrees or 90 degrees, without any generalization) and

142   across viewpoint, by training and testing on two different views (0 and 90 degrees). In this

143   experiment ten subjects were shown 50 two-second video clips (one for each of five actors,

144   actions, and two views, 0 and 90 degrees), each presented 20 times.

145       To examine whether form and motion information were necessary to construct invariant

146   action representations, in the second and third experiments we showed subjects limited "form"

147   (static image) or "motion" (point-light walkers) datasets. Specifically, in the second experiment,

148   ten subjects were shown 50 static images (one for each of five actors, actions, and two views, 0

149   and 90 degrees), which were single frames from the videos in Experiment 1, for 2 seconds

150   presented 20 times each. In the third experiment, ten subjects were shown 10 two-second video

151   clips, which consisted of point-light walkers traced along one actor's videos from two views in

152   experiment one (labelled by Mechanical Turk workers as described above), presented 100

153   times each.

154       In each experiment, subjects performed an action recognition task, where they were

155   asked after a random subset of videos or images (in a randomly interspersed 10% of the trials

156   for each video or image) what action was portrayed in the previous image or video. The purpose

157   of this behavioral task was to ensure subjects were attentive and assess behavioral

158   performance on the various datasets. The button order for each action was randomized across

159   trials to avoid systematic motor confounds in the decoding. Subjects were instructed to fixate

160   centrally. The videos were presented using Psychtoolbox to ensure accurate timing of stimulus

161 onset. Each video had a duration of 2s and a 2s inter-stimulus interval. The videos were shown

162 in grayscale at 3 x 5.4 degrees of visual angle on a projector with a 48 cm × 36 cm display, 140

163 cm away from the subject.

164

165 *MEG data acquisition and preprocessing*

166 The MEG data were collected using an Elekta Neuromag Triux scanner with 306 sensors, 102

167 magnetometers at 204 planar gradiometers, and were sampled at 1000 Hz. First the signals

168 were filtered using temporal Signal Space Separation (tSSS) with Elekta Neuromag software.

169 Next, Signal Space Projection (SSP) (Tesche et al. 1995) was applied to correct for movement

170 and sensor contamination. The MEG data were divided into epochs from -500 - 3500 ms,

171 relative to video onset, with the mean baseline activity removed from each epoch. The signals

172 were band-pass filtered from 0.1–100 Hz to remove external and irrelevant biological noise

173 (Acunzo et al. 2012; Rousselet 2012). The convolution between signals and bandpass filter was

174 implemented by wrapping signals in a way that may introduce edge effects at the beginning and

175 end of each trial. We mitigated this issue by using a large epoch window (-500-3500 ms) and

176 testing significance in a manner that takes into account temporal biases in the data (see

177 significance testing below). The above pre-processing steps were all implemented using the

178 Brainstorm software (Tadel et al. 2011).

179

180 *General MEG decoding methods*

181 MEG decoding analyses were performed with the Neural Decoding Toolbox (Meyers 2013), a

182 Matlab package implementing neural population decoding methods. In this decoding procedure,

183 a pattern classifier was trained to associate the patterns of MEG data with the identity of the

184 action in the presented image or video. The stimulus information in the MEG signal was

185 evaluated by testing the accuracy of the classifier on a separate set of test data. This procedure

186  was conducted separately for each subject and multiple re-splits of the data into training and
187  test data were utilized.

188      The time series data of the magnetic field measured in each sensor (including both the
189  magnetometers and gradiometers) were used as classifier features. We averaged the data in
190  each sensor into 100 ms overlapping bins with a 10 ms step size, and performed decoding
191  independently at each time point. Decoding analysis was performed using cross validation,
192  where the dataset was randomly divided into five cross validation splits. The classifier was then
193  trained on data from four splits (80% of the data), and tested on the fifth, held out split (20% of
194  the data) to assess the classifier's decoding accuracy.

195

196  *Decoding - feature pre-processing*

197      To improve signal to noise, we averaged together the ten different trials for each
198  semantic class (e.g. videos of run) in each given cross validation split of each subject's data so
199  there was one data point per stimulus per cross validation split. We next Z-score normalized
200  that data by calculating the mean and variance for each sensor using only the training data. We
201  then performed sensor selection using only the training data, by applying a five-way ANOVA to
202  each sensor's training data to test if the sensor was selective for the different actions. We use
203  sensors that were selective for action identity, i.e., show a significantly greater variation across
204  class than within class, with $p<0.05$ significance based on a F-test (if no sensors were deemed
205  significant, the one with the lowest p-value is selected). The selected sensors were then fixed
206  and used for testing. To avoid circularity in our feature pre-processing, the test data was never
207  used for the z-scoring or feature selection.

