

RESEARCH ARTICLE

Beyond core object recognition: Recurrent processes account for object recognition under occlusion

Karim Rajaei¹, Yalda Mohsenzadeh², Reza Ebrahimpour^{1,3*}, Seyed-Mahdi Khaligh-Razavi^{2,4*}

1 School of Cognitive Sciences (SCS), Institute for Research in Fundamental Sciences (IPM), Niavaran, Tehran, Iran, **2** Computer Science and AI Lab (CSAIL), MIT, Cambridge, Massachusetts, United States of America, **3** Department of Computer Engineering, Shahid Rajaee Teacher Training University, Tehran, Iran, **4** Department of Brain and Cognitive Sciences, Cell Science Research Center, Royan Institute for Stem Cell Biology and Technology, ACECR, Tehran, Iran

* rebrahimpour@strtu.edu (RE); skhaligh@mit.edu (S-MK-R)



OPEN ACCESS

Citation: Rajaei K, Mohsenzadeh Y, Ebrahimpour R, Khaligh-Razavi S-M (2019) Beyond core object recognition: Recurrent processes account for object recognition under occlusion. *PLoS Comput Biol* 15(5): e1007001. <https://doi.org/10.1371/journal.pcbi.1007001>

Editor: Leyla Isik, MIT, UNITED STATES

Received: October 17, 2018

Accepted: April 2, 2019

Published: May 15, 2019

Copyright: © 2019 Rajaei et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Dataset related to the manuscript are available at Megocclusion-vr3. RepOD (<http://dx.doi.org/10.18150/repod.2004402>).

Funding: MIT Libraries partially funded the costs for OA publication fees. SMK-R received a return-home fellowship grant from the Iranian National Elite Foundation. The funding was not specific to this study, and the funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Abstract

Core object recognition, the ability to rapidly recognize objects despite variations in their appearance, is largely solved through the feedforward processing of visual information. Deep neural networks are shown to achieve human-level performance in these tasks, and explain the primate brain representation. On the other hand, object recognition under more challenging conditions (i.e. beyond the core recognition problem) is less characterized. One such example is object recognition under occlusion. It is unclear to what extent feedforward and recurrent processes contribute in object recognition under occlusion. Furthermore, we do not know whether the **conventional** deep neural networks, such as AlexNet, which were shown to be successful in solving core object recognition, can perform similarly well in problems that go beyond the core recognition. Here, we characterize neural dynamics of object recognition under occlusion, using magnetoencephalography (MEG), while participants were presented with images of objects with various levels of occlusion. We provide evidence from multivariate analysis of MEG data, behavioral data, and computational modelling, demonstrating an essential role for recurrent processes in object recognition under occlusion. Furthermore, the computational model with local recurrent connections, used here, suggests a mechanistic explanation of how the human brain might be solving this problem.

Author summary

In recent years, deep-learning-based computer vision algorithms have been able to achieve human-level performance in several object recognition tasks. This has also contributed in our understanding of how our brain may be solving these recognition tasks. However, object recognition under more challenging conditions, such as occlusion, is less characterized. Temporal dynamics of object recognition under occlusion is largely unknown in the human brain. Furthermore, we do not know if the previously successful deep-learning algorithms can similarly achieve human-level performance in these more

Competing interests: The authors have declared that no competing interests exist.

challenging object recognition tasks. By linking brain data with behavior, and computational modeling, we characterized temporal dynamics of object recognition under occlusion, and proposed a computational mechanism that explains both behavioral and the neural data in humans. This provides a plausible mechanistic explanation for how our brain might be solving object recognition under more challenging conditions.

Introduction

There is abundance of feedforward, and recurrent connections in the primate visual cortex [1, 2]. The feedforward connections form a hierarchy of cortical areas along the visual pathway, playing a significant role in various aspects of visual object processing [3]. However, the role of recurrent connections in visual processing have remained poorly understood [1, 4–7].

Several complementary behavioral, neuronal, and computational modeling studies have confirmed that a large class of object recognition tasks called “core recognition” are largely solved through a single sweep of feedforward visual information processing [8–13]. Object recognition is defined as the ability to differentiate an object’s identity or category from many other objects having a range of identity-preserving changes [8]. Core recognition refers to the ability of visual system to rapidly recognize objects despite variations in their appearance, e.g. position, scale, and rotation [8].

Object recognition under challenging conditions, such as high variations [14, 15], object deletion and occlusion [16–24], and crowding [25–27] goes beyond the core recognition problem, which is thought to require more than the feedforward processes. Object recognition under occlusion is one of the key challenging conditions that occurs in many of the natural scenes we interact with every day. How our brain solves object recognition under such challenging condition is still an open question. Object deletion and object occlusion are shown to have different temporal dynamics [28]. While object deletion has been studied before in humans [17–20, 24, 29], we do not know much about the dynamics of object processing under the challenging condition of occlusion; in particular there has not been any previous MEG study of occluded object processing, with multivariate pattern analyses approach, linking models with both brain and behavior. Furthermore, as opposed to the core object recognition problem, where the **conventional** feedforward CNNs are shown to explain brain representations [10–13], we do not yet have computational models that successfully explain human brain representation and behavior under this challenging condition [Regarding feedforward CNNs, also see [30, 31] where CNNs cannot fully explain patterns of human behavioral decisions].

Few fMRI studies have investigated how and where occluded objects are represented in the human brain [32–36]. Hulme and Zeki [33] found that faces and houses in fusiform face area (FFA) and lateral occipital cortex (LOC) are represented similarly with and without occlusion. Ban et al. [35] used topographic mapping with simple geometric shapes (e.g. triangles), finding that the occluded portion of the shape is represented topographically in human V1 and V2, suggesting the involvement of early visual areas in object completion. A more recent study showed that the early visual areas may only code spatial information about occluded objects, but not their identity, and higher-order visual areas instead represent object-specific information, such as category or identity of occluded objects [36]. While these studies provide insights about object processing under occlusion, they do not provide any information about the temporal dynamics of these processes, and whether object recognition under occlusion requires recurrent processing.

Our focus in this study is understanding the temporal dynamics of object recognition under occlusion; and whether recurrent connections are critical in processing occluded objects? If yes, in what form are they engaged (e.g. long range feedback or local recurrent?), and how much is their contribution compared to the contribution of the feedforward visual information? We constructed a controlled image set of occluded objects, and used the combination of multivariate pattern analyses (MVPA) of MEG signals, computational modeling, backward masking, and behavioral experiments to characterize representational dynamics of object processing under occlusion, and to determine unique contributions of feedforward and recurrent processes.

Here, we provide five complementary evidence for the contribution of recurrent processes in recognizing occluded objects. *First*, MEG decoding time courses show that onset and peak for occluded objects—without backward masking—are significantly delayed compared to when the whole object is presented without occlusion. The timing of visual information plays an important role in discriminating the processing stages (i.e. feedforward vs. recurrent) with early brain responses reaching higher visual areas being dominantly feedforward and the delayed responses being mainly driven by recurrent processes [3, 4, 37–42]. *Second*, time-time decoding analysis (i.e. temporal generalization) suggests that occluded object processing goes through a relatively long sequence of stages that involve recurrent interaction—likely local recurrent. *Third*, the results of backward masking demonstrate that while the masking significantly impairs both human categorization performances and MEG decoding performances under occlusion, it has no significant effect on object recognition when objects are not occluded. *Fourth*, results from two computational models showed that a **conventional** feedforward CNN (AlexNet) that could achieve human-level performance in the no-occlusion condition, performed significantly worse than humans when objects were occluded. Additionally, the feedforward CNN could only explain the human MEG data when objects were presented without occlusion; but failed to explain the MEG data under the occlusion condition. In contrast, a hierarchical CNN with local recurrent connections (recurrent ResNet) achieved human-level performance and the representational geometry of the model was significantly correlated with that of the MEG neural data when objects were occluded. Finally, we quantified contributions of feedforward and recurrent processes in explaining the neural data, showing a significant unique contribution only for the recurrent processing under occlusion. These findings demonstrate significant involvement of recurrent processes in occluded object recognition, and improve our understand of object recognition beyond the core problem. To our knowledge this is the first MVPA study of MEG data linking feedforward and recurrent deep neural network architectures with both brain and behavior to investigate object recognition under occlusion in humans.

Results

We used multivariate pattern analysis (MVPA) of MEG data to characterize representational dynamics of object recognition under occlusion [43–46]. MEG along with MVPA allows for a fine-grained investigation of the underlying object recognition processes across time [46, 47]. Subjects ($N = 15$) were presented with images of objects with varying levels of occlusion (i.e., 0% = no-occlusion, 60% and 80% occlusion; see [Methods](#)). We also took advantage of the visual backward masking [48] as a tool to further control the feedforward and feedback flow of visual information processing. In the MEG experiment, each stimulus was presented for 34 ms [$2 \times$ screen frame rate (17ms) = 34ms], followed by a blank-screen ISI, and then in half of the trials followed by a dynamic mask ([S1 Fig](#)). We extracted and pre-processed MEG signals from -200 ms to 700 ms with regard to the stimulus onset. To calculate pairwise discriminability between objects, a support vector machine (SVM) classifier was trained and tested at each time

point (Fig 1a). MEG decoding time-courses show the pairwise discriminability of object images averaged across individuals. We first present the MEG results of the no-mask trials. After that in section “The neural basis of masking effect” we discuss the effect of backward masking.

Object recognition is significantly delayed under occlusion

We used pairwise decoding analysis of MEG signals to measure how object information evolves over time (Fig 1a). Significantly above-chance decoding accuracy means that objects can be discriminated using the information available from the brain data at that time-point. The decoding onset latency indicates the earliest time that the object-specific information becomes available and the peak decoding latency is the time-point wherein we have the highest object-discrimination performance.

We found that object information emerges significantly later under occlusion compared to the no-occlusion condition. Object decoding under no-occlusion had an early onset latency at 79ms [± 3 ms standard deviation (SD)] and was followed by a sharp increase reaching its maximum accuracy (i.e. peak latency) at 139 ± 1 ms (Fig 1b). This early and rapidly evolving dynamic is well consistent with the known time-course of the feedforward visual object processing [see S2 Fig and [38, 43, 44]].

However, when objects were partially occluded (i.e. 60% occlusion), decoding time-courses were significantly slower than the 0% occlusion condition: the onset for decoding accuracy was at 123 ± 15 ms followed by a gradual increase in decoding accuracy until it reached its peak decoding accuracy at 199 ± 3 ms (Fig 1b). The difference between onset latencies and peak latencies were both statistically significant with $p < 10^{-4}$ (two sided sign-rank test). Analysis of the behavioral response times was also consistent with the MEG object decoding curves, showing a significant delay in participants’ response times when increasing the occlusion level (S3 Fig). The slow temporal dynamics of object recognition under occlusion and the observed significant temporal delay in processing occluded objects compared to un-occluded (0% occlusion) objects do not match with a fully feedforward account of visual information processing. Previous studies have shown that first responses to visual stimuli that contain category-related information can reach higher visual areas in as little as 100 ms. [38, 49, 50]. Therefore, the observed late onset and the significant delay in peak and onset of the decoding accuracy for occluded objects, may be best explained by the engagement of recurrent processes.

Under 80% occlusion, the MEG decoding results did not reach significance [right-sided signrank test, FDR corrected across time, $p < 0.05$] (Fig 1b). However, behaviorally, human subjects still performed above-chance in object categorization even under 80% occlusion. This discrepancy might be due to MEG acquisition noise, whereas the behavioral categorization task is by definition free from that type of noise.

While the MEG and behavioral data have different levels of noise, we showed that within the MEG data itself, presented images with different levels of occlusion (0%, 60%, 80%) did not differ in terms of their level of MEG noise (S4 Fig). Thus, the difference in decoding performance between different levels of occlusion cannot be simply explained by the difference in noise. Furthermore, patterns of cross-decoding analyses (see the next section) demonstrate that the observed delay in peak latency under occlusion cannot be simply explained by a difference in signal strength.

