![JNeurosci — THE JOURNAL OF NEUROSCIENCE]

**Research Articles: Behavioral/Cognitive**

# "The spatiotemporal neural dynamics of object recognition for natural images and line drawings"

*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.

1 **Manuscript**

2 **Title: "The spatiotemporal neural dynamics of object recognition**
3 **for natural images and line drawings"**

4 **Abbreviated Title: "Object recognition for natural images and**
5 **drawings"**

6
7 **Authors:**
8 Johannes J.D. Singer[1,2*], Radoslaw M. Cichy[2¶], Martin N. Hebart[1,3¶]

9
10 **Affiliations:**
11 [1]*Vision and Computational Cognition Group, Max Planck Institute for Human*
12 *Cognitive and Brain Sciences, Leipzig, Germany*
13 [²]*Department of Education and Psychology, Freie Universität Berlin, Germany*
14 [3]*Department of Medicine, Justus-Liebig-Universität Gießen, Germany*
15 * Corresponding author
16 Email: johannes.singer@arcor.de
17 ¶ These authors contributed equally

18

25
26 **Conflict of interest statement:**
27 The authors declare no competing financial interests.

28

32
33
34 Number of pages: 56
35 Number of figures: 9
36 Abstract - Number of words: 242
37 Introduction - Number of words: 595
38 Discussion - Number of words: 1695
39

## 1. Abstract

Drawings offer a simple and efficient way to communicate meaning. While line drawings capture only coarsely how objects look in reality, we still perceive them as resembling real-world objects. Previous work has shown that this perceived similarity is mirrored by shared neural representations for drawings and natural images, which suggests that similar mechanisms underlie the recognition of both. However, other work has proposed that representations of drawings and natural images become similar only after substantial processing has taken place, suggesting distinct mechanisms. To arbitrate between those alternatives, we measured brain responses resolved in space and time using fMRI and MEG, respectively, while human participants (female and male) viewed images of objects depicted as photographs, line drawings, or sketch-like drawings. Using multivariate decoding, we demonstrate that object category information emerged similarly fast and across overlapping regions in occipital, ventral-temporal and posterior parietal cortex for all types of depiction, yet with smaller effects at higher levels of visual abstraction. In addition, cross-decoding between depiction types revealed strong generalization of object category information from early processing stages on. Finally, by combining fMRI and MEG data using representational similarity analysis, we found that visual information traversed similar processing stages for all types of depiction, yet with an overall stronger representation for photographs. Together our results demonstrate broad commonalities in the neural dynamics of object recognition across types of depiction, thus providing clear evidence for shared neural mechanisms underlying recognition of natural object images and abstract drawings.

**Keywords:**

2

65  object recognition - line drawings - fMRI - MEG - decoding - representational
66  similarity analysis

## 67  2. Significance Statement

68  When we see a line drawing, we effortlessly recognize it as an object in the world
69  despite its simple and abstract style. Here we asked to what extent this
70  correspondence in perception is reflected in the brain. To answer this question, we
71  measured how neural processing of objects depicted as photographs and line
72  drawings with varying levels of detail (from natural images to abstract line drawings)
73  evolves over space and time. We find broad commonalities in the spatiotemporal
74  dynamics and the neural representations underlying the perception of photographs
75  and even abstract drawings. These results indicate a shared basic mechanism
76  supporting recognition of drawings and natural images.

77

## 3. Introduction

Line drawings are universal in human culture and provide a simple and efficient tool for visualization. With just a few strokes we can depict the things that we encounter in everyday life in a way that is easily recognizable by others. Line drawings of objects can be recognized without any previous experience (Kennedy & Ross, 1975), by infants only a few months after birth (DeLoache et al., 1979), and across a large variation of styles and levels of detail of the drawing (Eitz et al., 2012). This ease of recognition raises the question as to how line drawings convey meaning so efficiently.

One possible explanation for our ability to recognize line drawings efficiently is that they resemble natural object images in terms of some core visual features that are central to object recognition (Fan et al., 2018). It has been suggested that these visual features correspond to the edges of an image (Biederman & Ju, 1988). Considering the architecture of visual cortex, it has been proposed that lines in drawings drive early visual brain areas in a similar fashion to edges in natural images and therefore lead to a similar representation of objects in the brain (Sayim & Cavanagh, 2011). This notion is supported by work demonstrating that the recognition of object drawings engages the same brain regions as photographs (Ishai et al., 2000; Kourtzi & Kanwisher, 2000) and that early and high-level visual brain regions similarly represent category information for drawings and natural object images (Haxby et al., 2001; Spiridon & Kanwisher, 2002). While these results indicate that drawings and natural object images share a representational format in some visually responsive brain regions, the exact spatial extent, the temporal dynamics, and the spatiotemporal evolution of the similarities in processing of natural object images and drawings remain largely unknown.

4

103    An alternative explanation for the recognition of drawings is that the visual

104    information retained in line drawings is too abstract and therefore insufficient to drive

105    visual recognition mechanisms tuned to natural images. According to this view,

106    additional processing steps are required to refine the representation of drawings,

107    making it more similar to the representation of natural images over time. For scenes,

108    there is evidence suggesting that similarities in processing of natural scene images

109    and scene drawings become progressively stronger with the depth of visual

110    processing (Walther et al., 2011) or even emerge only late in time (Lowe et al.,

111    2018). In addition, previous results in support of a shared representational format for

112    natural object images and drawings (Haxby et al., 2001; Spiridon & Kanwisher,

113    2002) used fMRI alone, making it impossible to infer whether the effects were driven

114    by the same or distinct underlying temporal dynamics for drawings and natural

115    images. This leaves open whether drawings and natural object images are similarly

116    processed from early on or whether the shared representational format is a result of

117    additional processing steps required for drawings.

118    To provide evidence in favor or against these explanations, here we resolved

119    the similarities and differences in processing of natural object images and drawings

120    across space and time. To this end, we measured fMRI and MEG in two sessions

121    while participants viewed object images depicted across three levels of visual

122    abstraction: colored photographs, detailed black-and-white line drawings, and

123    abstract sketch-like drawings. Using spatially and temporally-resolved multivariate

124    decoding and representational similarity analysis (RSA, Kriegeskorte et al., 2008),

125    we provide clear evidence in favor of common representational dynamics for objects

126    across levels of visual abstraction in visual cortex. These results elucidate the

127  representational nature of drawings in visual cortex and suggest common neural

128  mechanisms for object recognition across levels of visual abstraction.


129  **4. Materials and Methods**


130  **4.1. Participants**

131  Thirty-one healthy adults with normal or corrected-to-normal vision took part in the

132  study and provided their written informed consent before participating. In total, we

133  excluded 8 participants from the analysis of the fMRI data and 9 from the analysis of

134  the MEG data. We based the exclusion on withdrawn participation (one participant,

135  both fMRI and MEG sessions), low alertness (>20% missed catch trials, see

136  Experimental Task Paradigm, 2 fMRI sessions and 5 MEG sessions), missing data

137  (no structural image, one fMRI session), noisy data (>1% outlier volumes in

138  framewise intensity difference / excessive head motion, 4 fMRI sessions), and

139  excessive eye movements on the stimulus (>5% of experimental trials, see section

140  on eye movement recording and analysis, affecting 3 MEG sessions). Hence, for the

141  fMRI analyses, we included data of 23 participants (mean age=29.22, SD=3.97, 13

142  female, 10 male), while for the MEG analyses, we included 22 partly overlapping

143  participants (mean age=28.91, SD=4.02, 10 female, 12 male, 17 overlapping with

144  fMRI analysis). Please note that post-hoc analyses including all the subjects into the

145  analysis for whom data was available did not qualitatively change the pattern of

146  results, demonstrating that exclusion criteria did not alter the overall pattern of

147  results. The study was approved by the local ethics committee of the University

148  Medical Center Leipzig (012/20-ek) in accordance with the declaration of Helsinki,

149  and participants were reimbursed for their participation.

**4.2. Experimental stimuli**

We used object images of the same 48 categories in three different types of depiction (144 stimuli in total), each representing one level of visual abstraction (Fig. 1a). 24 of these object categories were natural objects (e.g. animals and plants), while the other 24 were man-made (e.g. food, tools and vehicles). For each category and type of depiction there was one exemplar. For the first type of depiction ("photos"), we used colored photographs of objects, cropped from their background. For the second type of depiction ("drawings"), we asked an artist to draw black and white line drawings based on the photos with a high level of detail. In these drawings, color and some texture features were abstracted while retaining most of the contours of the objects. Finally, in the third type of depiction ("sketches"), the artist was instructed to draw line drawings of the photos in a highly abstracted way. In comparison to the drawings and photos, the sketches distorted the contours and the size of some parts of the objects, and texture information was reduced to a minimum.

**4.2.1. Quantitative validation of the experimental stimuli**

To be able to meaningfully compare object recognition for photographs and drawings at different levels of visual abstraction we reasoned that our stimulus set is required to suffice two main criteria: stimuli in the three types of depiction should (1) differ in terms of their low-level visual features, reflecting a difference in the degree of visual abstraction and (2) be perceived similarly at a conceptual level by human participants.

First, to quantitatively validate that the stimuli in the three types of depiction differ in their level of visual abstraction, we extracted low-level visual features for

174    them using the deep convolutional neural network VGG16 (Simonyan & Zisserman,

175    2015). The network is widely used and prominent for its appearance at the ImageNet

176    Large Scale Visual Recognition Challenge (Russakovsky et al., 2015) in 2014, where

177    it reached a top-5 test accuracy of 92.7% on the ImageNet dataset. VGG16 contains

178    five convolutional blocks, each composed of a series of convolutional layers,

179    followed by a max pooling and a ReLU layer. After the convolutional layers there are

180    three fully connected layers. The last fully connected layer outputs class probability

181    values for all of the 1,000 classes in the ImageNet dataset (Deng et al., 2009) after

182    applying a softmax activation function. We used VGG16 as it has not only achieved

183    good performance in image recognition tasks but also repeatedly has been shown to

184    learn representations that resemble visual object representations in the human brain

185    (Güçlü & Gerven, 2015; Schrimpf et al., 2020; Storrs et al., 2021). As an

186    approximation for low-level visual feature representations, we extracted network

187    activations from pooling layer 2 in response to our object images (Bankson et al.,

188    2018; Greene & Hansen, 2020; Reddy et al., 2021; Xie et al., 2020). For feeding the

189    object images through the network, the objects were put on a square gray

190    background and resized to 224 x 224 pixels. Next, we computed representational

191    dissimilarity matrices (RDMs, Kriegeskorte et al., 2008) by correlating all activations

192    in a given type of depiction with each other and computing pairwise distances by

193    using 1-Pearson correlation as a distance measure. This yielded one low-level visual

194    RDM for each type of depiction. We finally compared these RDMs by correlating

195    their lower triangular parts to each other using Pearson correlation. This resulted in

196    one correlation value for a given comparison between two types of depiction (e.g.

197    photo-drawing), reflecting the degree of low-level feature similarity of the stimuli.

198    To ensure that human participants perceive the stimuli in the different types of

199    depiction similarly at a conceptual level, we used data from a previous study (Singer

200    et al., 2022) in which workers on Amazon Mechanical Turk had performed a triplet

201    odd-one out task (Hebart et al., 2020) on the same stimuli as used here. In this task

202    participants were instructed to find the odd-one out in triplets of object images

203    belonging to the same type of depiction. Based on these triplet judgments we

204    constructed human perceptual similarity matrices for each type of depiction

205    separately, describing the representational object space based on human behavior.

206    Subsequently, we correlated the lower triangular parts of the similarity matrices from

207    the different types of depiction using Pearson correlation, yielding a measure of

208    perceptual similarity between all types of depiction.


209    **4.3. Experimental design and procedure**

210    All participants first completed one fMRI experiment, followed by an MEG experiment

211    on a separate day, which took place on average 30.57 days after the first experiment

212    (range 7-85). Before the fMRI experiment, participants were familiarized with the

213    stimuli used in both experiments. This was done to ensure that every participant was

214    able to recognize the objects depicted in all of the images in order to rule out the

215    possibility of differences between types of depiction based on the recognizability of

216    the images.


217    **4.3.1. Experimental paradigm**

218    During both experiments (fMRI, MEG), subjects were presented with images of the

219    same object categories in three types of depiction (photos, drawings, sketches).

220    Depiction types were not mixed within runs but presented in separate runs to avoid

221    carry-over effects of consecutive presentation of the same object in different types of

222    depiction. Participants were instructed to maintain fixation at the center of the screen

223    indicated   by   a   fixation   cross   during   the   whole   experiment   (Fig.   1b-c).

224         Stimuli   were   presented   at   the   center   of   the   screen   overlaid   with   a   semi-

225    transparent crosshair fixation cross (Thaler et al., 2013), which subtended 0.63° in

226    the fMRI experiment and 0.5° of visual angle in the MEG experiment. The individual

227    stimulus size was manually adjusted before the experiment such that the area an

228    object image occupied on the screen was approximately equal for all object images.

229    Hence, the stimulus size could vary across object images, and one object image

230    subtended   on   average   4.34°   (Range   =   [2.97°,   5.85°])   in   the   fMRI   experiment   and

231    6.15° of visual angle (Range = [4.21°, 8.25°]) in the MEG experiment.

232         Stimulus   presentation   timings   were   adjusted   to   the   specifics   of   the   imaging

233    modality. In the MRI experiment, each stimulus was presented for 500ms followed by

234    an   interstimulus   interval   (ISI)   of   2500ms   (total   trial   duration   3s).   In   the   MEG

235    experiment,   each   stimulus   was   presented   for   450ms   followed   by   an   ISI   which   was

236    randomly   sampled   from   a   range   of   values   between   250ms   to   450ms   in   steps   of

237    50ms to reduce effects of phase synchronization (average total trial duration 800ms).

238         In   both   experiments,   stimulus   presentations   were   interleaved   with   catch   trials

239    in which participants were instructed to respond to a given stimulus, in order to keep

240    the   subjects   alert.   In   the   MRI   experiment,   participants   were   instructed   to   respond

241    with a button press when the fixation cross turned red. In the MEG experiment, they

242    were instructed to respond to a paperclip stimulus (which was presented in the type

243    of depiction of the corresponding run e.g., as a drawing) and to blink, in order to

244    reduce blinking artifacts during the experimental trials. In the MRI experiment, the ISI

245    for catch trials was equal to the ISI of experimental trials (total trial duration 3s).

