
Research Article: New Research | Sensory and Motor Systems

Sharpening of Hierarchical Visual Feature Representations of Blurred Images

Mohamed Abdelhack^{1,2} and Yukiyasu Kamitani^{1,2}

¹Graduate School of Informatics, Kyoto University, Yoshida-Honmachi, Sakyo-Ku, Kyoto 606-8501, Japan

²ATR Computational Neuroscience Laboratories, 2-2-2 Hikaridai, Seika, Soraku, Kyoto 619-0288, Japan

DOI: 10.1523/ENEURO.0443-17.2018

Received: 19 December 2017

Revised: 29 March 2018

Accepted: 10 April 2018

Published: 7 May 2018

Author contributions: MA and YK designed the study. MA performed experiments and analyses. MA and YK wrote the paper.

Funding: <http://doi.org/10.13039/501100001691>Japan Society for the Promotion of Science (JSPS)
JP15H05920
JP15H05710

Funding: ImPACT
NA

Funding: <http://doi.org/10.13039/501100001863>New Energy and Industrial Technology Development Organization (NEDO)
NA

Authors declare no competing financial interests.

Correspondence should be addressed to: Yukiyasu Kamitani, Ph.D. Graduate School of Informatics, Kyoto University, Yoshida-honmachi, Sakyo-ku, Kyoto 606-8501, Japan. Phone: +81-75-753-9133, Fax: +81-75-753-3145, E-mail: kamitani@i.kyoto-u.ac.jp

Cite as: eNeuro 2018; 10.1523/ENEURO.0443-17.2018

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

Accepted manuscripts are peer-reviewed but have not been through the copyediting, formatting, or proofreading process.

Copyright © 2018 Abdelhack and Kamitani

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

1 **1. Manuscript title:** Sharpening of hierarchical visual feature representations of blurred
2 images

3

4 **2. Abbreviated title:** Sharpening of visual features of blurred images

5

6 **3. Authors:**

7 Mohamed Abdelhack^{1,2} and Yukiyasu Kamitani^{1,2}

8

9 ¹Graduate School of Informatics, Kyoto University, Yoshida-honmachi, Sakyo-ku, Kyoto
10 606-8501, Japan

11 ²ATR Computational Neuroscience Laboratories, 2-2-2 Hikaridai, Seika, Soraku, Kyoto
12 619-0288, Japan

13

14 **4. Author contributions:**

15 MA and YK designed the study. MA performed experiments and analyses. MA and YK
16 wrote the paper.

17

18 **5. Correspondence should be addressed to:**

19 Yukiyasu Kamitani, Ph.D.

20 Graduate School of Informatics, Kyoto University, Yoshida-honmachi, Sakyo-ku, Kyoto
21 606-8501, Japan.

22 Phone: +81-75-753-9133, Fax: +81-75-753-3145

23 E-mail: kamitani@i.kyoto-u.ac.jp

24

25 **6. Number of figures:** 7

26 **7. Number of tables:** 0

27 **8. Number of multimedia:** 0

28 **9. Number of words for Abstract:** 148

29 **10. Number of words for Significance Statement:** 92

30 **11. Number of words for Introduction:** 695

31 **12. Number of words for Discussion:** 1669

32

33 **13. Acknowledgements:**

34 The authors thank Tomoyasu Horikawa, Kei Majima, Mitsuaki Tsukamoto, and Shuntaro
35 Aoki for help with data collection and analysis. We also thank the members of Kamitani
36 Lab for their valuable comments on this manuscript. This research was supported by

37 grants from JSPS KAKENHI Grant numbers JP15H05920, JP15H05710, and ImPACT
38 Program of Council for Science, Technology and Innovation (Cabinet Office,
39 Government of Japan), and the New Energy and Industrial Technology Development
40 Organization (NEDO), and Japanese Government Scholarship (MEXT).

41 We thank Karl Embleton, PhD, from Edanz Group (www.edanzediting.com/ac) for editing
42 a draft of this manuscript.

43 This study was conducted using the MRI scanner and related facilities of Kokoro
44 Research Center, Kyoto University

45

46 **14. Conflict of Interest:**

47 The authors declare no competing financial interests.

48

49 **15. Funding source:**

50 JSPS KAKENHI; JP15H05920; Mohamed Abdelhack, Yukiyasu Kamitani

51 JSPS KAKENHI; JP15H05710; Mohamed Abdelhack, Yukiyasu Kamitani

52 ImPACT Program of Council for Science, Technology and Innovation; N/A; Yukiyasu
53 Kamitani

54 The New Energy and Industrial Technology Development Organization (NEDO); N/A;
55 Yukiyasu Kamitani

56 Japanese Ministry of Education, Culture, Sports, Science and Technology Scholarship
57 (MEXT); N/A; Mohamed Abdelhack

58

59

60 **Abstract**

61 The robustness of the visual system lies in its ability to perceive degraded images. This
62 is achieved through interacting bottom-up, recurrent, and top-down pathways that
63 process the visual input in concordance with stored prior information. The interaction
64 mechanism by which they integrate visual input and prior information is still enigmatic.
65 We present a new approach using deep neural network (DNN) representation to reveal
66 the effects of such integration on degraded visual inputs. We transformed measured
67 human brain activity resulting from viewing blurred images to the hierarchical
68 representation space derived from a feedforward DNN. Transformed representations
69 were found to veer towards the original non-blurred image and away from the blurred
70 stimulus image. This indicated deblurring or sharpening in the neural representation,
71 and possibly in our perception. We anticipate these results will help unravel the interplay
72 mechanism between bottom-up, recurrent, and top-down pathways, leading to more
73 comprehensive models of vision.
74

75 **Significance statement**

76 One powerful characteristic of the visual system is its ability to complement visual
77 information for incomplete visual images. It operates by projecting information from
78 higher visual and semantic areas of the brain into the lower and mid-level
79 representations of the visual stimulus. We investigate the mechanism by which the
80 human brain represents blurred visual stimuli. By decoding fMRI activity into a
81 feedforward-only deep neural network reference space, we found that neural
82 representations of blurred images are biased towards their corresponding deblurred
83 images. This indicates a sharpening mechanism occurring in the visual cortex.

84 **Introduction**

85 Perception is the process by which humans and other animals make sense of the
86 environment around them. It involves integrating different sensory cues with prior
87 knowledge to arrive at a meaningful interpretation of the surroundings. This integration is
88 achieved by means of two neuronal pathways: a bottom-up stimulus driven pathway,
89 which processes sensory information hierarchically, and an intrinsic pathway, which
90 performs both recurrent processing through lateral pathways and projection of prior
91 information down the hierarchy (we call it top-down pathway for abbreviation; Friston,
92 2005; Arnal and Giraud, 2012; Clark, 2013; Summerfield and de Lange, 2014; Heeger,
93 2017). The interplay mechanism between these two pathways is still an open question.

94

95 Previous studies have given rise to two main hypotheses to explain the top-down
96 modulation process. The sharpening hypothesis states that top-down signals enhance
97 the neural representation in the lower visual areas, thus improving the quality of the
98 degraded sensory signal (Lee and Mumford, 2003; Hsieh et al., 2010; Kok et al., 2012;
99 Gayet et al., 2017). Conversely, the prediction error hypothesis (which originates from
100 computer science ideas; Shi et al. (2008) states that top-down signals provide expected
101 signal information that would be redundant if represented again in lower visual areas,
102 and therefore gets subtracted (Mumford, 1992; Rao and Ballard, 1999). This results in
103 an error signal that is repeatedly processed to update the prediction signal until the error
104 signal reaches zero, which corresponds with achieving a perceptual result (Murray et al.,
105 2002; den Ouden et al., 2009, 2012; Alink et al., 2010; Meyer and Olson, 2011;
106 Todorovic et al., 2011; Kok et al., 2012; Gordon et al., 2017). Most recently, these two
107 hypotheses were reconciled by a third hypothesis where prediction error is computed to

108 be later used to sharpen the neural representations (Kok et al., 2012). Models of vision
109 usually employ this mechanism to explain the interacting neural information processing
110 pathways (Lee and Mumford, 2003; Friston, 2005; Heeger, 2017).