Time-time decoding analysis for occluded objects suggests a neural architecture with recurrent interactions

We performed time-time decoding analysis measuring how information about object discrimination generalizes across time (Fig 2a). Time-time decoding matrices are constructed by

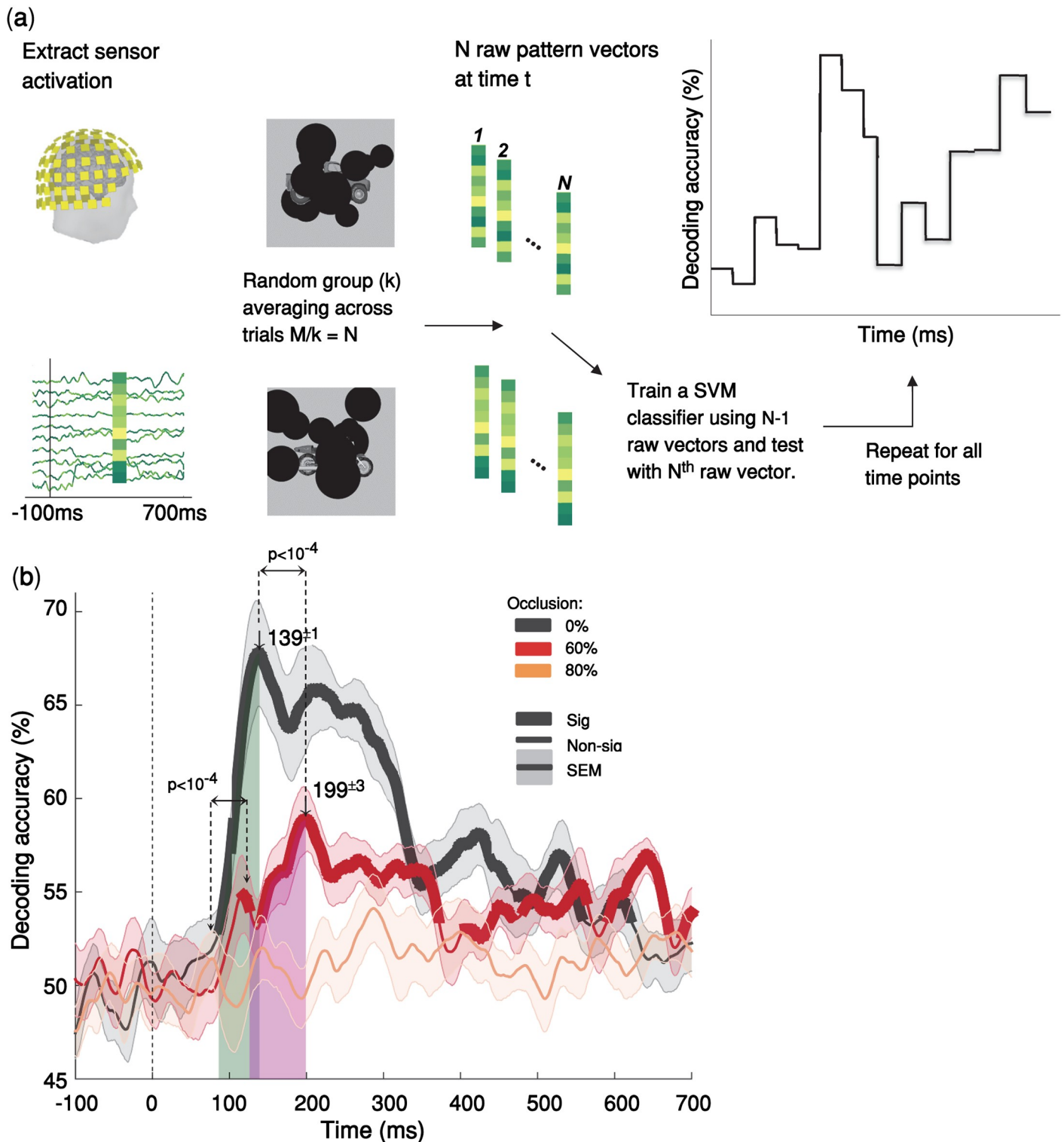


Fig 1. Temporal dynamics object recognition under various levels of occlusion. (a) Multivariate pattern classification of MEG data. We extracted MEG signals from -200 ms to 700 ms relative to the stimulus onset. At each time point (ms resolution), we computed average pairwise classification accuracy between all exemplars. (b) Time courses of pairwise decoding accuracies for the three different occlusion levels (without backward masking) averaged across 15 subjects. Thicker lines indicate a decoding accuracy significantly above chance (right-sided signrank test, FDR corrected across time, $p < 0.05$), showing that MEG signals can discriminate between object exemplars. Shaded error bars represent standard error of the mean (SEM). The two vertical shaded areas show the time from onset to peak, for 60% occluded and 0% occluded objects, which are largely non-overlapping. The onset latency is 79 ± 3 ms (mean \pm SD) in the no-occlusion condition; and 123 ± 15 ms in the

60% occlusion; the difference between onset latencies is significant ($p < 10^{-4}$, two-sided signrank test). Arrows above the curves indicate peak latencies. The peak latencies are 139 ± 1 ms and 199 ± 3 ms for the 0% occluded and partially occluded (60%) objects respectively. The difference between the peak latencies is also statistically significant ($p < 10^{-4}$). Images shown to participants are available from here: https://github.com/krajaei/Megocclusion/blob/master/Sample_occlusion_dataset.png.

<https://doi.org/10.1371/journal.pcbi.1007001.g001>

training a SVM classifier at a given time point and testing its generalization performance at all other time-points (see [Methods](#)). The pattern of temporal generalization provides useful information about the underlying processing architecture [51].

We were interested to see if there are differences between temporal generalization patterns of occluded and un-occluded objects. Different processing dynamics may lead to distinct patterns of generalization in the time-time decoding matrix [see [51] for a review]. For example, a narrow diagonal pattern suggests a hierarchical sequence of processing stages wherein information is sequentially transferred between neural stages. This hierarchical architecture is well consistent with the feedforward account of neural information processing across the ventral visual pathway. On the other hand, a time-time decoding pattern with off-diagonal generalization suggests a neural architecture with recurrent interactions between processing stages [see Fig 5 in [52]].

The temporal generalization pattern under no-occlusion (Fig 2b) indicated a sequential architecture, without an off-diagonal generalization until its early peak latency at 140 ms. This is consistent with a dominantly feedforward account of visual information processing. There was some off-diagonal generalization after 140 ms, however that was not of interest here, because the ongoing recurrent activity after the peak latency (as shown in Fig 1b) did not carry any information that further improves pairwise decoding performance of 0% occluded objects. On the other hand, when objects were occluded, the temporal generalization matrix (Fig 2c) indicated a significantly delayed peak latency at 199ms with extensive off-diagonal generalization before reaching its peak. In other words, for occluded objects, we see a discernible pattern of temporal generalization, which is characterized by 1) a relatively weak early diagonal pattern of the decoding accuracy during [100 150]ms with limited temporal generalization, which is in contrast with the high accuracy decoding of 0% occluded objects in the same time period. 2) A relatively late peak decoding accuracy with a wide generalization pattern around 200ms. This pattern of temporal generalization can be simulated by a hierarchical neural architecture with local recurrent interactions within the network [Fig 5 of [52]]

We also performed sensorwise decoding analysis to explore spatio-temporal dynamics of object information. To calculate sensorwise decoding, pairwise decoding analysis was conducted on 102 neighboring triplets of MEG sensors (2 gradiometers and 1 magnetometer in each location) yielding a decoding map of brain activity at each time-point. The sensorwise decoding patterns indicated the approximate locus of neural activity: in particular, we see that for both 0% occlusion (S2 Movie) and 60% occlusion (S1 Movie) conditions, during the onset of decoding as well as the peak decoding time, the main source of object decoding is in the left posterior-temporal sensors. From [110ms to 200ms], the peak of decoding accuracy remains locally around the same sensors, suggesting a sustained local recurrent activity.

Generalization across time and occlusion levels. Time-time decoding analyses can be further expanded by training the classifiers in one condition (e.g. occluded) and testing their ability to generalize to the other condition (e.g. un-occluded). The resulting pattern of generalization across time and occlusion level provides diagnostic information about the organization of brain processes (Fig 3). In particular, this can provide us with further evidence as to whether the observed decoding results under occlusion is due to changes in activation intensity (e.g. weakening of the signal), or a genuine difference in latency. As shown in ([51], Fig 4) each of these two come with distinct cross-condition generalization matrices. A change of signal

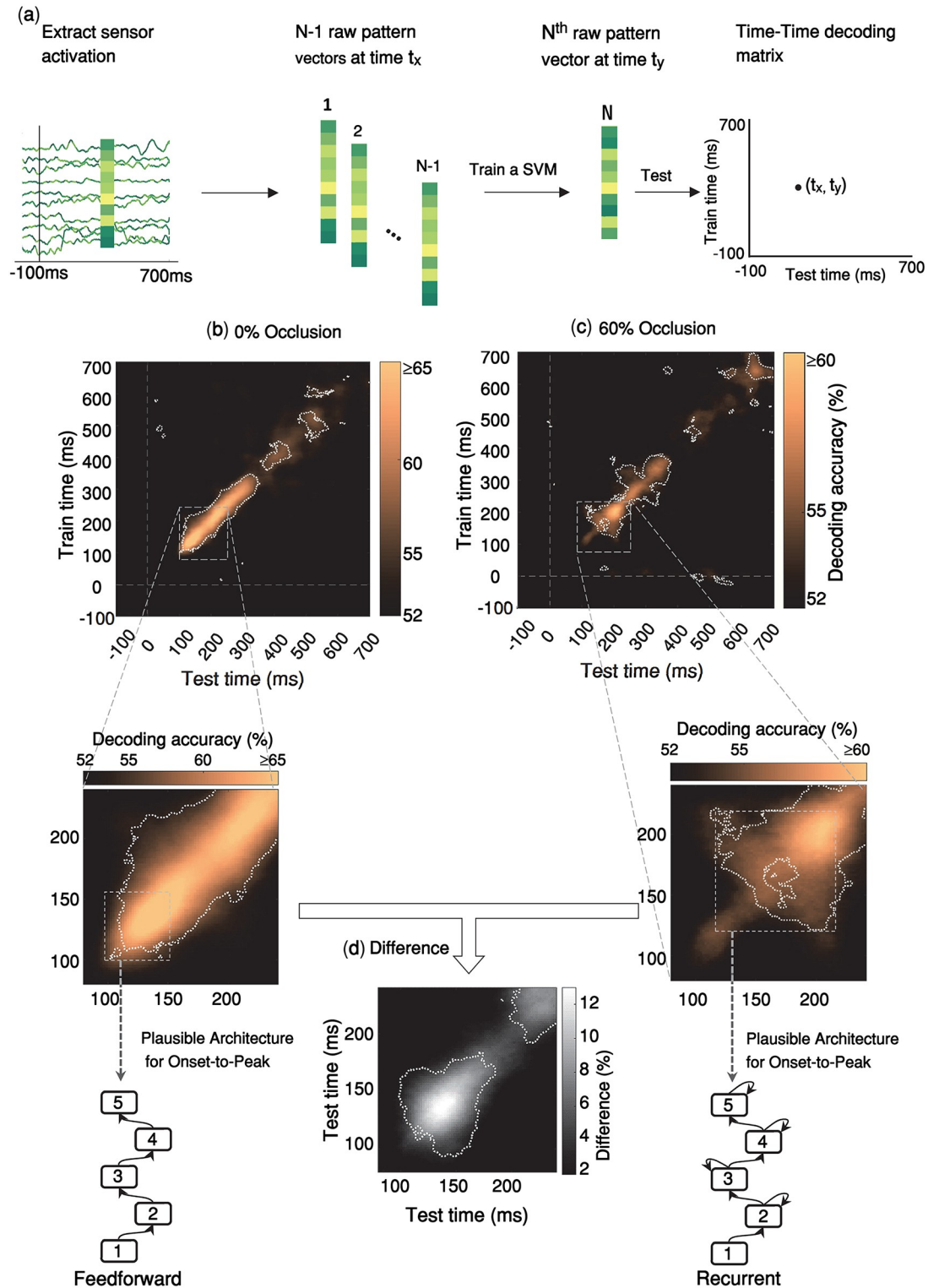


Fig 2. Temporal generalization patterns of object recognition with and without occlusion. (a) Time-time decoding analysis. The procedure is similar to the calculations of pairwise decoding accuracies explained in Fig 1, except that here the classifier is trained at a given time-point, and then tested at all time-points. In other words, for each pair of time points (t_x, t_y) , a SVM classifier is trained by N-1 MEG pattern vectors at time t_x and tested by the remaining one pattern vector at time t_y , resulting to an 801x801 time-time decoding matrix. (b-c) Time-time decoding accuracy and plausible processing architecture for no-occlusion and 60% occlusion. The results are for MEG trials without backward masking. Horizontal axis indicates testing times

and vertical axis indicates training times. Color bars represent percent of decoding accuracies (chance level = 50%); please note that in the time-time decoding matrices, the color bar ranges for 0% occlusion and 60% occlusion are different. Within the time-time decoding matrices, significantly above chance decoding accuracies, are surrounded by the white dashed contour lines (right-sided signrank test, FDR corrected across the whole 801x801 decoding matrix, $p < 0.05$). For each time-time decoding matrix, we also show the plausible processing architecture corresponding to it. These are derived from the observed patterns of temporal generalization from onset-to-peak decoding (shown by the gray dashed rectangles) [see Fig 5 of [52]]. Generalization patterns for the no-occlusion condition are consistent with a hierarchical feedforward architecture; whereas, for the occluded objects (60%) the temporal generalization patterns are consistent with a hierarchical architecture with recurrent connections. (d) Difference in time-time decoding accuracies between no-occlusion and occlusion conditions. Significantly above zero differences are surrounded by the white dashed contour lines (right-sided signrank test, FDR corrected across [80–240]ms matrix at $p < 0.05$).