246    Catch trials in the MEG experiment were followed by a longer ISI (range of values

247    between 1,050ms and 1,250ms in steps of 50ms) to give participants time to

248    respond and for the MEG signal to return back to baseline after the blink (average

249    total trial duration 1600ms).

250        In a given run, each object image of a given type of depiction was presented

251    twice in the MRI and eight times in the MEG experiment. Stimulus presentation order

252    was randomized while prohibiting immediate stimulus repetition. Catch trials

253    accounted for 20% of the trials in both experiments and were presented after every

254    4th to 6th object image presentation. In total, each participant completed 12 runs in

255    the MRI experiment (4 from each condition, randomized in order, total run duration

256    6min 16.5s) and 9 runs in the MEG experiment (3 from each condition, randomized

257    in order, total run duration 7min 44.8s), resulting in 8 stimulus presentations per

258    image and condition across runs in the MRI experiment and 24 stimulus

259    presentations per image and condition across runs in the MEG experiment.

260    **4.3.2. Functional localizer task**

261    Before the experimental task in the fMRI experiment, participants underwent one

262    functional localizer run independent from the experimental runs, which was later

263    used for defining regions of interest (ROIs). Subjects were presented with either fully

264    visible object images (objects), scrambled object images (scrambled), or a fixation

265    cross. Participants were instructed to fixate on the fixation cross and to respond with

266    a button press if the same object was presented in two consecutive trials. Objects

267    and scrambled objects were presented at the center of the screen for a duration of

268    400ms, followed by a presentation of a fixation cross for 350ms. Both types of

269    images were presented in blocks of 15s each and interleaved with blocks of 7.5s of

270    fixation. The localizer run comprised 12 blocks of fixation and 12 blocks of both

271    objects and scrambled objects with a total run duration of 7min 45s.

11

### 4.4. fMRI acquisition, preprocessing and univariate analysis

**4.4.1. fMRI acquisition**

We recorded fMRI data on a Siemens Magnetom Prisma Fit 3T system (Siemens, Erlangen, Germany) using a 32-channel head coil. Functional images were acquired using a multiband 3 sequence (TR=1.5s, TE=33.2ms, in-plane resolution: 2.49×2.49mm, matrix size=82×82, FOV=204mm, flip angle=70°, 57 slices, slice thickness=2.5mm) with whole brain coverage. Existing T1-weighted structural images obtained in previous studies were used that varied in exact sequence parameters (MPRAGE, voxel size = 1mm$^3$).

**4.4.2. fMRI preprocessing**

All preprocessing and univariate analyses of the fMRI data were conducted using SPM12 (https://www. l.ion.ucl.ac.uk/spm/) and custom scripts in Matlab R2021a (www.mathworks.com).

First, we screened functional data for outliers in image intensity difference and head motion. To this end, we carried out initial realignment and computed the difference in image intensity of each functional volume and its subsequent volume for each brain slice, excluding the eyeballs. Next, to determine outlier volumes, we scaled the differences between functional images relative to the overall mean of differences across all functional images. We excluded subjects for whom more than 1% of volumes showed a more than 30-fold increase in image intensity difference or a displacement of more than 0.5mm in any direction. For all other subjects, we removed and then linearly interpolated the images that exceeded the criteria.

Following outlier removal, functional images were realigned to the first image of the run, slice-time corrected, and coregistered to the anatomical image. The

12

296  functional images of the localizer task were smoothed with a Gaussian kernel

297  (FWHM = 5mm) while the images from the experimental runs were not smoothed.

298      Further, we estimated noise components for the functional images of the

299  experimental runs by using the aCompCor method (Behzadi et al., 2007)

300  implemented in the TAPAS PhysIO toolbox (Kasper et al., 2017). To this end, tissue-

301  probability maps for the gray matter, white matter, and cerebrospinal fluid (CSF)

302  were estimated based on the structural image of a participant, and noise

303  components were extracted based on the tissue-probability maps of the white matter

304  and CSF in combination with the fMRI time series.

305  **4.4.3. fMRI univariate analysis**

306  We modeled the fMRI responses to each object image in a given run with a general

307  linear model (GLM). The onsets and durations of each object image were entered as

308  regressors into the model and were convolved with a hemodynamic response

309  function (HRF) resulting in 48 regressors for the experimental conditions in each run.

310  As nuisance regressors, we included the noise components extracted from the white

311  matter and CSF maps as well as the movement parameters and their first and

312  second order derivatives. We repeated this GLM approach 20 times, each time

313  convolving with a different HRF obtained from an openly available library of HRFs

314  ([https://github.com/kendrickkay/GLMsingle](https://github.com/kendrickkay/GLMsingle)) which was derived from a large fMRI

315  dataset of participants viewing natural scenes (Allen et al., 2022). After fitting the

316  GLMs, for each voxel we extracted the beta values for the object image regressors

317  from the GLM with the HRF that had resulted in the minimum mean residual for that

318  given voxel. Since the true HRF is variable across subjects, tasks and even brain

319  regions (Polimeni & Lewis, 2021) this approach allows a closer approximation of the

320  true HRF in comparison to using the canonical HRF while it does not lead to

321  positively biased statistics at the group level. This procedure yielded 48 beta maps

322  (one for each object category) for each run and participant. For later searchlight

323  analyses, we normalized these beta maps to the MNI template brain.

324  The fMRI responses for the localizer experiment were modeled in a separate

325  GLM, with the onsets and durations of the blocks of objects and scrambled objects

326  convolved with the canonical HRF as regressors. Only movement parameters were

327  included as nuisance regressors in this GLM. From the resulting beta estimates we

328  computed three contrasts. The first contrast was used to localize activity in early

329  visual brain areas and was defined as scrambled > objects. The second contrast was

330  used to localize activity in object-selective cortex and was defined as objects >

331  scrambled. The third contrast was used to localize activity in posterior parietal cortex

332  and was defined as objects+scrambled > baseline. This way, we obtained three *t*-

333  maps for the three contrasts for each participant.

334  **4.4.4. Region-of-interest definition**

335  We focused on regions in early visual cortex (EVC) i.e., V1, V2, V3, the lateral

336  occipital complex (LOC), comprising object-selective regions LO and pFs in the

337  ventral stream, and on the posterior intraparietal sulcus (pIPS), comprising the

338  regions IPS0 and IPS1 in the dorsal stream.

339  To define EVC, we first transformed the subject-specific *t*-maps from the

340  scrambled > objects contrast from the localizer GLM into MNI-space. Based on these

341  transformed *t*-maps we computed a contrast comparing the group-level activation

342  against zero, which resulted in one *t*-map across subjects. We then thresholded this

343  *t*-map at the $p<0.001$ level and calculated the overlap between the thresholded *t*-map

344  and the combined anatomical definition of V1, V2 and V3 from the Glasser Brain

345  Atlas (Glasser et al., 2016). Finally, we transformed this overlap image back into the

14

346  native subject space for each subject resulting in subject-specific EVC masks.

347  Please note that a more fine-grained definition of the ROIs V1, V2, V3, and V4 based

348  on the Wang et al. (2015) atlas led to qualitatively similar results as the EVC

349  definition.

350      To define object-selective cortex, we manually identified the peaks in the

351  subject-specific *t*-maps of the objects > scrambled contrast from the localizer GLM

352  which corresponded anatomically to LO and pFS. We then defined spheres with a

353  radius of 6 voxels around both peaks, including only those voxels in the spheres that

354  had *t*-values corresponding to $p<0.0001$. This resulted in one ROI mask for LO and

355  pFS, respectively. Initial exploratory analyses revealed that LO and pFS yielded

356  highly comparable results. Therefore, we merged the two ROI masks into one

357  combined LOC mask. This resulted in one object-selective cortex mask for each

358  subject.

359      To define pIPS, we first combined the probability masks for IPS0 and IPS1

360  from the Wang et al. (2015) atlas and then thresholded this combined IPS0-1 mask

361  at a value of 20%. Next, we transformed the combined pIPS mask into the individual

362  subject space. Finally, we computed the overlap between the individual pIPS mask

363  and the subject-specific *t*-map of the contrast from the localizer GLM comparing all

364  objects and scrambled objects against baseline, thresholded at $p<0.0001$. This

365  procedure resulted in one pIPS ROI mask for each subject. In case the EVC, object-

366  selective cortex, or pIPS masks overlapped in a given subject, the overlapping

367  voxels were discarded from all masks.

### 4.5. MEG acquisition and preprocessing

**4.5.1. MEG acquisition**

Before the MEG measurement started, participants' head shape was digitized using a Polhemus FASTRAK device. Additionally, five coils were placed on the head of the participant which were later used to track the head position inside the MEG. During the experiment that took place inside a magnetically shielded room, we recorded neuromagnetic signals using a 306-channel NeuroMag VectorView MEG system (Elekta, Stockholm) with a sampling rate of 1,000Hz and an online filter between 0 and 330Hz.

**4.5.2. MEG preprocessing**

To remove external noise and correct for head movements during the MEG measurement, we applied temporal signal space separation (Taulu & Simola, 2006) and movement correction to the MEG data using the Maxfilter software (Elekta, Stockholm). All further preprocessing steps were implemented in Matlab R2021a (www.mathworks.com), using the utilities of the Fieldtrip toolbox (Oostenveld et al., 2011) and custom scripts.

First, Independent Component Analysis (ICA) was applied to the combined data from all blocks to identify components corresponding to eye movements, blinks, or heartbeat. The resulting ICA components were manually inspected in combination with their topographies and time courses, and only those components that could be clearly attributed to eye movements, blinks, or heartbeat were removed from the data. Using this procedure, for a given subject, we removed an average of 1.73 components (SD = 0.69). Please note that the removal of eye movement and blink-related independent components is only meant to clean the data from noise related

16

392  to these components and is unrelated to the exclusion criteria based on eye

393  movements on the stimulus. Next, the data were filtered with a 0.5Hz high pass filter

394  and a 40Hz low pass filter and segmented into trials starting 100ms prior to the onset

395  of a given stimulus and ending 1,001ms after the stimulus presentation. Importantly,

396  triggers indicating the beginning of the stimulus presentation were adjusted to match

397  the exact time of the onset of a given image presentation by aligning them to the

398  onset of the response of an optical sensor attached to the projection monitor in the

399  MEG. Following this step, data were baseline corrected with respect to the time

400  period -100ms to 0ms relative to stimulus onset and downsampled to 100Hz to

401  speed up later multivariate analyses. Finally, multivariate noise normalization was

402  applied to the data (following general guidelines for multivariate pattern analysis of

403  M/EEG data (Guggenmos et al., 2018)). In sum, this procedure resulted in trials of

404  111 timepoints across 306 channels for every participant.

405  **4.6. Eye movement recording and analysis**

406  During the MEG experiment, eye movements of the subject were recorded using an

407  SR Research EyeLink 1000 system (SR Research Ltd). These data were only

408  acquired for the purpose of identifying subjects that made a significant amount of eye

409  movements on the presented stimulus, which might be informative about the

410  stimulus identity and could therefore bias results of multivariate pattern analysis

411  (Mostert et al., 2018; Thielen et al., 2019). No reliable eye movement data could be

412  acquired for 4 subjects, so they were excluded from further eye movement data

413  analyses.

414       First, the data were filtered with a 0.1Hz high pass filter to remove slow drifts,

415  followed by segmentation into epochs beginning 100ms prior to and ending 500ms

17

416 after stimulus onset. Second, we removed epochs that contained estimated eye

417 movements with an amplitude greater than 3° of visual angle based on the

418 assumption that these movements could not have fallen on the presented stimulus

419 and thus could not constitute an eye movement on the stimulus but rather must

420 reflect noise or occasional non-informative eye movements beyond the stimulus.

421 Finally, we discarded the pupil diameter channel from the data and retained only the

422 horizontal and vertical position channels for further analyses.

423     As an index for an eye movement on the stimulus, we detected saccades and

424 microsaccades in the extracted clean epochs by using the microsaccade detection

425 algorithm by Engbert & Kliegl (2003). Subsequently, we computed the amplitude of

426 movement in a given detected micro-saccade and labeled only the micro-saccades

427 with an amplitude greater than 1.5° of visual angle as eye movements on the

428 stimulus, given that any smaller eye movements would be hard to distinguish from

429 noise. Finally, we computed the ratio of trials containing eye movements on the

430 stimulus to all experimental trials (excluding catch trials) to determine how many

431 experimental trials were contaminated by eye movements on the stimulus for a given

432 subject. Based on this estimate, we excluded three participants from the MEG

433 analysis because they showed such eye movements in more than 5% of the

434 remaining experimental trials.

435 **4.7. Multivariate decoding of object category information**

436 We used multivariate decoding on the preprocessed fMRI voxel patterns and MEG

437 channel patterns to determine where and when the category information of a

438 presented object can be read out from brain activity. To this end, separately for every

439 type of depiction, we trained and tested linear Support Vector Machine (SVM)

18

440 classifiers (Chang & Lin, 2011) to distinguish between the responses to two given

441 objects for every possible combination of objects, resulting in one accuracy value for

442 every pair of objects (50% chance level). Subsequently, we averaged all pairwise

443 accuracies to obtain a measure of overall object discriminability. This procedure was

444 repeated across ROIs or searchlights for the fMRI data and across time points for

445 the MEG data. All decoding analyses were performed separately for every

446 participant.

447 **4.7.1. Spatially-resolved multivariate fMRI decoding**

448 To ask where in the brain category information can be read out from fMRI voxel

449 activity patterns, we used both an ROI-based and a spatially unbiased searchlight

450 procedure (Haynes et al., 2007; Kriegeskorte et al., 2006).

451 For the ROI-based procedure, we arranged the beta values from the voxels in

452 a given ROI into pattern vectors for each object category and run. We then evaluated

453 classifiers using a leave-one-out cross-validation procedure, training on the pattern

454 vectors from three runs and testing on the pattern vector from the remaining run. We

455 repeated this procedure until every pattern vector had been used once for testing

456 (see Fig. 2a for visualization of the approach). This resulted in decoding accuracies

457 for every ROI, each type of depiction, and each participant.

