# Layer-specific spatiotemporal dynamics of feedforward and feedback in human visual object perception

Tony Carricarte[1,2,3*], Siying Xie[1+], Johannes Singer[1+], Robert Trampel[4], Laurentius Huber[5], Nikolaus Weiskopf[4,6,7] Radoslaw M. Cichy[1,2,3,8]

[1]Department of Education and Psychology, Freie Universität Berlin, 14195 Berlin, Germany

[2]Einstein Center for Neurosciences Berlin, Charité - Universitätsmedizin Berlin, 10117 Berlin, Germany

[3]Bernstein Center for Computational Neuroscience Berlin, Humboldt-Universität zu Berlin, 10117 Berlin, Germany

[4]Department of Neurophysics, Max Planck Institute for Human Cognitive and Brain Sciences, 04103 Leipzig, Germany

[5]Felix Bloch Institute for Solid State Physics, Faculty of Physics and Earth Sciences, Universität Leipzig, 04103 Leipzig, Germany

[7]Wellcome Centre for Human Neuroimaging, UCL Queen Square Institute of Neurology, University College London, WC1N 3AR London, United Kingdom

[8]Berlin School of Mind and Brain, Humboldt-Universität zu Berlin, 10117 Berlin, Germany

+ These authors contributed equally

* Correspondence: tcarricarte@gmail.com

**Abstract**

Visual object perception is mediated by information flow between regions of the ventral visual stream along feedforward and feedback anatomical connections. However, feedforward and feedback signals during naturalistic vision are rapid and overlapping, complicating their identification and precise functional specification. Here we recorded human layer-specific fMRI responses to naturalistic object images in early visual cortex (EVC) and lateral occipital complex (LOC) to isolate feedforward and feedback information signals spatially by their cortical layer specific termination pattern. We combined these layer-specific fMRI responses with electroencephalography (EEG) responses for the same images to segregate feedforward and feedback signals in both time and space. Feedforward signals emerge early in the middle layers of EVC and LOC, followed by feedback signals in the superficial layer of both regions, and the deep layer of EVC. Comparing the identified dynamics in LOC to a visual deep neural network (DNN), revealed that early feedforward signals in LOC encode medium complexity features, whereas later feedback signals increase the representational format to high complexity features. Together this specifies the spatiotemporal dynamics and functional role of feedforward and feedback information flow mediating visual object perception.

2

## Introduction

Human object vision relies on anatomical bidirectional connections along the ventral visual stream[1,2], spanning the visual hierarchy from early visual cortex (EVC) to the lateral occipital complex (LOC)[3]. These connections mediate visual computations via feedforward and feedback information flows, with complex overlapping spatiotemporal dynamics[4–6]. While the feedforward flow carries sensory information up the visual processing hierarchy[7,8], downstream feedback concurrently carries information down the hierarchy, refining and shaping feedforward neural dynamics[9]. Identifying the distinct contributions of these information flows to visual computation in space and time is therefore crucial for understanding the mechanism of human object vision.

Previous efforts to disentangle the specific signatures of feedforward from feedback have typically used one of two main approaches. The first involves invasive manipulations on the anatomically-[8] and functionally-defined[6,10,11] neural circuits in non-human primates, selectively targeting bottom-up vs top-down processes. While these interventions offer precise circuit-level control, their invasive nature makes them largely impractical for human research. The second approach comprises contrasting experimental conditions where the feedforward and feedback contributions vary, such as perception vs mental imagery[12], attention vs non-attention[13], and early vs late backward masking[5]. However, this method does not directly assess information flow during naturalistic vision and has conceptual limitations, as it relies on indirect comparisons of experimental conditions that may be confounded by incompletely controlled differences between them[14–16].
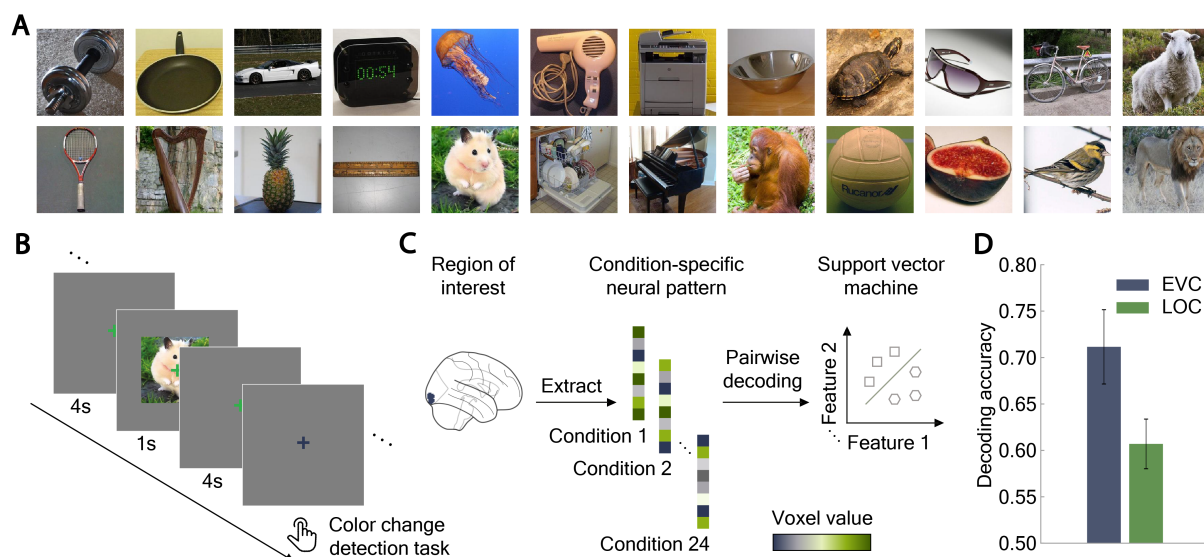
Instead, here we capitalized on the layer-specific anatomical connectivity found in the primate visual cortex[7,8,17], to dissect feedforward from feedback signals: while feedforward connections terminate primarily in the middle layer[18], feedback connections target superficial and deep layers[19,20].

Based on this three-compartment model of cortical depth, we characterized the layer-specific spatiotemporal dynamics of feedforward and feedback signals in object perception. Using sub-millimeter resolution fMRI at 7T, we recorded layer-specific brain activity in human EVC and LOC while participants viewed naturalistic object images. We then combined these layer-specific fMRI responses with time-resolved EEG responses for the same images within the framework of representational similarity analysis (RSA)[21–23] to resolve feedforward and feedback processing in millisecond resolution. Based on this, we then characterized the representational format in terms of visual feature complexity of feedforward and feedback processing using artificial deep neural networks.

## Results

76   Using sub-millimeter resolution 7T MRI, we recorded gradient-echo blood oxygenation level-dependent

77   (GE-BOLD) signals in human EVC and LOC in response to 24 different naturalistic object images (**Fig.**

78   **1A**). During the acquisition participants viewed naturalistic object images on real-world backgrounds in

79   random order (**Fig 1B**). We estimated the neural response to each condition, i.e., object image, by fitting

80   a general linear model.

81

82   Robust object information in both EVC and LOC is a precondition for further dissecting feedback and

83   feedforward-related aspects. To ascertain this, we extracted voxel values from EVC and LOC separately

84   to form condition-specific pattern vectors (**Fig. 1C**). Based on these pattern vectors we trained and tested

85   a support vector machine (SVM) to perform pairwise-object classification of objects. We found robust

86   decoding accuracy in both EVC (71,16 %, $P = 0.0039$, one-sample permutation test) and LOC (60,69

87   %, $P = 0.0092$), ascertaining reliable object information in those regions (**Fig. 1D**).
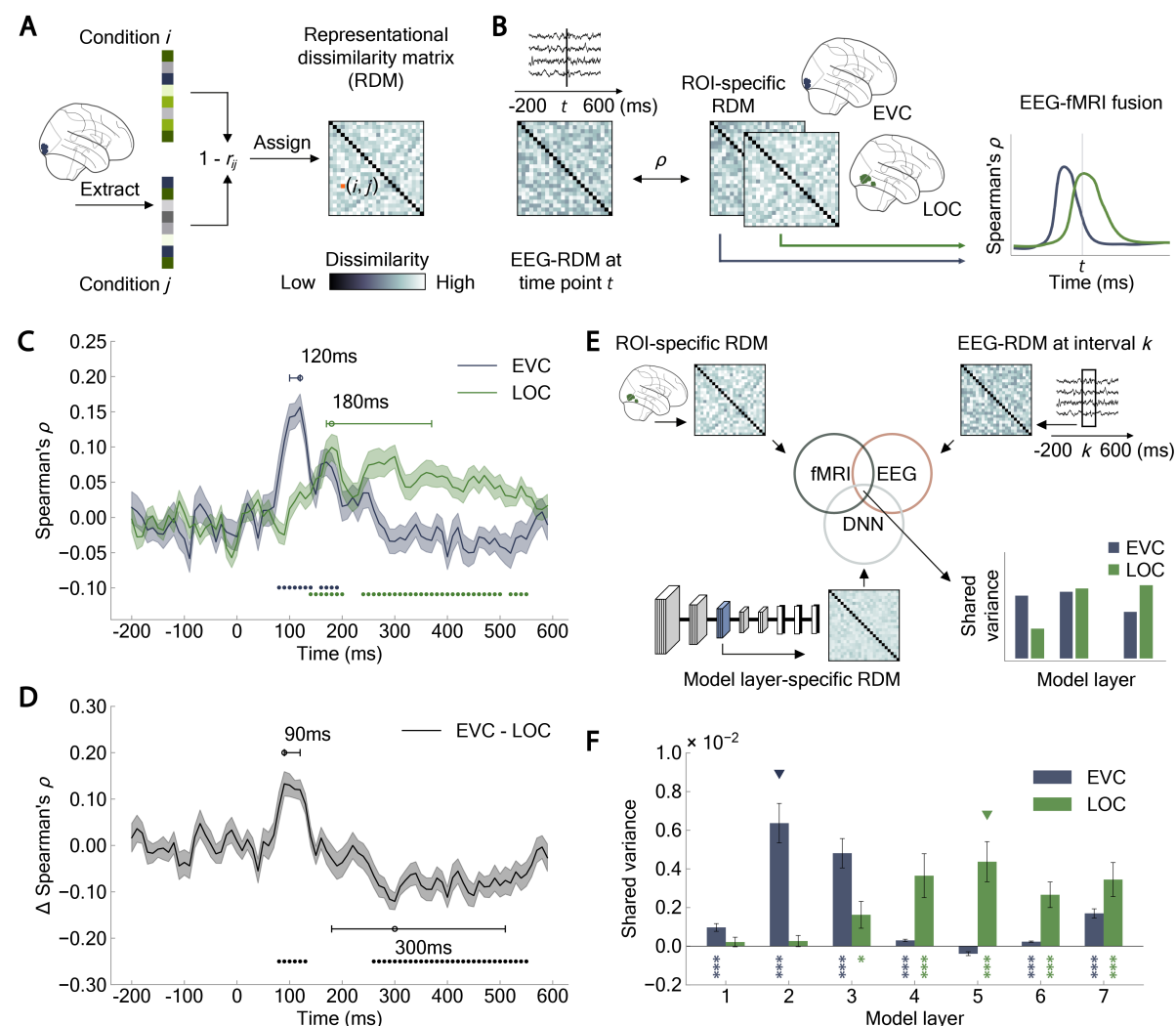
88



89

90   **Figure 1. Stimuli, experimental design and multivariate pattern classification. (A)** Stimulus set. The stimuli consisted of
91   24 different naturalistic object images. **(B)** fMRI experimental design. On each trial, participants viewed images for 1 s followed
92   by a 4-s baseline interval. Participants were required to perform a color-change detection task on the fixation cross that occurred
93   randomly throughout the experiment. **(C)** Extraction of voxel values to form condition-specific pattern vectors from region of
94   interest (example here: EVC) and pairwise-object classification using a support vector machine. ROIs are depicted for
95   visualization purposes only **(D)** Object-pairwise multivariate decoding output. Robust object-specific information was reliably
96   decoded from EVC (71,16 %, $P = 0.0039$) and LOC (60,69 %, $P = 0.0092$) using one-sample permutation tests. Error bars
97   indicate the standard error of the mean across participants.