208      Each sensor (including both magnetometers and gradiometers) was considered as an
209  independent sensor input into this algorithm, and the feature selection, like the other decoding
210  steps is performed separately at each 100ms time bin, and thus a different number of sensors

9

211  was selected for each subject at each time bin. The average number of sensors selected for

212  each subject across all significant decoding time bins is shown in Table 1.  These pre-

213  processing parameters have been shown to empirically improve MEG decoding signal to noise

214  in a previous MEG decoding study (Isik et al. 2014), however as we did not use absolute

215  decoding performance (rather significantly above chance decoding) as a metric for when

216  information is present in the MEG signals, we did not further optimize decoding performance

217  with the present data.

218

219  *Decoding - classification*

220      The pre-processed MEG data was then input into the classifier. Decoding analyses were

221  performed using a maximum correlation coefficient classifier, which computed the correlation

222  between each test vector and a mean training vector that is created from taking the mean of the

223  training data from a given class. Each test point was assigned the label of the class of the

224  training data with which it was maximally correlated. When we refer to classifier "training" this

225  could alternatively be thought of as learning to discriminate patterns of electrode activity

226  between the different classes in the training data, rather than a more involved training procedure

227  with a more complex classifier. We intentionally chose a very simple algorithm to see in the

228  simplest terms what information is coded in the MEG data. Prior work has also shown

229  empirically that results with a correlation coefficient classifier are very similar to standard linear

230  classifiers like support vector machines (SVMs) or regularized least squares (RLS) (Isik et al.

231  2014).

232      We repeated the above decoding procedure at each time bin to assess the decoding

233  accuracy versus time. We re-ran the above procedure 50 times for each subject. We measured

234  decoding accuracy as the average percent correct of the test set data across all decoding runs,

235 and reported decoding results for the average of ten subjects in each experiment. Plots and

236 latency measures were centered at the median value of each of the 100ms time bins.

237 For more details on these decoding methods see (Isik et al. 2014).

238

239 *Decoding invariant information*

240 To see if information in the MEG signals could generalize across a given transformation,

241 we trained the classifier on data from subjects viewing the stimuli under one condition (e.g. 0-

242 degree view) and tested the classifier on data from subjects viewing the stimuli under a

243 separate, held out condition (e.g. 90-degree view). This provided a strong test of invariance to a

244 given transformation. In all three experiments, we compared the within and across view

245 decoding. For the "within" view case, the classifier was trained on 80% of data from one view,

246 and tested on the remaining 20% of data from the same view. For the "across" view case, the

247 classifier was trained on 80% of data from one view, and tested on 20% of data from the

248 opposite view, so the same amount of training and test data was evaluated in each case.

249

250 *Significance testing*

251 We assessed action decoding significance using a permutation test. We ran the decoding

252 analysis for each subject with the labels randomly shuffled to create a null distribution. Shuffling

253 the labels breaks the relationship between the experimental conditions that occurred. We

254 repeated the procedure of shuffling the labels and running the decoding analysis 1000 times to

255 create a null distribution, and reported p-values as the percentage rank of the actual decoding

256 performance within the null distribution.

257 For each experiment and decoding condition, we averaged the null decoding data

258 across ten subjects and determined when the mean decoding across subjects was above the

259 mean null distribution. We define the decoding "onset time" as the first time the subject-

260  averaged decoding accuracy was greater than the subject-averaged null distribution, with $p <$

261  $0.05$. This provided a measure of when significant decodable information was first present in the

262  MEG signals, and is a standard metric to compare latencies between different conditions (Isik et

263  al. 2013; Cichy et al. 2016). Time of peak decoding accuracy for each condition, an alternative

264  established measure of decoding latency, was found to be much more variable (with 95%

265  confidence intervals that were on average over 400 ms larger than onset times), we therefore

266  restricted ourselves to using onset latency only.

267

268  *Assessing latency differences*

269  To compare when information arises in different decoding conditions (e.g. within versus

270  across view), we compared onset latency rather than raw decoding performance, because 1)

271  the raw magnitude of a classifier is difficult to interpret 2) we want to know *when* significant

272  information is present in each signal. To compare onset latencies for the within view versus

273  across view decoding, we performed 1000 bootstrap resamples of subjects and use the

274  resulting distribution to compute empirical 95%-confidence intervals (CI) for the onset latency of

275  each condition to estimate the temporal sensitivity of our measure (Hoenig and Heisey 2001),

276  as well as for the difference in onset latency between the two conditions. Specifically, in each

277  bootstrap run, we randomly selected a different subset of ten subjects with replacement,

278  computed onset latencies for each condition (as outlined above) and calculated the difference in

279  onset latency between the invariant and non-invariant conditions. We defined the onset

280  latencies for invariant and non-invariant decoding significantly different with $p<0.05$ if the

281  empirical 95% interval for their difference did not include 0 (Cichy et al. 2016).