<https://doi.org/10.1371/journal.pcbi.1007001.g002>

intensity leads to asymmetric generalizations across conditions because it can be more efficient to train a decoder with relatively high signal-to-noise ratio (SNR) data (e.g. 0% occlusion) and generalize it to low SNR data (e.g. 60% occlusion) rather than vice versa. However, in Fig 3, we do not see such asymmetry in decoding strengths, instead we observe different generalization patterns that are more consistent with changes in the speed of information processing (i.e. latency). More specifically, when the classifier is trained with 0% occlusion and tested on 60% occlusion (upper-left matrix in Fig 3a) the generalization pattern is shifted to the above diagonal and vice versa when trained with 60% occlusion and tested on 0% occlusion (the lower-right matrix).

Backward masking significantly impaired object recognition only under occlusion

Visual backward masking has been used as a tool to disrupt the flow of recurrent information processing, while feedforward processes are left relatively intact [4, 14, 48, 53–56]. Our time-time decoding results (Fig 4d 0% occluded) additionally supports the recurrent explanation of backward masking: off-diagonal generalization in time-time decoding matrices are representative of recurrent interactions; these off-diagonal components disappear when backward masking is present.

Considering the recurrent explanation of the masking effect, we further examined how the recurrent processes contribute in object processing under occlusion. We found that backward masking significantly reduced both MEG decoding accuracy time-course (Fig 4b) and subjects' behavioral performances (Fig 5d), only when objects were occluded. When occluded objects are masked, the MEG decoding time-course from 185ms to 237ms is significantly lower than the decoding time-course when in no-mask condition (Fig 4b, black horizontal lines; two-sided signrank test, FDR-corrected across time $p < 0.05$). On the other hand, for un-occluded objects, we did not find any significant difference between decoding time-courses of the mask and no-mask conditions (Fig 4a).

Consistent with the MEG decoding results, while the masking significantly reduced behavioral categorization performances when objects were occluded, it had no significant effect on the categorization performance for the un-occluded objects (Fig 5d) [two-sided signrank test]. Particularly, the backward masking removed the late MEG decoding peak (around 200ms) under occlusion (Fig 4f) likely due to disruption of later recurrent interactions.

Taken together, we demonstrated that visual backward masking, which is suggested to disrupt recurrent processes [4, 48, 53, 55, 57], significantly impairs object recognition only under occlusion. On the other hand, masking did not affect object processing under no occlusion, when information from the feedforward sweep is shown to be sufficient for object recognition. Thus, providing further evidence for the essential role of recurrent processes in object recognition under occlusion.

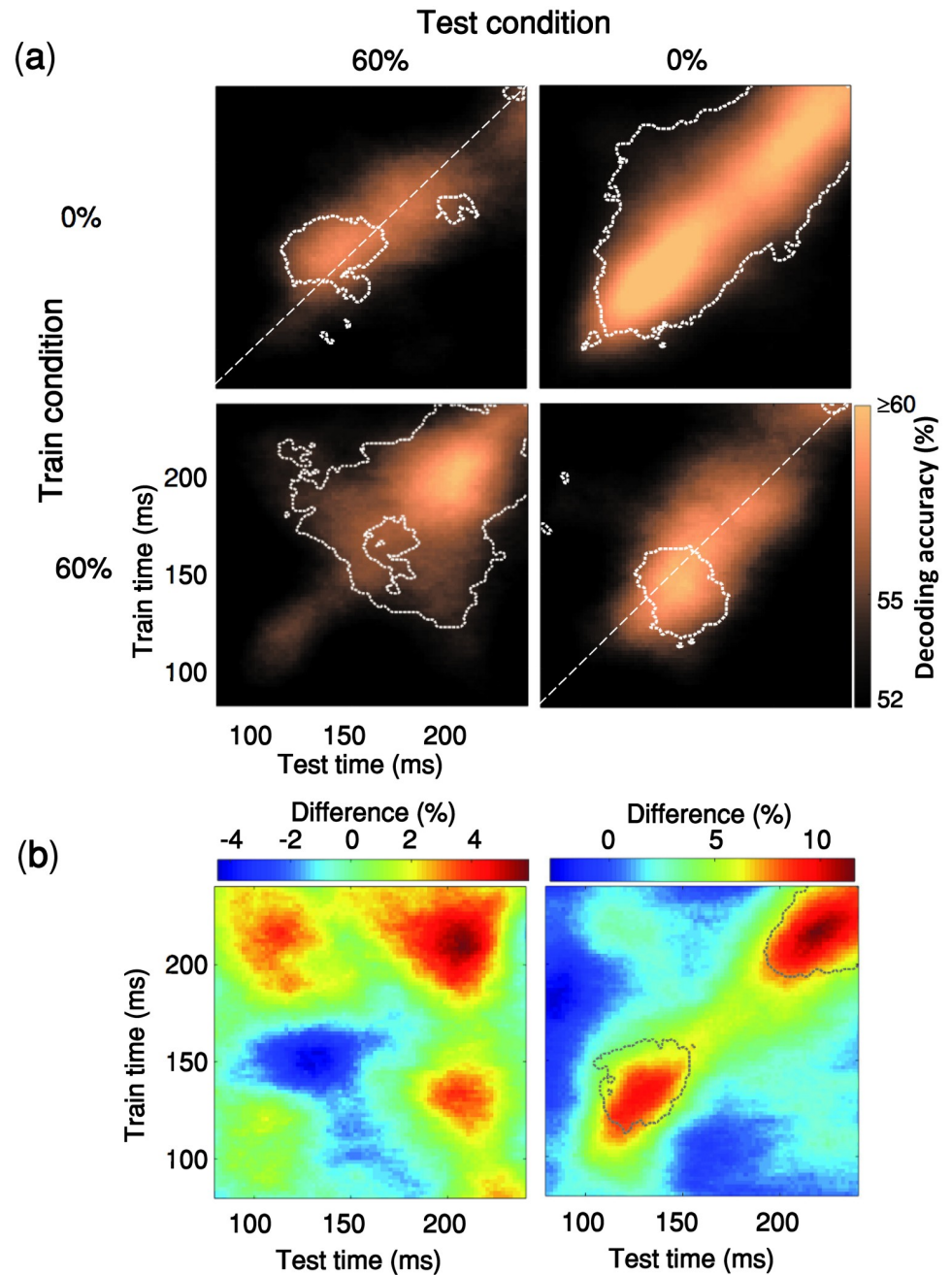


Fig 3. Generalization across time and occlusion levels. (a) The classifier is trained on an occlusion level (e.g. 0% occlusion) and tested on the other occlusion level (e.g. 60% occlusion). Time-points with significant decoding accuracy are shown inside the dashed contours (right-sided signrank test, FDR-corrected across time, $p < 0.05$). The contour of significant time-points has a shift towards the upper side of the diagonal when the classifier is trained with 0% occlusion and tested on 60% occlusion (i.e. 63% of significant time points are above the diagonal) whereas in the lower right matrix we see the opposite pattern (66% of significant time points are located below the diagonal). (b) The two color maps below the decoding matrices show the difference between the two decoding matrices located above them.

<https://doi.org/10.1371/journal.pcbi.1007001.g003>

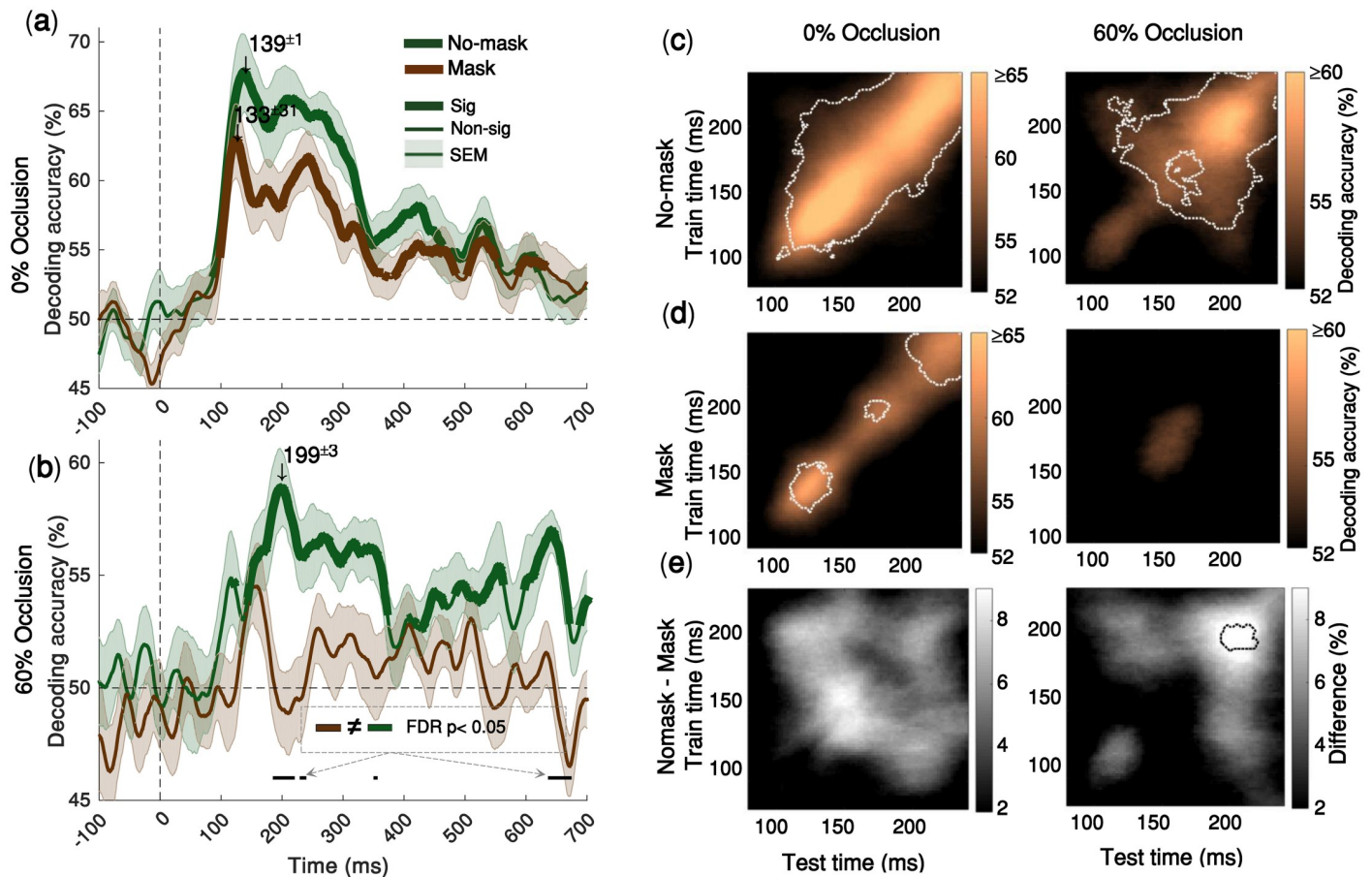


Fig 4. Backward masking significantly impairs object decoding under occlusion, but has no significant effect on object decoding under no occlusion. (a) Time-courses of the average pairwise decoding accuracies under no-occlusion. Thicker lines indicate significant time-points (right-sided signrank test, FDR corrected across time, $p < 0.05$). Shaded error bars indicate SEM (standard error of the mean). Downward pointing arrows indicate peak decoding accuracies. There is no significant difference between decoding time-courses for mask and no-mask trials, under no-occlusion (b) Time-courses of the average pairwise decoding under 60% occlusion (for 80% occlusion see S5 Fig). Under occlusion, the decoding onset latency for the no-mask trials is 123 ± 15 ms, with its peak decoding accuracy at 199 ± 3 ms; whereas the time-course for the masked trials does not reach statistical significance, demonstrating that backward masking significantly impairs object recognition under occlusion. Black horizontal lines below the curves show the time-points at which the two decoding curves are significantly different. This is particularly evident around the peak latency of the no-mask trials [from 185ms to 237ms]. (c, d) Time-time decoding matrices of 60% occluded and (0%) un-occluded objects with and without backward masking. Horizontal axes indicate testing times and the vertical axes indicate training times. Color bars show percent of decoding accuracies. Please note that in the time-time decoding matrices, the color bar ranges for 0% occlusion and 60% occlusion are different. Significantly above chance decoding accuracies, are surrounded by the white dashed contour lines (right-sided signrank test, FDR corrected across the whole 801×801 decoding matrix, $p < 0.05$). (e) Difference between time-time decoding matrices with and without backward masking. Statistically significant differences are surrounded by the black dotted contours (right-sided signrank test, FDR corrected across time at $p < 0.05$). There are significant differences between mask and no-mask only under occlusion.

<https://doi.org/10.1371/journal.pcbi.1007001.g004>

How well does a computational model with local recurrent interactions explain neural and behavioral data under occlusion?