111

112 Most of the top-down modulation studies utilized expectation of a previously-known
113 visual stimulus to drive the operation of top-down pathways, and have hence focused on
114 lower visual areas. Expectation-of-stimulus tasks facilitate comparison of a visualized
115 stimulus and an expected stimulus at the lower visual feature level. While such studies
116 have provided an empirical framework for the operation of top-down modulation driven
117 by expectation in the lower visual areas, they have not revealed its overall operation in
118 regular recognition-targeting visual tasks.

119

120 In this study, we tackle this question by investigating top-down pathway operation during
121 a natural-image visual recognition task throughout different levels of visual processing
122 ranging from lower visual areas (V1–3) to higher visual centers (LOC, FFA, and PPA).
123 We drive the operation of top-down modulation by applying degradation to natural
124 images by blurring them. When visual images are degraded, the visual sensory signal is
125 less reliable, and the visual cortex therefore depends more heavily on prior knowledge
126 driving the top-down pathway operation. To unmask the top-down effect, we investigate
127 how the neural representations of viewing blurred images deviate from a pure
128 feedforward representation leading to a sharpened representation along the visual
129 processing pathway.

130

131 To demonstrate such sharpening, we measured and analyzed brain activity from

132 functional magnetic resonance imaging (fMRI) brain data from different regions of the
133 lower and higher visual areas, to visualize the degradation effect on different levels of
134 neural processing. We utilized deep neural network (DNN) feature space as a proxy for
135 hierarchical representation. We used a feature decoding method devised by Horikawa
136 and Kamitani (2017a) to map brain activity into a DNN representation space. The
137 decoded features were analyzed for their similarity to the feedforward-only DNN features
138 of the stimulus images and original non-blurred images. These similarities are then
139 compared to their counterpart noisy DNN features, which account for decoding errors as
140 a baseline for pure-feedforward behavior, to find whether predicted features deviate from
141 the pure feedforward ones and how supplementing with prior knowledge about stimulus
142 categories would affect the sharpening behavior. We also compared the case where the
143 image content is successfully recognized with the one where it is not. If image
144 sharpening were in operation, it would be expected that the top-down effect would be
145 boosted due to successful perception.
146

147 **Methods and materials**

148 **Subjects**

149 Five healthy subjects (three males and two females, aged between 22 and 33) with normal
150 or corrected-to-normal vision took part in the fMRI experiments. The study protocol was
151 approved by the Ethics Committee of [author's institute]. All the subjects provided written
152 informed consent for their participation in the experiments.

153

154 **Visual stimuli**

155 Both original and blurred image stimuli were shown. The images were selected from the
156 ImageNet online database (Jia Deng et al., 2009), which is the database used for training,
157 testing, and validation of the pre-trained DNN model used in this study (see below). The
158 database contains images that are categorized by a semantic word hierarchy organized in
159 WordNet (Fellbaum, 2012). First, images with a resolution lower than 500 pixels were
160 excluded, then the remaining images were further filtered to select only those that showed
161 the main object at or close to the midpoint of the image. The selected images were then
162 cropped to a square that is centered on the midpoint. If no acceptable image remained after
163 this filtration process, another image was obtained from the worldwide web through an
164 image search.

165

166 We created three different levels of blurring for the blurred image stimuli. Blurring was
167 conducted by running a square-shaped averaging filter over the whole image. The size of the
168 filter relative to the image size dictated the degradation level. The three degradation filters
169 used had a side length of 6%, 12%, and 25% of the side length of the stimulus image. We
170 then added the original stimulus image represented by a level of 0% (Figure 1A).

171

172 **Experimental design**

173 Experiments were divided into: 1) the decoder training runs where natural undegraded
174 images were presented, and 2) the test image runs where the blurred images were
175 presented. Images included in the training and test datasets were mutually exclusive. In
176 the decoder training runs, we selected one stimulus image for each of the AlexNet
177 classification categories defined in the last layer, resulting in a total of 1000 stimulus
178 images. This training stimulus set selection was conducted to avoid any bias to certain
179 categories in the decoder. This dataset was divided into 20 runs of 50 images each. The
180 subject was instructed to press a button when the image was a repeat of the image
181 shown one-back. In each run, 5 of the 50 images were repeated in the following trial to
182 form the one-back task. Each image was shown once to the subject (except for the
183 one-back repetitions).

184

185 The test image runs consisted of two conditions. In the first condition, the subjects did
186 not have any prior information about the stimuli presented (no-prior condition). In the
187 second condition, the subjects were provided a semantic prior in the form of category
188 choices (category-prior condition). The stimuli in the category-prior condition consisted
189 of images from one of five object categories (airplane, bird, car, cat, or dog). The subject
190 was informed of these categories prior to the experiment, but not the order in which they
191 were to be presented.

192

193 The stimuli in both of the test conditions were presented in sequences of maximum
194 blurring to original image (25%, 12%, 6%, and 0% blurring). Each sequence consisted of

195 stimuli representing all four levels of blurring of the same original image. We selected
196 this order of presentation to avoid the subjects having a memory-prior of the less blurred
197 stimuli when viewing the more blurred ones. For each condition, the sequences for 80
198 images were randomly distributed across two runs (40 images each). The runs
199 belonging to the same test experimental condition were conducted in the same
200 experimental session. The training and test experiments were conducted over the
201 course of five months in total for all subjects.

202

203 All image presentation was performed using Psychtoolbox (Kleiner et al., 2007). Each
204 image (12 × 12 degrees) was presented in a flashing sequence for 8 s at 1 Hz (500 ms
205 on time). Images were displayed in the center of the display with a white central fixation
206 point. The fixation point changed from white to red 500 ms before each new stimulus
207 appeared. A 32-s pre-rest and 6-s post rest period were added at the beginning and end
208 of each run respectively. Subjects were required to fixate on the central point. For test
209 runs, subjects were required to provide voice feedback of their best guess of the
210 perceived content of the stimulus. They were also required to report the certainty level of
211 that guess by pressing one of two buttons, one indicating certainty and the other
212 indicating uncertainty. We checked if the vocal reports caused excessive motion by the
213 subject that leads to degradation in the data quality but found that the motion correction
214 results were comparable to runs without vocal response by the same subjects.

215

216 **MRI acquisition**

217 fMRI data was collected using a 3-Tesla MAGNETOM Verio (Siemens Medical Systems,
218 Erlangen, Germany) MRI scanner located in [institution where the scans were