282

283  *Temporal Cross Training*

284    Beyond decoding latency, we sought to examine the dynamics of the MEG decoding

285    using temporal-cross-training analysis (Meyers et al. 2008; Meyers 2013; Isik et al. 2014; King

286    and Dehaene 2014). In this analysis, rather than training and testing the classifier on the same

287    time point, a classifier was trained with data from one time point and then tested on data from all

288    other time points. Otherwise the decoding methods (including feature pre-processing, cross

289    validation and classification) were identical to the procedure outlined above. This method

290    yielded a matrix of decoding accuracies for each training and test time point, where the rows of

291    the matrix indicate the times when the classifier was trained, and the columns indicate the times

292    when the classifier was tested. The diagonal entries of this matrix contained the results from

293    when the classifier was trained and tested on data from the same time point (identical to the

294    procedure described above).

295

296    **Results**

297    *Readout of actions from MEG data is early and invariant*

298    Ten subjects viewed 2-second videos of five actions performed by five actors at two views (0

299    degrees and 90 degrees) (Figure 1, top row) while their neural activity was recorded in the

300    MEG. We then trained our decoding classifier on only on one view (0 degrees or 90 degrees),

301    and tested it on the second view (0 degrees or 90 degrees). We could read out action from the

302    subjects' MEG data in the case without any invariance ("within view" condition) at, on average,

303    250 ms (210-330 ms) (mean decoding onset latency across subjects based on $p < 0.05$

304    permutation test, 95% confidence intervals of onset latencies reported throughout in

305    parentheses, see Methods) post video onset (Figure 2a, blue trace). Each video began at a

306    random point in a given action sequence, suggesting that the brain can compute this

307    representation from different partial sequences of each action. We also observed a significant

308 rise in decoding after the video offset, consistent with offset responses that have been observed

309 in MEG decoding of static images (Carlson et al. 2011).

310       We next assessed if the MEG signals were invariant to changes in viewpoint by training

311 the classifier on data from subjects viewing actions performed at one view and testing on a

312 second held out view. This invariant "across-view" decoding arose on average at 230 ms (220-

313 270ms) (Figure 2a, red trace). The within and across view decoding were largely overlapping

314 (Figure 2a, insert), and their onset latencies were not significantly different (p = 0.13),

315 suggesting that the early action recognition signals are immediately view invariant. To ensure

316 that the lack of latency difference between the within and between view conditions was not due

317 to the fact that we are using 100ms overlapping time bins, we re-ran the decoding 10ms time

318 bins and 10ms step size (non-overlapping time bins). Although the overall decoding accuracy

319 was lower, the within and across view decoding onsets were still not significantly different (p =

320 0.62, Figure 2b).

321       We next examined which types of actions are decoding in both the within and across

322 decoding conditions. By analyzing the confusion matrices for the within- and across-view

323 decoding, we found that not only are coarse action distinctions made (e.g., between run/walk

324 and eat/drink), but so are fine action distinctions (e.g., between eat and drink) even at the

325 earliest decoding of 250 ms (Figure 3). Further, actions performed in a familiar context (i.e. run

326 and walk on a treadmill) were not better classified than those performed in an unfamiliar context

327 (i.e. eat and drink on a treadmill).

328

329 *The dynamics of invariant action recognition*

330 Given that the within- and across-view action decoding conditions had similar onset latencies,

331 we further compared the temporal profiles of the two conditions by asking if the neural codes for

332 each condition were stable over time. To test this, we trained our classifier with data at one time

14

333 point, and tested the classifier at all other time points. This yielded a matrix of decoding

334 accuracies for different train times by test times, referred to as a temporal cross training (TCT)

335 matrix (Meyers et al. 2008; Carlson et al. 2013; Meyers 2013; Isik et al. 2014). The diagonal of

336 this matrix shows when the classifier is trained and tested with data at the same time point, just

337 as the line plots in Figure 2a.

338       The within-view and across-view TCTs showed that the representations for actions, both

339 with and without view, are highly dynamic as there is little off-diagonal decoding that is

340 significantly above chance (Figure 4a-b). The window of significantly above chance decoding

341 performance from 200-400 ms, in particular, is highly dynamic and decoding only within a 50-

342 100 ms window is significantly above chance. At later time points, the above chance decoding

343 extends to a larger window that spans 300ms, suggesting the late representations for action are

344 more stable across time than the early representations. Further, we find that significant

345 decoding for the within and across view conditions were largely overlapping (Figure 4c) showing

346 that information for both conditions are represented at the same time scale in the MEG data.