98

99   We then proceeded to identify and examine the spatiotemporal neural dynamics of visual representations

100   in two steps that build upon each other. First, we assessed the macroscale of cortical regions (i.e., EVC

101   and LOC) to establish and thus validate representational EEG-fMRI fusion at 7T, and to characterize

102   the representational format of the identified dynamics. Based on this validation, we dissect feedforward

103   from feedback neural processing at the finer mesoscale level of cortical layers.

104

*Spatiotemporal neural dynamics of object representations in EVC and LOC at the macroscale of brain regions*



**Figure 2. Representational EEG-fMRI fusion at the macroscale.** (**A**) Representational similarity analysis. For each condition, we extracted the neural pattern from the region of interest (example here: EVC). We assessed the extent of pattern dissimilarity by calculating 1 – Pearson's correlation for all combinations of experimental conditions ($i, j$) and assigned the dissimilarity values to an fMRI representational dissimilarity matrix (RDM) indexed by the conditions in rows and columns, at entry ($i, j$). ROIs are depicted for visualization purposes only (**B**) Representational EEG-fMRI fusion. For each time point $t$, we correlated the EEG-RDM to the fMRI-RDMs of EVC and LOC using Spearman's rank order correlation. (**C**) Spatiotemporal neural dynamics at the macroscale level. Time course in EVC peaked earlier than in LOC. (**D**) Difference between EVC and LOC curves in **C**. EEG signals correlated first more with EVC than with LOC and later more with LOC than with EVC. Shaded area indicates the standard error of the mean across participants; colored circles indicate significant time points ($N = 32$, cluster-defining threshold $P < 0.05$, cluster threshold $P < 0.05$); uncolored circles and horizontal lines indicate peak latency means and 95% confidence intervals, respectively. (**E**) Commonality analysis. For each AlexNet layer, we correlated its RDM to each ROI-specific fMRI-RDM and the mean EEG-RDM at the time interval with significant ROI-specific temporal dynamics. (**F**) Format of representation ($\approx$ feature complexity) in EVC and LOC. Visual representations of low-complexity emerge early in EVC, while mid-to-high-level object representations emerge later in LOC. Error bars indicate the standard error of the mean across participants; colored asterisks indicate significant correlations ($N = 32$, right-tailed permutation tests, FDR-corrected; *$P < 0.05$; **$P < 0.01$; ***$P < 0.001$); colored triangles represent model layers with the highest occurrence proportion, determined through 1,000-iteration bootstraps.

To establish the time course of object processing in EVC and LOC at the macroscale, we employed representational EEG-fMRI fusion[22–24], which integrates time-resolved EEG with spatially resolved

128  fMRI measurements for a combined spatiotemporally resolved view. For MRI, we used the dataset

129  collected in this study, complemented with EEG responses from an existing dataset to the same set of

130  24 images[5] as in the MRI experiment. We assessed the representational geometry of EVC and LOC with

131  region-specific representational dissimilarity matrices (RDMs; **Fig. 2A**), and the representational

132  geometry of EEG signals with RDMs in a time-resolved way from -200 to 600 ms with respect to image

133  onset (**Fig. 2B**). We related EVC and LOC RDMs to the time-resolved EEG-RDMs, yielding a time

134  course of representational similarity for EVC and LOC each, for which we report peak latency and peak

135  latency differences with 95% confidence intervals in parenthesis as assessed by bootstrapping (1,000

136  iterations).

137

138  As expected from the well-established hierarchical organization of the visual system[7,8,25], object

139  processing in EVC preceded that in LOC (**Fig. 2C**): The EVC time course peaked at 120 ms (100 – 120

140  ms), whereas the LOC time course peaked later at 180 ms (170 – 370 ms), with a significant difference

141  between peak times of 60 ms (50 – 250 ms; $P < 0.001$). Furthermore, in direct comparison by subtraction

142  of the correlation curves of LOC from EVC, early representations correlated stronger with EVC than

143  with LOC at 90 ms (90 – 120 ms; **Fig. 2D**), whereas late representations correlated stronger with LOC

144  than with EVC at 300 ms (180 – 510 ms). A supplementary EEG-fMRI fusion analysis accounting for

145  similarities in representational geometry between LOC and EVC by partial correlations confirmed these

146  observations (**Suppl. Fig. 1A and B**). Our results extend EEG-fMRI fusion[5,22,24,26] from 3T to 7T fMRI

147  and provide additional support to the procedure, thus warranting a spatially finer investigation of the

148  spatiotemporal neural dynamics at the mesoscale level of cortical layers.
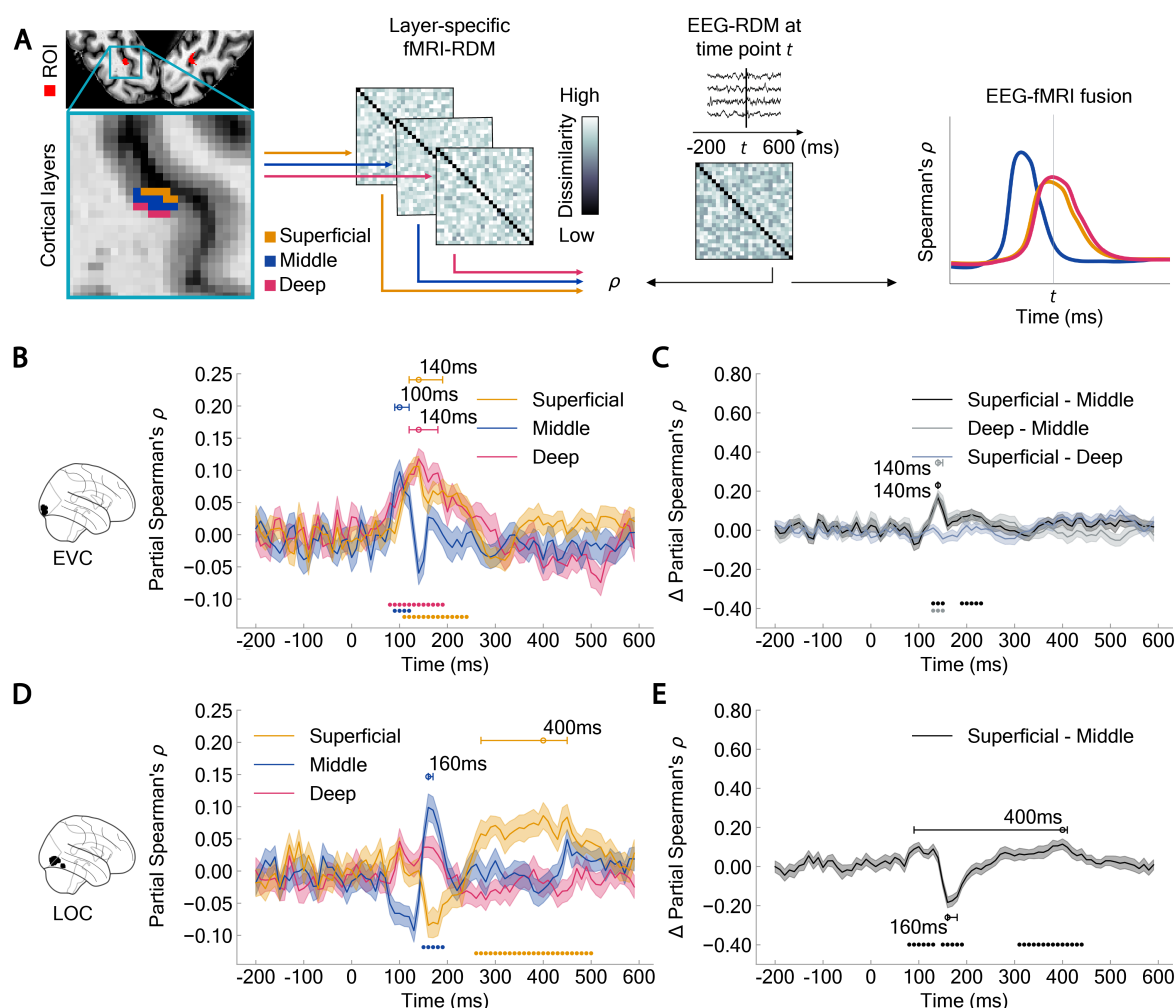
149

150  ***The format of object representations in EVC and LOC at the macroscale of brain regions***

151  The visual cortex represents objects in formats of increasing feature complexity from low to high along

152  the ventral visual stream[27–30]. To confirm this progression here, we related the neural dynamics of EVC

153  and LOC to the AlexNet deep neural network (DNN) trained on object categorization[31]. The underlying

154  rationale is that the feature complexity of visual representations progressively increases from lower to

155  higher layers across the network's layers, so that assessing the fit of the model layers to brain data reveals

156  feature complexity of the neural representations[30,32]. We conducted commonality analysis, linking each

157  of the 7 layer-specific DNN-RDMs to the ROI-specific fMRI-RDMs from EVC and LOC, and to the

158  mean EEG-RDM at corresponding time intervals, defined by the significant time points of the raw

159  curves in Fig. 2C. We report peak layers with 95% confidence intervals in parenthesis (1,000

160  bootstraps).

161

162  We found commonality with predominantly low model layers, with a peak at model layer 2 (2 – 2) in

163  EVC. and with predominantly middle to high model layers with a peak at model layer 5 (4 – 5) in LOC

164  (**Fig. 2F**). This indicates representations of primarily low complexity in EVC and mid- to high-

165    complexity in LOC. This pattern of results remained consistent when using an alternative experimental

166    choice based on significant time points from the difference between EVC and LOC curves in Fig. 2D

167    (**Suppl. Fig. 2**). Our results thus confirm a transition from representations of low complexity processed

168    early in EVC into representations of higher complexity processed later in LOC[28], further providing

169    additional support to the analytic approach at 7T and warranting a finer investigation of representational

170    format at the mesoscale level.

171

172    *Spatiotemporal neural dynamics of object representations in EVC and LOC at the mesoscale of*

173    *cortical layers*



174

**Figure 3. Representational EEG-fMRI fusion at the mesoscale. (A)**. For each cortical layer — deep, middle, and superficial — (example here: for EVC) we computed the partial Spearman's rank-order correlation between its layer-specific fMRI-RDM and the EEG-RDM at each time point *t*. (**B**) Layer-specific spatiotemporal neural dynamics in EVC. EEG signals correlated early across layers, with the time course in the middle layer peaking earlier than in the deep and superficial layers. (**C**) Difference between EVC layer curves in **B**. EEG signals correlated lately more with deep and superficial layers than with the middle layer (**D**) Layer-specific spatiotemporal neural dynamics in LOC. EEG signals correlated early in the middle layer and later in the superficial layer. (**E**) Difference between LOC layer curves in **D**. EEG signals correlated early more with the middle layer than with the superficial layer, and later more with the superficial layer than with the middle layer. Shaded area indicates the standard error of the mean across participants; colored circles indicate significant time points (*N* = 32, cluster-defining threshold *P* < 0.05, cluster threshold *P* < 0.05); uncolored circles and horizontal lines indicate peak latency means and 95% confidence intervals, respectively.

186

187    To investigate the spatiotemporal neural dynamics at the mesoscale level, we applied the research

188    procedure validated above to the finer level of cortical layers (**Fig. 3A**). To this end, we segmented the

189    cortical ribbon into three equidistant layers[33]: deep, middle and superficial. We then applied

190    representational EEG-fMRI fusion as established above, but now at the level of layers, to yield time-

191    resolved and layer-specific visual object processing time courses in EVC and LOC. To control for non-

192    specific macrovascular responses[34] that affect layer specificity[35], we conducted EEG-fMRI fusion

193    analysis partialing out for each layer the effect of the layers beneath.