219 conducted]. For image presentation experiments, an interleaved T2*-weighted multiband
220 accelerated EPI scan was performed to obtain images covering the whole brain. The
221 scanning parameters were TR = 2000 ms; TE = 43 ms; flip angle = 80°; FOV =
222 192 × 192 mm; voxel size = 2 × 2 × 2 mm; slice gap = 0 mm; number of slices = 76;
223 multiband factor = 4. For localizer experiments, an interleaved T2*-weighted
224 gradient-EPI scan was performed with the following parameters TR = 3000 ms; TE = 30
225 ms; flip angle = 80°; FOV = 192 × 192 mm; voxel size = 3 × 3 × 3 mm; slice gap = 0 mm;
226 number of slices = 46. For retinotopy experiments, an interleaved T2*-weighted
227 gradient-EPI scan was also performed where the scanning parameters were TR = 2000
228 ms; TE = 30 ms; flip angle = 80°; FOV = 192 × 192 mm; voxel size = 3 × 3 × 3 mm; slice
229 gap = 0 mm; number of slices = 30. T2-weighted turbo spin echo (TSE) images with the
230 same slice positions as the EPI images were also acquired, to act as high-resolution
231 anatomical images. The parameters for the anatomical sequences matching the image
232 presentation acquisition were TR = 11 000 ms; TE = 59 ms; flip angle = 160°; FOV =
233 192 × 192 mm; voxel size = 0.75 × 0.75 × 2.0 mm; slice gap = 0 mm; number of slices =
234 76. For the localizer experiment, the TSE parameters were TR = 7020 ms; TE = 69 ms
235 flip angle = 160°; FOV = 192 × 192 mm; voxel size = 0.75 × 0.75 × 3.0 mm; slice gap = 0
236 mm; number of slices = 48. For the retinotopy TSE acquisition the parameters were TR
237 = 6000 ms; TE = 58 ms; flip angle = 160°; FOV = 192 × 192 mm; voxel size =
238 0.75 × 0.75 × 3.0 mm. T1-weighted magnetization-prepared rapid acquisition
239 gradient-echo (MP-Rage) fine-structural images of the entire head were also acquired.
240 The scanning parameters for these were TR = 2250 ms; TE = 3.06 ms; TI = 900 ms; flip
241 angle = 9°; FOV = 256 × 256 mm; voxel size = 1 × 1 × 1 mm number of slices = 208.
242

243 **MRI data preprocessing**

244 After rejection of the first 8 seconds of each acquisition to avoid scanner instability
245 effects, the fMRI scans were preprocessed using SPM8 (<http://www.fil.ion.ucl.ac.uk/spm>,
246 RRID: SCR_007037), including 3D motion correction, slice-timing correction, and
247 co-registration to the appropriate high resolution anatomical images. Both scans were
248 then also co-registered to the T1 anatomical image. The EPI data were then interpolated
249 to $2 \times 2 \times 2$ mm voxels and further processed using Brain Decoder Toolbox 2
250 (<https://github.com/KamitaniLab/BrainDecoderToolbox2>, RRID: SCR_013150). Volumes
251 were shifted by 2 s (1 volume) to compensate for hemodynamic delays, then the linear
252 trend was removed from each run and the data were normalized. As each image was
253 presented for 8 s, it was represented by four fMRI volumes. These four volumes were
254 then averaged to provide a single image with increased signal to noise ratio for each
255 stimulus image. The averaged voxel values for each stimulus block were used as an
256 input feature vector for the decoding analysis.

257

258 **Region of interest construction**

259 Regions of interest (ROIs) were created for several regions in the visual cortex, including
260 the lower visual areas V1, V2, and V3, the intermediate area V4, and the higher visual
261 areas consisting of the lateral occipital complex (LOC), parahippocampal place area
262 (PPA), and fusiform face area (FFA).

263 First, anatomical 3D volumes and surfaces were reconstructed from T1 images using
264 the FreeSurfer reconstruction and segmentation tool
265 (<https://surfer.nmr.mgh.harvard.edu/>, RRID: SCR_001847). To delineate the areas V1–4,
266 a retinotopy experiment was conducted following a standard protocol (Engel et al., 1994;

267 Sereno et al., 1995) involving a rotating double wedge flickering checkerboard pattern.
268 The brain activity data for this experiment was analyzed using the FreeSurfer Fsfast
269 retinotopy analysis tool (<https://surfer.nmr.mgh.harvard.edu/fswiki/FsFast>, RRID:
270 SCR_001847). The analysis results were visually examined and ROIs were delineated
271 on a 3D inflated image of the cortical surface. Voxels comprising areas V1–3 were
272 selected to form the lower visual cortex (LVC) ROI.
273 Functional localizer experiments were conducted for the higher visual areas. Each
274 subject undertook eight runs of 12 stimulus blocks. For each block, intact and
275 pixel-scrambled images of face, object, and scene categories were presented in the
276 center of the screen (10 × 10 degrees). Each block contained 20 images from one of the
277 previous six categories. Each image was presented for 0.3 s followed by 0.45 s of blank
278 gray background. This led to each block having a duration of 15 seconds. Two blocks of
279 intact and scrambled images of the same category were always displayed consecutively
280 (the order of scrambled and intact images was randomly chosen), followed by a 15- s
281 rest period with a uniform gray background. Pre-rest and post-rest periods of 24 s and
282 6 s respectively were added to each run. The brain response to the localizer experiment
283 was analyzed using the FreeSurfer Fsfast event related analysis tool. Voxels showing
284 the highest activation response to intact images for each of the face, scene, and object
285 categories in comparison with their scrambled counterparts were visualized on a 3D
286 inflated image of the cortical surface and delineated to form FFA, PPA, and LOC regions
287 respectively. Voxels constituting the areas FFA, PPA, and LOC were then selected to
288 form the higher visual cortex (HVC) ROI and the aggregation of LVC, V4, and HVC was
289 used to form the visual cortex (VC) ROI. Selected ROIs for both retinotopy and localizer
290 experiments were transformed back into the original coordinates of the EPI images.

291

292 **Deep neural network model**

293 The neural representations were transformed into a deep neural network (DNN) feature
 294 proxy using the AlexNet DNN model (Krizhevsky et al., 2017). The Caffe implementation
 295 of the network packaged for the MatConvNet tool for MATLAB (Vedaldi and Lenc, 2015)
 296 was used for implementation. This network was trained to classify 1000 different image
 297 categories with images from the ImageNet database. The model consisted of 8 layers;
 298 the first five of which were convolutional layers, while the last three were fully-connected
 299 layers. The input to each layer is the output of the previous one as follows

$$y = f_8(f_7(\dots f_1(x_0) \dots))$$

300 where x is the input image and y is the resulting image classification vector and the

301 function f_n is the operation for each layer is

$$f_n(x_n) = r_n(z_n(x_n))$$

302 and

$$z_n(x_n) = c_n(w_n, x_n) + b_n$$

303 where r_n is a non-linearity function of the n^{th} layer (rectified linear operation for the first
 304 seven layers and softmax for the final layer), w_n is the n^{th} layer weight matrix that are
 305 pretrained in the model using the ImageNet dataset, c_n is the operation conducted at
 306 the n^{th} layer between its input and weights (convolution in the case of convolutional
 307 layers and matrix multiplication in the case of fully-connected layers), and b_n is the n^{th}
 308 layer bias term. We extract features from each layer n as the output of $z_n(x_n)$ before
 309 the application of the non-linearity.

310 One thousand features were extracted from each layer (out of 290400, 186624, 64896,
 311 64896, 43264, 4096, 4096, and 1000 features from DNN layer 1–8 respectively), with

312 the features with the highest feature decoding accuracy according to the mean accuracy
313 of the five subjects' data in Horikawa and Kamitani (2017a) being selected. All the
314 feature units in the last layer were selected, as this layer contained 1000 units in total.
315 The features from each layer were labelled as DNN1–DNN8.

316

317 **DNN feature decoding**

318 Multiple linear regression decoders were constructed to predict each feature extracted
319 from the voxels of each ROI from the layers of the DNN. The decoders were constructed
320 using sparse linear regression (SLR; Bishop, 2006). This algorithm assumes that each
321 feature can be predicted using a sparse number of voxels and selects the most
322 significant voxels for predicting the features (For details, see Horikawa and Kamitani,
323 2017a).

324

325 A decoder was constructed for each feature. Voxel selection was undertaken for each
326 ROI, to select the 500 voxels with the highest correlations with each feature value. fMRI
327 data and features of the training image dataset were first normalized to a zero mean with
328 one standard deviation. The mean and standard deviation values subtracted were also
329 recorded. The decoders were then trained on the normalized fMRI data and DNN
330 features. The recorded mean and standard deviation from the training fMRI data were
331 then used to normalize the test data before decoding the features. The resulting features
332 were denormalized only by multiplying it by the standard deviation but not the addition of
333 the mean to avoid the effect of baseline correlation in the subsequent data analysis. For
334 the correlation analysis, the feature vectors emanating from the DNN were normalized
335 by subtracting the mean of the training dataset, to match the predicted features. These

336 normalized feature vectors are referred to as “true” feature vectors in this study.