347

348 *Invariant action recognition is impaired in form- and motion-depleted stimuli*

349 To study the roles of two information streams, form and motion, in action recognition, subjects

350 viewed two limited stimulus-sets in the MEG. The first 'Form' stimulus set consisted of one static

351 frame from each video (containing no motion information). The second 'Motion' stimulus set,

352 consisted of point light figures that are comprised of dots on each actor's head, arm joints, torso,

353 and leg joints and move with the actor's joints (containing limited form information) (Johansson

354 1973). Ten subjects viewed each of the form and motion datasets in the MEG. We could decode

355 action from both datasets in the within view case without any invariance (Figure 5). The early

356 view-invariant decoding that was observed with full movies, however, was impaired for both the

357 form or motion datasets. In the form-only experiment, within view could be read out at 410 ms

358    (320- ms) and across view at 510ms (430- ms). The onset latencies of 410 ms and 510 ms are

359    the first significantly above chance time points for the average decoding across all ten subjects.

360    Although the average decoding across all ten subjects was significantly above chance, in more

361    than 5% of bootstrap runs (each randomly selecting a *different* subset of ten subjects with

362    replacement, see Methods), the decoding was not significantly above chance. Since we could

363    not calculate a significant onset time in the bootsrap runs that did not reach significantly above

364    chance decoding, the upper limit of the 95% CI for both the within and across view decoding is

365    missing and we did not detect a significant difference between the two conditions. In the motion-

366    only experiment, within view action information could be read out significantly earlier than

367    across view information: 210 ms (180-260 ms) versus 300 ms (300-510 ms), and was

368    significantly different between the two conditions (p = 0.013).

369

370    **Discussion**

371    We investigated the dynamics of invariant action recognition in the human brain and

372    found that action can be decoded from MEG signals as early as 200 ms post video onset,

373    considerably less than the 2s duration of each video and most action cycles (e.g., one drink

374    from a water bottle). This latency is similar to that found for biological motion detection in

375    evoked responses (Hirai et al. 2003; Jokisch et al. 2005; Hirai and Hiraki 2006; Pavlova et al.

376    2007). These results are also consistent with a recent MEG decoding study that classified two

377    actions, reaching and grasping, slightly after 200ms post video onset (Tucciarelli et al. 2015).

378    Crucially, we showed that these early neural signals are selective to a variety of full-body

379    actions as well as invariant to changes in 3-D viewpoint.

380    Interestingly we do not observe a difference in onset latency between invariant and non-

381    invariant action representations. While we cannot completely rule out differences at a finer scale

382    than we can resolve with our methods, this appears to be different than object recognition.

16

383    Invariant object information increases along subsequent layers of the ventral stream (Logothetis

384    and Sheinberg 1996; Rust and Dicarlo 2010) causing a delay in invariant decoding relative to

385    non-invariant decoding (Isik et al. 2014). Further, physiology data (Freiwald and Tsao 2010) and

386    computational models (Leibo et al. 2017) of static face recognition have shown that invariance

387    to 3D viewpoint, in particular, arises at a later processing stage than initial face recognition. One

388    possible account of this discrepancy is that even non-invariant ("within view") action

389    representations rely on higher-level visual features (that carry some degree of viewpoint

390    invariant information) than those used in basic object representations.

391        We characterized the dynamics of action representations using temporal cross training

392    and found that the decoding windows for within and across view decoding are largely

393    overlapping (Figure 4c), suggesting that the beyond onset latencies, the overall dynamics of

394    decoding are similar for non-invariant and view-invariant action representations. It has been

395    suggested that visual recognition, as studied with static object recognition, has a canonical

396    temporal representation that is demonstrated by highly diagonal TCT matrices (King and

397    Dehaene 2014). Our action results generally follow this pattern (Figure 4), but they are more

398    stable over time than previously reported for object decoding (Carlson et al. 2013a; Cichy et al.

399    2014; Isik et al. 2014).

400        As shown previously, we find that people can recognize and neural signals can

401    distinguish actions with either biological motion or form information removed from the stimulus

402    (Johansson 1973; Schindler and van Gool 2008; Singer and Sheinberg 2010). In particular,

403    decoding actions within-view is largely intact when form or motion cues are removed. This is

404    likely due to the fact that within-view decoding, unlike the across-view condition, requires little

405    generalization and can thus be performed using low-level cues in the form or motion stimuli. The

406    across-view decoding, on the other hand, requires substantially more generalization and cannot

407    be performed as well, or as quickly as the within-view decoding with form or motion depleted

408    stimuli. It is important to note, however, that the three experiments were completed separately

409    with different subjects, and therefore we cannot directly compare decoding with full videos to the

410    performance with form= or motion-depleted stimuli. Further, while our datasets are a best

411    attempt to isolate form and motion information, it is important to note that static images contain

412    implied motion and that point light figures contain some form information and have less motion

413    information than full movies. Nevertheless, the low-accuracy and delayed invariant decoding

414    with either limited stimulus set suggest that both form and motion information are necessary to

415    build a robust action representation.