Recent studies have shown that convolutional neural networks (CNNs) achieve human-level performance and explain neural data under non-challenging conditions—also referred to as the core object recognition [8, 10, 11]. Here, we first examined whether such feedforward CNNs (i.e. AlexNet) can explain the observed human neuronal and behavioral data in a challenging object recognition task when objects are occluded. The model accuracy was evaluated by the same object recognition task used to measure human behavioral performance (S6 Fig). A multiclass linear SVM classifier was trained on images from two occlusion levels and tested

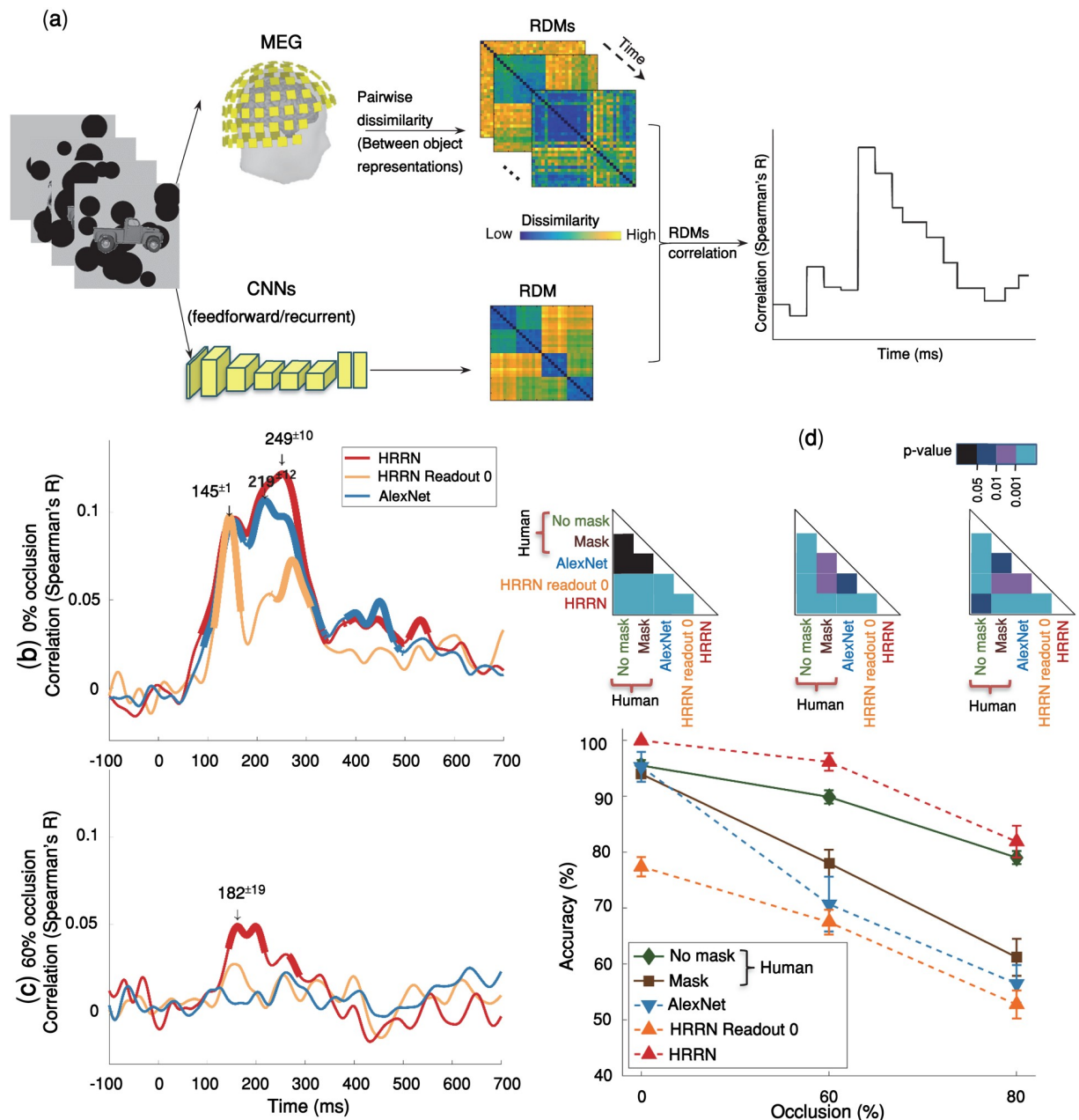


Fig 5. Comparing human MEG and behavioral data with feedforward and recurrent computational models of visual hierarchy. (a) Time-varying representational similarity analysis between human MEG data and the computational models. We, first, obtained representational dissimilarity matrices (RDM) for each computational model—using feature values of the layer before the softmax operation—, and for the MEG data at each time-point. For each subject, their MEG RDMs were correlated (Spearman's R) with the computational model RDMs (i.e. AlexNet & HRRN) across time; the results were then averaged across subjects. (b, c) Time-courses of RDM correlations between the models and the human MEG data. HRRN readout stage 0 represents the purely feedforward version of HRRN. Thicker lines show significant time points (right-sided signrank test, FDR-corrected across time, $p < = 0.05$). We indicate peak correlation latencies by numbers (mean \pm SD) above the downward pointing arrows. Under no-occlusion, AlexNet and HRRN demonstrate almost similar time-courses except that the peak latency for HRRN (249 ± 10 ms) is significantly later than the peak latency for AlexNet (219 ± 12 ms). However, under occlusion, only HRRN showed significant correlation with MEG data, with a peak latency of 182 ± 19 ms. (d) Object recognition performance of humans (mask and no-mask trials) and models [AlexNet and HRRN-ReadoutStage-0 (feedforward) and HRRN(recurrent)] across different levels of occlusion. We evaluated model accuracies on a multiclass recognition task similar to the multiclass behavioral experiment done in humans (S6 Fig). The models' performances were calculated by holding out an occlusion level for testing, and training a SVM classifier on the remaining levels of occlusion. Error bars are SEM.

<https://doi.org/10.1371/journal.pcbi.1007001.g005>

on the left-out occlusion level, using features from the penultimate layer of the model (e.g. ‘fc7’ in AlexNet). The classification task was to categorize images into car, motor, deer, or camel. This procedure was repeated for 15 times, and the mean categorization performance is reported here.

We, additionally, used representational similarity analysis (RSA) to assess model’s performance in explaining the human MEG data. RSA correlates time-resolved human MEG representations with that of the model, on the same set of stimuli. First, dissimilarity matrices (RDMs) were separately calculated for the MEG signals and the model. The model RDM were then correlated with MEG RDMs across time (Fig 5a; also see Methods).

We found that in the no-occlusion condition, the feedforward CNN achieved human-level performance and CNN representations were significantly correlated with the MEG data. Significant correlation between the model and MEG representational dissimilarity matrices (RDMs) started at ~90ms after the stimulus onset and remained significant for several hundred milliseconds with two peaks at 150ms and 220ms (Fig 5b). However, the feedforward CNNs (i.e. AlexNet and the purely feedforward version of HRRN (HRRN with readout stage 0)) failed to explain human MEG data when objects were occluded. And the model performance was significantly lower than that of human in the occluded object recognition task.

We were wondering if a model with local recurrent connections could account for object recognition under occlusion. Inspired by recent advancements in deep convolutional neural networks [58–61], we built a hierarchical recurrent ResNet (HRRN) that follows the hierarchy of the ventral visual pathway (Fig 6, also see Methods for more details about the model). The recurrent model (HRRN) could rival the human performance in the occluded object recognition task (Fig 5d), performing strikingly better than AlexNet in 60% and 80% occlusion. We also compared confusion matrices (patterns of errors) between the models and human (S7 Fig). Under the no-mask condition, HRNN had a significantly higher correlation with humans under 0% and 80% occlusion (the difference was not significant in 60% occlusion, S7b Fig).

Additionally, the HRRN model representation was significantly correlated with that of the human MEG data under occlusion [Fig 5c] (onset = 138±2ms; peak = 182±19ms). It is worth noting that the recurrent model here may be considered functionally equivalent to an arbitrarily deep feedforward CNN. We think the key difference between AlexNet and HRRN, is indeed in the number of non-linear transformations applied to the input image (please refer to section “Does a feedforward system with arbitrarily long depth work the same as a recurrent system with limited depth?” for a discussion about this).

The models here, the purely feedforward models (i.e. HRRN readout stage 0, and AlexNet) and the model with local recurrent connections, were both trained on the same object image dataset [ImageNet [62]] and had equal number of feedforward convolutional layers. Both models performed similarly in object recognition under no-occlusion, and achieved human-level performance. However, under occlusion, only the HRRN (i.e. the model with recurrent connections) could partially explain the human MEG data and achieved human-level performance, whereas the purely feedforward models failed to achieve human-level performance under occlusion—in both MEG and behavior.

Contribution of feedforward and recurrent processes in solving object recognition under occlusion

To quantify the contribution of feedforward and recurrent processes in solving object recognition under occlusion, we first correlated the feedforward and recurrent model RDMs with the

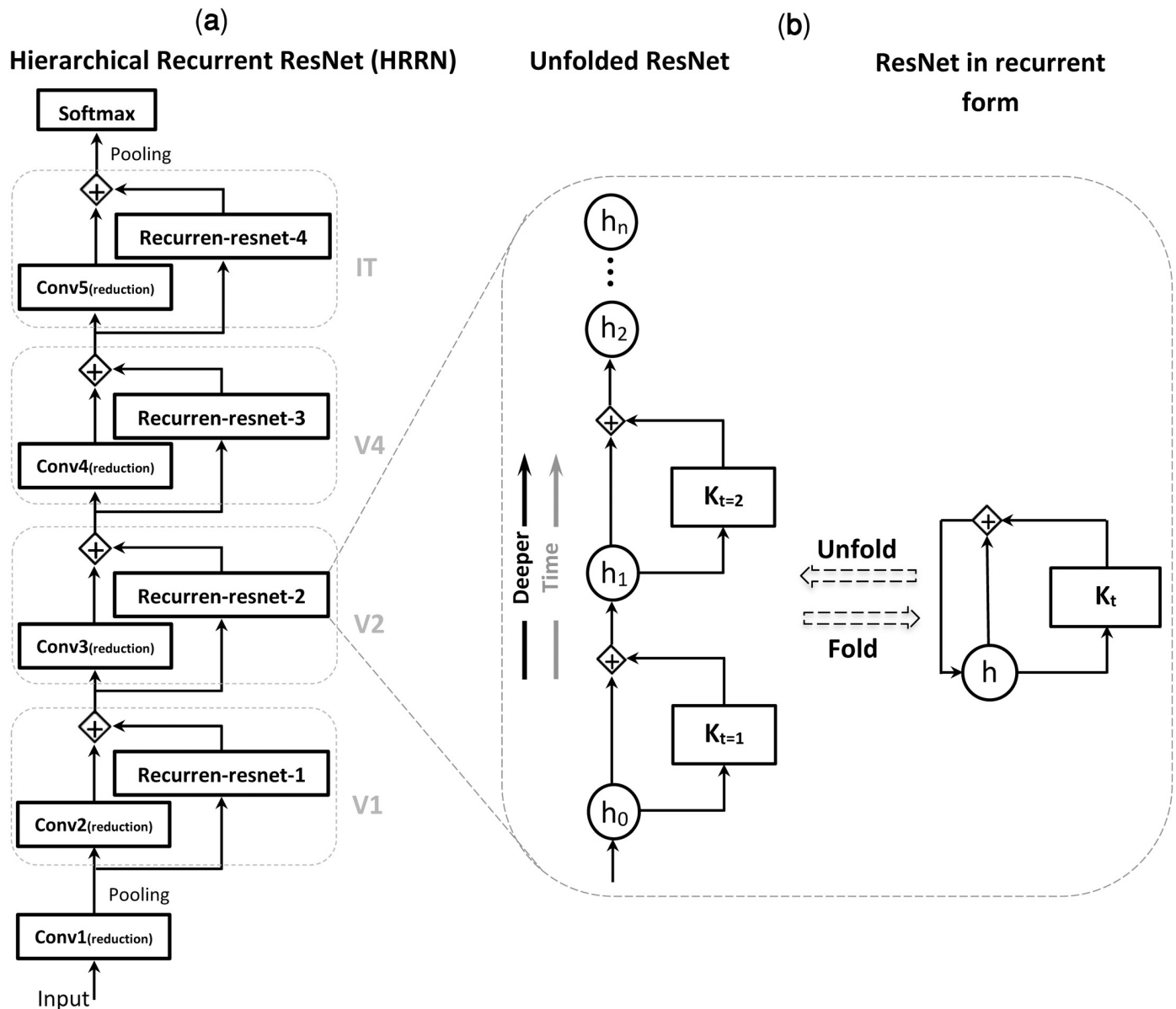


Fig 6. Hierarchical recurrent ResNet (HRRN) in unfolded form is equivalent to an ultra-deep ResNet. (a) A hierarchy of convolutional layers with local recurrent connections. This hierarchical structure models the feedforward and local recurrent connections along the hierarchy of ventral visual pathway (e.g. V1, V2, V4, IT). (b) Each recurrent unit is equivalent to a deep ResNet with arbitrary number of layers depending on the unfolding depth. h_t is the layer activity at a specific time (t) and K_t represents a sequence of nonlinear operations (e.g. convolution, batch normalization, and ReLU). [see [63] for more info].