194

195    In EVC we observed a correlation pattern suggesting two processing stages (**Fig. 3B, C**) indexed by

196    different profiles across layers and in timing. The first stage is marked by a peak in the middle layer at

197    100 ms (90 – 120 ms; **Fig. 3B**). The second stage is characterized by peaks at 140 ms in both the deep

198    (120 –180 ms) and superficial (120 – 190 ms) layers, with significant latency differences of 40 ms (10

199    – 50 ms; $P = 0.004$) between the middle layer and the deep layer and 40 ms (20 – 80 ms; $P = 0.002$)

200    between the middle layer and superficial layer, but not between the deep layer and the superficial layer

201    0 ms (-20 – 40 ms; $P = 0.814$). This pattern was further substantiated by subtracting the layer-specific

202    time courses which showed significant effects (**Fig. 3C**). We observed stronger correlations of late

203    representations at 140 ms with the deep (140–150 ms) and the superficial (140–140 ms) layers than with

204    the middle layer. This suggests initial feedforward processing in the middle layer and a later emerging

205    feedback processing with a distinctive layer profile in the deep and superficial layers in EVC.

206

207    In LOC, the results pattern also indicated two stages with a distinctive layer and temporal profile (**Fig.**

208    **3D, E**). We observe earlier processing in the middle layer, followed by processing in the superficial

209    layer later (**Fig. 3D**). In detail, time courses of object processing peaked first in the middle layer at 160

210    ms (160 – 170 ms), and later in the superficial layer at 400 ms (270 – 450 ms), with a significant

211    difference in peak latency of 240 ms (110 – 290 ms; $P < 0.012$). A direct comparison of time courses

212    by subtraction (**Fig. 3E**) substantiated the observation in that representations correlated stronger with

213    the middle layer early than the superficial layer at 160 ms (160 – 180 ms), and late representations

214    correlated stronger with the superficial layer than the middle layer later, at 400 ms (90 – 410 ms). This

215    suggests an initial feedforward processing in the middle layer and a later emerging feedback processing

216    with a distinctive layer profile in the superficial layer in LOC.

217

218    An analogous results pattern in EVC (**Suppl. Fig. 3A, B**) and LOC (**Suppl. Fig. 3C, D**) was observed

219    when assessing layers directly without partialing out the effect of deeper layers, confirming the

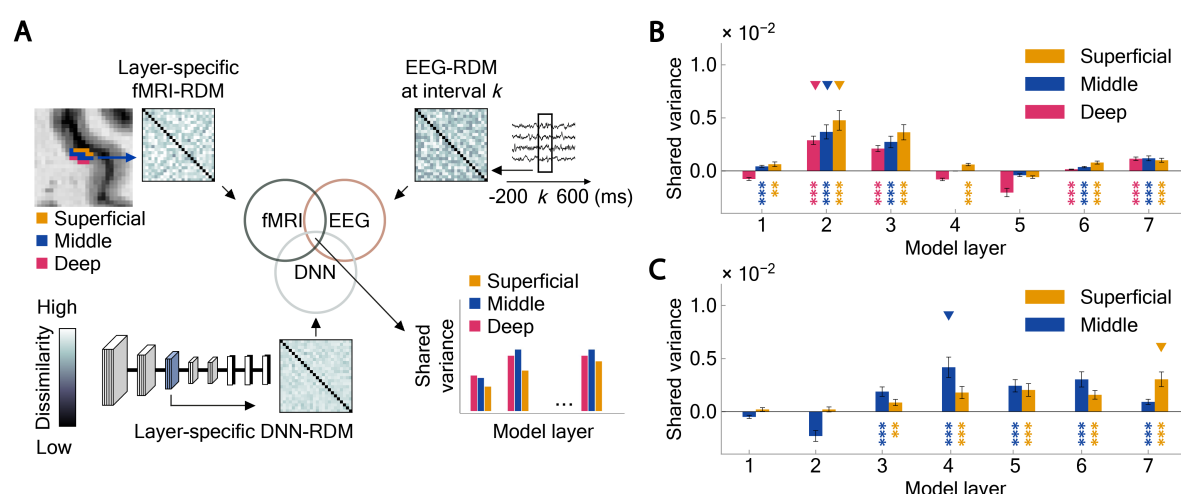220    robustness of the results to particular analysis choices.

221

222 Together, this resolves the spatiotemporal dynamics of feedforward and feedback information flow
223 during visual object processing across EVC and LOC through cortical layer-specific and temporally
224 distinct response profiles.

225

226 *The format of object representations in EVC and LOC at the mesoscale of cortical layers*

227 One interpretation of the observed layer-specific spatiotemporal dynamics, being guided by the
228 functional pattern of cortical layer connectivity[36,37], is that they indicate interareal communication via
229 feedforward and feedback connections[38,39]. However, an alternative explanation is that they arise intra-
230 areal communication via lateral connections[40] or from superficial bias in the fMRI measurements[41].

231

232 To disambiguate between these options, we characterized the format of representations by assessing
233 feature complexity as present in DNN layers, from low to high, applying the research approach as used
234 at the macroscale but now at the level of cortical layers. If the observed layer-specific spatiotemporal
235 dynamics reflect an interplay of feedforward and feedback information flow, we would expect the flow
236 to carry information of varying complexity[42,43] across different levels of the visual hierarchy[44,45],
237 resulting in distinct representational formats across cortical layers. In contrast, if the observed layer-
238 specific spatiotemporal dynamics reflect lateral connections modulating the neural gain[46–50] or a
239 superficial bias, we would expect uniform feature complexity across layers[51].

240

241 We conducted commonality analysis, linking each model layer-specific DNN-RDM to the deep, middle
242 and superficial layer-specific fMRI-RDM and mean EEG-RDM at corresponding time intervals, defined
243 by the significant time points of the raw curves in Fig. 3B and D indicating periods of layer-specific
244 neural dynamics (**Fig. 4A**).

245



246

**Figure 4. Commonality analysis between fMRI, EEG and AlexNet at the mesoscale.** (**A**) Procedure. For each AlexNet layer, we correlated its RDM to each layer-specific fMRI-RDM in EVC and LOC and the mean EEG-RDM at the time interval with significant layer-specific temporal dynamics. (**B**) Format of representation (≈ feature complexity) across cortical layers in EVC. Low model layers correlated strongly across layers in EVC. (**C**) Format of representation (≈ feature complexity) across cortical layers in LOC. Middle model layers correlated strongly with the middle layer in LOC, while high model layers

252  correlated primarily with the superficial layer. Colored asterisks indicate significant correlations ($N$ = 32, right-tailed
253  permutation tests, FDR-corrected; *$P$ < 0.05; **$P$ < 0.01; ***$P$ < 0.001); colored triangles represent model layers with the
254  highest occurrence proportion, determined through 1,000-iteration bootstraps.
255

256  In EVC (**Fig. 4B**) we observed a uniform results pattern across cortical layers: object representations

257  shared variance with all DNN layers, but most strongly with model layer 2 (2 – 3), indicating mainly

258  representations of low-to-mid level complexity. This suggests that the observed layer-specific

259  spatiotemporal dynamics in EVC reflect lateral connections modulating the neural gain[46–50] or

260  superficial bias[41] rather than feedback.

261

262  In contrast, in LOC (**Fig. 4C**) we observed a shift from mid-to high feature complexity across cortical

263  layers. The early emerging representations in the middle cortical layer (Fig. 3D) were of mid-

264  complexity, as indicated by a peak at model layer 4 (4 – 6). In contrast, the late emerging representations

265  in the superficial layers were of high complexity, with a peak at model layer 7 (4 – 7). This results

266  pattern was robust: making alternative experimental choices based on the average of the significant time

267  intervals from the difference between layer curves in Fig. 3C and E (**Suppl. Fig. 4 and 5**) and individual

268  time points spanning the full-time course (**Suppl. Fig. 6**) yielded equivalent results. Together, these

269  findings suggest a dynamic shift in representational format in LOC, transitioning from early

270  representations of mid-level feature complexity in the feedforward flow to later representations of high

271  feature complexity through interareal feedback.

272

273  **Discussion**

274  We leveraged layer-resolved fMRI, time-resolved EEG data and DNNs to identify and characterize the

275  spatiotemporal neural dynamics of feedforward and feedback information flow underlying object

276  perception. We validated our methods using 7T fMRI at the macroscale of cortical regions by replicating

277  the temporal dynamics and representational format in EVC and LOC observed in 3T fMRI

278  studies[5,22,26,30], allowing us to proceed to the mesoscale of cortical layers. There we made two key

279  observations. First, we observed distinct layer-specific temporal profiles for EVC and LOC. Visual

280  representations in the middle layers emerged earlier than in the deep and superficial layers of EVC and

281  the superficial layers of LOC. Second, the identified layer-specific dynamics in LOC had distinctive

282  visual feature complexity profiles: the early emerging middle layer representations were of mid-

283  complexity and the later emerging superficial layer representations were of high complexity.

284

285  ***Layer-specific EEG-fMRI fusion reveals sequential feedforward and feedback processing in visual***

286  ***object perception***

287  Visual object perception unfolds through intricate spatiotemporal neural dynamics[4–6,52,53], mediating

288  feedforward and feedback information flow[22,54] rapidly and through temporally overlapping responses.

289  Here, we disentangle the temporal dynamics of feedforward from feedback signals by leveraging the

290  anatomical canonical microcircuit of cortex[7,8] at the input and output stage of cortical object

291  processing[3,55,56] – EVC and LOC. We find that feedforward signals emerge early in the middle layer of

292  both EVC at 100 ms and in LOC at 160 ms, while feedback appears relatively later in deep and

293  superficial layers of EVC at 140 ms and superficial layer of LOC at 400 ms, supporting the idea of

294  sequential processing of visual information first through feedforward, then through feedback

295  processing[37,57]. Our results thus provide a functional temporal characterization for the anatomical

296  canonical cortical microcircuit model of feedforward and feedback connectivity in the human ventral

297  visual stream.

298

299  Our findings have theoretical implications. For example, they specify in spatiotemporal terms the

300  dynamical communication model in predictive coding theory[58] which posits distinct neural channels for

301  the transmission of sensory and predictive information[59]. Our results indicate that whereas sensory

302  signals are convened early in middle layers, subsequent predictive signals are transmitted later as

303  feedback in deep and superficial layers. This specification opens a new path for further testing

304  predictions of predictive coding by investigating the content and integration of predictive feedback and

305  sensory feedforward signals[60–62]. For example, our results invite investigating layer-specific neural

306  dynamics at distinct frequency bands for feedforward and feedback processing as predicted from

307  human[63,64] and non-human invasive studies[59,65–67].

308

309  Our approach facilitates resolving the interplay of feedforward and feedback processing in a variety of

310  visual research contexts. It may allow dissecting the distinctive roles of feedforward and feedback

311  information flow crucial to cognitive functions such as attention, expectation and memory by identifying

312  and characterizing their distinct spatiotemporal dynamics. Similarly, using multi- rather than univariate

313  methods[12,59,68] to assess layer-specific fMRI extends content-sensitive analysis[69–71] from the macro- to

314  the mesoscale of cortical organization, allowing for the assessment of the representational contents of

315  feedforward and feedback information flow underlying human cognition.

316

317  ***Layer-specific EEG-fMRI fusion clarifies neural dynamics of high-level visual cortex regions***

318  Neural dynamics in response to faces[22], scenes[72] and objects[73,74] assessed at the macroscale of cortical

319  regions commonly display a double peak pattern (see also **Fig. 2B**), with a sharp, early peak around 100-

320  130 ms followed by a wide second peak around 200-450 ms, suggesting distinct contributing neural

321  circuits and kinds of processing. Our layer-specific results clarify the neural circuits[22] and functional

322  nature of the components underlying this pattern: it results from mixing the early, narrow-peaked and

323  neural dynamics at the middle layer ~160 ms of feedforward information flow with the later, wide-

324  peaked and thus more persistent dynamics at the superficial layer ~400 ms of feedback information flow.

325  Whereas the early peak latency matches that of sensory feedforward signals[75], the late peak dynamics

326    match that of and attention-[76], consciousness-[77] or task-related[78] feedback, originating from higher-

327    order brain regions, such as the frontal eye fields[10] or the prefrontal cortex[79], outside the visual cortex.