337

338 The feature pattern correlation was computed for each stimulus image by aligning the
339 predicted 1000 features from each DNN layer and computing their Pearson correlation
340 coefficient with the corresponding true feature vector.

341

342 **Noise-matched features**

343 The decoded features of blurred stimulus images can be assumed to comprise the result
344 of both bottom-up and top-down processing in addition to fMRI noise, while those of the
345 true features from the DNN only contain the result of the bottom-up processing. To
346 isolate the effect of the top-down processing, we defined baseline features
347 (noise-matched feature) by adding noise to the true features. We could extract the
348 matching noise level from the decoded features of the non-blurred stimulus images
349 assuming that they do not elicit a sharpening top-down process and hence only contain
350 the bottom-up and fMRI noise components. Thus, we can add noise to the true DNN
351 features elicited from non-blurred stimulus images till their behavior matches that of the
352 decoded features of the same images. To perform this operation, Gaussian noise was
353 added to the true features so that the correlation between the noisy and the true features
354 equated the correlation between the decoded and the true features. This matching noise
355 level was calculated for each ROI/DNN layer pair in each subject.

356

357 **Feature gain**

358 The similarity of the decoded features to the original image features and that to the
359 stimulus image features were evaluated by the correlation coefficients, r_o and r_s ,

360 respectively, and the difference was calculated

$$361 \quad \Delta r_{\text{decode}} = r_o - r_s,$$

362 which indicates the bias toward the original features. To set a baseline, the same

363 difference of the correlation coefficients was calculated for the noise-matched features

364

$$365 \quad \Delta r_{\text{noise}} = r_{o_{\text{noise}}} - r_{s_{\text{noise}}}.$$

366 The feature gain was defined as the difference between these

$$\text{Feature Gain} = \Delta r_{\text{decode}} - \Delta r_{\text{noise}}.$$

367 A positive feature gain means that the decoded features are more biased towards the

368 original image features as compared to the noise-matched features.

369

370 **Content specificity**

371 To estimate the content specificity of the predicted features from the VC, their correlation

372 with the original image features was compared to that with the other original image

373 features. The correlation of predicted features for each stimulus was calculated for each

374 of the non-corresponding original images in the test dataset ($n = 39$), and the mean

375 correlation was then calculated. The mean over all the stimulus images from the stimuli

376 grouped by DNN layer with all blur levels pooled was calculated (Different image

377 correlation), and compared with the mean of the correlations with the corresponding

378 original images (Same image correlation). To compare this to the baseline correlation

379 between different images, the mean of the correlation between each stimulus image

380 vector and the feature vectors of other original images was calculated (True feature

381 correlation).

382

383 **Behavioral data extraction**

384 The subject vocal response was recorded manually from the voice recordings. The
385 written record was then revised with each subject to ensure accuracy. The record was
386 written as incorrect in the cases where the subject missed giving a voice response. In
387 the cases when the subject missed giving a button response, the previous button
388 response from the same sequence was used, except when the stimulus was the last
389 (original image) or the first (most degraded image) in the sequence, when the response
390 was set to certain and uncertain respectively. Correct responses were the ones identical
391 to the response of the last stimulus in each sequence (original image).

392

393 **Code accessibility**

394 The code described in the manuscript is freely available online at [URL redacted for
395 double-blind review]. The code is available as extended data (Extended Data 1). It was
396 created and run on MATLAB R2016b (RRID: SCR_001622) on a Linux Centos operating
397 system on a computer cluster for parallel computing. Data to reproduce our results are also
398 available at [URL redacted for double-blind review].

399 **Results**

400 **DNN feature decoding**

401 We posed a question on how top-down modulation in the visual cortex affects the neural
402 representation of blurred images. To address this question, we measured brain activity
403 while presenting blurred images. The protocol involved the presentation of stimuli in
404 blurred-to-original image sequences. Each sequence consisted of stimuli showing
405 different blur levels of the same image presented in the order of the most blurred to the
406 non-blurred original image (Figure 1A) so the subject is progressively receiving sharper
407 information about the stimulus. Subjects vocally reported the perceived object in each
408 stimulus, while also reporting their certainty of their perception. We conducted two
409 experiments using this protocol. In the first experiment, each image (stimulus sequence)
410 was chosen from a random object category and the subject had no prior information of
411 the object category (no-prior condition). In the second experiment, the stimulus
412 sequences were chosen from five predefined object categories (airplane, bird, car, cat,
413 and dog). The subjects were informed about the object categories of the set, but not of
414 each stimulus (category-prior condition). Using these two conditions we can analyze the
415 effect of adding prior information on the top-down effect in different visual areas.

416

417 To examine the effect of top-down modulation, we investigated the neural representation
418 of blurred images via the proxy of a hierarchical feedforward-only representational space
419 (Horikawa and Kamitani, 2017a). To transform brain data into the DNN feature space,
420 we trained multivoxel decoders to predict DNN features from brain activity data using a
421 separate training stimulus dataset consisting of 1000 natural non-blurred images. To

422 confirm that this choice of stimuli in training dataset did not cause the decoder to be
423 biased to non-blurred images, we conducted a content specificity analysis as will be
424 presented later.

425

426 Using the trained decoders, the brain activity pattern induced by each stimulus in the
427 blurred-to-original sequences was decoded (transformed) into the DNN feature space.
428 For each stimulus image, the Pearson correlation coefficient between its decoded
429 feature vector and the true features of the same stimulus image (r_s) at each layer was
430 computed. In addition, the correlation between the decoded feature vector and the true
431 features of the corresponding non-blurred original image (r_o) was computed (Figure 1B).
432 For non-blurred stimuli, r_s and r_o are identical.

433

434 **Feature gain computation**

435 The correlation with stimulus image features (r_s) reflects the degree to which image
436 features resulting from feedforward processing are faithfully decoded from brain activity,
437 while the correlation with original images (r_o) reflects the degree to which the decoded
438 features are “sharpened” by top-down processing, to be similar to those of the
439 non-blurred images. Figure 2A shows a scatter plot depicting a representative result for
440 prediction of the DNN layer 6 feature vector from the region of interest comprised of all
441 the visual areas (Visual cortex, VC) of Subject 4. DNN6 is a higher middle layer of
442 AlexNet where we can visualize the top-down effect on mid-level representations of
443 visual stimuli. It is also a fully connected layer that processes global stimulus information
444 rather than local information in the case of convolutional layers. This would lead to better
445 separated clusters that could show the transition from the most to the least blurred. Each

446 point represents a stimulus image pooling both category-prior and no-prior conditions,
447 and disregarding behavioral data while the mean points are also shown (white points
448 with black borders) to demonstrate how decreasing blur level leads to decoded features
449 veering towards original image features. Figure 2B shows the mean of the results in
450 Figure 2A, grouping stimuli by different blur levels. From the results of r_s and r_o , we
451 define Δr_{decode} as the difference between them. We notice from the representative data
452 that decoded features have higher correlation with the original image features than with
453 the stimulus image features, except when the blurring effect becomes too large, as in the
454 25% blur level. This suggests that a sharpening effect occurring in the visual cortex
455 causes the neural representations of viewing the blurred image to mimic those of a less
456 blurred version of it.