<https://doi.org/10.1371/journal.pcbi.1007001.g006>

average MEG RDMs extracted from two time spans: 80 to 150 ms, which is dominantly feed-forward [38, 42, 50], and 151 to 300 ms (significant involvement of recurrent processes). The results are shown in Fig 7a. HRRN and AlexNet both have a similar correlation with the MEG data at [80–150 ms]. However, the HRRN shows a significantly higher correlation with the MEG data at [151–300 ms] compared to the AlexNet.

We were further interested to determine the unique contribution of the models in explaining the neural data under occlusion. To this end, we calculated semipartial correlations between the model RDMs and the MEG RDMs (Fig 7b). We find that the HRRN and AlexNet perform

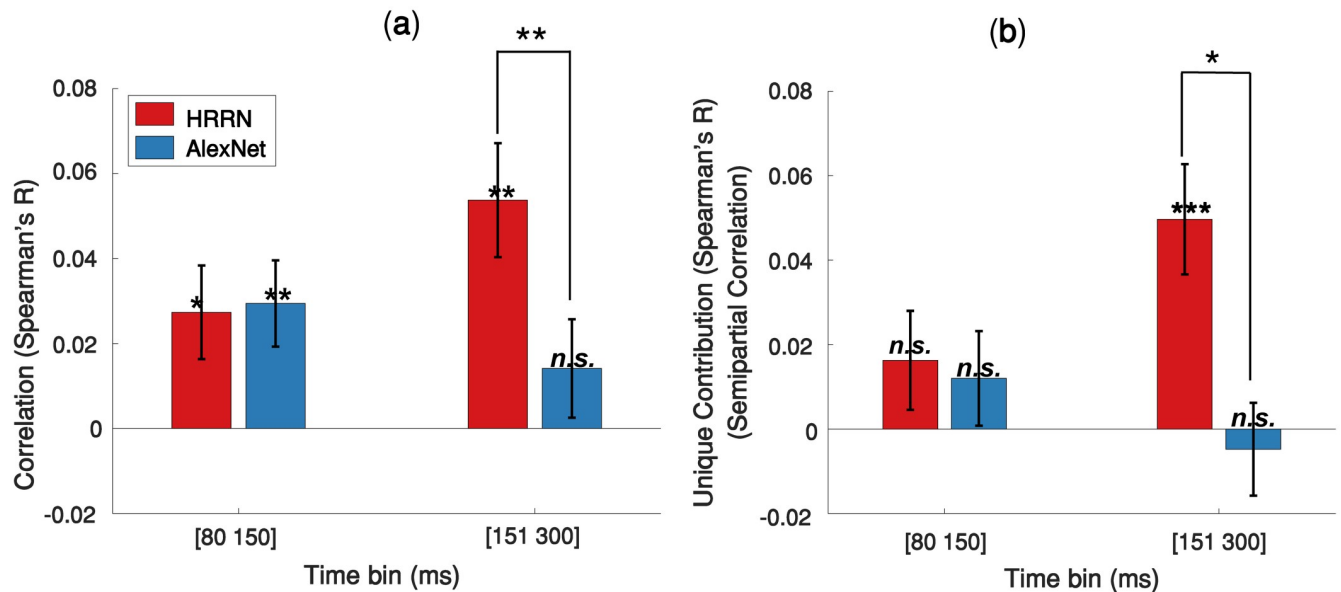


Fig 7. Contribution of the feedforward and recurrent models in explaining MEG data under 60% occlusion. (a) Correlation between the models RDMs and the average MEG RDM over two different time bins. (b) Unique contribution of each model (semipartial correlation) in explaining the MEG data. Error bars represent SEM (Standard Error of the Mean). Significantly above zero correlations/semipartial-correlations and significant differences between the two models are indicated by stars. * = $p < 0.05$; ** = $p < 0.01$; *** = $p < 0.001$.

<https://doi.org/10.1371/journal.pcbi.1007001.g007>

similarly in explaining the mainly feedforward component (i.e. 80-150ms) of the MEG data and they do not have a significant unique contribution. On the other hand, for the later component (150–300 ms) of the MEG data, only the HRRN model has a unique contribution.

Discussion

We investigated how the human brain processes visual objects under occlusion. Using multivariate MEG analysis, behavioral data, backward masking and computational modeling, we demonstrated that recurrent processing plays a major role in object recognition under occlusion.

One of the unique advantages of the current work in comparison to a number of previous studies that have investigated object recognition under partial visibility [20, 29] was bringing together neural data, behavioral data and computational models in a single study. This enabled us to study the link between brain and behavior and propose a neural architecture that is consistent with both. Whereas previous studies either related models with behavior (e.g. [20]), missing the link with brain data; or otherwise highlighting the role of recurrent processes using neural data [29] without suggesting a computational mechanism that could explain the data and the potential underlying mechanisms (however, please also see the recently published [24]).

Another unique advantage of the current study was the comparison of different deep neural network architectures and comparing their ability in explaining neural and behavioral data under occlusion. To our knowledge this was the first work that specifically compared deep convolutional networks with neural data in occluded object processing.

Beyond core object recognition

Several recent findings have indicated that a large class of object recognition tasks referred to as ‘core object recognition’ are mainly solved in the human brain within the first ~100 ms after

stimulus onset [38, 43, 45, 49, 64], largely associated with the feedforward path of visual information processing [4, 10–12]. More challenging tasks, such as object recognition under occlusion, go beyond the core recognition problem. So far it has been unclear whether the visual information from the feedforward sweep can fully account for this or otherwise recurrent information are essential to solve object recognition under occlusion.

Temporal dynamics. We found that under the no-occlusion condition, the MEG object-decoding time-course peaked at 140ms with an early onset at 79ms, consistent with findings from previous studies [8, 43–45]. Intracranial recordings in human visual cortex have shown that early responses to visual stimuli can reach higher visual areas in as little as 100 ms, with a peak decoding performance after 150 ms [38, 50]. These results suggest that category-related information around ~100 ms (± 20 ms) after stimulus onset are mainly driven by feedforward processes (see S2 Fig for time course of visual processing in human). M/EEG studies in humans have additionally shown that approximately after 150 ms, a very complex and dynamic phase of visual processing may strengthen the category-specific semantic representations [65, 66], which are likely to be driven by recurrent [49]. In our study, when objects were occluded, object decoding accuracy peaked at 200ms, with a late onset at 123ms—significantly later than the peak and onset under the no-occlusion condition (i.e. 140ms and 79ms)—suggesting the involvement of recurrent processes. Given the results from the temporal generalization analysis (Fig 2c), and the computational modeling, we argue for the engagement of mostly *local* recurrent connections as opposed to long-range top-down feedback in solving object recognition under occlusion for this image set. Previous studies also suggest that long-range top-down recurrent (e.g. PFC to IT) is prominently engaged after 200ms from stimulus onset [38, 67–69]. However, we do not rule out the possibility of involvement of long-range feedback in processing occluded objects, specifically when stimuli become more complicated (e.g. when objects are presented on natural backgrounds).

The additional time needed for processing occluded objects may facilitate object recognition by providing integrated semantic information from visible parts of the target objects, for example, via local recurrent in higher visual areas in the form of attractor networks [70]. In other words, partial semantic information (e.g. having wheels, having legs, etc.) may activate prior information associated with the category of the target object [66, 71]. Overall these suggest the observed temporal delay under 60% occlusion can be best explained by the engagement of recurrent processes—mostly local recurrent connections.

Computational modeling. Feedforward CNNs have been shown to be able to account for the core object recognition [10–12, 72–76]. The natural question to ask next is whether these models perform similarly well under more challenging conditions, beyond the core object recognition. To address this, we compared a conventional feedforward CNN with a recurrent convolutional network in terms of their object recognition performance, and their representational similarity with that of the human MEG data, under the challenging condition of occlusion. The feedforward CNN only achieved human-level performance when objects were not occluded; and performed significantly lower than the humans and the recurrent network when objects were occluded. The conventional feedforward CNN also failed to explain human neural data when objects were occluded. On the other hand, the convolutional network with local recurrent connections could achieve human-level performance in occluded object recognition and explained a significant variance of the human neural data. Thus, demonstrating that the conventional feedforward CNNs (such as AlexNet) do not account for object recognition under such challenging conditions, where recurrent computations have a prominent contribution.

Unique contribution of recurrent processes in solving object recognition under occlusion

AlexNet (i.e. feedforward model) and HRRN (i.e. recurrent model) both equally explained the early component of the MEG data (<150 ms) with and without occlusion. Consistent with this, the semipartial correlation analyses further revealed no unique variance for these models in explaining the early component of the MEG data. These results suggest that the early component of the MEG data under both conditions (with and without occlusion) are mainly feedforward and both AlexNet and HRRN share a common feedforward component that is significantly correlated with dominantly feedforward MEG representations before 150 ms (S8 Fig shows a plausible Venn diagram describing the relationship between the two models and the MEG data.).

On the other hand, the later component of the MEG data (>150 ms) under occlusion was only correlated with the recurrent model, which had a significant unique contribution in explaining the MEG data under this condition. Under no occlusion, while the later component of the MEG data is significantly correlated with both AlexNet and HRRN, only the HRRN model showed a significant unique contribution in explaining the data (S9 Fig). This shows that under no-occlusion the later component of the MEG data can still be partially explained by the common feedforward component of the two models, perhaps because object recognition under no-occlusion is primarily a feedforward process, however the recurrent model has some unique advantages in explaining later MEG components—even under no-occlusion.

Object occlusion vs. object deletion

Object recognition when part of an object is removed without an occluder is one of the challenging conditions that has been previously studied [17–20, 24, 29] and may partly look similar to occlusion. However, as shown by Johnson and Olshausen [28] deleting part of an object is different from occluding it with another object; even under short stimulus presentation times, there are meaningful differences between occlusion and deletion at the level of behavior (reaction time and accuracy of responses). Furthermore, Johnson and Olshausen report significant differences in ERP responses between occlusion and deletion, observed as early as ~130ms after stimulus onset. See S10 Fig for sample images of occlusion and deletion. Occlusion occurs when an object or shape appears in front of another one [28], in which case the occluding object might serve as an additional depth-base image cue for object completion. On the other hand, deletion occurs when part of an object is removed without additional cues about the shape or the place of the missing part. Occlusion is closely related with amodal completion which is an important visual process for perceptual filling-in of missing parts of an occluded object [28, 77]. Given the difference between these two phenomena at the level of stimulus set, we expect the dynamics of object processing (and the underlying computational mechanisms) to be also different when part of an object is occluded compared to when it is deleted. Consistent with this, Johnson and Olshausen [28] demonstrated that ERP responses in occipital and parietal electrodes are significantly different between object occlusion and deletion. Furthermore, there were significant behavioural differences between object occlusion and deletion, including differences in recognition memory and response confidence.

Object deletion has been previously studied in humans using a variety of techniques: Wyatte et. al. [18, 20] used human behaviour (in a backward masking paradigm), and computational modelling to show that object recognition, when part of the object is deleted, requires recurrent processing. Tang et. al. [24, 29] used intracranial field potential recording on epileptic patients to study temporal dynamics of object deletion; and proposed an attractor-based recurrent model that could explain the neural data. They found ~100 ms delay in processing

objects when part of the object was deleted, compared to when the whole object was present. In comparison, in our study we found ~60 ms delay in object processing when objects were occluded. This suggests that while object recognition under both occlusion and deletion requires recurrent processes, temporal dynamics of object deletion is slower, potentially due to the absence of the occluder, which can make the recognition task more difficult.

To summarize, while object deletion has been previously studied in humans, to our knowledge, temporal dynamics of object occlusion had not been studied before. In particular this was the first MVPA study in humans that characterized representational dynamics of object recognition under occlusion, and further provided a computational account of the underlying processes that explained both behavioral and neural data.

Does a feedforward system with arbitrarily long depth work the same as a recurrent system with limited depth?

While conventional CNNs could not account for object recognition beyond the core recognition problem, we do not rule out the possibility that much deeper CNNs could perform better under such challenging conditions.

Computational studies have shown that very deep CNNs outperform shallow ones on a variety of object recognition tasks [78–80]. Specifically, residual learning allows for a much deeper neural network with hundreds [58] and even thousands [59] of the layers providing better performance. This is due to the fact that the complex functions that can be represented by deeper architectures cannot be represented by shallow architectures [81]. Recent computational modeling studies have tried to clarify why increasing the depth of a network can improve its performance [60, 63]. These efforts have demonstrated that unfolding a recurrent architecture across time leads to a feedforward network with arbitrary depth, in which the weights are shared among the layers. Although such a recurrent network has far fewer parameters, Liao and Poggio [60] have empirically shown that it performs as well as a very deep feedforward network *without* shared weights. We also showed that a very deep ResNet (e.g. with 150 layers) can be reformulated into the form of a recurrent CNN with much fewer layers (e.g. five layers) (Fig 6). Thus, a compact architecture that resembles these very deep networks in terms of performance is a recurrent hierarchical network with much fewer layers. This compact architecture is probably what the human visual system has selected to be like [1, 2], given the biological constraints of having a very deep neural network inside the brain [82–86].