328    Thus, our results provide circuit-level mechanistic as well as functional interpretative guidance for

329    standard 3T human neuroimaging studies that typically cannot resolve visual information flow at this

330    fine spatiotemporal level.

331

332    ***Differential representational format across layer-specific dynamics in LOC implies interareal***

333    ***feedback mediating high-complexity visual features***

334    The analysis of the representational format of the neural dynamics in LOC revealed that early

335    feedforward signals primarily convey mid-complexity features, while feedback signals convey high-

336    complexity features. This indicates that feedback does not merely modulate the gain of pre-existing

337    features but is actively involved in the emergence of additional, more complex features that are absent

338    in feedforward processing[5]. This finding supports the notion that interareal feedback is integral to core

339    object recognition[80–82]. A potential source of this feedback might be the dorsolateral prefrontal cortex

340    (DLPFC)[83], whose silencing decreases feature complexity in monkey inferior temporal cortex[84], the

341    homologue of human LOC. Based on our results we predict that disrupting processing in human DLPFC,

342    e.g. through transcranial magnetic stimulation[85,86], will yield analogous effects specifically on

343    superficial, late cortical layer responses in LOC.

344

345    ***Limitations of the study***

346    We based our analyses on the GE-BOLD signal that is affected by locally nonspecific responses from

347    macrovasculature[41], compromising the estimation of the laminar response. To mitigate this effect, we

348    used Pearson's correlation as a scale-invariant measure and partialed out effects of lower cortical layers

349    when assessing the layers above. However, some residual bias may persist, particularly if the GE-BOLD

350    response across layers follows a point spread function[87] that a linear model cannot fully account for.

351    Human fMRI studies using contrasts with higher spatial specificity[88–90], optimized for condition-rich

352    experimental designs[21], are needed to confirm our findings.

353

354    ***Conclusion***

355    Understanding how the abundant feedforward and feedback connections in visual cortex mediate object

356    vision requires specifying their functional role. Here we provided two key advances in this regard. First,

357    leveraging layer-specific fMRI with EEG-fMRI fusion we dissociated temporally overlapping

358    feedforward and feedback processing in LOC and EVC, specifying their unique temporal profile.

359    Second, by assessing feature complexity, we showed that feedback in LOC actively increases the

360    complexity of LOC's feature format.

361

362    **Acknowledgments**

368

## Materials and Methods

370

### fMRI participants

372    10 adult volunteers (mean age 29.4 years; age range 20-37 years; 5 female) participated in the study and
373    provided written informed consent. The sample size was based on previous conventional layer fMRI
374    studies conducted at 7T[68,88,91]. All participants had normal or corrected-to-normal vision and no history
375    of neurological disorders. All participants received a monetary reward at the end of the study. The study
376    was approved by the Ethics Committee of the Faculty of Medicine at Leipzig University, Germany, and
377    conducted in accordance with the ethical principles of the Declaration of Helsinki, except for study
378    preregistration, which was not performed.

379

### Stimulus set

381    The stimulus set consisted of 24 naturalistic color images of everyday objects on real-world
382    backgrounds, each of a different category from the ImageNet image database (i.e. animals, tools,
383    vehicles, foods and others)[30].

384

### fMRI experimental design

386    The study was composed of a main experimental part and a localizer run.

387

388    The main experimental part consisted of 6 to 14 runs, each lasting 621 s. Each run started and ended
389    with a 10.5-s baseline period and included all 24 images, each repeated five times, for a total of 120
390    trials presented in random order. Each trial began with a 1-s stimulus-on interval, during which a
391    centrally displayed object image (5° visual angle) appeared on a grey background with a pink fixation
392    cross. This was followed by a 4-s stimulus-off interval, where only the fixation cross was shown.
393    Participants were instructed to maintain their gaze on the fixation cross throughout the experiment and
394    to perform a color-change detection task, pressing a button as soon as the fixation cross turned blue for
395    300 ms.

396

397    The functional localizer run was intended to define regions of interest (ROIs) EVC and LOC.
398    Participants completed it at the beginning of the recording session. The localizer consisted of 15-s blocks
399    displaying objects (not included in the main experiment) and scrambled objects overlayed with a fixation

400  cross, interleaved with 7.5-s baseline blocks showing only a fixation cross on a grey background. In
401  each block images were centrally presented (12° visual angle) for 400 ms, followed by a 350-ms display
402  of the fixation cross. Participants were instructed to maintain their gaze on the fixation cross and to press
403  a button if the same image appeared in consecutive trials. The localizer run included 12 blocks of each
404  image type, resulting in a total duration of 465 s. Block order was pseudo-randomized to avoid
405  immediate repetition of the same block type.
406

**MRI Procedure**

407
408  MRI data were acquired at the Max Planck Institute for Human Cognitive and Brain Sciences in Leipzig,
409  Germany. Four participants completed the experiment in one scanning session. Six participants
410  completed the experiment in two scanning sessions on two separate days. During the first scanning
411  session, we acquired a T1-weighted anatomical image, the functional localizer and 5 to 8 runs of the
412  main experiment. During the second scanning session, participants completed 6 to 8 runs of the main
413  experiment. Additionally, to enable distortion correction, five volumes with reversed phase-encoding
414  polarity were acquired following the first run of each main experimental session. To ensure that
415  participants were familiar with the experimental tasks, we provided them with verbal and written
416  instructions prior to the scanning, and the participants completed a 2-min training for both the localizer
417  and the main tasks.
418

**MRI acquisition parameters**

419
420  We acquired MR images on a Siemens Magnetom Terra 7T whole-body system (Siemens Healthineers,
421  Erlangen, Germany) with a single-channel-transmit and a 32-channel radio-frequency (RF) receive head
422  coil (Nova Medical Inc, Wilmington, USA). We acquired the functional data using a 2D Gradient-echo
423  (GE) echo planar imaging (EPI) sequence[92] (voxel size = 0.9 mm isotropic resolution, TE/TR =
424  26.2/3500 ms, in-plane field of view (FoV) 192 × 192 mm$^2$, 48 axial slices, flip angle = 75°, echo spacing
425  = 1.0 ms, GRAPPA factor = 3, partial Fourier = 6/8, phase encoding direction anterior-posterior). We
426  recorded anatomical data using an MP2RAGE sequence[93] (voxel size = 0.7 mm isotropic resolution,
427  TE/TR = 2.01/5590 ms, in-plane FoV 224 × 224 mm, GRAPPA factor = 2) yielding two inversion
428  contrasts (TI1 = 900 ms, flip angle 1 = 5°; TI2 = 2750 ms, flip angle 2 = 3°). The two inversion contrasts
429  were combined to produce T1-weighted MP2RAGE uniform (UNI) images with high contrast to noise
430  ratio.
431

**MRI preprocessing**

432
433  For each recording session, we spatially realigned the functional volumes to their mean volume using
434  SPM12 (http://www.fil.ion.ucl.ac.uk/spm). To correct for geometric distortions in the phase encoding
435  direction, we calculated a deformation field based on reverse gradient estimation, using the Advanced
436  Normalization Tools (ANTs) software package (http://stnava.github.io/ANTs/). In detail, we combined

the mean functional volume (forward image) with the mean volume acquired with opposing phase encoding direction to generate a distortion-corrected template. Next, we estimated the deformation map by registering the forward image to the corrected template reference using non-linear (SyN) transformations. Finally, the deformation map was used to produce a distortion-corrected mean functional volume.

To co-register the anatomical and the distortion-corrected mean functional volumes, we initially referred to the Glasser's atlas[94] and the Kanwisher's atlas[95] to identify the approximate location of EVC and LOC, respectively. Then we outlined a volume containing these regions in the occipito-temporal cortex of both hemispheres on the individual participant's native space (from here on referred to as manual mask) using ITK-SNAP with the 3D paintbrush tool (v.3.8)[96]. We then estimated a second deformation map in ANTs by registering the distortion-corrected mean volume to the T1-weighted volume, applying nonlinear (SyN) transformations within the manual mask. We visually inspected the fixed and registered volumes in ITK-SNAP for each participant. If the volume was not correctly registered within the region of the manual mask, we repeated the registration, adding linear transformations (rigid and affine) with stricter convergence criteria and increased iterations, until an accurate alignment was achieved. To minimize spatial resolution loss during resampling, we combined both deformation maps and resampled the functional images to the anatomical reference in a single interpolation step using a fifth-order spline function.

Finally, we spatially smoothed the functional localizer images using a 6-mm full width at half maximum (FWHM) Gaussian kernel. Functional images of the main runs were not smoothed to preserve spatial specificity.

Before segmenting the T1-weighted UNI volume into grey matter, white matter and cerebrospinal fluid following the procedure in ITK-SNAP outlined here (https://www.youtube.com/watch?v=tSA77mFTwcg&t=1042s), we corrected for bias field effects with a customized script[97]. Next, we divided the cortical ribbon into laminar and columnar compartments using LAYNII (v2.2.1)[98]. In detail, applying the equi-distant model[33], we segmented the gray matter into three cortical depths: deep, middle and superficial (**Fig. 3A**). Here, we used the term cortical layers to refer to the depth-dependent compartments along the cortical ribbon, distinct from the actual anatomical layers found in cortex. Additionally, we segmented the gray matter into columnar compartments within the manual masks.

**fMRI univariate analysis**

To estimate neural responses, we ran separate General Linear Model (GLM) analyses in SPM12 for each pre-processed functional run, i.e., all main experimental runs and the localizer run. All analyses

474  were conducted in each participant's native anatomical space. Specifically, we modelled 25 regressors

475  (i.e., 24 object images + baseline) for each main experimental run, and 3 regressors (i.e., objects,

476  scrambled objects, and baseline) for the localizer run. We created the regressors by convolving a boxcar

477  function representing the onsets and durations of the corresponding condition with the canonical (2

478  Gamma) hemodynamic response function (HRF). We incorporated the motion estimates into the model

479  as nuisance regressors. By fitting a GLM, we obtained beta weight estimates for each condition (i.e.,

480  regressor) for each run, which were subsequently used in further analyses.

481

482  **Definition of regions of interest**

483  At the macroscale, we defined two regions of interest (ROIs) for each participant: EVC, comprising

484  areas V1, V2, and V3, and LOC in a two-step procedure. First, we used anatomical masks from the

485  above-mentioned brain atlases for EVC[94] and for LOC[95]. These masks were resampled from the

486  MNI152 space into each participant's individual space. Second, we identified the overlap between the

487  participant-specific anatomical masks and the corresponding functional contrast T-statistic map from

488  the localizer experiment, retaining the top 2,000 voxels. Specifically, we ranked the voxels according to

489  the objects + scrambled > baseline contrast T-statistic for EVC, and the objects > scrambled contrast T-

490  statistic for LOC. Due to lack of activation in the lateral occipitotemporal cortex for one participant (3),

491  we used the objects + scrambled > baseline contrast to define LOC. Any voxels overlapping across ROIs

492  were excluded. This process resulted in one final EVC and LOC mask for each participant.

493

494  At the mesoscale, we defined six ROIs based on the two brain regions – EVC and LOC – and the three

495  cortical layers – deep, middle and superficial (see above for definition of cortical layers). Here, the ROI

496  definition was analogous, with the following deviations to address issues arising specifically at the

497  mesoscale level. Voxels with higher spatial resolution exhibit a reduced SNR and increased

498  susceptibility to noise sources[41], compromising the quality of the recorded signal. Additionally, signal

499  reliability gradually decreases along the ventral visual cortex from lower to higher visual areas[99], likely

500  due to signal loss and susceptibility-induced distortions[100]. To address this, we chose a higher number

501  of voxels for EVC compared to LOC. This resulted in 3000 and 1500 voxels with the highest T-statistic

502  within the cortical ribbon of EVC and LOC, respectively. We then assigned each voxel to one of the

503  three layers and selected only those sharing full overlap of columnar compartments.

504

505  **EEG data – paradigm, acquisition and analysis**

506  We used a subset of the EEG data ($N = 32$) collected by Xie et al., 2024[5] for the same images as in the

507  fMRI experiment. Below is a summary of the relevant experimental procedures and acquisition steps.