416        Importantly these invariant action representations cannot be explained by low-level

417    stimulus features, such as motion energy as the output of a standard motion energy model

418    (Simoncelli and Heeger 1998) cannot significantly discriminate action across viewpoint

419    (Tacchetti et al. 2016). While we cannot fully rule out the effects of eye movements or shifts in

420    covert attention, eye movement patterns cannot be accounting for our early MEG decoding

421    accuracy, because we do not observe a significant shift in the eye positions between different

422    actions until after 600 ms post video onset and further the same decoder applied to MEG

423    signals does not successfully decode action information using raw eye position data (Figure 2c).

424        The five actions tested in this study comprise only a small subset of the wide variety of

425    familiar actions we recognize in our daily lives. The five-way classification shows similar

426    decoding across between all five actions, including both coarse and fine action distinctions

427    (Figure 3a-d). These five actions were selected to be highly familiar, and thus we do not know to

428    what extent familiarity is necessary for the immediate invariance we observed. Indeed, modeling

429    and theoretical work suggest that in order to build templates to be invariant to non-affine

430    transformations such as changes in 3-D viewpoint, one must learn templates from different

431    views of each given category (Leibo et al. 2015). It remains an open question how this

18

432  invariance would translate to unfamiliar actions and how many examples would be needed to

433  learn invariant representations of new actions.

434      Finally, the longer latency and greater cross-temporal stability of action decoding raises

435  the question of whether recurrent and feedback connections are used to form invariant action

436  representations. This is difficult to test explicitly without high spatiotemporal resolution data. It is

437  indeed likely that feedback and recurrent connections occur within the 200 ms of our earliest

438  decoding (Lamme and Roelfsema 2000). However, further studies have shown that purely

439  feedforward computational models can discriminate actions invariant to viewpoint, and produce

440  representations that explain a significant amount of variance in the human MEG data (Tacchetti

441  et al. 2016).

442      Taken as a whole, our results show that the brain computes action selective

443  representations remarkably quickly and, unlike in the recognition of static faces and objects, at

444  the same time that it computes invariance to non-affine transformations that are orthogonal to

445  the recognition task. This may represent a key difference between action and object visual

446  processing. Moreover, our findings suggest that both form and motion information are

447  necessary to construct these fast invariant representations of human action sequences. The

448  methods and results presented here provide a framework to study the dynamic neural

449  representations evoked by natural videos, and open the door to probing neural representations

450  for higher level visual and social information conveyed by video stimuli.

451

452  **Acknowledgements**

459 **References**

460

461 Acunzo DJ, Mackenzie G, van Rossum MCW. 2012. Systematic biases in early ERP and ERF

462     components as a result of high-pass filtering. J Neurosci Methods. 209:212–218.

463 Beauchamp MS, Lee KE, Haxby J V, Martin A. 2003. FMRI responses to video and point-light

464     displays of moving humans and manipulable objects. J Cogn Neurosci. 15:991–1001.

465 Carlson TA, Hogendoorn H, Kanai R, Mesik J, Turret J. 2011. High temporal resolution

466     decoding of object position and category. J Vis. 11.

467 Carlson T, Tovar DA, Alink A, Kriegeskorte N. 2013. Representational dynamics of object vision:

468     the first 1000 ms. J Vis. 13:1-.

469 Cichy RM, Pantazis D, Oliva A. 2014. Resolving human object recognition in space and time.

470     Nat Neurosci. 17:455–462.

471 Cichy RM, Pantazis D, Oliva A. 2016. Similarity-Based Fusion of MEG and fMRI Reveals

472     Spatio-Temporal Dynamics in Human Cortex During Visual Object Recognition. Cereb

473     Cortex. bhw135.

474 DiCarlo JJ, Cox DD. 2007. Untangling invariant object recognition. Trends Cogn Sci. 11:333–

475     341.

476 Dinstein I, Gardner JL, Jazayeri M, Heeger DJ. 2008. Executed and Observed Movements

477     Have Different Distributed Representations in Human aIPS. J Neurosci. 28:11231–11239.

478 Dinstein I, Thomas C, Behrmann M, Heeger DJ. 2008. A mirror up to nature. Curr Biol. 18:R13-

479     8.

480 Downing PE, Jiang Y, Shuman M, Kanwisher N, Downing PE, Jiang Y, Jiang Y, Shuman M,

481     Shuman M, Kanwisher N, Kanwisher N. 2001. A cortical area selective for visual

482      processing of the human body. Science. 293:2470–2473.

483    Freeman J, Ziemba CM, Heeger DJ, Simoncelli EP, Movshon JA. 2013. A functional and

484      perceptual signature of the second visual area in primates. Nat Neurosci. 16:974–981.

485    Freiwald WA, Tsao DY. 2010. Functional compartmentalization and viewpoint generalization

486      within the macaque face-processing system. Science. 330:845–851.