From a computational viewpoint, recognition of complex images might require more processing efforts; in other words, they might need to go through more layers of processing to be prepared for the final readout. Similarly, in a recurrent architecture, more processing means more iterations. Our modeling results supports this assumption, showing that under more challenging recognition tasks, more iterations are required to reach human-level performance. For example, under 60% and 80% occlusion, the HRRN model reached human level performance, respectively after going through 13 recurrent stages, and 43 recurrent stages (S11 Fig). With more iterations, the HRRN model tends to achieve a performance slightly better than the average categorization performance in humans.

Our choice of a recurrent architecture, as opposed to an arbitrarily deep neural network, is mainly driven by the plausibility of such architecture with the hierarchy of vision, where there is only a limited number of processing layers. However, in terms of performance in real-world object recognition tasks (e.g. object recognition under occlusion), the key in achieving a good performance is the number of non-linear operations, which can come either in the form of deeper networks in a feedforward architecture or otherwise more iterations in a recurrent architecture.

The neural basis of masking effect

Backward masking is a useful tool for studying temporal dynamics of visual object processing [48, 53]. It can impair recognition of the target object and reduce or eliminate perceptual visibility through the presentation of a second stimulus (mask) immediately or with an interval after the target stimulus, e.g. 50 ms after the target's onset. While the origin of masking effect was not the focus of the current study, our MEG results could provide some insights about the underlying mechanisms of backward masking.

There are several accounts of backward masking in the literature: Breitmeyer and Ganz [87] provided a purely feedforward explanation (two-channel model), arguing that the mask travels rapidly through the fast channel disrupting recognition of the target object traveling through the slow channel. A number of other studies, however, suggest that the masking mainly interferes with the top-down feedback processes [4, 48, 53, 55]. And finally, Macknik and Martinez-Conde [57] explain the masking effect by the lateral inhibition mechanism of neural circuits within different levels of the visual hierarchy; arguing that the mask interferes with the recognition of the target object through lateral inhibition (i.e. inhibitory interactions between target and mask).

The last two accounts of masking, while being different, both argue for the disruption of recurrent processes by the mask: either the top-down recurrent processes, or the local recurrent processes (e.g. lateral interactions). With a short interval between the target and mask, the mask may interfere with the fast recurrent processes (i.e. local recurrent) while with a relatively long interval it may interfere with the slow recurrent processes (i.e. top-down feedback).

Our results of MEG decoding time-courses, time-time decoding and behavioral performances under the no-occlusion condition does not support the purely feedforward account of visual backward masking. We showed that the backward masking did not have a significant effect on disrupting the fast feedforward processes of object recognition under no occlusion (MEG: Fig 4a; behaviorally: Fig 5d). On the other hand, when objects were occluded the backward masking significantly impaired object recognition both behaviorally (Fig 5d) and neurally (Fig 4b). Additionally, the time-time decoding results (Fig 4c, 4d and 4f) showed that backward masking, under no occlusion, had no significant effect on disrupting the diagonal component of the temporal generalization matrix that is mainly associated with the feedforward path of visual processing. On the other hand, the masking removed the off-diagonal components and the late peak (>200ms) observed in the temporal generalization matrix of the occluded objects.

Taken together, our MEG and behavioral results are in favor of a recurrent account for backward masking. Particularly in our experiment with a short stimulus onset asynchrony (SOA = time from stimulus onset to the mask onset), the mask seems to have affected mostly the local recurrent connections.

Methods

Ethics statement

The study was conducted according to the Declaration of Helsinki. The study involved human participants. The experiment protocol was approved by the local ethics committee at Massachusetts institute of technology. Volunteers completed a consent form before participating in the experiment and were financially compensated after finishing the experiment.

Occluded objects image set

Images of four different object categories (i.e. camel, deer, car, and motorcycle) were used as the stimulus set (see https://github.com/krajaei/Megocclusion/blob/master/Sample_occlusion_dataset.png for sample images of occlusion). Object images were transformed to be similar in

terms of size and contrast level. To generate an occluded image, in an iterative process we added several black circles (as artificial occluders) of different sizes in random positions on the image. The configuration of black circles (i.e. number, size, and their positions on the images) were randomly selected as such that a V1-like model could not discriminate between images with 0%, 60% and 80% occlusion. To simulate the type of occlusion that occurs in natural scenes, the black circles are positioned in both front and back of the target object. The percent of object occlusion is defined as the percent of the target object covered by the black occluders. We defined three levels of occlusion: 0% (no occlusion), 60% occlusion and 80% occlusion. Black circles also existed in the 0% occlusion, but not covering the target object; this was to make sure that the difference observed between occluded and un-occluded objects cannot be solely explained by the presence of these circles. The generated image set is comprised of 12 conditions: four objects \times three occlusion levels. For each condition, we generated $M = 64$ sample images varying by the occlusion pattern and the target object position. To remove the potential effect of low-level visual features in object discrimination—objects positions were slightly changed around the center of the images (by $\Delta x \leq 15$, $\Delta y \leq 15$ pixels). Overall, we controlled for low-level image statistics, as such that images of different levels of occlusion could not be discriminated simply by using low-level visual features (i.e. Gist and V1 model).

Participants and MEG experimental design

Fifteen young volunteers (22–38 year-old, all right-handed; 7 female) participated in the experiment. During the experiment, participants completed eight runs; each run consisted of 192 trials and lasted for approximately eight minutes (total experiment time for each participant = ~ 70 min). Each trial started with 1sec fixation followed by 34ms ($2 \times$ screen frame rate (17ms) = 34ms) presentation of an object image (6° visual angle). In half the trials, we employed backward masking in which a dynamic mask was presented for 102ms shortly after the stimulus offset—inter-stimulus-interval (ISI) of 17ms—(S1 Fig). In each run, each object image (i.e. camel, deer, car, motor) was repeated 8 times under different levels of occlusions without backward masking; and another 8 repetitions with backward masking. In other words, each condition (i.e. combination of object-image, occlusion-level, mask or no-mask) was repeated 64 times over the duration of the whole experiment.

Every 1–3 trials, a question mark appeared on the screen (lasted for 1.5 sec) prompting participants to choose animate if the last stimulus was deer/camel and inanimate if the last stimulus was car/motor (S1 Fig; see S12 Fig for behavioral performance of animate/inanimate task). Participants were instructed to only respond and blink during the question trials to prevent contamination of MEG signals with motor activity and the eye-blink artifact. The question trials were excluded from further MEG analyses.

The dynamic mask was a sequence of random images ($n = 6$ images; each presented for 17ms) selected from a pool of the synthesized mask images. They were generated by using a texture synthesis toolbox that is available at: <http://www.cns.nyu.edu/~lcv/texture/> [88]. The synthesized images have low-level feature statistics similar to the original stimuli.

MEG acquisition

To acquire brain signals with millisecond temporal resolution, we used 306-sensors MEG system (Elekta Neuromag, Stockholm). The sampling rate was 1000Hz and band-pass filtered online between 0.03 and 330 Hz. To reduce noise and correct for head movements, raw data were cleaned by spatiotemporal filters [Maxfilter software, Elekta, Stockholm; [89]]. Further pre-processing was conducted by Brainstorm toolbox [90]. Trials were extracted -200ms to

700ms relative to the stimulus onset. The signals were then normalized by their baseline (-200ms to 0ms), and were temporally smoothed by low-pass filtering at 20Hz.

Behavioral task of multiclass object recognition

We ran a psychophysical experiment, outside MEG, to evaluate human performance on a multi-class occluded object recognition task. Sixteen subjects participated in a session lasting about 40 minutes. The experiment was a combination of mask and no-mask trials that were randomly distributed across the experiment. Each trial, started by a fixation point presented for 0.5s followed by a stimulus presentation of 34ms. In the masked trials, a dynamic mask of 102ms was presented after a short ISI of 17ms (S5 Fig). Subjects were instructed to respond accurately and as soon as possible after detecting the target stimulus. They were asked to categorize the object images by pressing one of the pre-assigned four keys on a keyboard corresponding to the four object categories: camel, deer, car, and motorcycle.

Overall, 16 human subjects (25 to 40 years-old) participated in this experiment. Before the experiment, participants performed a short training phase on a totally different image-set to learn the task and reach a predefined performance level in the multi-class object recognition task. The main experiment consisted of 768 trials that were randomly distributed into four blocks of 192 trials (32 repetitions of object images with small variations in position and the pattern of occlusion \times three occlusion levels \times two masking conditions \times four object categories = 768). Images of 256x256 pixels size were presented at a distance of 70 cm at the center of a CRT monitor with the frame rate of 60 Hz and a resolution of 1024x768. We used the MATLAB based psychophysics toolbox of [91].

Multivariate pattern analyses (MVPA)

Pairwise decoding analysis. To measure temporal dynamics of object information processing, we used pairwise decoding analysis on the MEG data [44, 45, 92]. For each subject, at each time-point, we created a data matrix of 64-trials \times 306-sensors per condition. We used a support vector machine (SVM) to pairwise decode any two conditions, with a leave-one-out cross-validation approach. At each time-point, for each condition, $N-1$ pattern vectors were used to train the linear classifier [SVM; LIBSVM, [93], software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>], and the remaining N^{th} vector was used for evaluation. The above procedure was repeated 100 times with random reassignment of the data to training and testing sets. This was then averaged over the six pairwise decoding accuracies. The SVM decoding analysis is done independently for each subject and then we report the average decoding performance over these individuals (Fig 1a).

Time-time decoding analysis. We also reported time-time decoding accuracies, obtained by cross-decoding across time. For each pair of objects, a SVM classifier is trained at a given time and tested at all other time-points, thus showing the generalization of the classifier across time. The results are then averaged across all the pairwise classifications. This yields an 801x801 ($t = -100$ to 700 ms) matrix of average pairwise decoding accuracies for each subject. Fig 2 shows the time-time decoding matrices averaged across 15 participants. To test for statistical significance, we did one-sided signrank test against the chance-level and then corrected for multiple comparison using FDR.

Sensorwise decoding analysis. We also examined a sensorwise visualization of pairwise object decoding across time (S1 and S2 Movies). To this end, we trained and tested the SVM classifier at the level of sensors (i.e. combination of three neighboring sensors) across the whole 306 sensors. First, 306 MEG sensors were grouped into 102 triplets (Elekta Triux system; 2 gradiometers and 1 magnetometer in each location). At each time-point, we applied the

same pairwise decoding procedure as previously explained in 4.5.1, this time at the level of groups of 3 adjacent sensors (instead of taking all the 306 MEG sensors together). Average pairwise decoding accuracies across subjects, at each time point, are color-coded across the head surface. We used black dots to indicate channels with significantly above chance accuracy (FDR-corrected across both time and sensors), and gray dots to show accuracies with $p < 0.05$, before correcting for multiple comparison. At each time-point, we also specify the channel with peak decoding accuracy by a red dot.

Representational similarity analysis (RSA) over time. We used representational similarity analysis (RSA) [43, 44, 94–96], to compare representations of computational models with time-resolved representations derived from MEG data.

For the MEG data, representational dissimilarity matrices (RDM) were calculated at each time-point by computing the dissimilarity ($1 - \text{Spearman's R}$) between all pairs of the MEG patterns elicited by object images. Time-resolved MEG RDMs were then correlated (Spearman's R) with the computational model RDMs, yielding a correlation vector over time (Fig 5a).

Semipartial correlation. We additionally calculated semipartial correlations between the MEG RDMs and the computational model RDMs (Fig 7b). Semipartial correlation indicates the unique relationship between a model representation (e.g. AlexNet RDM) and the MEG data, by taking out the shared contribution of other models (e.g. HRRN RDM) in explaining the MEG data. The semipartial correlation is computed as follows:

$$r_{1(2.3)} = \frac{r_{12} - r_{13}r_{23}}{\sqrt{1 - r_{23}^2}}$$

' $r_{1(2.3)}$ ' is the correlation between *model* X_1 and *model* X_2 controlling for the effect of *model* X_3 (i.e. removing X_3 only from X_2) [97].

To construct CNN model RDMs, we used the extracted features from the penultimate layer of the networks (i.e. the layer before softmax operation). Significant correlations were determined by one-sided signrank test ($p < 0.5$, FDR-corrected across time).

Significance testing

We used the non-parametric Wilcoxon signrank test [98] for random effect analysis. To determine time-points with significantly above chance decoding accuracy (or significant RDM correlations), we used a right-sided signrank test across $n = 15$ participants. To adjust p-values for multiple comparisons (e.g. across time), we further applied the false discovery rate (FDR) correction [99] [RSA-Toolbox: is available from <https://github.com/rsagroup/rsatoolbox> [100]].

To determine whether two time-courses (e.g. correlation or decoding) are significantly different at any time interval, we used a two-sided signrank test, FDR corrected across time.