508

509  The EEG experiment employed and backward masking paradigm consisting of 2,544 trials. In each trial,

510  an object image was briefly presented for 17 ms and followed by a dynamic mask under two conditions:

16

511   an early mask condition with inter-stimulus interval (ISI) of 17 ms or a late mask condition with an ISI

512   of 600 ms. Object images and masks were randomly paired on each trial. All stimuli were centrally

513   displayed on a gray background, with a size of 5° visual angle, and overlaid with a bull's-eye fixation

514   symbol. Participants were instructed to maintain fixation and refrain from blinking during trials, except

515   during designated task trials (i.e., two-alternative forced-choice task, 25% of total trials) where blinking

516   was allowed after making a response. The inter-trial interval (ITI) ranged from 900 – 1,100 ms;

517   following the task trials, the ITI was extended to 2,000 ms to reduce motor artifacts.

518

519   For the present analyses, we focused on data from the late mask condition (ISI = 600 ms), yielding

520   approximately 53 trials per object image. We analyzed the time window from 200 ms pre-stimulus to

521   600 ms post-stimulus, before mask onset.

522

523   EEG data were recorded with a 64-electrode ActiCap system and a Brainvision actiChamp amplifier.

524   Electrodes were positioned according to the international 10-10 system, with a ground electrode and a

525   reference electrode placed on the scalp. The signals were sampled at 1,000 Hz and online filtered

526   between 0.03 and 100 Hz. EEG data were preprocessed using Brainstorm-3[101]. Noisy channels (mean =

527   2.2, SD = 1.8) were removed, and the data were low-pass filtered at 40 Hz. Independent component

528   analysis was applied to remove eye movement and other artifact components (mean = 2.7, SD = 0.9).

529   Data were then segmented into epochs from -200 ms to 600 ms relative to stimulus onset, baseline-

530   corrected, and multivariate noise normalization[102] was applied for decoding analyses.

531

532   Multivariate analysis was performed on a participant-specific basis using SVMs as implemented in the

533   LIBSVM toolbox in MATLAB (2021a). To determine when the brain processes object information,

534   time-resolved decoding analysis was conducted from -200 ms to 600 ms relative to target image onset

535   in 10 ms intervals. At each time point, trial-specific EEG channel activations were extracted and

536   arranged into 64-dimensional pattern vectors for each of the 24 object image conditions. For each

537   condition, trials were randomly grouped into four equally sized bins and averaged to create four pseudo-

538   trials, repeated over 100 permutations. Pseudo-trials were then divided into a training set (three pseudo-

539   trials) and a testing set (one pseudo-trial) for pairwise object identity decoding. This process was

540   repeated for all pairwise combinations of object conditions.

541

542   **Multivariate pattern analysis of fMRI data**

543   To decode object information from voxel activation patterns in EVC and LOC, we used SVMs as

544   implemented in the scikit-learn library in Python. Analyses were conducted separately for each

545   participant and for EVC and LOC at the macroscale. We aggregated run-specific voxel-wise beta

546   estimates derived from the GLM into pattern vectors for each of the 24 conditions. To enhance the signal

547   to noise ratio, we averaged beta weights across two randomly selected subsets of trials into two pseudo-

548  trials for each condition. We repeated this process 300 times, each time performing pairwise object

549  decoding, for all object condition combinations, using a two-fold cross-validation approach. We

550  averaged the results across all iterations and all object condition combinations to yield one grand-

551  average decoding accuracy for each ROI and participant.

552

553  **Representational similarity analysis**

554  To relate object representations across different signal spaces (i.e., voxel activation patterns in fMRI,

555  sensor activation patterns in EEG, embeddings in DNNs), we used representational similarity analysis[21].

556  This approach is based on the rationale that if two images elicit similar neural representations, their

557  corresponding fMRI, EEG or DNN signal patterns should also be similar. We used two variants of RSA:

558  (i) representational similarity analysis-based fusion[21,23,24], relating spatially localized object

559  representations (from fMRI) to specific temporal dynamics (from EEG) to resolve the spatiotemporal

560  dynamics with which visual representations emerge, and ii) representational similarity analysis-based

561  commonality analysis[103,104] to determine the visual feature complexity (from layer-specific DNN

562  embeddings) of spatiotemporally identified dynamics. For this we calculated the common variance

563  between spatially localized object representations (from fMRI), specific temporal dynamics (from

564  EEG), and layer embeddings (from DNNs). We detail both approaches below after describing the

565  specifics of summarizing representational similarity for each signal space in representational

566  dissimilarity matrices (RDMs).

567

568  *Construction of fMRI representational dissimilarity matrices*

569  To construct the fMRI-RDMs, we first extracted and vectorized beta weight estimates for each of the

570  24 conditions to form fMRI neural patterns for a given ROI (area or cortical layer) and applied

571  multivariate noise normalization in order to enhance the SNR. We then quantified the extent of pattern

572  similarity by calculating the Pearson's correlation for all pairwise combinations of experimental

573  conditions. Output correlations were transformed into dissimilarity values using 1 – Pearson's

574  correlation, and organized into a 24×24 fMRI representational dissimilarity matrix (RDM), indexed in

575  rows and columns by the compared conditions. For each participant, this yielded two fMRI-RDMs at

576  the macroscale level of brain areas: the EVC RDM and LOC RDM; and six layer-specific fMRI-RDMs

577  at the mesoscale level of cortical layers: Deep EVC RDM, Middle EVC RDM, Superficial EVC RDM,

578  and Deep LOC RDM, Middle LOC RDM and Superficial LOC RDM. For all further analyses we

579  averaged the MRI-RDMs across participants.

580

581  *Construction of EEG representational dissimilarity matrices*

582  A detailed description of the creation of the EEG-RDMs is provided by Xie et. al., (2024)[5]. Briefly, for

583  each participant and each time point in the epoch from -200 to 600 ms, we pairwise decoded object

584  information using EEG channel activation patterns, via SVMs as implemented in the LIBSVM

585 toolbox[105] in MATLAB (2021a). The decoding accuracies for each pair of object images were assembled

586 into 24×24 EEG-RDM for each time point (801 in total), with rows and columns indexed by the

587 conditions. Each matrix was symmetric across the diagonal, with diagonal entries left undefined.

588

589 *Construction of DNN representational dissimilarity matrices*

590 To compute DNN-RDMs, we used the AlexNet architecture[31] trained on visual object categorization on

591 the ImageNet dataset[106]. In detail, we fed the pre-trained network with the same set of 24 object images

592 as input and extracted the activation pattern for each object condition from each of the five convolutional

593 layers and the two subsequent fully-connected layers. Next, we quantified the pairwise pattern

594 dissimilarity using 1 – Pearson's correlation for all pairs of image combinations and organized the

595 outputs into 24×24 DNN-RDMs, with rows and columns representing Pearson's based dissimilarity

596 measure between object conditions. This resulted in a total of 7 layer-specific DNN RDMs.

597

598 *Representational similarity analysis-based EEG-fMRI fusion*

599 We used representational similarity analysis-based fusion[21,23,24], relating fMRI-RDMs to EEG-RDMs

600 using Spearman rank order correlation. We applied this analysis at the macroscale and at the mesoscale,

601 for each participant separately. At the macroscale, we related region-specific RDMs (i.e., EEG, LOC)

602 to EEG-RDMs, yielding one time course for each ROI for each participant. At the mesoscale we related

603 region- and layers-specific ROIs to EEG-RDMs, yielding one time course for each of the six ROI-layer

604 combinations for each participant.

605

606 *Representational similarity analysis-based commonality analysis*

607 To characterize the level of feature complexity of the neural representations identified above in space

608 and time, we related them to DNN embeddings across model layers. Feature complexity increases

609 progressively across DNN layers: early model layers are more sensitive to low complexity features,

610 whereas later layers are tuned to high-complexity features[107,108], analogous to the ventral visual

611 hierarchy in humans and non-human primates[30,32]. Thus, by linking the neural representations to those

612 from the DNN layers, we can determine the feature complexity of spatially localized object

613 representations to particular time points. We did this using representational similarity analysis-based

614 commonality analysis[109]. In detail, we calculated the commonality coefficient corresponding to the

615 shared variance among fMRI-RDMs for each brain region and cortical layer, EEG-RDMs at time points

616 with significant effects in EEG-fMRI fusion (averaged across time), and DNN RDMs for each model

617 layer.

618

619 **Quantification and statistical analysis**

620 To assess statistical significance, we performed non-parametric statistical analyses that do not rely on

621 assumptions about the data distribution. The empirical estimations, including decoding accuracy as well

as correlation values from the representational similarity analysis-based fusion and from the commonality analysis, were tested against a null distribution created by sign permutation (1,000 permutations). We randomly multiplied the participant-specific data by ±1 to generate permutation samples, recomputed the statistic for each sample, and derived $P$ values.

We controlled the familywise error rate across time points using cluster-based statistics. First, $P$-value maps were thresholded at $P < 0.05$ to define temporally contiguous suprathreshold clusters. These clusters were then used to construct an empirical null distribution of maximum cluster weights, and a corrected threshold was determined at the 95th percentile of the right tail of this distribution. For the commonality analyses we corrected $P$ values for multiple comparisons by FDR-correction.

We estimated the 95% confidence intervals for the peak latencies in the RSA-derived time courses using bootstrapping (1,000 paired bootstrap resamples with replacement). For each bootstrap we calculated the statistic, resulting in a bootstrap estimate of peak latencies from which we derived the confidence intervals.

To estimate confidence intervals for peak latency differences, we used an analogous bootstrapping approach, resampling the mean peak-to-peak latency difference for each resample. This generated a distribution of mean differences, from which we derived the 95% confidence interval. We set $P < 0.05$, i.e., if the 95% confidence interval did not include 0, we rejected the null hypothesis of no peak-to-peak latency differences.

## References

1. Mishkin, M., Ungerleider, L. G. & Macko, K. A. Object vision and spatial vision: two cortical pathways. *Trends in Neurosciences* **6**, 414–417 (1983).

2. Kravitz, D. J., Saleem, K. S., Baker, C. I., Ungerleider, L. G. & Mishkin, M. The ventral visual pathway: an expanded neural framework for the processing of object quality. *Trends in Cognitive Sciences* **17**, 26–49 (2013).

3. Grill-Spector, K., Kourtzi, Z. & Kanwisher, N. The lateral occipital complex and its role in object recognition. *Vision Research* **41**, 1409–1422 (2001).

4. Lamme, V. A., Supèr, H. & Spekreijse, H. Feedforward, horizontal, and feedback processing in the visual cortex. *Current Opinion in Neurobiology* **8**, 529–535 (1998).

5. Xie, S., Singer, J., Yilmaz, B., Kaiser, D. & Cichy, R. M. The representational nature of spatio-temporal recurrent processing in visual object recognition. 2024.07.30.605751 Preprint at https://doi.org/10.1101/2024.07.30.605751 (2024).

6.   Lamme, V. A. F., Zipser, K. & Spekreijse, H. Figure-ground activity in primary visual cortex is suppressed by anesthesia. *Proceedings of the National Academy of Sciences* **95**, 3263–3268 (1998).

7.   Markov, N. T. *et al.* Anatomy of hierarchy: Feedforward and feedback pathways in macaque visual cortex. *Journal of Comparative Neurology* **522**, 225–259 (2014).

8.   Felleman, D. & Van Essen, D. Distributed hierarchical processing in the primate cerebral cortex. *Cereb Cortex* **1**, 1–47 (1991).

9.   Kreiman, G. & Serre, T. Beyond the feedforward sweep: feedback computations in the visual cortex. *Annals of the New York Academy of Sciences* **1464**, 222–241 (2020).

10.  Hüer, J., Saxena, P. & Treue, S. Pathway-selective optogenetics reveals the functional anatomy of top–down attentional modulation in the macaque visual cortex. *Proc. Natl. Acad. Sci. U.S.A.* **121**, e2304511121 (2024).

11.  Nurminen, L., Merlin, S., Bijanzadeh, M., Federer, F. & Angelucci, A. Top-down feedback controls spatial summation and response amplitude in primate visual cortex. *Nature Communications* **9**, 2281 (2018).