458 For the spatially unbiased searchlight analysis, we defined a sphere with a

459 radius of 4 voxels around a given voxel and formed pattern vectors based on all the

460 beta values within this sphere. Analogous to the ROI-based procedure, we then

461 evaluated classifiers using a leave-one-out cross-validation procedure. This

462 evaluation procedure was iterated over all possible searchlights, yielding accuracy

463 values across the whole brain for each type of depiction and each participant

19

464 separately. The resulting searchlight maps were subsequently smoothed with a

465 Gaussian kernel (FWHM = 5mm).

466 **4.7.2. Temporally-resolved multivariate MEG decoding**

467 For the temporally-resolved decoding analyses, we arranged the preprocessed MEG

468 data into pattern vectors containing the MEG data across channels for every object

469 category, trial and time point. Subsequently, to improve the signal-to-noise ratio, we

470 averaged data from two trials of the same object category into one supertrial,

471 resulting in twelve supertrials per object category and time point. We then evaluated

472 SVM classifiers using a leave-one-out cross-validation framework, training the

473 classifiers on eleven supertrials and testing on the left out supertrial and repeating

474 this procedure until every supertrial had been used once for testing (see Fig. 5a for

475 visualization of the approach). To increase the robustness of the results, we

476 repeated the whole cross-validation procedure and the averaging of trials into

477 supertrials five times while randomizing the assignment from trials to supertrials.

478 Accuracies were subsequently averaged across repetitions. This procedure was

479 repeated for every time point and for each type of depiction separately, which

480 resulted in object decoding time courses for every type of depiction and every

481 participant.

482 **4.7.3. fMRI and MEG cross-decoding of category information between types of**

483 **depiction**

484 To determine where and when object category information generalizes between

485 types of depictions, we used cross-decoding. This approach was analogous to the

486 regular decoding procedure, but instead of training and testing on data from the

487 same type of depiction, we trained a classifier on data from one type of depiction

20

488 (e.g., photos) and tested on data from another type (e.g., drawings). We carried out

489 cross-decoding for three types of comparisons: photo-drawing, photo-sketch and

490 drawing-sketch. Further, we computed the cross-decoding accuracies for both train-

491 test directions and averaged the accuracies subsequently. This way, data from both

492 types of depiction was used once for training and once for testing the classifier.

493 Analogous to the regular decoding procedure, we repeated this procedure across

494 ROIs and searchlights for fMRI and across time points for MEG data, resulting in

495 cross-decoding accuracies across space and time for the three comparisons and for

496 each participant separately.

497 **4.7.4. MEG temporal generalization analysis**

498 To investigate at which points in time the object category MEG pattern information

499 generalized to other points in time, we used the temporal generalization method

500 (King & Dehaene, 2014). For a given time point, we trained a classifier analogous to

501 the temporally-resolved decoding procedure. To determine the generalization of this

502 classifier across time, we tested the classifier on patterns not only at the matching

503 time point but at all timepoints. Then, we repeated this training-generalization

504 approach for every time point, yielding a time × time temporal generalization matrix

505 of decoding accuracies for each type of depiction and each participant.

506 **4.8. RSA-based MEG-fMRI fusion**

507 For combining the information about visual processing in the spatial dimension from

508 fMRI data with the temporal dimension from MEG data, we used RSA-based MEG-

509 fMRI fusion (Cichy et al., 2014; Cichy & Oliva, 2020; Hebart et al., 2018; see Fig.8a

510 for a visualization of the approach). The basic idea behind RSA is to characterize the

511 representational space in a given measurement component (e.g. an fMRI ROI) with

512    an RDM. An RDM describes the representational space in terms of pairwise

513    distances between responses to all of the conditions of interest, thereby abstracting

514    from the activity patterns of measurement channels (e.g. fMRI voxels or MEG

515    sensors). RDMs can be obtained e.g. across different regions in the brain or points in

516    time and can subsequently be compared by correlating them. If two RDMs exhibit a

517    positive correlation, it is assumed that the underlying representational geometry is

518    similar. Following this rationale, we computed RDMs for each time point, ROI, type of

519    depiction and each subject separately. For this, we first averaged all run-wise fMRI

520    or trial-wise MEG pattern vectors for a given object category extracted at different

521    ROIs or time points. Subsequently, we computed the pairwise dissimilarities between

522    pattern vectors as 1 - Pearson correlation and stored these dissimilarities in one

523    RDM for a given ROI or time point. Then, we correlated the lower triangular parts of

524    the ROI-specific and temporally-resolved RDMs with each other using Pearson

525    correlation, resulting in MEG-fMRI fusion time courses for each ROI, each type of

526    depiction and each participant separately.

527    **4.9. Statistical analyses**

528    To assess the statistical significance of the decoding accuracies as well as RDM

529    correlations, we used non-parametric sign-permutation tests (Nichols & Holmes,

530    2002). To this end, we obtained null distributions by randomly permuting the sign of

531    the results at the participant level a total number of 10,000 times. Based on these

532    null distributions, we obtained $p$-values for the empirical results and thresholded

533    these $p$-values at the $p<0.001$ level. $P$-values obtained for decoding accuracies were

534    based on one-sided tests, while $p$-values for RDM-correlations as well as differences

535    of decoding accuracies were based on two-sided tests. Uncorrected $p$-values were

536    only used for inference when testing decoding accuracies against chance in

537    individual ROIs since results for photos, drawings, and sketches were treated as

538    testing separate hypotheses. However, when testing for pairwise differences

539    between conditions (i.e. photo vs. drawing, photo vs. sketch, drawing vs. sketch) or

540    when testing cross-decoding accuracies for multiple combinations of depiction types

541    (i.e. photo-drawing, photo-sketch, drawing-sketch) against chance within a given

542    ROI, we corrected the $p$-values with the Benjamini-Hochberg FDR correction

543    (Benjamini & Hochberg, 1995).

544        For statistical tests across voxels or time involving a large number of multiple

545    comparisons, we applied cluster correction to control the alpha-error rate (Maris &

546    Oostenveld, 2007). The data points that exceeded the $p<0.001$ threshold were

547    clustered based on temporal or spatial adjacency, and the maximum cluster size was

548    computed for each permutation. This way, we obtained a null-distribution of the

549    maximum cluster size statistic. Finally, the clusters in the empirical results were then

550    thresholded based on the null-distribution of the maximum cluster size statistic at the

551    $p<0.05$ level. To correct for multiple tests of significance of pairwise differences

552    between conditions (e.g., photo vs. drawing, photo vs. sketch, drawing vs. sketch) or

553    for testing cross-decoding accuracies for multiple combinations of depiction types

554    (i.e. photo-drawing, photo-sketch, drawing-sketch), we obtained the cluster-size

555    statistic which corresponded to the given statistical threshold ($p<0.05$) for all of the

556    multiple tests and used the maximum cluster-size statistic computed across tests as

557    the threshold for all clusters from all tests.

558        In order to estimate confidence intervals for the decoding accuracy and RDM

559    correlation peak latencies, we used a bootstrapping procedure. For this, we

560    randomly sampled participant specific time series with replacement for a total

561  number of 100,000 times. Next, we averaged the results across participants for

562  every bootstrap sample and then estimated the peak latency by finding the maximum

563  of the average time series. Based on the mean and standard deviation of the

564  resulting distribution of peak latencies we computed the 95% confidence intervals of

565  the peak latency.

566      For comparing decoding accuracy and RDM correlation peak latencies we

567  used a bootstrapping procedure analogous to the approach described above.

568  However, instead of estimating confidence intervals of peak latencies of one

569  condition we estimated the confidence intervals of the difference between conditions

570  by subtracting the peak latencies for two given conditions estimated for each

571  bootstrap sample. This yielded a distribution of peak latency differences from which

572  we obtained the 95% confidence interval of the difference. We regarded a given

573  difference between peak latencies as significant if the confidence interval of the

574  difference did not include zero.

575      Finally, to test for the statistical equivalence of decoding accuracy or RDM

576  correlation peak latencies we used a two one-sided tests procedure (TOST) (Lakens,

577  2017).

578  **4.10. Data and code availability**

579  All results of the decoding and RSA analyses are publicly available via

580  https://osf.io/vsc6y/ along with preprocessed fMRI and MEG data from an exemplary

581  subject. The raw MEG and fMRI data can be accessed on OpenNeuro via

582  https://openneuro.org/datasets/ds004330                                    and

583  https://openneuro.org/datasets/ds004331. Code to reproduce the results and figures

584 in          the          paper          is          provided          via

585 https://github.com/Singerjohannes/object_drawing_dynamics.

## 5. Results

### 5.1. Natural object images and line drawings differ in low-level visual features but are perceived similarly

589 To ensure that our stimulus set is well suited for comparing object recognition across

590 different levels of visual abstraction, we aimed to quantitatively validate that objects

591 are perceived similarly by human subjects at a conceptual level despite differences

592 at the visual level. As a proxy for low-level visual features, we first extracted features

593 from pooling layer 2 of the deep convolutional neural network VGG16 (Simonyan &

594 Zisserman, 2015) for all of the object images, in line with previous work (Bankson et

595 al., 2018; Greene & Hansen, 2020; Reddy et al., 2021; Xie et al., 2020). We then

596 computed RDMs based on the extracted features separately for the different types of

597 depiction and correlated the lower triangular parts of the RDMs between types of

598 depiction. As expected, photos and drawings showed the highest RDM correlation

599 ($r$=0.79) while the correlation for photos and sketches ($r$=0.41) as well as the

600 correlation between drawings and sketches ($r$=0.45) were lower. Next, to confirm that

601 human subjects perceive the object images in the different types of depiction

602 similarly at a conceptual level, we used previously acquired data (Singer et al., 2022)

603 where workers on Amazon Mechanical Turk indicated which of three object images

604 they thought was the odd-one out (Hebart et al., 2020). These triplet judgments were

605 used to construct perceptual similarity matrices for each type of depiction separately,

606 which we subsequently correlated to each other to estimate their representational

607  similarity. As expected, human subjects perceived all types of depictions highly

608  similarly (all pairwise correlations *r*=0.97). In sum, these analyses quantitatively

609  confirm that while there is a gradual difference in low-level visual features across the

610  three types of depiction reflecting the degree of visual abstraction, there is also a

611  correspondence in how human participants perceive these images at a conceptual

612  level.

613  **5.2. Object category information can be decoded and generalizes**

614  **across types of depiction in early and high-level visual cortex**

615  Based on previous findings (Haxby et al., 2001; Spiridon & Kanwisher, 2002; Walther

616  et al., 2011), we hypothesized that information about the category of a presented

617  object is represented in early and high-level visual cortex for natural images as well

618  as for line drawings and that this information generalizes across levels of visual

619  abstraction. To test this hypothesis, we trained and tested SVM classifiers on the

620  fMRI data to decode the category of a presented object for each ROI and for every

621  type of depiction separately. We focused on EVC and LOC as proxies for early and

622  high-level visual processing, respectively. In addition, we explored the region pIPS in

623  the dorsal stream since a growing body of evidence supports an important role of

624  regions in the dorsal visual pathway for object recognition (see Freud et al. (2016)

625  and Ayzenberg & Behrmann (2022) for a review) and has shown a selectivity for

626  object format in these regions (Freud et al. 2018; Snow et al. 2011).

627       The category decoding results for EVC, LOC, and pIPS are shown in Fig. 2b.

628  Category information could be decoded with accuracies significantly above chance

629  from the voxel activity patterns from EVC, LOC, as well as pIPS for all types of

630  depiction (*p*<0.001, sign-permutation test). When directly comparing decoding

631 accuracies between types of depiction within an ROI, we found that there were no

632 significant differences between any of the types of depiction in either EVC (all

633 $p$>0.205, sign-permutation test, FDR-corrected), LOC (all $p$>0.083, sign-permutation

634 test, FDR-corrected) or pIPS (all $p$>0.364, sign-permutation test, FDR-corrected).

635 Finally, decoding accuracies for all types of depiction were higher in EVC than in

636 both LOC and pIPS (all $p$<0.001, sign-permutation test, FDR-corrected) and higher in

637 LOC than in pIPS (all $p$<0.001, sign-permutation test, FDR-corrected), which is

638 expected given the strong visual differences between object categories in a given

639 type of depiction. To control that these differences in decoding accuracies between

640 ROIs are not simply driven by a larger number of voxels for any of the ROIs, we

641 carried out the same decoding analysis after equating the number of voxels included

642 in all ROI masks by randomly subsampling the bigger ROI masks. This control

643 analysis led to comparable results, demonstrating that the differences between ROIs

644 are not driven by a larger ROI size of any of the ROIs. In sum, this suggests that

645 information about the category of a presented object is represented in early and

646 high-level visual brain regions for all levels of visual abstraction.

647 To identify the degree to which category information generalizes between

648 photos, drawings and sketches, we carried out cross-decoding. The rationale behind

649 this approach is that if the classifier trained on data from one type of depiction (e.g.

650 photos) can be used for data from another type of depiction (e.g. drawings), it is

651 concluded that the underlying representational format is similar. We evaluated three

652 different comparisons - photo-drawing, photo-sketch and drawing-sketch, resulting in

653 three values for each ROI. We found significant cross-decoding accuracies between

654 all types of depiction already in EVC but also in LOC and pIPS (all $p$<0.001, sign-

655 permutation test; FDR-corrected, Fig. 2c). To evaluate the robustness of these

656  findings, we also correlated the lower triangular parts of RDMs of different types of

657  depiction with each other for each ROI separately. This led to qualitatively similar

658  results, confirming the cross-decoding results.

659      Next, to further examine the degree of generalization between types of

660  depiction, we asked if the decoding accuracies within types of depiction were

661  different compared to the cross-decoding accuracies across types of depiction in

662  each ROI. If these accuracies are not significantly different, this would indicate that

663  the representation of object category is invariant to the type of depiction. If, however,

664  the cross-decoding accuracies are lower than the decoding accuracies within type of

665  depiction, this would indicate that the representation is tolerant but not invariant to

666  the type of depiction (Hebart & Baker, 2018). The results for all comparisons are

667  shown in Fig. 3. Accuracies across types of depiction were significantly lower than

668  the corresponding accuracies within types of depiction for all comparisons in both

669  EVC and LOC (all $p<0.002$, sign-permutation test, FDR-corrected). In pIPS, only the

670  comparisons "Photo minus Photo-Sketch" and "Drawing minus Drawing-Sketch"

671  reached significance (all $p<0.003$, sign-permutation test, FDR-corrected), while the

672  other comparisons were only marginally significant (all $p=0.051$, sign-permutation

673  test, FDR-corrected). The fact that these differences were less pronounced in pIPS

674  might be explained by the overall smaller decoding accuracies in pIPS. Moreover,

675  the differences in LOC and pIPS were significantly smaller than the ones in EVC,

676  and the differences in pIPS were smaller than the ones in LOC (all $p<0.004$, sign-

677  permutation test, FDR-corrected). These smaller effects in LOC and pIPS are

678  consistent with the idea of gradually increasing tolerance to the type of depiction with

679  depth of visual processing, yet, could also be explained by overall smaller decoding

680  accuracies in LOC and pIPS. Overall, these results suggest that while the

681  representation of object category in EVC, LOC and to some extent in pIPS is not

682  invariant to the type of depiction, it exhibits tolerance to the type of depiction.