Onset latency. We defined onset latency as the earliest time where performance became significantly above chance for at least ten consecutive milliseconds. Mean and standard deviation (SD) for onset latencies were calculated by leave-one-subject-out repeated for $N = 15$ times.

Peak latency. The time for peak decoding accuracy was defined as the time where the decoding accuracy was the maximum value. The mean and SD for peak latencies were calculated similar to the onset latencies.

Computational modeling

Feedforward computational model (AlexNet). We used a well-known CNN (AlexNet) [101] that is shown to account for the core object recognition [10, 12, 74, 75]. CNNs are

cascades of hierarchically organized feature extraction layers. Each layer has several hundred convolutional filters and each convolutional filter scans various places on the input generating a feature map at its output. A convolutional layer may be followed by a local or global pooling layer merging outputs of a group of units. The pooling layers make the feature maps invariant to small variations [81]. AlexNet has eight cascading layers: five convolutional layers, some of which followed by pooling layers, and three fully-connected layers [101]. The last fully-connected layer is a 1000-way softmax that corresponds to the 1000 class labels. The network has 60 million free parameters. A pre-trained version of the model, which is trained on 1.2 million images from ImageNet dataset [102] is used for the experiments here. We used the features extracted by the fc7 layer (before softmax operation) as the model output.

Hierarchical recurrent ResNet (HRRN). In convolutional neural networks, performance in visual recognition tasks can be substantially improved by adding to the depth of the network [78, 79, 103]. However, this comes at a cost: deeper networks of simply stacking layers (plain nets) have higher training errors due to the vanishing gradients (degradation) [104] problem that prevents convergence in the training phase. To address this problem, He et al. [58] introduced a deep residual learning framework. Residual networks can overcome the vanishing gradient problem during learning by employing *identity shortcut connections* that allow bypassing residual layers. This framework enables training ultra-deep networks, e.g. with 1202 layers, leading to much better performances compared to the shallower networks [58, 59].

Residual connections give ResNet an interesting characteristic of having several possible pathways with different lengths from the network's input to the output instead of a single deep pathway [61]. For example, the ultra-deep 152-layers ResNet in its simplest form—by skipping all the residual layers—is a hierarchy of five convolutional layers. By including additional residual layers, more complex networks with various depths are constructed [see table 1 in [58]]. Each series of the residual modules can be reformulated into the form of a convolutional layer with recurrent connections [63]. Additionally, Liao and Poggio [60] show that a ResNet with shared weights can retain most of the performance of the corresponding network with non-shared weights.

In this study, we proposed a generalization of this convolutional neural network by redefining residual layers as local recurrent connections. As shown in Fig 6, we reformulated the 152-layers ResNet of [58] into the form of a five-layer convolutional network with folded residual layers as its local recurrent connections. The unfolded HRRN (= 152 layers ResNet) is deeper than the AlexNet and has different normalization (i.e. batch normalization) and filter sizes, however, they still have a similar number of free parameters (60M). Comparing AlexNet with a purely feedforward version of HRRN (readout stage 0, with five layers), AlexNet performs slightly better than HRRN in readout stage 0, and gradually after 4 iterations HRRN reaches a categorization performance similar to that of AlexNet (S13 Fig). The model is pre-trained on ImageNet 2012 dataset with a training set similar to that of Alexnet (1.2 million training images). It is shown experimentally that an unfolded recurrent CNN (with shared weights) is similar to a very deep feedforward network with non-shared weights [60]. In our analyses, we used the extracted features of the penultimate layer (i.e. layer pool5, which is before the softmax layer) as the model output.

Supporting information

S1 Fig. MEG experimental paradigm. The experiment was divided into two types of trials: mask and no mask trials, shown in random order. Each trial started by a fixation of 1 sec followed by a target stimulus presented for 34ms. In the mask trials, after a short inter-stimulus-interval (ISI) of 17ms, a dynamic mask of 102ms duration was presented. Every 1–3 trials

(average = 2) a question mark appeared on the screen. Subjects were asked to select whether the last image was animate or inanimate. They were also instructed to restrict their blinking (and swallowing) to the question-mark trials.

(TIF)

S2 Fig. Time course of visual processing and masking in humans. Earliest responses reaching each of the visual areas from V1 to IT are indicated by the oblique lines, when a stimulus is on for 34ms, followed by an ISI of 17 ms, followed by a mask. The grey shaded area indicates the effect of mask when it disrupts the information that is being fed back from higher visual areas to lower visual areas. The approximate timings are set according to human [38, 50, 105] and non-human (i.e. macaque) studies [4] controlling for the fact that the macaque cortex is smaller, with a shorter neural distance and therefore faster transmission of visual information [37].

(TIF)

S3 Fig. Average response times of behavioral experiment across three occlusion levels. The results are averaged over $n = 15$ human participants. Error bars represent SEM. Significant difference between occlusion levels are indicated by stars (signrank test). *** = $p < 0.001$.

(TIF)

S4 Fig. Split-half replicability for different conditions. The MEG trials for each condition (i.e. 0% occlusion, 60% occlusion, and 80% occlusion) were divided into two halves, the replicability is measured as the correlation between these two halves. Thicker lines indicate significantly above chance correlations (right sided sign-rank test, FDR corrected across time, $p < 0.05$). No significant difference was observed between the replicability of different conditions (two-sided signrank test, FDR-corrected across time), thus indicating that different conditions do not differ in their level of noise. In more details, for each condition, we randomly split $M = 64$ trial repetitions into two groups of 32 trials. Distance matrices were then calculated for average raw pattern vectors of each group by computing pairwise dissimilarity (1-correlation) between the patterns (12x12 matrices; 12 experimental stimuli). Spearman's R was used as the replicability measure between these two split-half matrices across time.

(TIF)

S5 Fig. Time-course of average pairwise decoding accuracy for mask and no mask trials under 80 percent occlusion. Shaded error bars represent standard error of the mean (SEM). Decoding accuracy was not significantly above chance at any time-point for both mask and no mask (right-sided signrank test, FDR-corrected across time, $p < 0.5$).

(TIF)

S6 Fig. Experimental design of the multiclass behavioral task. The behavioral experiment had two types of trials: mask and no mask trials (in random order). Each trial started by 0.5sec fixation, followed by a short presentation of stimulus for 34ms. In the masked trials, 17ms after the stimulus offset (short ISI) a dynamic mask of 100ms was presented. The dynamic mask was a sequence of synthesized images. The subjects were instructed to respond as soon and accurate as possible. Subject's response was to categorize the presented image by pressing one of the four pre-specified keys on a keyboard corresponding to the four object categories (camel, deer, car, and motorcycle).

(TIF)

S7 Fig. Confusion matrices of the human (mask/no-mask) and models across three occlusion levels. To compare patterns of errors in the models and humans, we computed confusion matrices. To obtain a confusion matrix, we first trained a SVM classifier on a multiclass object

recognition task similar to the behavioral experiment. Then, we calculated the percentage of predicted labels assigned to a category. We display these percentages using color-codes in the matrix. Elements in the main diagonal of the confusion matrix show classification performances and off-diagonal elements show errors made in the classification. **(a)** Confusion matrices for the three levels of occlusion. The color bar, at the bottom-right corner, indicates the percentage of labels assigned to a category. **(b)** Bars indicate correlations between confusion matrices of the models with that of humans (mask and no-mask). Stars show significant differences between HRRN and AlexNet (signrank test, across subjects). * = $p < 0.05$; ** = $p < 0.01$. (TIF)

S8 Fig. Venn diagram of MEG and the models. Red area indicates the unique contribution of HRRN in explaining MEG data. AlexNet has no unique contribution likely due to a component shared between the two models (i.e. feedforward component). (TIF)

S9 Fig. Contribution of the feedforward and recurrent models in explaining MEG data under 0% occlusion. (a) Correlation between the models RDMs and the average MEG RDM over two different time bins. (b) Unique contribution of each model (semipartial correlation) in explaining the MEG data. Error bars represent SEM (Standard Error of the Mean). Significantly above zero correlations/semipartial-correlations and significant differences between the two models are indicated by stars. * = $p < 0.05$; ** = $p < 0.01$; *** = $p < 0.001$. (TIF)

S10 Fig. Sample images of object occlusion versus deletion. (TIF)

S11 Fig. HRRN accuracy across readout stages for different levels of occlusion. Shaded error bars indicate SD. Black circles are average accuracies across $n = 16$ human participants. Readout stage: readout stage refers to the number of local recurrent iterations involved in processing the input image throughout the hierarchy of the network. Readout stage 0 is when the model is fully feedforward (no local recurrent is active). And readout stage 1 is when only one recurrent iteration is engaged and readout stage n is when the network has gone through n recurrent iterations. (TIF)

S12 Fig. Behavioral performance of animate/inanimate categorization task of the MEG experiment. Stars indicate significant differences between mask and no-mask trials. The results are averaged over $N = 15$ human participants. (TIF)

S13 Fig. Categorization accuracies of HRRN across different readout stages compared with Alexnet for different occlusion levels. Shaded error bars indicate SD. Readout stage: readout stage refers to the number of local recurrent iterations involved in processing the input image throughout the hierarchy of the network. Readout stage 0 is when the model is fully feedforward (no local recurrent is active). And readout stage 1 is when only one recurrent iteration is engaged and readout stage n is when the network has gone through n recurrent iterations. Black circles are average accuracies for Alexnet, which are shown around the approximate corresponding HRRN readout stages. (TIF)

S1 Movie. Sensorwise visualization of pairwise object decoding across time for no-occlusion condition. Color map represents percent of decoding accuracies across head surface

(chance level = 50%). Circles indicate neighboring triplets (102 triplets) of MEG sensors (2 gradiometers and 1 magnetometer in each location). Significantly above chance decoding accuracies after correction for multiple comparison are shown by black dots (FDR-corrected across 102 triplets and 801 time-points). Gray dots indicate decoding accuracies with $p < 0.05$ (right sided signed rank test) that did not remain significant after FDR-correction. At each time point, the peak decoding accuracy is indicated by a red dot.

(MP4)

S2 Movie. Sensorwise visualization of pairwise object decoding across time for occlusion condition. Color map represents percent of decoding accuracies across head surface (chance level = 50%). Circles indicate neighboring triplets (102 triplets) of MEG sensors (2 gradiometers and 1 magnetometer in each location). Significantly above chance decoding accuracies after correction for multiple comparison are shown by black dots (FDR-corrected across 102 triplets and 801 time-points). Gray dots indicate decoding accuracies with $p < 0.05$ (right sided signed rank test) that did not remain significant after FDR-correction. At each time point, the peak decoding accuracy is indicated by a red dot.

(MP4)

Acknowledgments

The study was conducted at the Athinoula A. Martinos Imaging Center at the McGovern Institute for Brain Research, Massachusetts Institute of Technology. We would like to thank Aude Oliva and Dimitrios Pantazis for their help and support in conducting this study. We would also like to thank Radoslaw Martin Cichy for helpful comments.

Author Contributions

Data curation: Karim Rajaei, Yalda Mohsenzadeh.

Formal analysis: Karim Rajaei.

Investigation: Karim Rajaei, Reza Ebrahimpour, Seyed-Mahdi Khaligh-Razavi.

Methodology: Karim Rajaei, Yalda Mohsenzadeh, Reza Ebrahimpour, Seyed-Mahdi Khaligh-Razavi.

Project administration: Seyed-Mahdi Khaligh-Razavi.

Software: Karim Rajaei.

Supervision: Reza Ebrahimpour, Seyed-Mahdi Khaligh-Razavi.

Validation: Seyed-Mahdi Khaligh-Razavi.

Visualization: Karim Rajaei, Yalda Mohsenzadeh, Seyed-Mahdi Khaligh-Razavi.

Writing – original draft: Karim Rajaei, Yalda Mohsenzadeh, Seyed-Mahdi Khaligh-Razavi.

Writing – review & editing: Karim Rajaei, Yalda Mohsenzadeh, Seyed-Mahdi Khaligh-Razavi.