12.  Bergmann, J. *et al.* Cortical depth profiles in primary visual cortex for illusory and imaginary experiences. *Nature Communications* **15**, 1002 (2024).

13.  Lawrence, S. J., Norris, D. G. & de Lange, F. P. Dissociable laminar profiles of concurrent bottom-up and top-down modulation in the human visual cortex. *eLife* **8**, e44422 (2019).

14.  Carricarte, T. *et al.* Laminar dissociation of feedforward and feedback in high-level ventral visual cortex during imagery and perception. *iScience* **27**, 110229 (2024).

15.  Dowdle, L. *et al.* Characterizing top-down microcircuitry of complex human behavior across different levels of the visual hierarchy. Preprint at https://doi.org/10.1101/2022.12.03.518973 (2023).

16.  Liu, C. *et al.* Layer-dependent multiplicative effects of spatial attention on contrast responses in human early visual cortex. *Progress in Neurobiology* **207**, 101897 (2021).

17.  Rockland, K. S. & Pandya, D. N. Laminar origins and terminations of cortical connections of the occipital lobe in the rhesus monkey. *Brain Research* **179**, 3–20 (1979).

18.  Rockland, K. S. What do we know about laminar connectivity? *NeuroImage* **197**, 772–784 (2019).

19.  Siu, C., Balsor, J., Merlin, S., Federer, F. & Angelucci, A. A direct interareal feedback-to-feedforward circuit in primate visual cortex. *Nat Commun* **12**, 4911 (2021).

20.  Federer, F., Ta'afua, S., Merlin, S., Hassanpour, M. S. & Angelucci, A. Stream-specific feedback inputs to the primate primary visual cortex. *Nature Communications* **12**, 228 (2021).

21. Kriegeskorte, N., Mur, M. & Bandettini, P. A. Representational similarity analysis - connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience* **2**, (2008).

22. Cichy, R. M., Pantazis, D. & Oliva, A. Resolving human object recognition in space and time. *Nature Neuroscience* **17**, 455–462 (2014).

23. Cichy, R. M. & Oliva, A. A M/EEG-fMRI Fusion Primer: Resolving Human Brain Responses in Space and Time. *Neuron* **107**, 772–781 (2020).

24. Cichy, R. M., Pantazis, D. & Oliva, A. Similarity-Based Fusion of MEG and fMRI Reveals Spatio-Temporal Dynamics in Human Cortex During Visual Object Recognition. *Cerebral Cortex* **26**, 3563–3579 (2016).

25. Essen, D. C. V. & Maunsell, J. H. R. Hierarchical organization and functional streams in the visual cortex. *Trends in Neurosciences* **6**, 370–375 (1983).

26. Singer, J. J. D., Karapetian, A., Hebart, M. N. & Cichy, R. M. Identifying and characterizing scene representations relevant for categorization behavior. *Imaging Neuroscience* **3**, imag_a_00449 (2025).

27. DiCarlo, J. J., Zoccolan, D. & Rust, N. C. How Does the Brain Solve Visual Object Recognition? *Neuron* **73**, 415–434 (2012).

28. Kourtzi, Z. & Connor, C. E. Neural Representations for Object Perception: Structure, Category, and Adaptive Coding. *Annual Review of Neuroscience* vol. 34 45–67 (2011).

29. Muukkonen, I., Ölander, K., Numminen, J. & Salmela, V. R. Spatio-temporal dynamics of face perception. *NeuroImage* **209**, 116531 (2020).

30. Cichy, R. M., Khosla, A., Pantazis, D., Torralba, A. & Oliva, A. Comparison of deep neural networks to spatio-temporal cortical dynamics of human visual object recognition reveals hierarchical correspondence. *Scientific Reports* **6**, 27755 (2016).

31. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet Classification with Deep Convolutional Neural Networks. in *Advances in Neural Information Processing Systems* (eds. Pereira, F., Burges, C. J., Bottou, L. & Weinberger, K. Q.) vol. 25 (Curran Associates, Inc., 2012).

32. Güçlü, U. & van Gerven, M. A. J. Deep Neural Networks Reveal a Gradient in the Complexity of Neural Representations across the Ventral Stream. *J. Neurosci.* **35**, 10005 (2015).

33. Waehnert, M. D. *et al.* Anatomically motivated modeling of cortical laminae. *NeuroImage* **93**, 210–220 (2014).

34. Barth, M. & Norris, D. G. Very high-resolution three-dimensional functional MRI of the human visual cortex with elimination of large venous vessels. *NMR in Biomedicine* **20**, 477–484 (2007).

721  35.  Markuerkiaga, I., Barth, M. & Norris, D. G. A cortical vascular model for examining the specificity of the laminar BOLD signal. *NeuroImage* **132**, 491–498 (2016).

723  36.  van Kerkoerle, T., Self, M. W. & Roelfsema, P. R. Layer-specificity in the effects of attention and working memory on activity in primary visual cortex. *Nature Communications* **8**, 13804 (2017).

725  37.  Self, M. W., van Kerkoerle, T., Supèr, H. & Roelfsema, P. R. Distinct Roles of the Cortical Layers of Area V1 in Figure-Ground Segregation. *Current Biology* **23**, 2121–2129 (2013).

727  38.  Barzegaran, E. & Plomp, G. Four concurrent feedforward and feedback networks with different roles in the visual cortical hierarchy. *PLOS Biology* **20**, e3001534 (2022).

729  39.  Grossberg, S. Linking the laminar circuits of visual cortex to visual perception: Development, grouping, and attention. *Neuroscience & Biobehavioral Reviews* **25**, 513–526 (2001).

731  40.  Smith, M. A., Jia, X., Zandvakili, A. & Kohn, A. Laminar dependence of neuronal correlations in visual cortex. *Journal of Neurophysiology* **109**, 940–947 (2013).

733  41.  Polimeni, J. R., Fischl, B., Greve, D. N. & Wald, L. L. Laminar analysis of 7T BOLD using an imposed spatial activation pattern in human V1. *NeuroImage* **52**, 1334–1346 (2010).

735  42.  Schwiedrzik, C. M. & Freiwald, W. A. High-Level Prediction Signals in a Low-Level Area of the Macaque Face-Processing Hierarchy. *Neuron* **96**, 89-97.e4 (2017).

737  43.  Stecher, R. & Kaiser, D. Representations of imaginary scenes and their properties in cortical alpha activity. *Scientific Reports* **14**, 12796 (2024).

739  44.  Morgan, A. T., Petro, L. S. & Muckli, L. Scene Representations Conveyed by Cortical Feedback to Early Visual Cortex Can Be Described by Line Drawings. *J. Neurosci.* **39**, 9410 (2019).

741  45.  Papale, P. *et al.* The representation of occluded image regions in area V1 of monkeys and humans. *Current Biology* **33**, 3865-3871.e3 (2023).

743  46.  Del Rosario, J. *et al.* Lateral inhibition in V1 controls neural and perceptual contrast sensitivity. *Nature Neuroscience* (2025) doi:10.1038/s41593-025-01888-4.

745  47.  Gilbert, C. D. & Wiesel, T. N. Morphology and intracortical projections of functionally characterised neurones in the cat visual cortex. *Nature* **280**, 120–125 (1979).

747  48.  Rockland, K. S. & Lund, J. S. Widespread Periodic Intrinsic Connections in the Tree Shrew Visual Cortex. *Science* **215**, 1532–1534 (1982).

749  49.  Gilbert, C. & Wiesel, T. Clustered intrinsic connections in cat visual cortex. *J. Neurosci.* **3**, 1116 (1983).

751  50.  Stettler, D. D., Das, A., Bennett, J. & Gilbert, C. D. Lateral Connectivity and Contextual Interactions in Macaque Primary Visual Cortex. *Neuron* **36**, 739–750 (2002).

753  51. Fujita, I. & Fujita, T. Intrinsic connections in the macaque inferior temporal cortex. *Journal of*
754  *Comparative Neurology* **368**, 467–486 (1996).

755  52. Grootswagers, T., Robinson, A. K., Shatek, S. M. & Carlson, T. A. Mapping the dynamics of
756  visual feature coding: Insights into perception and integration. *PLOS Computational Biology* **20**,
757  e1011760 (2024).

758  53. Gallimore, C. G., Ricci, D. A. & Hamm, J. P. Spatiotemporal dynamics across visual cortical
759  laminae support a predictive coding framework for interpreting mismatch responses. *Cerebral*
760  *Cortex* **33**, 9417–9428 (2023).

761  54. Dobs, K., Isik, L., Pantazis, D. & Kanwisher, N. How face perception unfolds over time. *Nature*
762  *Communications* **10**, 1258 (2019).

763  55. Grill-Spector, K. & Malach, R. The human visual cortex. *Annual Review of Neuroscience* vol. 27
764  649–677 (2004).

765  56. Grill-Spector, K. *et al.* A sequence of object-processing stages revealed by fMRI in the human
766  occipital lobe. *Human Brain Mapping* **6**, 316–328 (1998).

767  57. Callaway, E. M. Feedforward, feedback and inhibitory connections in primate visual cortex.
768  *Neural Networks* **17**, 625–632 (2004).

769  58. Rao, R. P. N. & Ballard, D. H. Predictive coding in the visual cortex: a functional interpretation
770  of some extra-classical receptive-field effects. *Nat Neurosci* **2**, 79–87 (1999).

771  59. Scheeringa, R. & Fries, P. Cortical layers, rhythms and BOLD signals. *NeuroImage* **197**, 689–698
772  (2019).

773  60. Larkum, M. E., Zhu, J. J. & Sakmann, B. A new cellular mechanism for coupling inputs arriving
774  at different cortical layers. *Nature* **398**, 338–341 (1999).

775  61. Aru, J., Suzuki, M. & Larkum, M. E. Cellular Mechanisms of Conscious Processing. *Trends in*
776  *Cognitive Sciences* **24**, 814–825 (2020).

777  62. Schuman, B., Dellal, S., Prönneke, A., Machold, R. & Rudy, B. Neocortical Layer 1: An Elegant
778  Solution to Top-Down and Bottom-Up Integration. *Annual Review of Neuroscience* **44**, 221–252
779  (2021).

780  63. Chen, L., Cichy, R. M. & Kaiser, D. Alpha-frequency feedback to early visual cortex orchestrates
781  coherent naturalistic vision. *Science Advances* **9**, eadi2321 (2023).

782  64. Xie, S., Kaiser, D. & Cichy, R. M. Visual Imagery and Perception Share Neural Representations
783  in the Alpha Frequency Band. *Current Biology* **30**, 2621-2627.e5 (2020).

784  65. Bastos, A. M. *et al.* Visual Areas Exert Feedforward and Feedback Influences through Distinct
785  Frequency Channels. *Neuron* **85**, 390–401 (2015).

66. van Kerkoerle, T. *et al.* Alpha and gamma oscillations characterize feedback and feedforward processing in monkey visual cortex. *Proceedings of the National Academy of Sciences* **111**, 14332–14341 (2014).

67. Michalareas, G. *et al.* Alpha-Beta and Gamma Rhythms Subserve Feedback and Feedforward Influences among Human Visual Cortical Areas. *Neuron* **89**, 384–397 (2016).

68. Muckli, L. *et al.* Contextual Feedback to Superficial Layers of V1. *Current Biology* **25**, 2690–2695 (2015).

69. Haynes, J.-D. & Rees, G. Decoding mental states from brain activity in humans. *Nature Reviews Neuroscience* **7**, 523–534 (2006).

70. Kriegeskorte, N. & Kievit, R. A. Representational geometry: integrating cognition, computation, and the brain. *Trends in Cognitive Sciences* **17**, 401–412 (2013).

71. Mur, M., Bandettini, P. A. & Kriegeskorte, N. Revealing representational content with pattern-information fMRI—an introductory guide. *Social Cognitive and Affective Neuroscience* **4**, 101–109 (2009).

72. Henriksson, L., Mur, M. & Kriegeskorte, N. Rapid Invariant Encoding of Scene Layout in Human OPA. *Neuron* **103**, 161-171.e3 (2019).