683      Together, these findings corroborate earlier studies showing that category

684  information can be decoded and is similarly represented in early and high-level

685  visual cortex for natural object images and abstract drawings.


686  **5.3. Large parts of occipital and ventral temporal cortex conjointly**

687  **carry object category information which generalizes across levels**

688  **of visual abstraction**

689  While the results from the ROI analyses suggest a shared representational format of

690  object category information across types of depiction in these ROIs, they leave open

691  the spatial extent of this shared representation beyond these ROIs. To identify where

692  category information is reflected in the brain across levels of visual abstraction and

693  where it generalizes between types of depiction, we carried out a spatially unbiased

694  searchlight analysis, iterating the decoding procedure over all possible searchlight

695  locations in the brain.

696      The searchlight maps for decoding within types of depiction are shown in Fig.

697  4a. We found significant accuracies across large parts of occipital, ventral-temporal,

698  and to some extent also posterior parietal cortex ($p$<0.05, cluster-based permutation

699  test), with a strong overlap in the significance maps across types of depiction. Yet,

700  significant voxels for photos extended more into anterior parts of ventral-temporal

701  cortex than for drawings and sketches. To quantify the overlap between types of

702  depiction, we conducted a conjunction analysis based on the intersection between all

703  voxels that were significant for all three types of depiction (Nichols et al., 2005). The

704  resulting conjunction map shows where category information was conjointly found

29

705 across levels of visual abstraction (Fig. 4a). Confirming our initial observation, the

706 conjunction map covered large parts of the occipital and ventral-temporal cortex, as

707 well as a part of posterior parietal cortex. Beyond these similarities, no significant

708 differences in decoding accuracies were found between different types of depictions

709 (all *p*>0.05, cluster-based permutation test).

710 The results for the searchlight cross-decoding between different types of

711 depiction can be seen in Fig. 4b. We found significant cross-decoding accuracies

712 between all types of depiction in large regions in occipital and ventral-temporal

713 cortex, and to a smaller extent in posterior parietal cortex (*p*<0.05, cluster-based

714 permutation test). The conjunction map for all three types of comparisons showed a

715 broad overlap for all three comparisons mirroring the results from the within-type

716 decoding.

717 In sum, this suggests that - beyond localized regions in early and high-level

718 visual cortex - a large part of the ventral visual stream as well as parts of the dorsal

719 visual stream reflect information about the object in a format that can be generalized

720 across different levels of visual abstraction of the image.

721 **5.4. Category information can be decoded rapidly from MEG activity**

722 **patterns and generalizes early across types of depiction**

723 Having established where category information can be decoded and where it

724 generalizes across types of depiction, we investigated when information about the

725 category of a presented object can be read out and when this information

726 generalizes across types of depiction. Assuming that drawings and natural object

727 images are similarly processed from early on in the visual system (Sayim &

728 Cavanagh, 2011), we expected (1) that object category information should emerge

30

729 with similar temporal dynamics for all types of depiction and (2) that category

730 information should generalize early. In contrast, if additional processing is required to

731 resolve the abstract visual information in drawings, we expected delayed emergence

732 of category information for drawings and sketches in comparison to photos and

733 generalization of category information only late in time. To distinguish between these

734 alternatives, we trained and tested SVM classifiers either on the MEG channel

735 patterns from the same or different types of depiction to decode the category of a

736 presented object for each time point analogous to the fMRI decoding procedure.

737 The results of the temporally-resolved MEG decoding analyses within photos,

738 drawings and sketches are shown in Figure 5b. Irrespective of the type of depiction,

739 there was a rapid early rise in decoding accuracy, followed by a steady decline that

740 continued into the end of the trial and that remained significant for all three levels of

741 depiction ($p<0.05$, cluster-based permutation test). Overall, time courses were very

742 similar for the three conditions, peaking at 100ms for all conditions (photo peak 95%

743 confidence interval (CI) = [99.91ms 100.09ms], drawing peak CI = [98.61ms

744 101.45ms], sketch peak CI = [86.15ms 105.49ms]), with no significant differences

745 between peak latencies (all $p>0.05$, based on bootstrap CI). A two one-sided tests

746 procedure (TOST) testing for statistical equivalence of the peak latencies revealed

747 significant results for all comparisons (Photo vs. Drawing, $MD$=-0.5ms, $p<0.001$;

748 Photo vs. Sketch, $MD$=5ms, $p=0.004$; Drawing vs. Sketch, $MD$=5.5ms, $p=0.007$,

749 FDR-corrected). Despite these similarities, the overall accuracy for the three

750 conditions was different, as highlighted in the difference time courses (Fig 5c). There

751 were significantly higher decoding accuracies for photos than for both drawings and

752 sketches and significantly higher decoding accuracies for drawings than for sketches

753 (all p<0.05, cluster-based permutation test). These differences suggest that there

31

754   was a gradual decrease in the strength of the representation of category information

755   with an increasing level of visual abstraction potentially related to the additional

756   visual information (e.g. color, texture) contained in photos and drawings.

757   　　　　The cross-decoding time courses, which are depicted in Fig. 6a, showed a

758   similar pattern for all comparisons, with a sharp increase shortly after stimulus

759   presentation leading up to a peak after which accuracies declined slowly and

760   remained significant for all three comparisons up until the end of the trial ($p$<0.05

761   cluster-based permutation test). Accuracies for all three comparisons peaked at

762   100ms (photo-drawing 95% peak CI = [93.40ms 104.78ms], photo-sketch CI =

763   [85.54ms 105.21ms], drawing-sketch CI = [86.18ms 105.58ms]) with no significant

764   differences between any of the peak latencies (all $p$>0.05, based on bootstrap CI).

765   Testing for equivalence of the peak latencies revealed significant results for the

766   comparison of photo-drawing and drawing-sketch peaks ($p$=0.008, TOST, FDR-

767   corrected) but non-significant results for the other two comparisons (both $p$=0.24,

768   TOST, FDR-corrected). Please note that an analysis correlating the lower triangular

769   parts of RDMs of different types of depiction with each other for each time point led

770   to comparable results, corroborating these findings.

771   　　　　Moreover, to further assess the generalization between types of depiction, we

772   compared decoding accuracies within types of depiction with accuracies across

773   types of depiction in a time-resolved fashion. We found significantly higher decoding

774   accuracies within types of depiction for all comparisons, which remained significant

775   throughout most of the trial (p<0.05, cluster-based permutation test; Fig. 6b).

776   Differences increased rapidly, peaked early around 100ms, and declined afterwards.

777   In line with the results from the fMRI data, these results suggest that the

778    representation of object category is tolerant rather than invariant to the type of

779    depiction.

780        Together, these results show that object category can be decoded from

781    stimulus evoked brain activity for natural images and drawings regardless of the level

782    of visual abstraction, with similar temporal dynamics but a larger effect for natural

783    object images which decreased across levels of visual abstraction. Furthermore,

784    object category information generalized strongly across all levels of visual

785    abstraction beginning already in early stages of visual processing and persisting into

786    late stages of visual processing. This suggests that recognition of drawings and

787    natural object images share strong similarities from early on in visual processing.

788    **5.5. Comparable generalization of category information across time**

789    **for all levels of visual abstraction**

790    The temporally-resolved decoding analyses suggest that object category information

791    emerges similarly fast for all types of depiction and generalizes early across

792    depiction types. Yet, there might be differences in the dynamics and the stability of

793    the representations between levels of visual abstraction. Such differences in the

794    temporal dynamics between types of depiction would indicate differences in the

795    underlying neural mechanisms for recognition of natural object images and line

796    drawings. To investigate how the representation of category information for photos,

797    drawings and sketches generalizes across time, we used temporal generalization

798    analysis (King & Dehaene, 2014; Meyers et al., 2008), training a classifier on one

799    time point and testing on all other time points for every type of depiction separately.

800        The resulting time × time generalization matrices for photos, drawings and

801    sketches are shown in Fig. 7a. We found a similar pattern for all three types of

33

802 depiction with strong generalization of the representation of category information

803 beginning shortly after stimulus onset and continuing across the whole trial period

804 (*p*<0.05, cluster-based permutation test). For all types of depiction there was a

805 strong on-diagonal pattern from 50ms to ~200ms with comparatively weak off-

806 diagonal accuracies early on. Later on, there was a stronger off-diagonal component

807 after ~200ms until ~500ms. The overall pattern observed in the temporal

808 generalization matrices was qualitatively similar across types of depiction indicating

809 that the representation of the category of a presented object underwent comparable

810 representational transformations in time for all types of depiction.

811        The direct comparison of the pattern of generalization between photos,

812 drawings and sketches, shown in Fig. 7b, revealed significant differences between

813 all depiction types (*p*<0.05, cluster-based permutation test). Accuracies for photos

814 were overall higher than for sketches, with the strongest differences spanning on-

815 diagonal elements. Moreover, there were significantly higher decoding accuracies for

816 drawings than for sketches. The strongest differences again mostly covered on-

817 diagonal elements, yet with some distributed off-diagonal differences. For the

818 comparison of photos and drawings the differences were less strong and did not

819 show a clear pattern as for the other comparisons. Significant differences were more

820 distributed with higher values for photos mostly on the diagonal and also some off-

821 diagonal elements showing higher values for drawings.

822        In sum, these results demonstrate similarities in the overall pattern of

823 generalization of category information across time but also differences in the

824 strength of generalization. These differences were strongest for on-diagonal

825 elements for the photo-sketch and drawing-sketch comparison, suggesting

826 differences in the overall representation of category information between types of

34

827 depiction but less so for the generalization across time. Differences between photos

828 and drawings were less pronounced and scattered, limiting a strong interpretation of

829 these results.

830 **5.6. Similarities and differences in the combined spatiotemporal**

831 **dynamics of object recognition for different levels of visual**

832 **abstraction**

833 Our results so far suggest that there are broad commonalities in the spatial and

834 temporal dynamics of the representation of object category across levels of visual

835 abstraction. Further, object category information generalized strongly from early

836 visual processing stages on. Yet, the temporally-resolved decoding results and

837 temporal generalization results indicate that there were differences in the strength of

838 representation between types of depiction while the spatially-resolved decoding

839 results did not show such differences. Hence, the question remains where

840 differences in the neural dynamics between photos, drawings and sketches arise

841 and at what time they arise in a given region. To combine the temporal and spatial

842 information from MEG and fMRI data, we used RSA-based MEG-fMRI fusion (Cichy

843 et al., 2014; Cichy & Oliva, 2020; Hebart et al., 2018). We computed RDMs for each

844 ROI for the fMRI data and for each time point for the MEG data and and correlated

845 the lower triangular parts of the temporally-resolved and ROI-specific RDMs (see

846 Fig. 8a for visualization of the approach). This way, we could ask in what ROI and at

847 what point in time the representation of objects was similar, revealing the

848 spatiotemporal dynamics of object processing for photos, drawing and sketches. For

849 visualization purposes we also carried out the MEG-fMRI fusion analysis using a

850 spatially unbiased searchlight approach (Cichy et al., 2016), iterating the RDM

35

851 correlation across searchlights and timepoints. The resulting MEG-fMRI fusion

852 movies are publicly available via https://osf.io/vsc6y/.

853 The fusion time courses for all types of depiction in EVC and LOC are shown

854 in Fig. 8b and d respectively. In EVC we found an early increase in MEG-fMRI

855 correlation for all types of depiction leading up to peaks, followed by a sharp

856 decrease and another rise. After this second rise in correlation the MEG-fMRI

857 correlations slowly decayed for drawings and sketches while for photos there was

858 another increase. Finally, there was a last spike in correlation for all types of

859 depiction around 500 to 540ms likely reflecting effects induced by the offset of the

860 stimulus. Peak latencies for all types of depiction were found in the time from 90ms

861 to 100ms (95% CI photo = [89.47ms 106.55ms], drawing = [83.32ms 95.29ms],

862 sketch = [81.63ms 103.06ms]) with no significant differences between any types of

863 depiction (all $p$>0.05, based on bootstrap CI of difference). Moreover, we tested for

864 equivalence of the peak latencies which revealed non-significant results for all

865 comparisons (all $p$>0.626, FDR-corrected). The comparison of fusion time courses

866 between photos and both drawings and sketches in EVC, shown in Fig. 8c, revealed

867 that there were significantly higher correlations for photos than for both drawings and

868 sketches ($p$<0.05, cluster-based permutation test). These differences were strongest

869 in the time around 100ms to 200ms and the time from ~300ms to ~500ms.

870 Differences between drawings and sketches in EVC were only small but significant

871 ($p$<0.05, cluster-based permutation test). In sum, object information regardless of the

872 level of visual abstraction first peaked in EVC at around 100ms and then re-emerged

873 later around 200ms, after which the representation slowly decayed for drawings and

874 sketches, while for photos there was another late rise.

875     In LOC we found a rise in correlation for all types of depiction with peaks at

876     150ms for all three types of depiction (95% CI photo = [141.60ms 155.65ms],

877     drawing = [136.50ms 156.39ms], sketch = [146.03ms 154.92ms]), significantly later

878     than the peak latencies in EVC for all types of depictions (all $p<0.05$, based on

879     bootstrap CI of difference; Fig. 8d). After these peaks, the correlation decayed up

880     until the end of the trial only interrupted by a small rise shortly after the offset of the

881     stimulus. There were no significant differences between peak latencies of different

882     types of depiction in LOC (all $p>0.05$, based on CI of difference). Testing for

883     statistical equivalence revealed non-significant results for all comparisons of peak

884     latencies (all $p=0.841$, TOST, FDR-corrected). Furthermore, MEG-fMRI correlations

885     were stronger for photos than for both drawings and sketches in LOC while there

886     were no significant differences between drawings and sketches ($p<0.05$ cluster-

887     based permutation test; Fig. 8e). Significant differences between photos and both

888     drawings and sketches in LOC were mostly confined to early time points before

889     150ms.