457

458 One shortcoming of this measure (Δr_{decode}) is that it does not have an appropriate
459 baseline for sharpening. A value of Δr_{decode} equal to zero implies that decoded features
460 are equally similar to stimulus and original image features, but it does not mean that
461 there is no sharpening. Thus, we defined a baseline for no sharpening according to the
462 behavior of feedforward-only processing. Decoded features from feedforward-only
463 processing were modeled by stimulus image features plus Gaussian noise. The noise
464 level was determined to match the decoding errors with the non-blurred images used as
465 stimuli, in which no sharpening was assumed to be involved. Noise was added to the
466 point where the decoded and noise-added features had nearly identical correlations to
467 the original image features (Figure 2B and C; 0% blur level in each). The same level of
468 noise (the mean across images in each subject and DNN layer) was added to the
469 stimulus features of the blurred images. We then computed r_s and r_o for the

470 noise-matched features, from which we could obtain the noise-matched baseline Δr_{noise}
471 (Figure 2C).

472

473 By comparing B and C in Figure 2, it is possible to note an opposite trend in how the
474 features are correlated. As the decoder was trained to predict image features, the
475 natural trend for Δr_{decode} would be negative, similar to Δr_{noise} . This indicates a level of
476 alteration in the neural representation of the blurred images, to improve the match with
477 the original images.

478

479 By subtracting the noise baseline from Δr_{decode} , we obtained the “feature gain” incurred
480 by top-down processing (Figure 2D). The value of the feature gain indicates how the
481 top-down pathways affect the predicted features in comparison with pure feedforward
482 behavior. Figure 2E shows the results of the mean feature gain for different subjects for
483 each layer. We can observe positive significant feature gains for most of the DNN layers
484 and blur levels (17 out of 24 DNN layer/blur level combinations; *t*-test across subjects
485 with Bonferroni correction, $p < 0.002$, Bonferroni correction factor = 24). This suggests
486 that top-down processing modulates neural representations to bias them towards the
487 original images. We also noticed that the fully-connected layers DNN6–8 had more
488 pronounced positive feature gains than the convolutional layers. Another notable issue is
489 that the 12% blur level shows better feature gain relative to both 6% and 25% blur levels
490 in higher visual areas. One possible explanation is that at 6% blur level the local
491 information start to unravel leading to sharpening at the shallower layers only.

492

493 One possible cause for this result is the training scheme of the decoders that only used

494 natural non-blurred images. This could have biased the output features to those
495 resembling natural images. In this case, the features could be correlated to any natural
496 image features. We investigated this possibility by measuring the content specificity of
497 the predicted features. We computed the correlation of predicted features (excluding
498 those with a 0% blur level) with the corresponding original image feature (r_o). This was
499 then compared with the mean correlation of the same predicted features, but with the
500 original features of different images. This measure provided information on how tightly
501 the predicted features were associated with the presented stimulus content, as opposed
502 to natural images in general. Figure 3 shows the result of such a content specificity
503 analysis. There are significant differences between correlations with the same image
504 features and mean correlations with different image features in all layers (t -test across
505 subjects, $p < 0.05$, uncorrected), indicating a tight association of the predicted features
506 with the stimulus image content, and ruling out a decoder bias explanation.

507

508 As mentioned before, the DNN model used in this study implements hierarchical
509 processing that is synonymous with that happening in the visual cortex. Previous studies
510 have shown homology between the features of the DNNs and the representations in the
511 visual cortex (Cadieu et al., 2014; Khaligh-Razavi and Kriegeskorte, 2014; Yamins et al.,
512 2014; Güçlü and van Gerven, 2015; Horikawa and Kamitani, 2017a). To this point, we
513 have shown the results of features predicted from the collection of all denoted visual
514 areas (VC). We further investigated the separate visual areas of the lower, intermediate,
515 and higher visual areas, to examine the homology between the feature gain and the
516 visual cortex hierarchy (Figure 4). We showed that the feature gain also shows similar
517 homologies to the visual hierarchy, in that we could observe that shallower DNN layers

518 showed larger feature gain from the lower visual areas (V1–3), while deeper DNN layers
519 showed larger feature gain from the higher visual areas (LOC, FFA, PPA). The results
520 are significantly positive for most of the layers and regions of interest (ROI), especially in
521 the higher visual areas and fully connected layers (t -test, $p < 0.05$, uncorrected).
522 However, DNN1 did not show significantly positive feature gains. These results imply
523 that feature gain also follows the same visual homology in the visual cortex areas, and
524 that the top-down effect is more pronounced in higher visual areas.

525

526 **Effects of prior knowledge and recognition**

527 In the previous analyses, the data from different experimental conditions were pooled
528 together. We then further investigated the difference between the category-prior and
529 no-prior conditions. We compared the feature gain means grouped according to the
530 experimental condition (category-prior vs. no-prior) while pooling all the behavioral
531 responses (Figure 5). We performed two-way ANOVA on the feature gain data using the
532 ROI and the experimental conditions as the independent variables. The addition of a
533 prior caused significant enhancement to the feature gain in layers DNN4, 7, 8 ($p < 0.006$,
534 Bonferroni correction factor = 8). The difference was most pronounced in DNN8 ($p =$
535 0.0000026). This result indicates that addition of prior information enhances top-down
536 modulation, thereby causing an increase in feature gain. This implies augmented
537 sharpening of neural representations.

538

539 This result, however, pooled both correctly and incorrectly reported results. When
540 considering behavioral data, there are considerable differences between category-prior
541 and no-prior conditions. The category-prior condition was characterized by a higher

542 number of correct responses (235 out of 300 total instances for five subjects) compared
543 with the no-prior condition (92 out of 300 total instances for five subjects). However, in
544 the category prior condition, the task was to choose one of five categories. This could
545 lead to false positives, as if a subject responded in a random manner, 20% of the
546 responses would be likely to be correct. In some cases when the stimulus was highly
547 degraded, the best guess response by the subjects could be random. To attempt to curb
548 this problem, we could use the certainty level as an indicator of correctness, especially
549 for the category prior. We found from the behavioral results that nearly all the trials
550 labelled as certain were also correctly recognized (category prior: 138 out of 139 certain
551 trials were correct; no-prior: 57 out of 70 certain trials were correct). This further supports
552 the observation that adding priors aids recognition.

553

554 We further analyzed our data by grouping it according to both experimental condition
555 (category-prior and no-prior) and recognition performance (correct and incorrect). We
556 show the results of the mean feature gain over subject means for each DNN layer in
557 Figure 6. For each experimental condition, we performed a three-way ANOVA test using
558 ROI, recognition performance, and blur level as independent variables. For the
559 category-prior condition, we found significant enhancement in feature gain when an
560 image was correctly recognized in DNN6, while for the no-prior condition no significant
561 enhancement was found ($p < 0.003$, Bonferroni correction factor = 16). From these
562 results, we notice that the effect of recognition leads to a very faint enhancement in the
563 feature gain.

564

565 We also analyzed our data by grouping it according to certainty level (certain and

566 uncertain). We show the results of mean feature gain over subject means for each DNN
567 layer of this analysis in Figure 7. For the category-prior condition, when an image was
568 recognized with certainty we found significant enhancement in feature gain in DNN5,
569 while for the no-prior condition significant enhancement was found in DNN1 and 7.

570

571 From the results of Figures 6 and 7, we can observe that in some layers and conditions,
572 recognition has a significant boosting effect on feature gain. However, we also found a
573 considerable feature gain even without recognition that indicates a sharpening effect not
574 guided by subjective recognition. This could be caused by a lower-level sharpening
575 associated with local similarity or object component sharpening that could be common
576 across different objects (like body parts in animals).