487    Grossman E, Donnelly M, Price R, Pickens D, Morgan V, Neighbor G, Blake R. 2000. Brain

488      Areas Involved in Perception of Biological Motion. J Cogn Neurosci. 12:711–720.

489    Grossman ED, Blake R. 2002. Brain Areas Active during Visual Perception of Biological Motion.

490      Neuron. 35:1167–1175.

491    Grossman ED, Jardine NL, Pyles JA. 2010. fMR-Adaptation Reveals Invariant Coding of

492      Biological Motion on the Human STS. Front Hum Neurosci. 4:15.

493    Hamilton AF d. C, Grafton ST. 2006. Goal Representation in Human Anterior Intraparietal

494      Sulcus. J Neurosci. 26:1133–1137.

495    He K, Zhang X, Ren S, Sun J. 2015. Delving Deep into Rectifiers: Surpassing Human-Level

496      Performance on ImageNet Classification.

497    Hirai M, Fukushima H, Hiraki. 2003. An event-related potentials study of biological motion

498      perception in humans. Neurosci Lett. 344:41–44.

499    Hirai M, Hiraki K. 2006. The relative importance of spatial versus temporal structure in the

500      perception of biological motion: An event-related potential study. Cognition. 99:B15–B29.

501    Hoenig JM, Heisey DM. 2001. The Abuse of Power. Am Stat. 55:19–24.

502    Isik L, Meyers EM, Leibo JZ, Poggio T. 2014. The dynamics of invariant object recognition in the

503      human visual system. J Neurophysiol. 111:91–102.

504  Isik L, Meyers EM, Leibo JZ, Poggio TA. 2013. The dynamics of invariant object recognition in

505      the human visual system. J Neurophysiol. in press.

506  Johansson G. 1973. Visual perception of biological motion and a model for its analysis. Percept

507      Psychophys. 14:201–211.

508  Jokisch D, Daum I, Suchan B, Troje NF. 2005. Structural encoding and recognition of biological

509      motion: evidence from event-related potentials and source analysis. Behav Brain Res.

510      157:195–204.

511  Karpathy A, Toderici G, Shetty S, Leung T, Sukthankar R, Fei-Fei L. 2014. Large-Scale Video

512      Classification with Convolutional Neural Networks. In: 2014 IEEE Conference on Computer

513      Vision and Pattern Recognition. IEEE. p. 1725–1732.

514  King J-R, Dehaene S. 2014. Characterizing the dynamics of mental representations: the

515      temporal generalization method. Trends Cogn Sci. 18:203–210.

516  Lamme VAF, Roelfsema PR. 2000. The distinct modes of vision offered by feedforward and

517      recurrent processing. Trends Neurosci. 23:571–579.

518  Le Q V., Zou WY, Yeung SY, Ng AY. 2011. Learning hierarchical invariant spatio-temporal

519      features for action recognition with independent subspace analysis. In: CVPR 2011. IEEE.

520      p. 3361–3368.

521  Leibo JZ, Liao Q, Anselmi F, Freiwald WA, Poggio T. 2017. View-Tolerant Face Recognition

522      and Hebbian Learning Imply Mirror-Symmetric Neural Tuning to Head Orientation. Curr

523      Biol. 27:62–67.

524  Leibo JZ, Liao Q, Anselmi F, Poggio TA. 2015. The Invariance Hypothesis Implies Domain-

525      Specific Regions in Visual Cortex. PLOS Comput Biol. 11:e1004390.

526  Lingnau A, Downing PE. 2015. The lateral occipitotemporal cortex in action. Trends Cogn Sci.

527   Logothetis NK, Sheinberg DL. 1996. Visual object recognition. Annu Rev Neurosci. 19:577–621.

528   Meyers EM. 2013. The neural decoding toolbox. Front Neuroinform. 7.

529   Meyers EM, Freedman DJ, Kreiman G, Miller EK, Poggio T. 2008. Dynamic population coding

530       of category information in inferior temporal and prefrontal cortex. J Neurophysiol.

531       100:1407–1419.

532   Michels L, Lappe M, Vaina LM. 2005. Visual areas involved in the perception of human

533       movement from dynamic form analysis. Neuroreport. 16:1037–1041.

534   Moeslund TB, Granum E. 2001. A Survey of Computer Vision-Based Human Motion Capture.

535       Comput Vis Image Underst. 81:231–268.

536   Oosterhof NN, Tipper SP, Downing PE. 2012. Viewpoint (in)dependence of action

537       representations: an MVPA study. J Cogn Neurosci. 24:975–989.

538   Oosterhof NN, Tipper SP, Downing PE. 2013. Crossmodal and action-specific: neuroimaging

539       the human mirror neuron system. Trends Cogn Sci. 17:311–318.

540   Oosterhof NN, Wiggett AJ, Diedrichsen J, Tipper SP, Downing PE. 2010. Surface-Based

541       Information Mapping Reveals Crossmodal Vision–Action Representations in Human

542       Parietal and Occipitotemporal Cortex. J Neurophysiol. 104.

543   Oram MW, Perrett DI. 1996. Integration of form and motion in the anterior superior temporal

544       polysensory area (STPa) of the macaque monkey. J Neurophysiol. 76:109–129.