890     Finally, we also explored the spatiotemporal dynamics of visual processing for

891     photos, drawings and sketches in the region pIPS in the dorsal visual pathway. The

892     MEG-fMRI correlation and MEG-fMRI correlation difference time courses for pIPS

893     are shown in Fig. 9a and 9b, respectively. In pIPS, correlations increased up to a

894     peak at 130ms for photos and drawings and at 150ms for sketches (95% CI photo =

895     [57.15ms 207.54], drawing = [6.95ms 312.74ms], sketch = [91.27ms 204.49ms]),

896     followed by a sharp decrease and a late rise in correlation after the offset of the

897     stimulus. No significant differences between peak latencies were found (all $p>0.05$,

898     based on CI of difference) and equivalence tests revealed non-significant results for

899     all comparisons (all $p=0.99$, TOST, FDR-corrected). Please note that due to overall

37

900 weaker effects in pIPS, the first peak in the time series could not be reliably detected

901 using the whole trial period for detection. Therefore, we restricted the peak detection

902 to the time period from the beginning of the trial up to the time the stimulus

903 presentation ended (-100ms to 450ms). Moreover, we found significant differences

904 between the correlations for all types of depiction ($p$<0.05 cluster-based permutation

905 test, Fig. 9b). However, these differences were overall rather small and did not follow

906 a clear pattern, limiting strong interpretation of these effects.

907 Taken together, the spatiotemporal dynamics of object recognition followed a

908 comparable pattern across levels of visual abstraction. For all types of depiction,

909 object information first peaked in EVC and later in LOC. This was followed by re-

910 emergence of object information in EVC and a late phase of object processing with a

911 sustained response in EVC. In addition, even in high-level visual regions in the

912 dorsal visual pathway there was no difference in the emergence of object information

913 between types of depiction. Despite these similarities, photos were distinctive in

914 terms of the strength of representational similarity between fMRI and MEG data and

915 showed both early and late differences in EVC as well as particularly early

916 differences in LOC. In pIPS, differences were overall less pronounced and less

917 stable, making the interpretation of these effects more challenging. In sum, these

918 results indicate additional processing at multiple stages for photos in comparison to

919 both drawings and sketches.

920 **6. Discussion**

921 In this study, we sought to identify the spatiotemporal neural dynamics underlying

922 the processing of object drawings and to determine the similarities and differences to

923 the processing of natural object images. Specifically, we used fMRI and MEG to

924    distinguish between two alternative predictions: That photos, drawings, and sketches

925    share the same representational format in both space and time, or that, alternatively,

926    additional, potentially time-consuming processes would be required for the

927    recognition of drawings and sketches. While these two predictions are not mutually

928    exclusive, our findings only confirm the former prediction in four ways. First, we

929    demonstrated that information about the category of a presented object could be

930    read out from brain activity similarly fast and across large parts of the ventral visual

931    stream as well as posterior parietal cortex, regardless of the type of depiction of the

932    image. Second, the representation of object category information generalized

933    beginning early in visual processing. Third, results from temporal generalization

934    analyses suggest that there were qualitatively similar temporal dynamics for photos,

935    drawings and sketches. Finally, the MEG-fMRI fusion results indicate that visual

936    information processing follows similar stages, first peaking in EVC and then later in

937    LOC for all types of depiction, with similar dynamics even in pIPS outside the ventral

938    visual stream. In sum, this demonstrates that there are broad temporal and spatial

939    commonalities in the neural dynamics as well as similar underlying representations

940    for natural images and drawings from early on in visual processing.

941        In addition, we did not find evidence confirming the latter prediction proposing

942    additional processing for drawings and sketches. Rather, our results suggest the

943    opposite, that is, enhanced processing for photos at multiple stages. We found a

944    gradual decline in the strength of category representations across levels of visual

945    abstraction in the MEG data, with photos showing the strongest representation,

946    followed by drawings and sketches. Moreover, the comparison of the spatiotemporal

947    dynamics between types of depiction showed that photos exhibited a stronger

948    representation both early and late in time in early visual brain regions, and

949      exclusively early on in high-level visual cortex as compared to both drawings and

950      sketches.

951            Collectively, our findings substantiate the hypothesis that line drawings

952      resemble natural object images in terms of some core visual features (Fan et al.,

953      2018), leading to a similar representation of drawings and natural images in the brain

954      (Sayim & Cavanagh, 2011). Contrary to the hypothesis of additional processing for

955      the recognition of line drawings, our results suggest that more in-depth processing is

956      elicited by natural object images at multiple stages. Finally, these results indicate

957      that the same neural mechanisms that support natural object recognition might also

958      hold for drawings across different levels of visual abstraction.

959            Despite the abstraction of substantial amounts of visual information in line

960      drawings, we found broad commonalities in the neural dynamics of object

961      recognition for natural object images and line drawings. In combination with earlier

962      findings (Haxby et al., 2001; Lowe et al., 2018; Spiridon & Kanwisher, 2002; Walther

963      et al., 2011), these results provide evidence for the hypothesis that the information

964      retained in line drawings serves as a basis for visual recognition, consistent with an

965      edge-based account of recognition (Biederman & Ju, 1988). However, our results

966      also show that object representations are stronger for photos compared to drawings

967      or sketches. This is consistent with the theory on the role of surface information in

968      object recognition (Tanaka et al., 2001) and empirical evidence (for a review see:

969      Bramão et al., 2011) which propose that visual information only contained in natural

970      images such as color and texture exerts influence on object recognition. Our findings

971      substantiate this notion and suggest that while edge-based information in drawings

972      might be sufficient to elicit qualitatively similar spatiotemporal representational

973   dynamics as for natural images, surface information significantly contributes to object

974   recognition at multiple processing stages.

975       Previous work has suggested a shared representational format for objects

976   depicted as natural photographs or line drawings in early and high-level visual cortex

977   (Haxby et al., 2001; Spiridon & Kanwisher, 2002), while for scenes such similarities

978   have been shown to become stronger or to only arise late in the visual hierarchy

979   (Lowe et al., 2018; Walther et al., 2011). Our results corroborate and extend earlier

980   findings in object recognition by demonstrating that commonalities between natural

981   object images and line drawings emerge early in time and early in the visual

982   hierarchy. Yet, these results conflict to some part with previous work in scene

983   recognition. This discrepancy might be explained by the fact that our stimulus set

984   comprised a single exemplar instead of multiple exemplars per category (Lowe et al.,

985   2018; Walther et al., 2011), which emphasizes low-level visual feature differences in

986   decoding category information. Yet, it is possible that these partly conflicting findings

987   point to a distinction in the representation and relevance of low-level visual features

988   such as edges in object and scene recognition (Groen et al., 2017), which invites

989   further exploration.

990       We demonstrated that object category information emerges similarly fast in

991   the brain for abstract drawings as compared to color photographs. This suggests that

992   object recognition can be resolved with the same amount of processing resources for

993   different levels of visual abstraction of the image. This is consistent with previous

994   computational work showing that representations for photographs and drawings at

995   different levels of visual abstraction become highly similar when being processed in

996   feedforward deep convolutional neural networks trained to categorize natural object

997   images (Fan et al., 2018; Singer et al., 2022). While other work has demonstrated

998 that additional recurrent processing is necessary for resolving degraded (Wyatte et

999 al., 2012), occluded (Rajaei et al., 2019; Tang et al., 2018) or otherwise challenging

1000 images (Kar et al., 2019), our findings indicate that no additional mechanisms are

1001 needed for the robust recognition of abstract drawings. Future research could

1002 identify precisely in which cases visual recognition can be resolved with or without

1003 the need for additional processing which might serve as an important constraint for

1004 future efforts in modeling object recognition.

1005 One difference in visual processing between photos and both drawings and

1006 sketches was found with MEG-fMRI fusion very early on in high-level visual cortex.

1007 LOC exhibited a faster rise of object representations for photos than for the other

1008 depiction types. While the source of this specific difference is unclear and not found

1009 for MEG and fMRI data alone, one possible explanation for this finding is the marked

1010 difference in the spatial frequency spectrum between drawings and photographs.

1011 While drawings and sketches contain mainly high spatial frequency information,

1012 photos additionally contain low spatial frequency information (Walther et al., 2011).

1013 This increased presence of low spatial frequency information may have contributed

1014 to an earlier rise of information related to rapid extraction of coarse information (Bar,

1015 2003; Bar et al., 2006; Kveraga et al., 2007; Musel et al., 2014; Petras et al., 2019;

1016 Peyrin et al., 2010; Schyns & Oliva, 1994; Sugase et al., 1999). Future studies that

1017 carefully control spatial frequency in an image might reveal to what extent the

1018 spatiotemporal dynamics of object recognition are influenced by different spatial

1019 frequency patterns (Perfetto et al., 2020).

1020 Previous work has shown that regions in the dorsal visual stream respond

1021 differently to real objects and images of objects (Freud et al. 2018; Snow et al.

1022 2011). Therefore, we explored if a similar sensitivity for the type of depiction of an

1023 object (i.e. photo vs. drawing) can be observed in these regions. We found category

1024 information that could be generalized across types of depiction and no evidence for

1025 differences in the emergence of category representations between types of depiction

1026 in pIPS. However, we also found differences in the strength of the spatiotemporal

1027 dynamics of object recognition in pIPS, which may lend support to a differentiation of

1028 drawings and natural images in the dorsal stream. However, the effects in pIPS were

1029 overall smaller and less pronounced as compared to the results in occipito-temporal

1030 cortex, which might point to low SNR in this region, potentially constraining the

1031 conclusions that can be drawn from our results. Future investigations focusing

1032 specifically on the dorsal visual pathway could use experimental designs and

1033 imaging protocols tailored to these regions to be able to more clearly contribute to

1034 the growing evidence for the involvement of the dorsal visual pathway in object

1035 recognition (Freud et al. 2016; Ayzenberg & Behrmann 2022).

1036 It should be noted that while we ensured that stimuli in different types of

1037 depiction are perceptually different, we did not explicitly control which details were

1038 included in the drawings and sketches and which not. Some details such as

1039 junctions and curvatures have been shown to crucially contribute to the

1040 recognizability of drawings (Walther et al., 2011; Walther & Shen, 2014). To rule out

1041 that differences in the representation of drawings and natural images can simply be

1042 explained by differences in the recognizability, we ensured that the participants were

1043 able to recognize all stimuli. Yet, the presence of some features in a drawing might

1044 determine if it is processed similarly as natural images in the visual system or not.

1045 While our results do not answer the question what features exactly allow for the

1046 recognition of drawings, they demonstrate that the features that are commonly

1047 retained lead to a similar representation of drawings and natural images in the brain.

1048     Disentangling what types of features contribute to the representations in the visual

1049     system is an overarching goal in visual neuroscience and ongoing efforts as well as

1050     future investigations might reveal distinct contributions of different features (Bankson

1051     et al., 2018; Groen et al., 2018).

1052     **6.1. Conclusion**

1053     In conclusion, our results show that the set of core visual features retained in line

1054     drawings is sufficient to elicit a processing cascade in the visual system that is

1055     remarkably similar to the one of natural images. This suggests that the same neural

1056     mechanisms that support natural object recognition might also hold for abstract line

1057     drawings. While we did not find any evidence for the involvement of additional

1058     processing for drawings, our findings indicate that visual information unique to

1059     natural images modulates visual processing at multiple stages. These results

1060     contribute to the understanding of how drawings convey meaning efficiently on the

1061     one hand and provide important insights into the neural mechanisms that underlie

1062     robust object recognition on the other hand.

1063

## 7. References

Allen, E. J., St-Yves, G., Wu, Y., Breedlove, J. L., Prince, J. S., Dowdle, L. T., Nau, M.,

Caron, B., Pestilli, F., Charest, I., Hutchinson, J. B., Naselaris, T., & Kay, K. (2022). A

massive 7T fMRI dataset to bridge cognitive neuroscience and artificial intelligence.

*Nature Neuroscience*, *25*(1), 116–126. https://doi.org/10.1038/s41593-021-00962-x

Ayzenberg, V., & Behrmann, M. (2022). Does the brain's ventral visual pathway

compute object shape? *Trends in Cognitive Sciences*.

https://doi.org/10.1016/j.tics.2022.09.019

Bankson, B. B., Hebart, M. N., Groen, I. I., & Baker, C. I. (2018). The temporal evolution of

conceptual object representations revealed through models of behavior, semantics

and deep neural networks. *NeuroImage*, *178*, 172–182.

Bar, M. (2003). A Cortical Mechanism for Triggering Top-Down Facilitation in Visual Object

Recognition. *Journal of Cognitive Neuroscience*, *15*(4), 600–609.

https://doi.org/10.1162/089892903321662976

Bar, M., Kassam, K. S., Ghuman, A. S., Boshyan, J., Schmid, A. M., Dale, A. M.,

Hämäläinen, M. S., Marinkovic, K., Schacter, D. L., Rosen, B. R., & Halgren, E.

(2006). Top-down facilitation of visual recognition. *Proceedings of the National

Academy of Sciences of the United States of America*, *103*(2), 449.

https://doi.org/10.1073/pnas.0507062103

Behzadi, Y., Restom, K., Liau, J., & Liu, T. T. (2007). A component based noise correction

method (CompCor) for BOLD and perfusion based fMRI. *NeuroImage*, *37*(1), 90–

101. https://doi.org/10.1016/j.neuroimage.2007.04.042

Benjamini, Y., & Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and

Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society.

Series B (Methodological)*, *57*(1), 289–300.

Biederman, I., & Ju, G. (1988). Surface versus edge-based determinants of visual

recognition. *Cognitive Psychology*, *20*(1), 38–64. https://doi.org/10.1016/0010-

1091    0285(88)90024-2

1092    Bramão, I., Reis, A., Petersson, K. M., & Faísca, L. (2011). The role of color information on

1093        object recognition: A review and meta-analysis. *Acta Psychologica*, *138*(1), 244–253.

1094        https://doi.org/10.1016/j.actpsy.2011.06.010

1095    Chang, C.-C., & Lin, C.-J. (2011). LIBSVM: A library for support vector machines. *ACM*

1096        *Transactions on Intelligent Systems and Technology*, *2*(3), 27:1-27:27.