577

578 **Discussion**

579 In this study, we have demonstrated sharpening of the neural representations of blurred
580 visual stimuli. This sharpening can assist the visual system in achieving successful
581 prediction. It originates from endogenous processing elicited by top-down projections or
582 recurrent connections (or both) in the visual cortex. Compared with pure-feedforward
583 behavior, the neural representations of blurred images tended to be biased towards
584 those of corresponding original images, even though the original images had not yet
585 been viewed. This sharpening effect was also found to follow a visual hierarchy similar to
586 that in the visual cortex. We found that this sharpening was content-specific, and not just
587 due to a natural image bias. It was also shown to be boosted by giving category
588 information to the subject prior to stimulus viewing. This indicated that adding a more
589 specific prior leads to further sharpening of the neural representations. However, we did
590 not find that recognition had a strong role in boosting the enhancement process.

591

592 In our experimental protocol, the subjects viewed blurred stimuli in randomly organized
593 sequences. In each sequence, different levels of blur of the same image were shown,
594 ordered from the most blurred to the non-blurred stimulus (Figure 1A). This ensured that
595 subjects did not have pixel level information. Nonetheless, the results show a tendency
596 for the blurred images' neural representations to correlate with the original images
597 (Figure 2A and B). Conversely, the feedforward behavior demonstrated by the noisy
598 DNN output showed an opposite tendency (Figure 2C). We computed the feature gain to
599 investigate how the predicted DNN features deviated from pure feedforward behavior.
600 Feature gain analysis showed that the predicted features are rather correlated with the
601 original image features (Figure 2E). This indicates that a sharpening effect happens

602 across the visual cortex, leading to a more natural-image-like neural representation.

603

604 We notice also that feature gain is relatively higher in deeper layers of the DNN (DNN6–
605 8, Figure 2E). This effect could be caused by the nature of image degradation. Image
606 blurring tends to conceal localized details in favor of the global shape information. This
607 could lead to the subject attempting to recognize the global object while ignoring
608 localized details. Another observation is that feature gain in deeper layers drops
609 between the 12% and 6% blur levels. At the 6% blur level, localized details start to
610 unravel in most of the stimulus images. This could cause the lower layers' feature gain to
611 increase at the expense of the deeper layers' feature gain (Figure 2E). If we consider the
612 stimulus sequence from the most blurred stimulus to the original image stimulus, we
613 could visualize the time scale of the top-down effect where deeper layers peak earlier
614 than shallower ones. This could be one effect of our image presentation protocol where
615 the subject is accumulating evidence at each level starting from global shape evidence
616 followed by localized details to confirm their concordance with the global shape
617 evidence.

618

619 This representation was also confirmed as not being due to a natural-image bias caused
620 by the decoder training dataset, which consisted of natural unaltered images (Figure 3).
621 These results are in line with previous studies showing that neural representations are
622 improved due to a top-down effect (Lee and Mumford, 2003; Hsieh et al., 2010; Kok et
623 al., 2012; Gayet et al., 2017). Kok et al. (2012) demonstrated that even though the
624 overall neural activation weakens, the neural representations improve when a stimulus
625 is in agreement with the expectation. Gayet et al. (2017) also showed that visual working

626 memory enhances the neural representations of viewed stimuli. The reverse hierarchy
627 theory also suggests that top-down modulation serves to fine-tune sensory signals by
628 means of predictions initially made using lower spatial frequency features (Hochstein
629 and Ahissar, 2002; Ahissar and Hochstein, 2004). Furthermore, Revina et al. (2017)
630 showed that blurred stimuli can generate top-down processes that generalize to higher
631 spatial frequencies. Modelling studies that incorporated top-down and recurrent
632 connections have also shown a sharpening-like effect under an image degradation
633 scheme visible in text-based CAPTCHAs (George et al., 2017).

634

635 Our analysis also shows that feature gain follows a similar hierarchy to the visual cortex
636 (Figure 4). This indicates that the sharpening process occurs in the same hierarchical
637 processing localization as normal processing where low level sharpening occurs in lower
638 visual areas and enhancement of higher level features occurs in the higher visual areas.
639 It could be indicative of a convergent mechanism by which bottom-up and top-down
640 pathways are integrated into a single neural representation of the stimulus. This
641 suggestion could be supported by previous reports on the prediction of visual features.
642 Horikawa and Kamitani (2017a) demonstrated that visual perception and mental
643 imagery yielded feature prediction that was homologous with that of the visual cortex
644 hierarchy. Horikawa and Kamitani (2017b) also showed similar results from dream
645 induced brain activity. Earlier studies showed strong representational similarities
646 between the deeper layers of DNN and the brain activity in the inferior temporal cortex
647 (IT; Cadieu et al., 2014; Khaligh-Razavi and Kriegeskorte, 2014; Yamins et al., 2014;
648 Güçlü and van Gerven, 2015). This could be further investigated by high resolution
649 imaging to reveal the layered structures in the visual cortex and analyze the neural

650 representations in each layer. We demonstrate that top-down effects also show similar
651 homology, thus suggesting that DNN-based methods are useful for studying visual
652 top-down pathways since it can reveal the localization of the sharpening by means of
653 DNN layer feature gain.

654

655 When we added a category-prior to the task, the number of competing categories for
656 recognition decreased, thus the subjects tended to have a more directed top-down effect,
657 due to the fewer number of competing stimuli (Bar and Aminoff, 2003). This led to a
658 higher feature gain that was especially noted in the higher layers (Figure 5). This further
659 supports the idea of neural representation sharpening when given a prior describing the
660 stimulus content, as the top-down signal would be more correlated with the correct
661 recognition results, thus leading to a stronger feature gain.

662

663 We also found that when subjects successfully recognized the image content, the
664 feature gain in some layers predicted from lower visual areas was significantly improved.
665 However, this was not salient as a general trend (Figures 6, 7). It was expected that
666 recognition would lead to a boost in feature gain from the sensory competition
667 perspective, as the subject would attend to the successfully recognized object in the
668 stimulus image, leading to a directed top-down effect (Moran and Desimone, 1985;
669 Kastner et al., 1998, 1999). Hsieh et al. (2010) also showed that successful recognition
670 of binary images leads to a neural representation that is more correlated with the natural
671 image, although in their study, recognition was driven by ground truth image viewing
672 before watching the degraded stimulus at a later time.

673

674 Our results could be justified by the findings in de-Wit et al. (2012), where top-down
675 prediction was shown to be topologically inaccurate as it led to activity reduction in the
676 whole of V1, rather than at the predicted location. Revina et al. (2017) also found that
677 brain patterns due to top-down modulation did not share information with the
678 corresponding bottom-up signals. In the prediction error realm, successful prediction
679 would lead to zero error in the higher visual areas, and thus feature gain would decrease.
680 Our results show an opposite, albeit weak effect, which nonetheless supports the
681 representation-sharpening rather than prediction error hypothesis. Thus, prediction error
682 mechanisms do not appear to be in operation when stimuli are blurred or they could be
683 calculated and used as the sharpening signal as proposed in Kok et al. (2012). The
684 sharpening effect without recognition may be driven by more localized and lower-level
685 feature mechanisms. These mechanisms would enhance features corresponding to
686 local components of the main objects that were common across many objects (i.e. eyes
687 in animals). These local enhancement effects could lead to different recognition results.
688 This was shown to be true for computer vision DNN-based deblurring algorithms, where
689 the results of the enhancement process can lead to different results, according to the
690 desired object (Bansal et al., 2017).

691

692 From these results, we can deduce that top-down modulation is in operation when visual
693 input is degraded, even in the absence of a memory or expectation prior. Previous
694 studies have proposed that the brain makes an initial processing step using low spatial
695 frequency information. This step generates predictions of the content of the image in the
696 orbitofrontal cortex; these predictions are then used to drive the top-down modulation
697 effect (Bar and Aminoff, 2003; Bar et al., 2006; Kveraga et al., 2007; Breitmeyer, 2014).