73. Mohsenzadeh, Y., Qin, S., Cichy, R. M. & Pantazis, D. Ultra-Rapid serial visual presentation reveals dynamics of feedforward and feedback processes in the ventral visual pathway. *eLife* **7**, e36329 (2018).

74. Motlagh, S. C., Joanisse, M., Wang, B. & Mohsenzadeh, Y. Unveiling the neural dynamics of conscious perception in rapid object recognition. *NeuroImage* **296**, 120668 (2024).

75. Hung, C. P., Kreiman, G., Poggio, T. & DiCarlo, J. J. Fast readout of object identity from macaque inferior temporal cortex. *Science* **310**, 863–866 (2005).

76. Hopf, J.-M. *et al.* Neural Sources of Focused Attention in Visual Search. *Cerebral Cortex* **10**, 1233–1241 (2000).

77. Cul, A. D., Baillet, S. & Dehaene, S. Brain Dynamics Underlying the Nonlinear Threshold for Access to Consciousness. *PLOS Biology* **5**, e260 (2007).

78. Duan, Y., Zhan, J., Gross, J., Ince, R. A. A. & Schyns, P. G. Pre-frontal cortex guides dimension-reducing transformations in the occipito-ventral pathway for categorization behaviors. *Current Biology* **34**, 3392-3404.e5 (2024).

79. Bar, M. A Cortical Mechanism for Triggering Top-Down Facilitation in Visual Object Recognition. *Journal of Cognitive Neuroscience* **15**, 600–609 (2003).
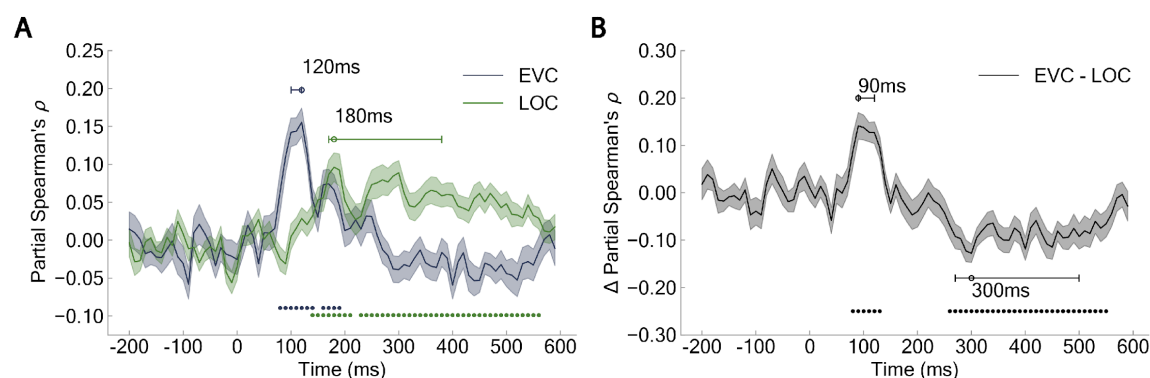
80. Kietzmann, T. C. *et al.* Recurrence is required to capture the representational dynamics of the human visual system. *Proceedings of the National Academy of Sciences* **116**, 21854–21863 (2019).

81. Kar, K., Kubilius, J., Schmidt, K., Issa, E. B. & DiCarlo, J. J. Evidence that recurrent circuits are critical to the ventral stream's execution of core object recognition behavior. *Nature Neuroscience* **22**, 974–983 (2019).

82. von Seth, J., Nicholls, V. I., Tyler, L. K. & Clarke, A. Recurrent connectivity supports higher-level visual and semantic object representations in the brain. *Commun Biol* **6**, 1–15 (2023).

83. Hamm, J. P., Shymkiv, Y., Han, S., Yang, W. & Yuste, R. Cortical ensembles selective for context. *Proceedings of the National Academy of Sciences* **118**, e2026179118 (2021).

84. Kar, K. & DiCarlo, J. J. Fast Recurrent Processing via Ventrolateral Prefrontal Cortex Is Needed by the Primate Ventral Stream for Robust Core Visual Object Recognition. *Neuron* **109**, 164-176.e5 (2021).

85. Hallett, M. Transcranial Magnetic Stimulation: A Primer. *Neuron* **55**, 187–199 (2007).

86. Bergmann, T. O. *et al.* Concurrent TMS-fMRI for causal network perturbation and proof of target engagement. *NeuroImage* **237**, 118093 (2021).

87. Markuerkiaga, I., Marques, J. P., Gallagher, T. E. & Norris, D. G. Estimation of laminar BOLD activation profiles using deconvolution with a physiological point spread function. *J Neurosci Methods* **353**, 109095 (2021).

88. Huber, L. *et al.* High-Resolution CBV-fMRI Allows Mapping of Laminar Activity and Connectivity of Cortical Input and Output in Human M1. *Neuron* **96**, 1253-1263.e7 (2017).

89. Huber, L. *et al.* Techniques for blood volume fMRI with VASO: From low-resolution mapping towards sub-millimeter layer-dependent applications. *NeuroImage* **164**, 131–143 (2018).

90. Pizzuti, A. *et al.* Imaging the columnar functional organization of human area MT+ to axis-of-motion stimuli using VASO at 7 Tesla. *Cerebral Cortex* **33**, 8693–8711 (2023).

91. Kok, P., Bains, L. J., van Mourik, T., Norris, D. G. & de Lange, F. P. Selective Activation of the Deep Layers of the Human Primary Visual Cortex by Top-Down Feedback. *Current Biology* **26**, 371–376 (2016).

92. Moeller, S. *et al.* Multiband multislice GE-EPI at 7 tesla, with 16-fold acceleration using partial parallel imaging with application to high spatial and temporal whole-brain fMRI. *Magnetic Resonance in Medicine* **63**, 1144–1153 (2010).

93. Marques, J. P. *et al.* MP2RAGE, a self bias-field corrected sequence for improved segmentation and T1-mapping at high field. *NeuroImage* **49**, 1271–1281 (2010).

94. Glasser, M. F. *et al.* A multi-modal parcellation of human cerebral cortex. *Nature* **536**, 171–178 (2016).

95. Julian, J. B., Fedorenko, E., Webster, J. & Kanwisher, N. An algorithmic method for functionally defining regions of interest in the ventral visual pathway. *NeuroImage* **60**, 2357–2364 (2012).

96. Yushkevich, P. A. *et al.* User-Guided Segmentation of Multi-modality Medical Imaging Datasets with ITK-SNAP. *Neuroinformatics* **17**, 83–102 (2019).

97. Lüsebrink, F., Sciarra, A., Mattern, H., Yakupov, R. & Speck, O. T1-weighted in vivo human whole brain MRI dataset with an ultrahigh isotropic resolution of 250 μm. *Scientific Data* **4**, 170032 (2017).

98. Huber, L. (Renzo) *et al.* LayNii: A software suite for layer-fMRI. *NeuroImage* **237**, 118091 (2021).

99. Badwal, M. W., Bergmann, J., Roth, J., Doeller, C. F. & Hebart, M. N. The scope and limits of fine-grained image and category information in the ventral visual pathway. *J. Neurosci.* **45**, e0936242024 (2024).

100. Malekian, V. *et al.* Mitigating susceptibility-induced distortions in high-resolution 3DEPI fMRI at 7T. *NeuroImage* **279**, 120294 (2023).

101. Tadel, F., Baillet, S., Mosher, J. C., Pantazis, D. & Leahy, R. M. Brainstorm: A User-Friendly Application for MEG/EEG Analysis. *Computational Intelligence and Neuroscience* **2011**, 879716 (2011).

102. Guggenmos, M., Sterzer, P. & Cichy, R. M. Multivariate pattern analysis for MEG: A comparison of dissimilarity measures. *NeuroImage* **173**, 434–447 (2018).

103. Mood, A. M. Partitioning Variance in Multiple Regression Analyses as a Tool For Developing Learning Models. *American Educational Research Journal* **8**, 191–202 (1971).

104. Reichwein Zientek, L. & Thompson, B. Commonality Analysis: Partitioning Variance to Facilitate Better Understanding of Data. *Journal of Early Intervention* **28**, 299–307 (2006).

105. Chang, C.-C. & Lin, C.-J. LIBSVM: A library for support vector machines. *ACM Trans. Intell. Syst. Technol.* **2**, 27:1-27:27 (2011).

106. J. Deng *et al.* ImageNet: A large-scale hierarchical image database. in *2009 IEEE Conference on Computer Vision and Pattern Recognition* 248–255 (2009). doi:10.1109/CVPR.2009.5206848.

107. Zeiler, M. D. & Fergus, R. Visualizing and Understanding Convolutional Networks. Preprint at https://doi.org/10.48550/arXiv.1311.2901 (2013).

108. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A. & Torralba, A. Object Detectors Emerge in Deep Scene CNNs. Preprint at https://doi.org/10.48550/arXiv.1412.6856 (2015).
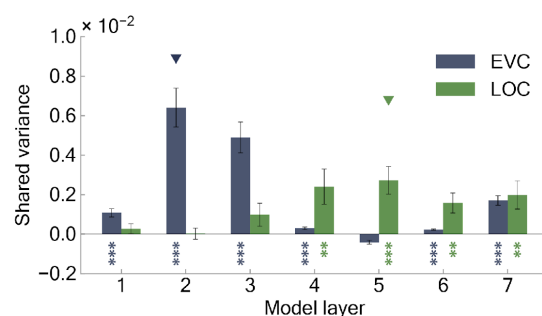
884    109. Hebart, M. N., Bankson, B. B., Harel, A., Baker, C. I. & Cichy, R. M. The representational
885        dynamics of task and object processing in humans. *Elife* **7**, e32816 (2018).
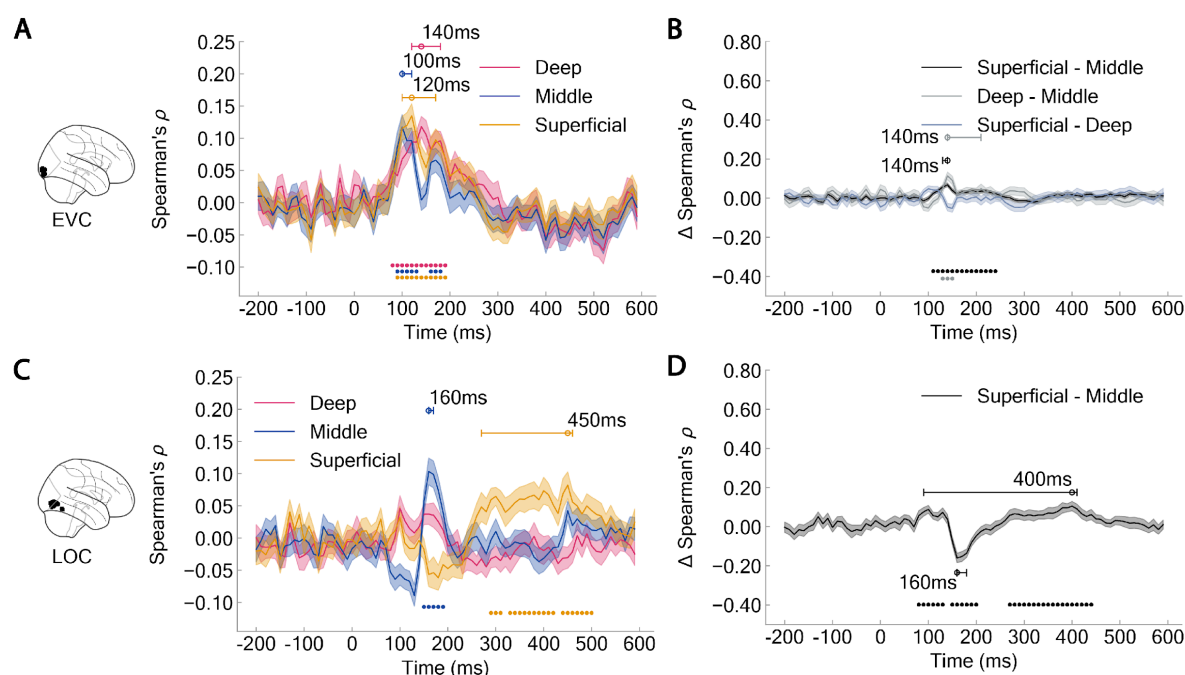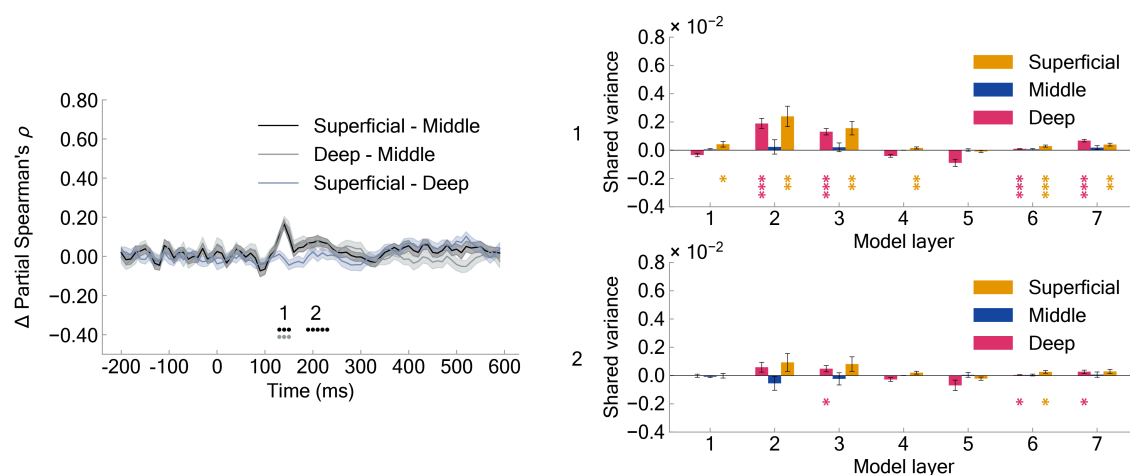886