1097        https://doi.org/10.1145/1961189.1961199

1098    Cichy, R. M., & Oliva, A. (2020). A M/EEG-fMRI Fusion Primer: Resolving Human Brain

1099        Responses in Space and Time. *Neuron*, *107*(5), 772–781.

1100        https://doi.org/10.1016/j.neuron.2020.07.001

1101    Cichy, R. M., Pantazis, D., & Oliva, A. (2014). Resolving human object recognition in space

1102        and time. *Nature Neuroscience*, *17*(3), 455–462. https://doi.org/10.1038/nn.3635

1103    Cichy, R. M., Pantazis, D., & Oliva, A. (2016). Similarity-Based Fusion of MEG and fMRI

1104        Reveals Spatio-Temporal Dynamics in Human Cortex During Visual Object

1105        Recognition. Cerebral Cortex, 26(8), 3563–3579.

1106        https://doi.org/10.1093/cercor/bhw135

1107    DeLoache, J. S., Strauss, M. S., & Maynard, J. (1979). Picture perception in infancy. *Infant*

1108        *Behavior and Development*, *2*, 77–89. https://doi.org/10.1016/S0163-6383(79)80010-

1109        7

1110    Deng, J., Dong, W., Socher, R., Li, L.-J., Kai Li, & Li Fei-Fei. (2009). ImageNet: A large-scale

1111        hierarchical image database. *2009 IEEE Conference on Computer Vision and*

1112        *Pattern Recognition*, 248–255. https://doi.org/10.1109/CVPR.2009.5206848

1113    Eitz, M., Hays, J., & Alexa, M. (2012). How Do Humans Sketch Objects? *ACM Transactions*

1114        *on Graphics - TOG*, *31*. https://doi.org/10.1145/2185520.2185540

1115    Engbert, R., & Kliegl, R. (2003). Microsaccades uncover the orientation of covert attention.

1116        *Vision Research*, *43*(9), 1035–1045. https://doi.org/10.1016/S0042-6989(03)00084-1

1117    Fan, J. E., Yamins, D. L. K., & Turk-Browne, N. B. (2018). Common Object Representations

1118        for Visual Production and Recognition. *Cognitive Science*, *42*(8), 2670–2698.

1119   https://doi.org/10.1111/cogs.12676

1120 Freud, E., Macdonald, S. N., Chen, J., Quinlan, D. J., Goodale, M. A., & Culham, J. C.

1121   (2018). Getting a grip on reality: Grasping movements directed to real objects and

1122   images rely on dissociable neural representations. *Cortex*, *98*, 34-48.

1123 Freud, E., Plaut, D. C., & Behrmann, M. (2016). 'What' is happening in the dorsal visual

1124   pathway. *Trends in Cognitive Sciences*, *20*(10), 773-784.

1125 Glasser, M. F., Coalson, T. S., Robinson, E. C., Hacker, C. D., Harwell, J., Yacoub, E.,

1126   Ugurbil, K., Andersson, J., Beckmann, C. F., Jenkinson, M., Smith, S. M., & Van

1127   Essen, D. C. (2016). A multi-modal parcellation of human cerebral cortex. *Nature*,

1128   *536*(7615), 171–178. https://doi.org/10.1038/nature18933

1129 Greene, M. R., & Hansen, B. C. (2020). Disentangling the Independent Contributions of

1130   Visual and Conceptual Features to the Spatiotemporal Dynamics of Scene

1131   Categorization. *The Journal of Neuroscience*, *40*(27), 5283.

1132   https://doi.org/10.1523/JNEUROSCI.2088-19.2020

1133 Groen, I. I., Greene, M. R., Baldassano, C., Fei-Fei, L., Beck, D. M., & Baker, C. I. (2018).

1134   Distinct contributions of functional and deep neural network features to

1135   representational similarity of scenes in human brain and behavior. *Elife*, *7*, e32962.

1136   https://doi.org/10.7554/eLife.32962

1137 Groen, I. I., Silson, E. H., & Baker, C. I. (2017). Contributions of low- and high-level

1138   properties to neural processing of visual scenes in the human brain. *Philosophical*

1139   *Transactions of the Royal Society of London. Series B, Biological Sciences*,

1140   *372*(1714), 20160102. https://doi.org/10.1098/rstb.2016.0102

1141 Güçlü, U., & Gerven, M. A. J. van. (2015). Deep Neural Networks Reveal a Gradient in the

1142   Complexity of Neural Representations across the Ventral Stream. *Journal of*

1143   *Neuroscience*, *35*(27), 10005–10014. https://doi.org/10.1523/JNEUROSCI.5023-

1144   14.2015

1145 Guggenmos, M., Sterzer, P., & Cichy, R. M. (2018). Multivariate pattern analysis for MEG: A

1146   comparison of dissimilarity measures. *NeuroImage*, *173*, 434–447.

47

1147        https://doi.org/10.1016/j.neuroimage.2018.02.044

1148    Haxby, J. V., Gobbini, M. I., Furey, M. L., Ishai, A., Schouten, J. L., & Pietrini, P. (2001).

1149        Distributed and overlapping representations of faces and objects in ventral temporal

1150        cortex. *Science*, *293*(5539), 2425–2430.

1151    Haynes, J.-D., Sakai, K., Rees, G., Gilbert, S., Frith, C., & Passingham, R. E. (2007).

1152        Reading Hidden Intentions in the Human Brain. *Current Biology*, *17*(4), 323–328.

1153        https://doi.org/10.1016/j.cub.2006.11.072

1154    Hebart, M. N., & Baker, C. I. (2018). Deconstructing multivariate decoding for the study of

1155        brain function. *NeuroImage*, *180*(Pt A), 4–18.

1156        https://doi.org/10.1016/j.neuroimage.2017.08.005

1157    Hebart, M. N., Bankson, B. B., Harel, A., Baker, C. I., & Cichy, R. M. (2018). The

1158        representational dynamics of task and object processing in humans. *ELife*, *7*,

1159        e32816. https://doi.org/10.7554/eLife.32816

1160    Hebart, M. N., Zheng, C. Y., Pereira, F., & Baker, C. I. (2020). Revealing the

1161        multidimensional mental representations of natural objects underlying human

1162        similarity judgements. *Nature Human Behaviour*, *4*(11), 1173–1185.

1163        https://doi.org/10.1038/s41562-020-00951-3

1164    Ishai, A., Ungerleider, L. G., Martin, A., & Haxby, J. V. (2000). The representation of objects

1165        in the human occipital and temporal cortex. *Journal of Cognitive Neuroscience, 12*

1166        *Suppl 2*, 35–51. https://doi.org/10.1162/089892900564055

1167    Kar, K., Kubilius, J., Schmidt, K., Issa, E. B., & DiCarlo, J. J. (2019). Evidence that recurrent

1168        circuits are critical to the ventral stream's execution of core object recognition

1169        behavior. *Nature Neuroscience*, *22*(6), 974–983. https://doi.org/10.1038/s41593-019-

1170        0392-5

1171    Kasper, L., Bollmann, S., Diaconescu, A. O., Hutton, C., Heinzle, J., Iglesias, S., Hauser, T.

1172        U., Sebold, M., Manjaly, Z.-M., Pruessmann, K. P., & Stephan, K. E. (2017). The

1173        PhysIO Toolbox for Modeling Physiological Noise in fMRI Data. *Journal of*

1174        *Neuroscience Methods*, *276*, 56–72. https://doi.org/10.1016/j.jneumeth.2016.10.019

1175    Kennedy, J. M., & Ross, A. S. (1975). Outline Picture Perception by the Songe of Papua.

1176        *Perception*, *4*(4), 391–406. https://doi.org/10.1068/p040391

1177    King, J.-R., & Dehaene, S. (2014). Characterizing the dynamics of mental representations:

1178        The temporal generalization method. *Trends in Cognitive Sciences*, *18*(4), 203–210.

1179        https://doi.org/10.1016/j.tics.2014.01.002

1180    Kourtzi, Z., & Kanwisher, N. (2000). Cortical Regions Involved in Perceiving Object Shape.

1181        *The Journal of Neuroscience*, *20*(9), 3310. https://doi.org/10.1523/JNEUROSCI.20-

1182        09-03310.2000

1183    Kriegeskorte, N., Goebel, R., & Bandettini, P. (2006). Information-based functional brain

1184        mapping. *Proceedings of the National Academy of Sciences*, *103*(10), 3863–3868.

1185        https://doi.org/10.1073/pnas.0600244103

1186    Kriegeskorte, N., Mur, M., & Bandettini, P. (2008). Representational similarity analysis—

1187        Connecting the branches of systems neuroscience. *Frontiers in Systems*

1188        *Neuroscience*, *2*, 4–4.

1189    Kveraga, K., Boshyan, J., & Bar, M. (2007). Magnocellular Projections as the Trigger of Top-

1190        Down Facilitation in Recognition. *The Journal of Neuroscience*, *27*(48), 13232.

1191        https://doi.org/10.1523/JNEUROSCI.3481-07.2007

1192    Lakens, D. (2017). Equivalence Tests: A Practical Primer for t Tests, Correlations, and Meta-

1193        Analyses. *Social Psychological and Personality Science*, *8*(4), 355–362.

1194        https://doi.org/10.1177/1948550617697177

1195    Lowe, M. X., Rajsic, J., Ferber, S., & Walther, D. B. (2018). Discriminating scene categories

1196        from brain activity within 100 milliseconds. *Cortex*, *106*, 275–287.

1197        https://doi.org/10.1016/j.cortex.2018.06.006

1198    Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data.

1199        *Journal of Neuroscience Methods*, *164*(1), 177–190.

1200        https://doi.org/10.1016/j.jneumeth.2007.03.024

1201    Meyers, E. M., Freedman, D. J., Kreiman, G., Miller, E. K., & Poggio, T. (2008). Dynamic

1202        Population Coding of Category Information in Inferior Temporal and Prefrontal

49

1203        Cortex. *Journal of Neurophysiology*, *100*(3), 1407–1419.

1204        https://doi.org/10.1152/jn.90248.2008

1205    Mostert, P., Albers, A. M., Brinkman, L., Todorova, L., Kok, P., & de Lange, F. P. (2018). Eye

1206        Movement-Related Confounds in Neural Decoding of Visual Working Memory

1207        Representations. *ENeuro*, *5*(4), ENEURO.0401-17.2018.

1208        https://doi.org/10.1523/ENEURO.0401-17.2018

1209    Musel, B., Kauffmann, L., Ramanoël, S., Giavarini, C., Guyader, N., Chauvin, A., & Peyrin,

1210        C. (2014). Coarse-to-fine categorization of visual scenes in scene-selective cortex.

1211        *Journal of Cognitive Neuroscience*, *26*(10), 2287–2297.

1212        https://doi.org/10.1162/jocn_a_00643

1213    Nichols, T., Brett, M., Andersson, J., Wager, T., & Poline, J.-B. (2005). Valid conjunction

1214        inference with the minimum statistic. *NeuroImage*, *25*(3), 653–660.

1215        https://doi.org/10.1016/j.neuroimage.2004.12.005

1216    Nichols, T. E., & Holmes, A. P. (2002). Nonparametric permutation tests for functional

1217        neuroimaging: A primer with examples. *Human Brain Mapping*, *15*(1), 1–25.

1218        https://doi.org/10.1002/hbm.1058

1219    Oostenveld, R., Fries, P., Maris, E., & Schoffelen, J.-M. (2011). FieldTrip: Open Source

1220        Software for Advanced Analysis of MEG, EEG, and Invasive Electrophysiological

1221        Data. *Computational Intelligence and Neuroscience*, *2011*, 1–9.

1222        https://doi.org/10.1155/2011/156869

1223    Perfetto, S., Wilder, J., & Walther, D. B. (2020). Effects of Spatial Frequency Filtering

1224        Choices on the Perception of Filtered Images. *Vision*, *4*(2).

1225        https://doi.org/10.3390/vision4020029

1226    Petras, K., ten Oever, S., Jacobs, C., & Goffaux, V. (2019). Coarse-to-fine information

1227        integration in human vision. *NeuroImage*, *186*, 103–112.

1228        https://doi.org/10.1016/j.neuroimage.2018.10.086

1229    Peyrin, C., Michel, C. M., Schwartz, S., Thut, G., Seghier, M., Landis, T., Marendaz, C., &

1230        Vuilleumier, P. (2010). The Neural Substrates and Timing of Top–Down Processes

1231    during Coarse-to-Fine Categorization of Visual Scenes: A Combined fMRI and ERP

1232    Study. *Journal of Cognitive Neuroscience*, *22*(12), 2768–2780.

1233    https://doi.org/10.1162/jocn.2010.21424

1234    Polimeni, J. R., & Lewis, L. D. (2021). Imaging faster neural dynamics with fast fMRI: A need

1235    for updated models of the hemodynamic response. *Progress in neurobiology*, *207*,

1236    102174. https://doi.org/10.1016/j.pneurobio.2021.102174

1237    Rajaei, K., Mohsenzadeh, Y., Ebrahimpour, R., & Khaligh-Razavi, S.-M. (2019). Beyond core

1238    object recognition: Recurrent processes account for object recognition under

1239    occlusion. *PLOS Computational Biology*, *15*(5), e1007001–e1007001.

1240    https://doi.org/10.1371/journal.pcbi.1007001

1241    Reddy, L., Cichy, R. M., & VanRullen, R. (2021). Representational Content of Oscillatory

1242    Brain Activity during Object Recognition: Contrasting Cortical and Deep Neural

1243    Network Hierarchies. *ENeuro*, *8*(3), ENEURO.0362-20.2021.

1244    https://doi.org/10.1523/ENEURO.0362-20.2021

1245    Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy,

1246    A., Khosla, A., & Bernstein, M. (2015). Imagenet large scale visual recognition

1247    challenge. *International Journal of Computer Vision*, *115*(3), 211–252.

1248    Sayim, B., & Cavanagh, P. (2011). What Line Drawings Reveal About the Visual Brain.

1249    *Frontiers in Human Neuroscience*, *5*, 118. https://doi.org/10.3389/fnhum.2011.00118

1250    Schrimpf, M., Kubilius, J., Hong, H., Majaj, N. J., Rajalingham, R., Issa, E. B., Kar, K.,

1251    Bashivan, P., Prescott-Roy, J., Geiger, F., Schmidt, K., Yamins, D. L. K., & DiCarlo,

1252    J. J. (2020). Brain-Score: Which Artificial Neural Network for Object Recognition is

1253    most Brain-Like? *BioRxiv*, 407007. https://doi.org/10.1101/407007

1254    Schyns, P. G., & Oliva, A. (1994). From blobs to boundary edges: Evidence for time- and

1255    spatial-scale-dependent scene recognition. *Psychological Science*, *5*(4), 195–200.