698 This top-down effect comes about in the form of sharpening of neural representations
699 resulting from viewing degraded images. The mechanisms by which this effect
700 materializes have been mostly overlooked in previous literature, due to the difficulty in
701 finding a baseline for measurement. There has been more focus on the source of this
702 top-down modulation effect than on how it materializes in the visual cortex (Bar et al.,
703 2006; Chiou and Lambon Ralph, 2016). As we demonstrate here, the DNN
704 representations could offer a plausible proxy for representing brain activity and for
705 attaining a pure-feedforward baseline that can be used for measuring top-down effects.
706 The illustrated enhancement was shown to be affected by the presence of prior
707 semantic information, leading to a boost in the enhancement effect that was more visible
708 in higher-level features. To the contrary, successful recognition did not also cause an
709 overall boost in neural representation enhancement. Our results contribute to the
710 long-standing question of how top-down and recurrent pathways affect bottom-up
711 signals to achieve successful perception, which is believed to cause the hallucinatory
712 symptoms associated with psychological disorders such as schizophrenia when their
713 balance is disrupted (for review, see Friston et al., 2016; Jardri et al., 2016). Moreover,
714 our stimulus presentation protocol could be used to test more comprehensive models of
715 decision making under accumulation of evidence tasks (Platt and Glimcher, 1999). We
716 have examined the question from a more general perspective of vision, which has
717 allowed us to achieve a more comprehensive understanding of the vision process.

References

- 718 Ahissar M, Hochstein S (2004) The reverse hierarchy theory of visual perceptual
719 learning. *Trends Cogn Sci (Regul Ed)* 8:457–464.
- 720 Alink A, Schwiedrzik CM, Kohler A, Singer W, Muckli L (2010) Stimulus predictability
721 reduces responses in primary visual cortex. *J Neurosci* 30:2960–2966.
- 722 Arnal LH, Giraud A-L (2012) Cortical oscillations and sensory predictions. *Trends Cogn
723 Sci (Regul Ed)* 16:390–398.
- 724 Bansal A, Sheikh Y, Ramanan D (2017) PixelINN: Example-based Image Synthesis.
725 arXiv preprint arXiv:170805349.
- 726 Bar M, Aminoff E (2003) Cortical analysis of visual context. *Neuron* 38:347–358.
- 727 Bar M, Kassam KS, Ghuman AS, Boshyan J, Schmid AM, Dale AM, Hämäläinen MS,
728 Marinkovic K, Schacter DL, Rosen BR, Halgren E (2006) Top-down facilitation
729 of visual recognition. *Proc Natl Acad Sci U S A* 103:449–454.
- 730 Bishop C (2006) *Pattern Recognition and Machine Learning*. Springer.
- 731 Breitmeyer BG (2014) Contributions of magno- and parvocellular channels to conscious
732 and non-conscious vision. *Philos Trans R Soc Lond, B, Biol Sci* 369:20130213.
- 733 Cadieu CF, Hong H, Yamins DLK, Pinto N, Ardila D, Solomon EA, Majaj NJ, DiCarlo JJ
734 (2014) Deep neural networks rival the representation of primate IT cortex for
735 core visual object recognition. *PLoS Comput Biol* 10:e1003963.
- 736 Chiou R, Lambon Ralph MA (2016) The anterior temporal cortex is a primary semantic
737 source of top-down influences on object recognition. *Cortex* 79:75–86.
- 738 Clark A (2013) Whatever next? Predictive brains, situated agents, and the future of
739 cognitive science. *Behav Brain Sci* 36:181–204.
- 740 de-Wit LH, Kubilius J, Wagemans J, Op de Beeck HP (2012) Bistable Gestalts reduce

- 741 activity in the whole of V1, not just the retinotopically predicted parts. *J Vis* 12.
- 742 Den Ouden HEM, Friston KJ, Daw ND, McIntosh AR, Stephan KE (2009) A dual role for
743 prediction error in associative learning. *Cereb Cortex* 19:1175–1185.
- 744 Den Ouden HEM, Kok P, de Lange FP (2012) How prediction errors shape perception,
745 attention, and motivation. *Front Psychol* 3:548.
- 746 Engel SA, Rumelhart DE, Wandell BA, Lee AT, Glover GH, Chichilnisky EJ, Shadlen MN
747 (1994) fMRI of human visual cortex. *Nature* 369:525.
- 748 Fellbaum C (2012) WordNet. The encyclopedia of applied linguistics (Chapelle C, ed).
749 Hoboken, NJ, USA: John Wiley & Sons, Inc.
- 750 Friston K (2005) A theory of cortical responses. *Philos Trans R Soc Lond, B, Biol Sci*
751 360:815–836.
- 752 Friston K, Brown HR, Siemerikus J, Stephan KE (2016) The dysconnection hypothesis
753 (2016). *Schizophr Res* 176:83–94.
- 754 Gayet S, Guggenmos M, Christophel TB, Haynes J-D, Paffen CLE, Van der Stigchel S,
755 Sterzer P (2017) Visual working memory enhances the neural response to
756 matching visual input. *J Neurosci* 37:6638–6647.
- 757 George D, Lehrach W, Kansky K, Lázaro-Gredilla M, Laan C, Marthi B, Lou X, Meng Z,
758 Liu Y, Wang H, Lavin A, Phoenix DS (2017) A generative vision model that
759 trains with high data efficiency and breaks text-based CAPTCHAs. *Science*
760 358.
- 761 Gordon N, Koenig-Robert R, Tsuchiya N, van Boxtel JJ, Hohwy J (2017) Neural markers
762 of predictive coding under perceptual uncertainty revealed with Hierarchical
763 Frequency Tagging. *elife* 6.
- 764 Güçlü U, van Gerven MAJ (2015) Deep Neural Networks Reveal a Gradient in the

- 765 Complexity of Neural Representations across the Ventral Stream. *J Neurosci*
766 35:10005–10014.
- 767 Heeger DJ (2017) Theory of cortical function. *Proc Natl Acad Sci U S A* 114:1773–1782.
- 768 Hochstein S, Ahissar M (2002) View from the top: hierarchies and reverse hierarchies in
769 the visual system. *Neuron* 36:791–804.
- 770 Horikawa T, Kamitani Y (2017a) Generic decoding of seen and imagined objects using
771 hierarchical visual features. *Nat Commun* 8:15037.
- 772 Horikawa T, Kamitani Y (2017b) Hierarchical Neural Representation of Dreamed Objects
773 Revealed by Brain Decoding with Deep Neural Network Features. *Front*
774 *Comput Neurosci* 11:4.
- 775 Hsieh PJ, Vul E, Kanwisher N (2010) Recognition alters the spatial pattern of fMRI
776 activation in early retinotopic cortex. *J Neurophysiol* 103:1501–1507.
- 777 Jardri R, Hugdahl K, Hughes M, Brunelin J, Waters F, Alderson-Day B, Smailes D,
778 Sterzer P, Corlett PR, Leptourgos P, Debbané M, Cacia A, Denève S (2016)
779 Are hallucinations due to an imbalance between excitatory and inhibitory
780 influences on the brain? *Schizophr Bull* 42:1124–1134.
- 781 Jia Deng, Wei Dong, Socher R, Li-Jia Li, Kai Li, Li Fei-Fei (2009) ImageNet: A
782 large-scale hierarchical image database. In: 2009 IEEE Conference on
783 Computer Vision and Pattern Recognition, pp 248–255. IEEE.
- 784 Kastner S, De Weerd P, Desimone R, Ungerleider LG (1998) Mechanisms of directed
785 attention in the human extrastriate cortex as revealed by functional MRI.
786 *Science* 282:108–111.
- 787 Kastner S, Pinsk MA, De Weerd P, Desimone R, Ungerleider LG (1999) Increased
788 activity in human visual cortex during directed attention in the absence of visual