545   Pavlova M, Lutzenberger W, Sokolov AN, Birbaumer N, Krägeloh-Mann I. 2007. Oscillatory

546       MEG response to human locomotion is modulated by periventricular lesions. Neuroimage.

547       35:1256–1263.

548   Peelen M V, Downing PE. 2005. Selectivity for the human body in the fusiform gyrus. J

549       Neurophysiol. 93:603–608.

550 Perrett DI, Smith PAJ, Mistlin AJ, Chitty AJ, Head AS, Potter DD, Broennimann R, Milner AD,

551     Jeeves MA. 1985. Visual analysis of body movements by neurones in the temporal cortex

552     of the macaque monkey: A preliminary report. Behav Brain Res. 16:153–170.

553 Rousselet GA. 2012. Does Filtering Preclude Us from Studying ERP Time-Courses? Front

554     Psychol. 3:131.

555 Rust NC, Dicarlo JJ. 2010. Selectivity and tolerance ("invariance") both increase as visual

556     information propagates from cortical area V4 to IT. J Neurosci. 30:12978–12995.

557 Schindler K, van Gool L. 2008. Action snippets: How many frames does human action

558     recognition require? In: 2008 IEEE Conference on Computer Vision and Pattern

559     Recognition. IEEE. p. 1–8.

560 Simoncelli EP, Heeger DJ. 1998. A model of neuronal responses in visual area MT. Vision Res.

561     38:743–761.

562 Singer JM, Sheinberg DL. 2010. Temporal cortex neurons encode articulated actions as slow

563     sequences of integrated poses. J Neurosci. 30:3133–3145.

564 Tacchetti A, Isik L, Poggio T. 2016. Spatio-temporal convolutional neural networks explain

565     human neural representations of action recognition.

566 Tadel F, Baillet S, Mosher JC, Pantazis D, Leahy RM. 2011. Brainstorm: a user-friendly

567     application for MEG/EEG analysis. Comput Intell Neurosci. 2011:879716.

568 Tesche CD, Uusitalo MA, Ilmoniemi RJ, Huotilainen M, Kajola M, Salonen O. 1995. Signal-

569     space projections of MEG data characterize both distributed and well-localized neuronal

570     sources. Electroencephalogr Clin Neurophysiol. 95:189–200.

571 Tucciarelli R, Turella L, Oosterhof NN, Weisz N, Lingnau A. 2015. MEG Multivariate Analysis

572     Reveals Early Abstract Action Representations in the Lateral Occipitotemporal Cortex. J

573       Neurosci. 35:16034–16045.

574    Vaina LM, Solomon J, Chowdhury S, Sinha P, Belliveau JW. 2001. Functional neuroanatomy of

575       biological motion perception in humans. Proc Natl Acad Sci U S A. 98:11656–11661.

576    Vangeneugden J, Peelen M V, Tadin D, Battelli L. 2014. Distinct neural mechanisms for body

577       form and body motion discriminations. J Neurosci. 34:574–585.

578    Vangeneugden J, Pollick F, Vogels R. 2009. Functional differentiation of macaque visual

579       temporal cortical neurons using a parametric action space. Cereb Cortex. 19:593–611.

580    Wurm MF, Ariani G, Greenlee MW, Lingnau A. 2016. Decoding Concrete and Abstract Action

581       Representations During Explicit and Implicit Conceptual Processing. Cereb Cortex.

582       26:3390–3401.

583    Wurm MF, Lingnau A, Wurm XMF, Lingnau A. 2015. Decoding Actions at Different Levels of

584       Abstraction. J Neurosci. 35:7727–7735.

585

586

587 **Figure legends**

588 ***Figure 1 – Action recognition dataset***

589 *(a) We used a dataset of two-second videos depicting five actors performing five actions from*

590 *five viewpoints. Frames from one example walk video at 90 degrees (top) and one example*

591 *drink video at 0 degrees (bottom) are shown. We extended this dataset to (b) a "Form only"*

592 *dataset, containing single (action informative) frames from each two-second movie, and (c) a*

593 *"Motion only" dataset of point light videos created by labeling joints on actors in each video (a,*

594 *bottom).*

595

596 ***Figure 2 – Action decoding from video data***

597 ***(a-b) Within and across view action decoding from MEG data.*** *We can decode action by*

598 *training and testing a simple on the same view ('within-view' condition), or, to assess viewpoint*

599 *invariance, training on one view (0 degrees or 90 degrees) and testing on second view ('across*

600 *view' condition), in (a) 100 ms overlapping bins (10 ms step size), or (b) 10 ms non-overlapping*

601 *bins. Results are from the average of ten subjects. Error bars represent standard error across*