1256    https://doi.org/10.1111/j.1467-9280.1994.tb00500.x

1257    Simonyan, K., & Zisserman, A. (2015). Very Deep Convolutional Networks for Large-Scale

1258    Image Recognition. *ArXiv:1409.1556 [Cs]*. http://arxiv.org/abs/1409.1556

1259   Singer, J. J. D., Seeliger, K., Kietzmann, T. C., & Hebart, M. N. (2022). From photos to

1260        sketches—How humans and deep neural networks process objects across different

1261        levels of visual abstraction. *Journal of Vision*, *22*(2), 4–4.

1262        https://doi.org/10.1167/jov.22.2.4

1263   Snow, J. C., Pettypiece, C. E., McAdam, T. D., McLean, A. D., Stroman, P. W., Goodale, M.

1264        A., & Culham, J. C. (2011). Bringing the real world into the fMRI scanner: Repetition

1265        effects for pictures versus real objects. *Scientific reports*, *1*(1), 1-10.

1266   Spiridon, M., & Kanwisher, N. (2002). How distributed is visual category information in

1267        human occipito-temporal cortex? An fMRI study. *Neuron*, *35*(6), 1157–1165.

1268   Storrs, K. R., Kietzmann, T. C., Walther, A., Mehrer, J., & Kriegeskorte, N. (2021). Diverse

1269        Deep Neural Networks All Predict Human Inferior Temporal Cortex Well, After

1270        Training and Fitting. *Journal of Cognitive Neuroscience*, *33*(10), 2044–2064.

1271        https://doi.org/10.1162/jocn_a_01755

1272   Sugase, Y., Yamane, S., Ueno, S., & Kawano, K. (1999). Global and fine information coded

1273        by single neurons in the temporal visual cortex. *Nature*, *400*(6747), 869–873.

1274        https://doi.org/10.1038/23703

1275   Tanaka, J., Weiskopf, D., & Williams, P. (2001). The role of color in high-level vision. *Trends*

1276        *in Cognitive Sciences*, *5*(5), 211–215. https://doi.org/10.1016/s1364-6613(00)01626-

1277        0

1278   Tang, H., Schrimpf, M., Lotter, W., Moerman, C., Paredes, A., Ortega Caro, J., Hardesty, W.,

1279        Cox, D., & Kreiman, G. (2018). Recurrent computations for visual pattern completion.

1280        *Proceedings of the National Academy of Sciences*, *115*(35), 8835.

1281        https://doi.org/10.1073/pnas.1719397115

1282   Taulu, S., & Simola, J. (2006). Spatiotemporal signal space separation method for rejecting

1283        nearby interference in MEG measurements. *Physics in Medicine and Biology*, *51*(7),

1284        1759–1768. https://doi.org/10.1088/0031-9155/51/7/008

1285   Thaler, L., Schütz, A. C., Goodale, M. A., & Gegenfurtner, K. R. (2013). What is the best

1286        fixation target? The effect of target shape on stability of fixational eye movements.

1287          *Vision Research*, *76*, 31–42. https://doi.org/10.1016/j.visres.2012.10.012

1288 Thielen, J., Bosch, S. E., van Leeuwen, T. M., van Gerven, M. A. J., & van Lier, R. (2019).

1289          Evidence for confounding eye movements under attempted fixation and active

1290          viewing in cognitive neuroscience. *Scientific Reports*, *9*(1), 17456–17456.

1291          https://doi.org/10.1038/s41598-019-54018-z

1292 Wang, L., Mruczek, R. E., Arcaro, M. J., & Kastner, S. (2015). Probabilistic Maps of Visual

1293          Topography in Human Cortex. *Cerebral cortex (New York, N.Y. : 1991)*, *25*(10),

1294          3911–3931. https://doi.org/10.1093/cercor/bhu277

1295 Walther, D. B., Chai, B., Caddigan, E., Beck, D. M., & Fei-Fei, L. (2011). Simple line

1296          drawings suffice for functional MRI decoding of natural scene categories.

1297          *Proceedings of the National Academy of Sciences of the United States of America*,

1298          *108*(23), 9661–9666. https://doi.org/10.1073/pnas.1015666108

1299 Walther, D. B., & Shen, D. (2014). Nonaccidental properties underlie human categorization

1300          of complex natural scenes. *Psychological Science*, *25*(4), 851–860.

1301          https://doi.org/10.1177/0956797613512662

1302 Wyatte, D., Curran, T., & O'Reilly, R. (2012). The Limits of Feedforward Vision: Recurrent

1303          Processing Promotes Robust Object Recognition when Objects Are Degraded.

1304          *Journal of Cognitive Neuroscience*, *24*, 2248–2261.

1305          https://doi.org/10.1162/jocn_a_00282

1306 Xie, S., Kaiser, D., & Cichy, R. M. (2020). Visual Imagery and Perception Share Neural

1307          Representations in the Alpha Frequency Band. *Current Biology*, *30*(13), 2621-

1308          2627.e5. https://doi.org/10.1016/j.cub.2020.04.074

1309

## 8. Figure legends

**Figure 1. Stimuli and experimental paradigm. a) Stimulus set used in the experiment.** We used images of the same 48 object categories in three types of depiction (144 stimuli in total). Objects were depicted as either photos, drawings, or sketches, with each type of depiction reflecting a different level of visual abstraction. **b) MEG paradigm.** In the MEG experiment, participants viewed sequences of object images in random order while fixating on a central fixation cross. Their task was to respond to rare catch trials by pressing a button and blinking. **c) fMRI paradigm.** Analogous to the MEG experiment, in the fMRI experiment participants viewed sequences of object images in random order while fixating on the central fixation cross. Object sequences were interspersed with catch trials in which participants were instructed to respond with a button press. Stimulus presentation timing and ISIs were adjusted according to the modality-specific requirements.

**Figure 2. Representation and generalization of category information in early and high-level visual cortex at different levels of visual abstraction. a) Spatially resolved decoding procedure.** We trained SVM classifiers on the voxel activity patterns of a given ROI or searchlight to classify if a given pattern belonged to object category *i* or *j* for all possible pairs of objects using a leave-one-out cross-validation framework. Subsequently, we averaged the pairwise decoding accuracies for all object pairs, resulting in decoding accuracies across ROIs or searchlights for each participant and type of depiction separately. **b) Category decoding accuracies in early and high-level visual cortex across levels of visual abstraction.** We found above chance decoding accuracies for all types of depiction in EVC, LOC and pIPS. There were no significant differences in decoding accuracies between types of depiction in any of the ROIs. **c) Cross-decoding accuracies between types of depiction across ROIs.** We found significant cross-decoding accuracies between all types of depictions in EVC, LOC as well as pIPS. Error bars reflect the standard error of the mean across participants.

**Figure 3. Tolerance for the type of depiction in low- and high-level visual cortex.** Decoding accuracies within type of depiction were significantly higher than the decoding accuracies across types of depiction in EVC and LOC but only to some extent in pIPS. These differences were smaller for LOC and pIPS compared to EVC and smaller in pIPS than in LOC. Error bars reflect the standard error of the mean across participants.

1339 **Figure 4. Representation and generalization of category information at different levels of visual**

1340 **abstraction across the whole brain. a) Searchlight significance maps of the decoding of object**

1341 **category across levels of visual abstraction.** The color-coded masks indicate individually

1342 significant voxels for the decoding accuracies for each type of depiction separately. The conjunction

1343 map (color-coded in white) indicates conjointly significant voxels for all types of depiction. While

1344 significant areas for photos spanned more anterior parts of ventral temporal cortex than the ones for

1345 drawings and sketches, overall we found large parts of occipital and ventral-temporal and a part of

1346 posterior parietal cortex that conjointly reflected category information regardless of the level of visual

1347 abstraction. **b) Significance maps of the cross-decoding of object category between types of**

1348 **depiction.** Searchlight cross-decoding across the whole brain resulted in significant accuracies

1349 between all types of depiction in large parts of occipital and ventral-temporal cortex as well as a part

1350 of posterior parietal cortex. The conjunction map revealed a broad overlap in the locus of the

1351 generalizable information between all types of depiction.

1352 **Figure 5. MEG-based category information resolved in time across levels of visual abstraction.**

1353 **a) Temporally resolved decoding procedure.** For each time point, an SVM classifier was applied to

1354 MEG channel pattern supertrials in a repeated leave-one-out cross-validation scheme for all object

1355 categories $i$ and $j$. The resulting decoding accuracies were then averaged over all possible object

1356 pairs and all repetitions for every time point. This yielded decoding accuracy time-courses for every

1357 participant and every type of depiction separately. **b) Temporally resolved decoding accuracies**

1358 **across levels of visual abstraction.** For all types of depiction, category information emerged rapidly

1359 after stimulus presentation, peaking at 100ms and gradually declining afterwards. The decline was

1360 interrupted by a small increase in accuracies shortly after stimulus offset around 530ms. **c)**

1361 **Differences in decoding accuracies between types of depiction.** When directly comparing

1362 decoding accuracies between types of depiction across time we found that accuracies were

1363 significantly higher for photos compared to both drawings and sketches. Additionally, drawing

1364 accuracies were higher than sketch accuracies. Shaded areas represent the standard error of the

1365 mean across participants for each time point. Colored lines below the accuracy plots indicate

1366 significant time points ($p<0.05$, cluster-based permutation test).

1367 **Figure 6. Generalization of object category information between types of depiction across**

1368 **time. a) Cross-decoding between types of depiction across time.** Between all types of depiction,
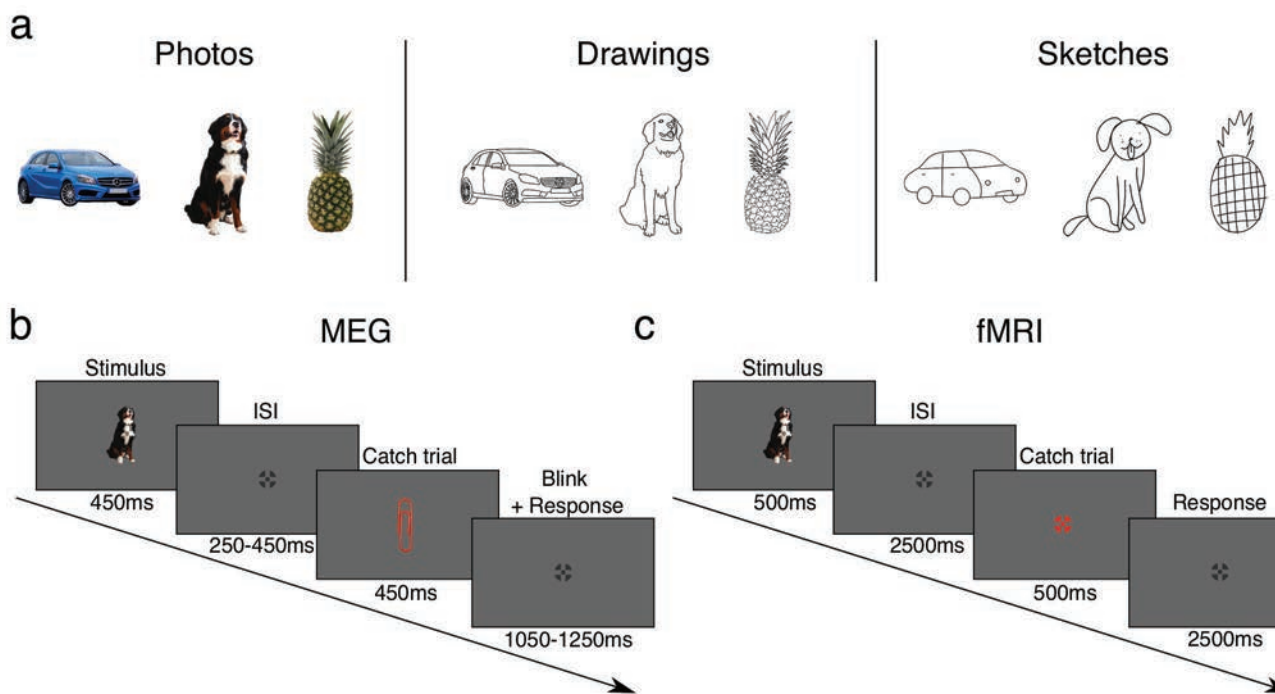
1369 we found high cross-decoding accuracies based on the MEG data already early on, peaking around

1370 100ms, and remaining high until shortly after the offset of the stimulus. **b) Differences between**

1371 **decoding within and across types of depiction across time.** Beginning early, peaking around

1372 100ms, and remaining throughout most of the trial, decoding accuracies within type of depiction were

1373 significantly higher than decoding accuracies across types of depiction. These differences declined

1374 after the peak. Shaded areas reflect the standard error of the mean across participants for each time