- 789 stimulation. *Neuron* 22:751–761.
- 790 Khaligh-Razavi S-M, Kriegeskorte N (2014) Deep supervised, but not unsupervised,
791 models may explain IT cortical representation. *PLoS Comput Biol* 10:e1003915.
- 792 Kleiner M, Brainard D, Pelli D, Ingling A, Murray R, Broussard C (2007) What's new in
793 psychtoolbox-3. *Perception* 36.
- 794 Kok P, Jehee JFM, de Lange FP (2012) Less is more: expectation sharpens
795 representations in the primary visual cortex. *Neuron* 75:265–270.
- 796 Krizhevsky A, Sutskever I, Hinton GE (2017) ImageNet classification with deep
797 convolutional neural networks. *Commun ACM* 60:84–90.
- 798 Kveraga K, Boshyan J, Bar M (2007) Magnocellular projections as the trigger of
799 top-down facilitation in recognition. *J Neurosci* 27:13232–13240.
- 800 Lee TS, Mumford D (2003) Hierarchical Bayesian inference in the visual cortex. *J Opt*
801 *Soc Am A Opt Image Sci Vis* 20:1434–1448.
- 802 Meyer T, Olson CR (2011) Statistical learning of visual transitions in monkey
803 inferotemporal cortex. *Proc Natl Acad Sci U S A* 108:19401–19406.
- 804 Moran J, Desimone R (1985) Selective attention gates visual processing in the
805 extrastriate cortex. *Science* 229:782–784.
- 806 Mumford D (1992) On the computational architecture of the neocortex. *Biol Cybern*
807 66:241–251.
- 808 Murray SO, Kersten D, Olshausen BA, Schrater P, Woods DL (2002) Shape perception
809 reduces activity in human primary visual cortex. *Proc Natl Acad Sci U S A*
810 99:15164–15169.
- 811 Platt ML, Glimcher PW (1999) Neural correlates of decision variables in parietal cortex.
812 *Nature* 400:233–238.

- 813 Rao RP, Ballard DH (1999) Predictive coding in the visual cortex: a functional
814 interpretation of some extra-classical receptive-field effects. *Nat Neurosci* 2:79–
815 87.
- 816 Revina Y, Petro LS, Muckli L (2017) Cortical feedback signals generalise across different
817 spatial frequencies of feedforward inputs. *Neuroimage*.
- 818 Sereno MI, Dale AM, Reppas JB, Kwong KK, Belliveau JW, Brady TJ, Rosen BR, Tootell
819 RB (1995) Borders of multiple visual areas in humans revealed by functional
820 magnetic resonance imaging. *Science* 268:889–893.
- 821 Shi, Sun YQ, Huifang (2008) *Image and Video Compression for Multimedia Engineering:*
822 *Fundamentals, Algorithms, and Standards*. CRC Press, Inc.
- 823 Summerfield C, de Lange FP (2014) Expectation in perceptual decision making: neural
824 and computational mechanisms. *Nat Rev Neurosci* 15:745–756.
- 825 Todorovic A, van Ede F, Maris E, de Lange FP (2011) Prior expectation mediates neural
826 adaptation to repeated sounds in the auditory cortex: an MEG study. *J Neurosci*
827 31:9118–9123.
- 828 Vedaldi A, Lenc K (2015) *Matconvnet: convolutional neural networks for MATLAB*. In:
829 *Proceedings of the 23rd ACM international conference on Multimedia*, pp 689–
830 692. New York, New York, USA: ACM Press.
- 831 Yamins DLK, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ (2014)
832 Performance-optimized hierarchical models predict neural responses in higher
833 visual cortex. *Proc Natl Acad Sci U S A* 111:8619–8624.
- 834

Legends

Figure 1 | Study design: (A) The stimulus sequence was divided into sequences of four stimuli each. Stimuli in the same sequence contained different blur levels of the same image organized from the highest blur level (25%) to the lowest (0%). Each stimulus was presented for 8 seconds. (B) Overview of the feature decoding analysis protocol; fMRI activity was measured as the subjects viewed the stimulus images presented, described in A. Trained decoders were used to predict DNN features from fMRI activity patterns. The decoded features were then analyzed for their similarity with the true DNN features of both the original image (r_o) and stimulus image (r_s). The same procedure was also conducted for noise matched DNN features that are composed of true DNN features with additional Gaussian noise to match predicted features from fMRI.

Figure 2 | Correlation of decoded features with original and stimulus image

features: (A) Scatter plot showing feature correlation of DNN6 features decoded from the whole visual cortex (VC) of subject 4, with original image features (r_s ; x-axis) and stimulus image features (r_o ; y-axis). Each point represents a stimulus image for all blurring levels except 0% while the white points with black borders show the mean of all points of the same blur level. Diagonal dotted line represents the line of equal correlation ($\Delta r_{\text{decode}} = 0$). (B) Representative result from DNN6 features decoded from the whole visual cortex (VC) of subject 4. Solid lines represent the mean correlation at different blur levels while pooling different experimental conditions and behavioral response data. The difference between r_o and r_s is labelled as Δr_{decode} . (C) Representative result showing mean noise-matched feature correlation with the original and stimulus image features for different blur levels. Noise-matching was performed to match the correlation of the DNN6 predicted features of the 0% blur stimuli decoded from VC of subject 4 (thus obtaining equal values with the decoded features at the 0% level). The difference between r_o and r_s yields the noise baseline (Δr_{noise}). (D) Feature gain is defined as the difference between Δr_{decode} and Δr_{noise} . Δr could be defined as the displacement along the r_o axis of the point on the plot from the line of equal correlation. So by subtracting the vector representing noise matched feature correlations from decoded feature correlation, we can calculate feature gain. (E) Mean feature gain is indicated for each DNN layer for

features decoded from VC at different stimulus blur levels (excluding the 0% level). Error bars indicate 95% confidence interval (CI) across five subjects.

Figure 3 | Content specificity of decoded features with blurred images: Same image correlation indicates correlation of predicted features (blur levels pooled, excluding 0%) with corresponding original image features. Different images correlation indicates the mean of correlations of the same predicted features with original image features of different images. The mean correlation is shown for different DNN layers. Error bars indicate 95% CI across five subjects.

Figure 4 | Feature gain across visual areas: Feature gain for features predicted from different visual areas. Mean feature gain is indicated for each DNN layer (blur levels pooled, 0% excluded). Error bars indicate 95% CI across five subjects.

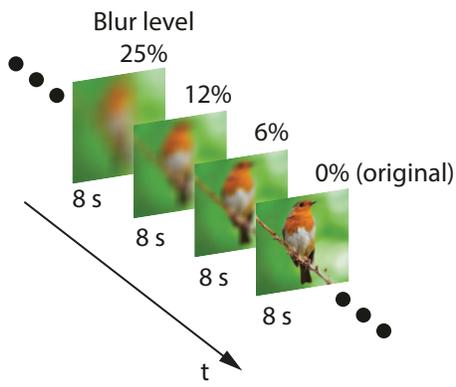
Figure 5 | Effect of category prior: Feature gain for features predicted from different visual areas grouped by experimental condition (category-prior vs. no-prior). Mean feature gain is indicated for each DNN layer (blur levels pooled, 0% excluded). Error bars indicate 95% CI across five subjects.

Figure 6 | Effect of behavioral performance: Feature gain for features predicted from different visual areas grouped by experimental condition (category-prior vs. no-prior) and recognition (correct vs. incorrect). Legends include the total number of occurrences of each response across subjects. Mean feature gain is indicated for each DNN layer (blur levels pooled, 0% excluded). Error bars indicate 95% CI across five subjects.

Figure 7 | Effect of confidence level: Feature gain for features predicted from different visual areas grouped by experimental condition (category-prior vs. no-prior) and confidence level (certain vs. uncertain). Legends include the total number of occurrences of each response across subjects. Mean feature gain is indicated for each DNN layer (blur levels pooled, 0% excluded). Error bars indicate 95% CI across five subjects.

Extended Data 1 | Code for replicating the study results

A



B