602 *subjects. Horizontal line indicates chance decoding accuracy. Line at bottom of plot indicates*

603 *group-level significance with p<0.05 permutation test, for the average null distribution across the*

604 *ten subjects. The first time point in this line is the onset time for each condition, reported in the*

605 *main text. Inset shows a zoom of decoding time courses from 175-525 ms post-video onset. (c)*

606 ***Action decoding from eye tracking data.*** *We trained a linear classifier on the output of*

607 *eyetracking data from a separate experiment. We trained the classifier with 80% of the data*

608 *from all views, and tested on the 20% of held out data. Decoding methods are otherwise*

609 *analgous to the MEG decoding procedure Results are from the average of five different*

610 *subjects. Error bars represent standard error across subjects. Horizontal line indicates chance*

611   *decoding (20%). Decoding does not pass the group-level significance threshold of p<0.05 as*

612   *determined by a permutation test.*

613

614   ***Figure 3 – Confusion matrices for action video dataset.*** *Confusion matrices for the within*

615   *and across view decoding conditions in the video dataset for **(a)** within view decoding at 250 ms*

616   *post-video onset, **(b)** across view decoding at 250ms post-video onset, **(c)** within view decoding*

617   *at 500ms post-video onset, **(d)** across view decoding at 500 ms post-video onset, **(e)** subjects'*

618   *average behavioral accuracy in Experiment 2. Y-axis shows true action labels and X-axis shows*

619   *the classifier's prediction (a-d) or subjects' mean response (e). Colorbar indicates the fraction of*

620   *videos a given action (Y-axis) that was labeled by the classifier or subject as another action (X-*

621   *axis).*

622

623   ***Figure 4 – Dynamics of action representations.*** *A temporal cross training matrix showing the*

624   *decoding results for training a classifier at each point in time (y-axis) and testing the classifier at*

625   *all other times (x-axis), zoomed in to the time period from 0-1500ms post-video onset, for **(a)***

626   *within-view decoding, and **(b)** across-view decoding for subjects watching the 2-view video*

627   *dataset (Experiment 1). Colorbar indicates mean decoding accuracy for ten subjects. Black dots*

628   *indicate points when decoding is significantly above chance at group level based on p<0.05*

629   *significance test. Results along the diagonal for the within and across view decoding are the*

630   *same as shown in the line plots in Figure 3. **(c)** Significantly above chance decoding time points,*

631   *based on a p<0.05 permutation test, for the within view (blue) and across view (red) conditions*

632   *overlaid on the same plot for the entire time window (-500-3500 ms post video onset).*

633

634   ***Figure 5 – The effects of form and motion on invariant action recognition. (a)*** *Action can*

635   *also be decoded invariantly to view from form information alone (static images) **(b)** Action can*

636    *be decoded from biological motion only (point light walker stimuli). Results are each from the*

637    *average of ten subjects. Error bars represent standard error across subjects. Horizontal line*

638    *indicates chance decoding (20%). Line at bottom of plot indicates group-level significance with*

639    *p<0.05 permutation test, for the average null distribution across the ten subjects. The first time*

640    *point in this line is the onset time for each condition, reported in the main text.*

641

| Experiment | Subject | Num. sensors selected (within view) | Num. sensors selected (across view) |
|---|---|---|---|
| video | 1 | 9 | 11 |
| video | 2 | 7 | 7 |
| video | 3 | 7 | 9 |
| video | 4 | 13 | 18 |
| video | 5 | 6 | 7 |
| video | 6 | 6 | 6 |
| video | 7 | 9 | 10 |
| video | 8 | 7 | 10 |
| video | 9 | 8 | 10 |
| video | 10 | 9 | 12 |
| | | | |
| frame | 11 | 11 | 11 |
| frame | 12 | 10 | 4 |
| frame | 13 | 20 | 27 |
| frame | 14 | 4 | 5 |
| frame | 15 | 6 | 6 |
| frame | 16 | 8 | 9 |
| frame | 17 | 12 | 19 |
| frame | 18 | 16 | 23 |
| frame | 19 | 7 | 7 |
| frame | 20 | 8 | 10 |
| | | | |
| point light | 21 | 44 | 62 |
| point light | 22 | 28 | 20 |
| point light | 23 | 26 | 36 |
| point light | 24 | 29 | 32 |
| point light | 25 | 16 | 24 |
| point light | 26 | 3 | 3 |
| point light | 27 | 8 | 11 |
| point light | 28 | 24 | 25 |
| point light | 29 | 24 | 21 |
| point light | 30 | 10 | 15 |

642

643 **Table 1** - The average number of sensors selected for decoding (based on a ANOVA on the
644 training data, see Methods) for each of the 10 subjects in each experiment. The entire decoding
645 procedure, including sensor selection is repeated at each time bin. Here we report the average
646 number of sensors selected during the peak decoding time point for each subject.