1375 point. Colored lines below the accuracy plots indicate significant time points ($p$<0.05, cluster-based
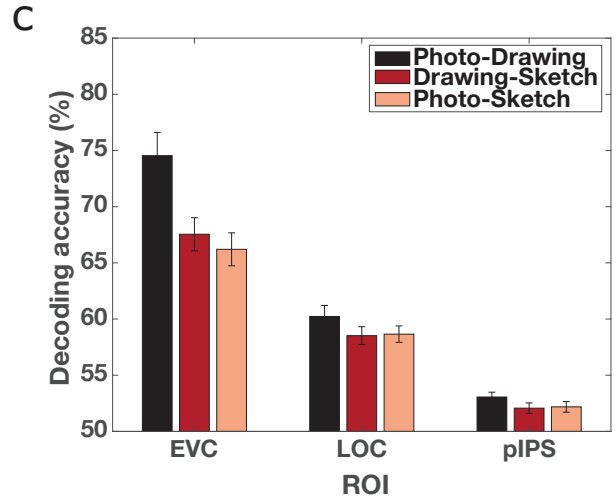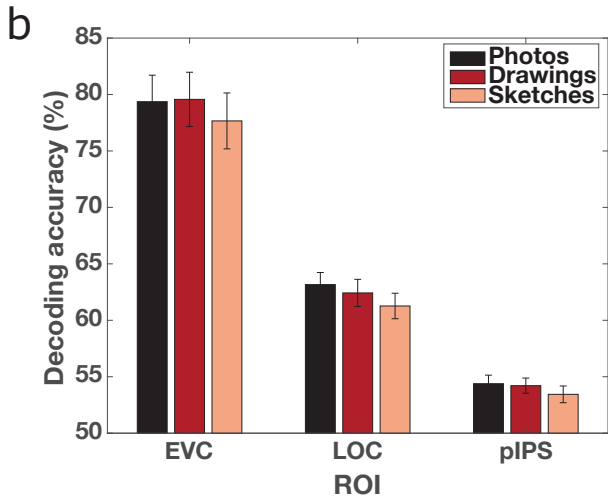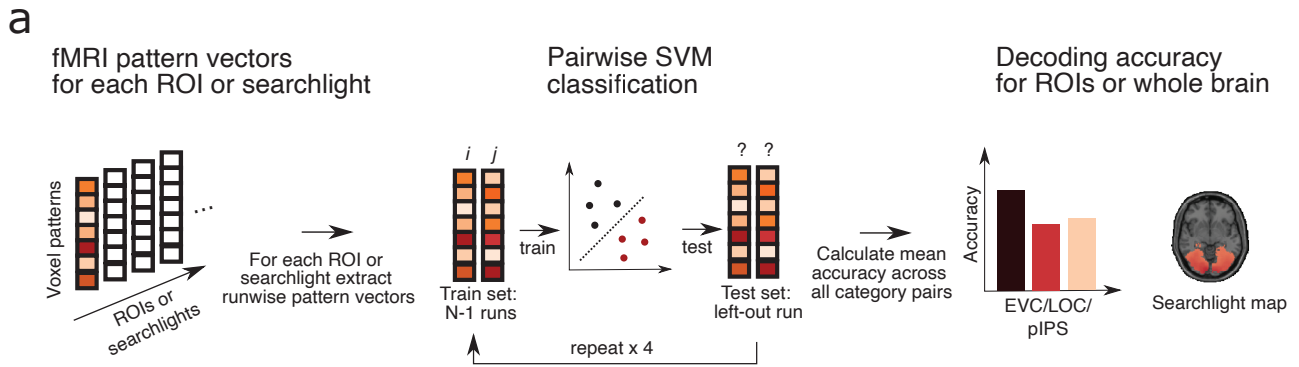
1376 permutation test).

1377 **Figure 7. Generalization of the representation of category information across time for all types**

1378 **of depiction. a) Temporal generalization matrices for the three types of depiction.** For all three

1379 types of depiction, we found strong generalization across time covering large parts of the trial. The

1380 pattern for the three types of depiction was qualitatively similar with a strong early on-diagonal

1381 component followed by a later component which showed additional strong off-diagonal elements. **b)**

1382 **Differences in temporal generalization between types of depiction.** The direct comparison of the

1383 temporal generalization between types of depiction revealed that there were differences in the

1384 strength of generalization. These differences between photo and sketches as well as drawings and

1385 sketches were most pronounced for on-diagonal elements. In addition, we found differences between

1386 photos and drawings which were less pronounced and without a clear pattern. Significant time points

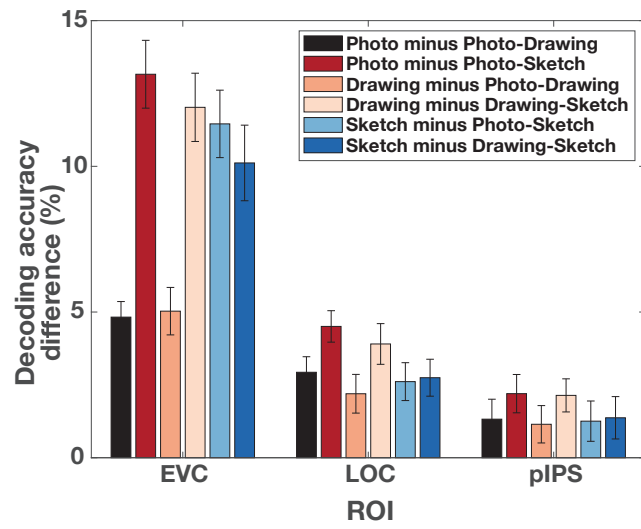1387 are indicated by the outlined areas ($p$<0.05, cluster-based permutation test).

1388 **Figure 8. Spatiotemporal dynamics of object recognition for different levels of visual**

1389 **abstraction in EVC and LOC. a) RSA-based MEG-fMRI fusion procedure.** For combining the

1390 spatial and temporal information from fMRI and MEG we first computed RDMs by calculating the

1391 pairwise dissimilarities (1- Pearson correlation) between all object-specific pattern vectors for every

1392 ROI or time point and every type of depiction separately. Then we correlated the lower triangular parts

1393 of the ROI-wise and time-wise RDMs for each ROI and every type of depiction separately. This

1394 yielded MEG-fMRI fusion time courses for each ROI reflecting the spatiotemporal dynamics of object

1395 recognition for photos, drawings and sketches. **b) MEG-fMRI fusion time courses in EVC.** In EVC

1396 we found an early peak in MEG-fMRI correlation around 100ms for all types of depiction, with no

1397 differences in peak latencies. **c) MEG-fMRI fusion difference time courses between types of**

1398 **depiction in EVC.** Photos showed a stronger correlation than both drawings and sketches in EVC,

56

1399    particularly around 100ms to 200ms and around 300ms to 500ms after stimulus presentation. The

1400    differences between drawings and sketches in EVC were small but significant. **d) MEG-fMRI fusion**

1401    **time courses in LOC.** MEG-fMRI correlations in LOC peaked significantly later than in EVC around

1402    150ms. There were no differences between peak latencies of different types of depiction. **e) MEG-**

1403    **fMRI fusion difference time courses between types of depiction in LOC.** In LOC there were no

1404    significant differences between drawings and sketches while photos showed a stronger correlation

1405    than both drawings and sketches particularly early on before 150ms. Shaded areas represent the

1406    standard error of the mean across participants for each time point. Colored lines below the accuracy

1407    plots indicate significant time points ($p$<0.05, cluster-based permutation test).
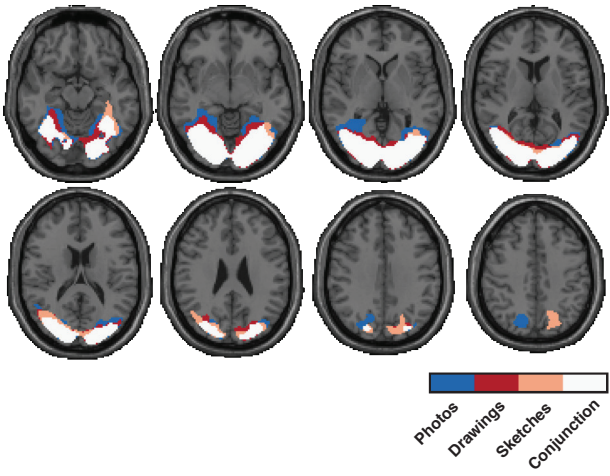
1408    **Figure 9. Spatiotemporal dynamics of object processing for different levels of visual**

1409    **abstraction in pIPS. a) MEG-fMRI fusion time courses in pIPS.** For all types of depiction, the

1410    MEG-fMRI correlation increased early and peaked at 130ms for photos and drawings and for 150ms

1411    for sketches. There were no significant differences between peak latencies. **b) MEG-fMRI fusion**

1412    **differences in pIPS.** Differences in MEG-fMRI correlations between types of depiction in pIPS were

1413    small but significant. Overall these effects were rather unstable changing directionality over time.

1414    Shaded areas represent the standard error of the mean across participants for each time point.

1415    Colored lines above the accuracy plots indicate significant time points ($p$<0.05, cluster-based

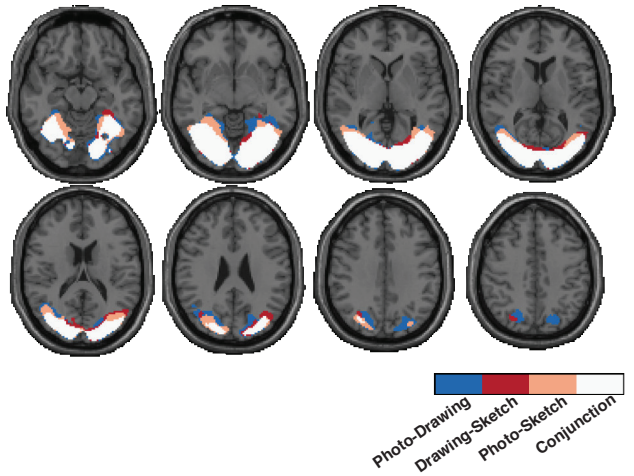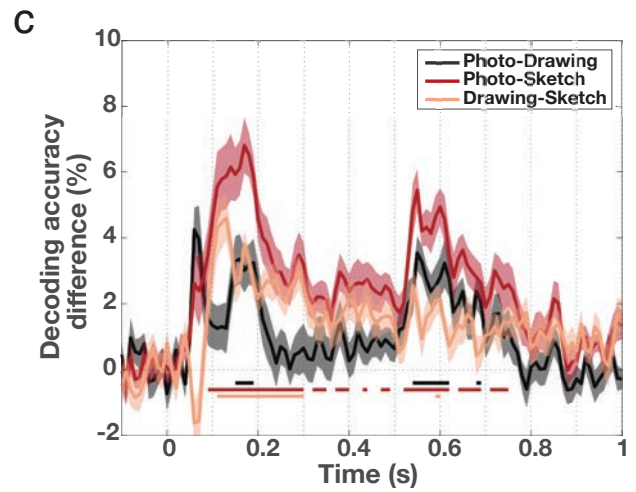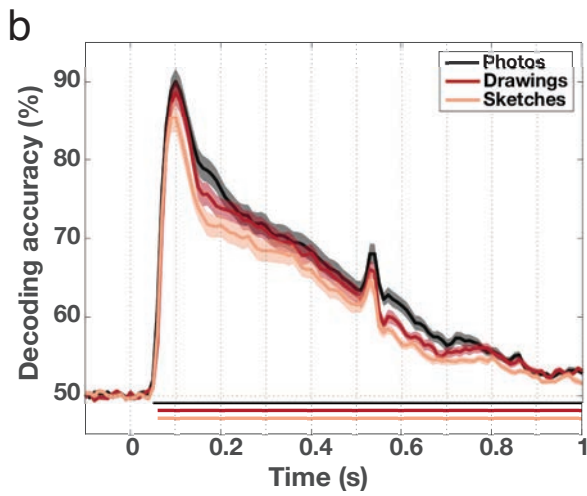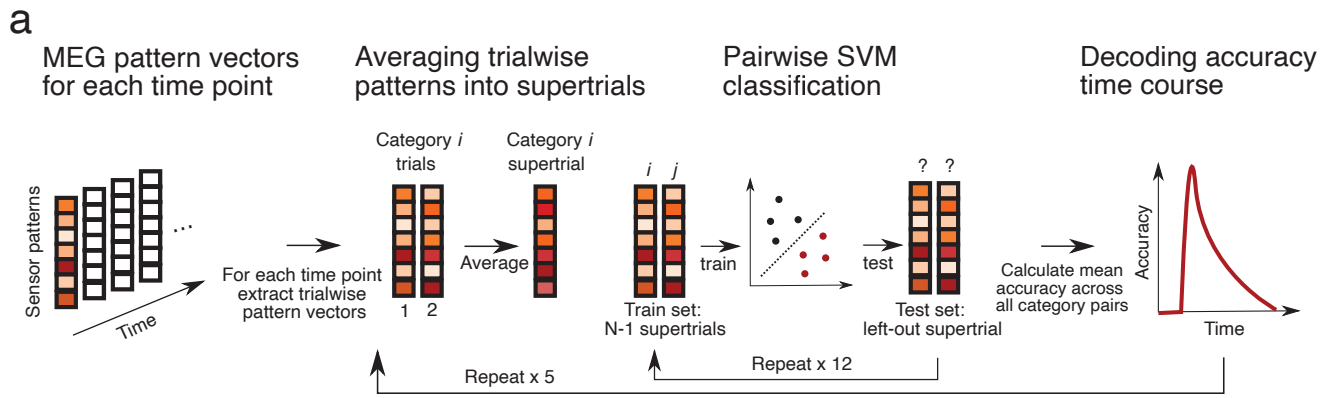1416    permutation test). The y-axes were scaled to be consistent with Figure 8b-e.

a

Photos | Drawings | Sketches

b MEG

Stimulus

450ms

ISI

250-450ms

Catch trial

450ms

Blink + Response

1050-1250ms

c fMRI

Stimulus

500ms

ISI

2500ms

Catch trial

500ms

Response

2500ms

a

fMRI pattern vectors
for each ROI or searchlight

Pairwise SVM
classification

Decoding accuracy
for ROIs or whole brain

Voxel patterns

ROIs or
searchlights

For each ROI or
searchlight extract
runwise pattern vectors

*i* *j*

train

test

? ?

Train set:
N-1 runs

Test set:
left-out run

Calculate mean
accuracy across
all category pairs

Accuracy

EVC/LOC/
pIPS

Searchlight map

repeat x 4

b



c

a



b

Photos  Drawings  Sketches  Conjunction

Photo-Drawing  Drawing-Sketch  Photo-Sketch  Conjunction

a

MEG pattern vectors
for each time point

Averaging trialwise
patterns into supertrials

Pairwise SVM
classification

Decoding accuracy
time course



b



c

a

b

a



b

a

Object specific fMRI pattern vectors

fMRI RDM

fMRI RDMs for each ROI

MEG-fMRI fusion time course for each ROI

Object A   Object B

Pairwise dissimilarity (1- Pearson R)

A
B

Object specific MEG pattern vectors

MEG RDM

MEG RDMs for each time point

Object A   Object B

Pairwise dissimilarity (1- Pearson R)

A
B

RDM similarity (Pearson R)

Similarity

Time

b

EVC

Pearson correlation

— Photo
— Drawing
— Sketch

Time (s)

c

EVC

Pearson correlation diff.

— Photo-Drawing
— Photo-Sketch
— Drawing-Sketch

Time (s)

d

LOC

Pearson correlation

— Photo
— Drawing
— Sketch

e

LOC

— Photo-Drawing
— Photo-Sketch
— Drawing-Sketch

a

**pIPS**



b

**pIPS**