## Supplementary information



**Supplementary Figure 1. Representational EEG-fMRI fusion at the macroscale using partial correlations.** For EVC we partialed out the effect of LOC and for LOC we partialed out the effect of EVC. Early representations correlated stronger with EVC than with LOC and later representations stronger with LOC than with EVC. Shaded area indicates the standard error of the mean across participants; colored circles indicate significant time points ($N = 32$, cluster-defining threshold $P < 0.05$, cluster threshold $P < 0.05$); uncolored circles and horizontal lines indicate peak latency means and 95% confidence intervals, respectively.

**Supplementary Figure 2. Commonality analysis at the macroscale based on significant time intervals of curve EVC – LOC differences.** Visual representations of low-complexity emerge primarily in EVC, while mid-to-high-level object representations emerge in LOC. Colored asterisks indicate significant correlations ($N = 32$, right-tailed permutation tests, FDR-corrected; $*P < 0.05$; $**P < 0.01$; $***P < 0.001$); colored triangles represent model layers with the highest occurrence proportion, determined through 1,000-iteration bootstraps.
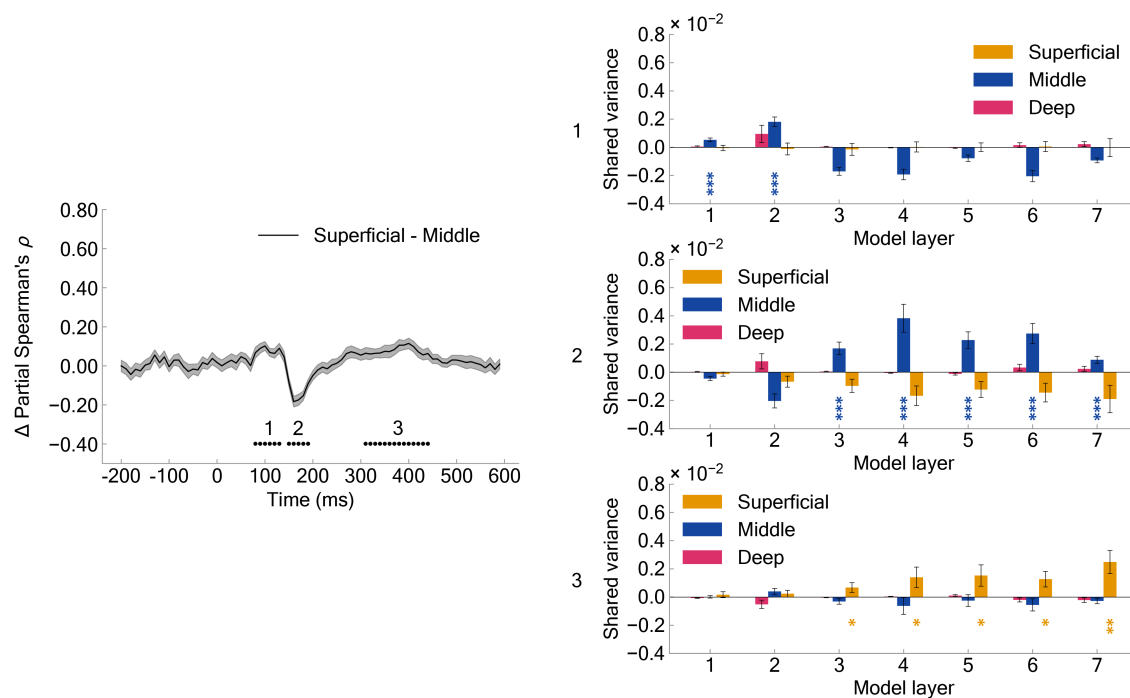
**Supplementary Figure 3. Representational EEG-fMRI fusion at the mesoscale without partial correlations.** Object representations derived from GE-BOLD signals are strongly influenced by non-specific macrovascular responses, which can compromise layer-specificity. To address this, we applied partial rank-order Spearman correlation, reducing the laminar influence of the layers beneath. Specifically, for the superficial layer, we partialed out the effect of the middle layer. Similarly, for the middle layer, we partialed out the effect of the deep layer. The deep layer remained unaffected by this approach. (**A**) Early representations emerged across layers in EVC, while late representations appeared in middle and superficial layers. (**B**) Early representations emerged in middle layers of LOC, whereas late representations appeared in superficial layers. Shaded area indicates the standard error of the mean across participants; colored circles indicate significant time points ($N = 32$, cluster-defining threshold $P < 0.05$, cluster threshold $P < 0.05$); uncolored circles and horizontal lines indicate peak latency means and 95% confidence intervals, respectively.
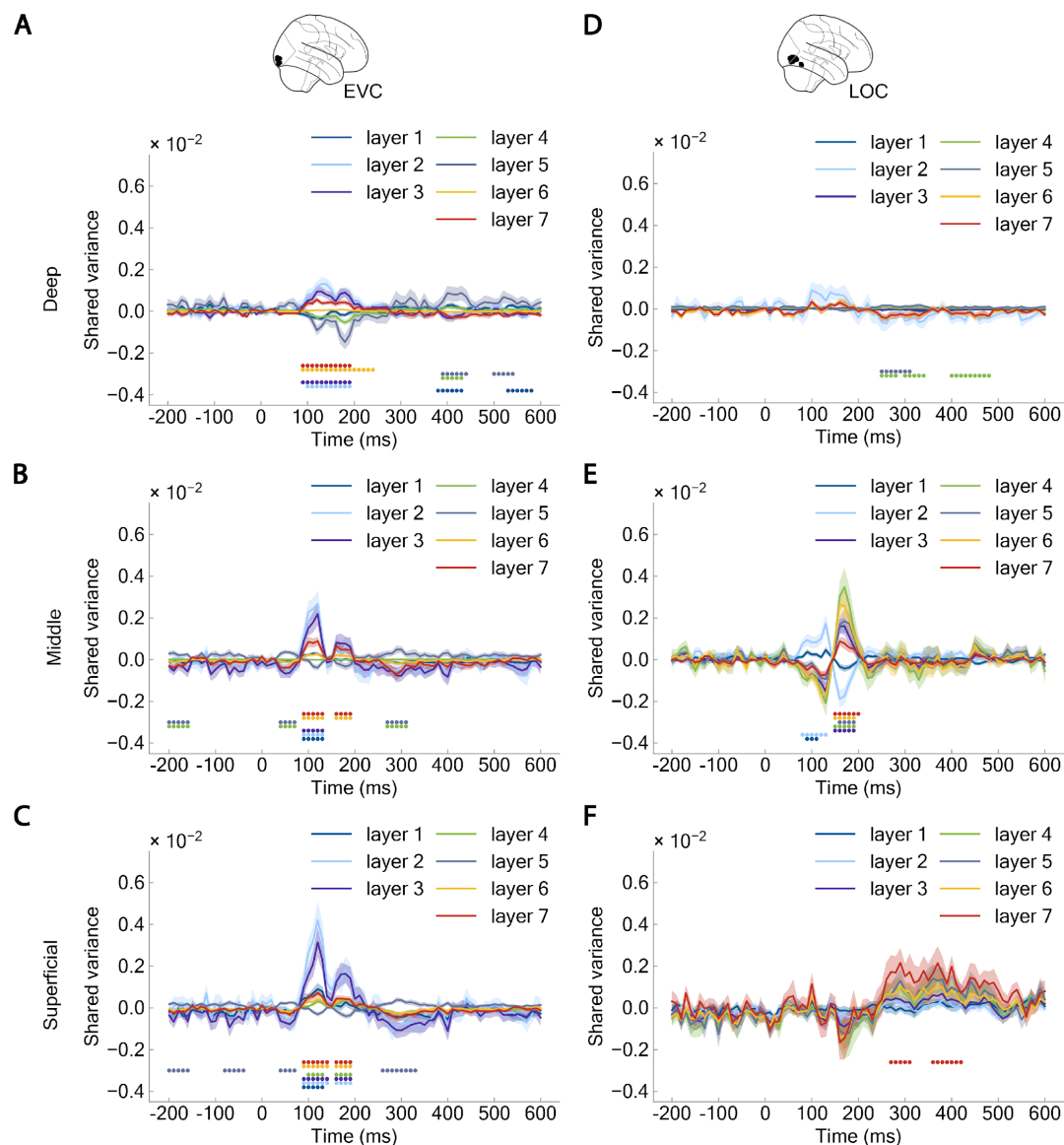
**Supplementary Figure 4. Commonality analysis in EVC at the mesoscale based on significant time intervals of curve differences across cortical layers.** We defined two clusters (labeled 1 and 2) from the significant time intervals. For each cluster, we performed commonality analysis by linking each model layer-specific DNN-RDM to each layer-specific fMRI-RDM in EVC and LOC and the mean EEG-RDM at the time interval with significant layer-specific temporal dynamics. Error bars indicate the standard error of the mean across participants; colored asterisks indicate significant correlations ($N = 32$, right-tailed permutation tests, FDR-corrected; *$P < 0.05$; **$P < 0.01$; ***$P < 0.001$).

**Supplementary Figure 5. Commonality analysis in LOC at the mesoscale based on significant time intervals of curve differences across cortical layers.** We defined three clusters (labeled 1, 2 and 3) from the significant time intervals. For each cluster, we performed commonality analysis by linking each model layer-specific DNN-RDM to each layer-specific fMRI-RDM in EVC and LOC and the mean EEG-RDM at the time interval with significant layer-specific temporal dynamics. Error bars indicate the standard error of the mean across participants; colored asterisks indicate significant correlations ($N = 32$, right-tailed permutation tests, FDR-corrected; *$P < 0.05$; **$P < 0.01$; ***$P < 0.001$).

**Supplementary Figure 6. Commonality analysis at the mesoscale based on all individual time points.** (**A, D**) Deep layer, (**B, E**) Middle layer and (**C, F**) Superficial layer in (**A-C**) EVC and (**D-F**) LOC. Shaded area indicates the standard error of the mean across participants; colored circles indicate significant time points ($N = 32$, cluster-defining threshold $P < 0.05$, cluster threshold $P < 0.05$).