References

1. Lamme VA, Super H, Spekreijse H. Feedforward, horizontal, and feedback processing in the visual cortex. *Current opinion in neurobiology*. 1998; 8(4):529–35. PMID: [9751656](https://pubmed.ncbi.nlm.nih.gov/9751656/)
2. Sporns O, Zwi JD. The small world of the cerebral cortex. *Neuroinformatics*. 2004; 2(2):145–62. <https://doi.org/10.1385/NI:2:2:145> PMID: [15319512](https://pubmed.ncbi.nlm.nih.gov/15319512/)

3. Felleman DJ, Van Essen DC. Distributed hierarchical processing in the primate cerebral cortex. *Cerebral cortex*. 1991; 1(1):1–47. <https://doi.org/10.1093/cercor/1.1.1> PMID: 1822724
4. Lamme VA, Roelfsema PR. The distinct modes of vision offered by feedforward and recurrent processing. *Trends in neurosciences*. 2000; 23(11):571–9. PMID: 11074267
5. Gilbert CD, Li W. Top-down influences on visual processing. *Nature Reviews Neuroscience*. 2013; 14(5):350–63. <https://doi.org/10.1038/nrn3476> PMID: 23595013
6. Kafaligonul H, Breitmeyer BG, Ögmen H. Feedforward and feedback processes in vision. *Frontiers in psychology*. 2015; 6.
7. Klink PC, Dagnino B, Gariel-Mathis M-A, Roelfsema PR. Distinct Feedforward and Feedback Effects of Microstimulation in Visual Cortex Reveal Neural Mechanisms of Texture Segregation. *Neuron*. 2017.
8. DiCarlo JJ, Cox DD. Untangling invariant object recognition. *Trends in cognitive sciences*. 2007; 11(8):333–41. <https://doi.org/10.1016/j.tics.2007.06.010> PMID: 17631409
9. DiCarlo JJ, Zoccolan D, Rust NC. How does the brain solve visual object recognition? *Neuron*. 2012; 73(3):415–34. <https://doi.org/10.1016/j.neuron.2012.01.010> PMID: 22325196
10. Khaligh-Razavi S-M, Kriegeskorte N. Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS Comput Biol*. 2014; 10(11):e1003915. <https://doi.org/10.1371/journal.pcbi.1003915> PMID: 25375136
11. Yamins DL, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences*. 2014; 111(23):8619–24.
12. Cadieu CF, Hong H, Yamins DL, Pinto N, Ardila D, Solomon EA, et al. Deep neural networks rival the representation of primate IT cortex for core visual object recognition. *PLoS Comput Biol*. 2014; 10(12):e1003963. <https://doi.org/10.1371/journal.pcbi.1003963> PMID: 25521294
13. Wen H, Shi J, Chen W, Liu Z. Deep Residual Network Predicts Cortical Representation and Organization of Visual Features for Rapid Categorization. *Scientific Reports*. 2018; 8(1):3752. <https://doi.org/10.1038/s41598-018-22160-9> PMID: 29491405
14. Ghodrati M, Farzmaahi A, Rajaei K, Ebrahimpour R, Khaligh-Razavi S-M. Feedforward object-vision models only tolerate small image variations compared to human. *Frontiers in computational neuroscience*. 2014; 8:74. <https://doi.org/10.3389/fncom.2014.00074> PMID: 25100986
15. Karimi-Rouzbahani H, Bagheri N, Ebrahimpour R. Hard-wired feed-forward visual mechanisms of the brain compensate for affine variations in object recognition. *Neuroscience*. 2017; 349:48–63. <https://doi.org/10.1016/j.neuroscience.2017.02.050> PMID: 28245990
16. Rensink RA, Enns JT. Early completion of occluded objects. *Vision research*. 1998; 38(15):2489–505.
17. Nielsen KJ, Logothetis NK, Rainer G. Dissociation between local field potentials and spiking activity in macaque inferior temporal cortex reveals diagnosticity-based encoding of complex objects. *Journal of Neuroscience*. 2006; 26(38):9639–45. <https://doi.org/10.1523/JNEUROSCI.2273-06.2006> PMID: 16988034
18. Wyatt D, Curran T, O'Reilly R. The limits of feedforward vision: Recurrent processing promotes robust object recognition when objects are degraded. *Journal of Cognitive Neuroscience*. 2012; 24(11):2248–61. https://doi.org/10.1162/jocn_a_00282 PMID: 22905822
19. O'Reilly RC, Wyatt D, Herd S, Mingus B, Jilk DJ. Recurrent processing during object recognition. *Frontiers in psychology*. 2013; 4:124. <https://doi.org/10.3389/fpsyg.2013.00124> PMID: 23554596
20. Wyatt D, Jilk DJ, O'Reilly RC. Early recurrent feedback facilitates visual object recognition under challenging conditions. 2014.
21. Kosai Y, El-Shamayleh Y, Fyall AM, Pasupathy A. The role of visual area V4 in the discrimination of partially occluded shapes. *Journal of Neuroscience*. 2014; 34(25):8570–84. <https://doi.org/10.1523/JNEUROSCI.1375-14.2014> PMID: 24948811
22. Choi H, Pasupathy A, Shea-Brown E. Predictive coding in area V4: dynamic shape discrimination under partial occlusion. *arXiv preprint arXiv:161205321*. 2016.
23. Spoerer C, McClure P, Kriegeskorte N. Recurrent Convolutional Neural Networks: A Better Model Of Biological Object Recognition Under Occlusion. *bioRxiv*. 2017:133330.
24. Tang H, Schrimpf M, Lotter W, Moerman C, Paredes A, Caro JO, et al. Recurrent computations for visual pattern completion. *Proceedings of the National Academy of Sciences*. 2018:201719397.
25. Livne T, Sagi D. Multiple levels of orientation anisotropy in crowding with Gabor flankers. *Journal of vision*. 2011; 11(13):18-. <https://doi.org/10.1167/11.13.18> PMID: 22101017
26. Manassi M, Herzog M, editors. Crowding and grouping: how much time is needed to process good Gestalt? *Perception*; 2013.

27. Clarke AM, Herzog MH, Francis G. Visual crowding illustrates the inadequacy of local vs. global and feedforward vs. feedback distinctions in modeling visual perception. *Frontiers in psychology*. 2014; 5.
28. Johnson JS, Olshausen BA. The recognition of partially visible natural objects in the presence and absence of their occluders. *Vision research*. 2005; 45(25):3262–76.
29. Tang H, Buia C, Madhavan R, Crone NE, Madsen JR, Anderson WS, et al. Spatiotemporal dynamics underlying object completion in human ventral visual cortex. *Neuron*. 2014; 83(3):736–48. <https://doi.org/10.1016/j.neuron.2014.06.017> PMID: 25043420
30. Eberhardt S, Cader JG, Serre T, editors. How deep is the feature analysis underlying rapid visual categorization? *Advances in neural information processing systems*; 2016.
31. Rajalingham R, Issa EB, Bashivan P, Kar K, Schmidt K, DiCarlo JJ. Large-scale, high-resolution comparison of the core visual object recognition behavior of humans, monkeys, and state-of-the-art deep artificial neural networks. *bioRxiv*. 2018:240614.
32. Rauschenberger R, Liu T, Slotnick SD, Yantis S. Temporally unfolding neural representation of pictorial occlusion. *Psychological Science*. 2006; 17(4):358–64. <https://doi.org/10.1111/j.1467-9280.2006.01711.x> PMID: 16623695
33. Hulme OJ, Zeki S. The sightless view: neural correlates of occluded objects. *Cerebral Cortex*. 2007; 17(5):1197–205. <https://doi.org/10.1093/cercor/bhl031> PMID: 16844722
34. Hegdé J, Fang F, Murray SO, Kersten D. Preferential responses to occluded objects in the human visual cortex. *Journal of vision*. 2008; 8(4):16-. <https://doi.org/10.1167/8.4.16> PMID: 18484855
35. Ban H, Yamamoto H, Hanakawa T, Urayama S-i, Aso T, Fukuyama H, et al. Topographic representation of an occluded object and the effects of spatiotemporal context in human early visual areas. *Journal of Neuroscience*. 2013; 33(43):16992–7007. <https://doi.org/10.1523/JNEUROSCI.1455-12.2013> PMID: 24155304
36. Erlikhman G, Caplovitz GP. Decoding information about dynamically occluded objects in visual cortex. *NeuroImage*. 2017; 146:778–88. <https://doi.org/10.1016/j.neuroimage.2016.09.024> PMID: 27663987
37. Thorpe SJ, Fabre-Thorpe M. Seeking categories in the brain. *Science*. 2001; 291(5502):260–3. PMID: 11253215
38. Liu H, Agam Y, Madsen JR, Kreiman G. Timing, timing, timing: fast decoding of object information from intracranial field potentials in human visual cortex. *Neuron*. 2009; 62(2):281–90. <https://doi.org/10.1016/j.neuron.2009.02.025> PMID: 19409272
39. Khaligh-Razavi S-M, Carlin J, Martin CR, Kriegeskorte N. The effects of recurrent dynamics on ventral-stream representational geometry. *Journal of vision*. 2015; 15(12):1089-.
40. Grootswagers T, Carlson T. Decoding the emerging representation of degraded visual objects in the human brain. *Journal of vision*. 2015; 15(12):1087-.
41. Kaneshiro B, Guimaraes MP, Kim H-S, Norcia AM, Suppes P. A Representational Similarity Analysis of the Dynamics of Object Processing Using Single-Trial EEG Classification. *PloS one*. 2015; 10(8): e0135697. <https://doi.org/10.1371/journal.pone.0135697> PMID: 26295970
42. Mohsenzadeh Y, Qin S, Cichy RM, Pantazis D. Ultra-Rapid serial visual presentation reveals dynamics of feedforward and feedback processes in the ventral visual pathway. *Elife*. 2018; 7:e36329. <https://doi.org/10.7554/eLife.36329> PMID: 29927384
43. Carlson T, Tovar DA, Alink A, Kriegeskorte N. Representational dynamics of object vision: the first 1000 ms. *Journal of vision*. 2013; 13(10):1-. <https://doi.org/10.1167/13.10.1> PMID: 23908380
44. Cichy RM, Pantazis D, Oliva A. Resolving human object recognition in space and time. *Nature neuroscience*. 2014; 17(3):455–62. <https://doi.org/10.1038/nn.3635> PMID: 24464044
45. Isik L, Meyers EM, Leibo JZ, Poggio T. The dynamics of invariant object recognition in the human visual system. *Journal of neurophysiology*. 2014; 111(1):91–102. <https://doi.org/10.1152/jn.00394.2013> PMID: 24089402
46. Grootswagers T, Wardle SG, Carlson TA. Decoding dynamic brain patterns from evoked responses: A tutorial on multivariate pattern analysis applied to time series neuroimaging data. *Journal of cognitive neuroscience*. 2017.
47. Contini EW, Wardle SG, Carlson TA. Decoding the time-course of object recognition in the human brain: From visual features to categorical decisions. *Neuropsychologia*. 2017.
48. Breitmeyer B, Ögmen H. *Visual masking: Time slices through conscious and unconscious vision*; 2006. Oxford University Press; 2006.
49. Thorpe SJ. The speed of categorization in the human visual system. *Neuron*. 2009; 62(2):168–70. <https://doi.org/10.1016/j.neuron.2009.04.012> PMID: 19409262

50. Cichy RM, Pantazis D, Oliva A. Similarity-based fusion of MEG and fMRI reveals spatio-temporal dynamics in human cortex during visual object recognition. *Cerebral Cortex*. 2016; 26(8):3563–79. <https://doi.org/10.1093/cercor/bhw135> PMID: 27235099
51. King J, Dehaene S. Characterizing the dynamics of mental representations: the temporal generalization method. *Trends in cognitive sciences*. 2014; 18(4):203–10. <https://doi.org/10.1016/j.tics.2014.01.002> PMID: 24593982
52. King J-R, Pescetelli N, Dehaene S. Brain mechanisms underlying the brief maintenance of seen and unseen sensory information. *Neuron*. 2016; 92(5):1122–34. <https://doi.org/10.1016/j.neuron.2016.10.051> PMID: 27930903
53. Lamme VA, Zipser K, Spekreijse H. Masking interrupts figure-ground signals in V1. *Journal of cognitive neuroscience*. 2002; 14(7):1044–53. <https://doi.org/10.1162/089892902320474490> PMID: 12419127
54. Bacon-Macé N, Macé MJ-M, Fabre-Thorpe M, Thorpe SJ. The time course of visual processing: Backward masking and natural scene categorisation. *Vision research*. 2005; 45(11):1459–69. <https://doi.org/10.1016/j.visres.2005.01.004> PMID: 15743615
55. Fahrenfort JJ, Scholte HS, Lamme VA. Masking disrupts reentrant processing in human visual cortex. *Journal of cognitive neuroscience*. 2007; 19(9):1488–97. <https://doi.org/10.1162/jocn.2007.19.9.1488> PMID: 17714010
56. Serre T, Oliva A, Poggio T. A feedforward architecture accounts for rapid categorization. *Proceedings of the National Academy of Sciences*. 2007; 104(15):6424–9.
57. Macknik SL, Martinez-Conde S. The role of feedback in visual masking and visual processing. *Advances in cognitive psychology*. 2007; 3(1–2):125–52.
58. He K, Zhang X, Ren S, Sun J, editors. Deep residual learning for image recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition*; 2016.
59. He K, Zhang X, Ren S, Sun J, editors. Identity mappings in deep residual networks. *European Conference on Computer Vision*; 2016: Springer.
60. Liao Q, Poggio T. Bridging the gaps between residual learning, recurrent neural networks and visual cortex. *arXiv preprint arXiv:160403640*. 2016.
61. Veit A, Wilber MJ, Belongie S, editors. Residual networks behave like ensembles of relatively shallow networks. *Advances in Neural Information Processing Systems*; 2016.
62. Deng J, Dong W, Socher R, Li L-J, Li K, Fei-Fei L, editors. Imagenet: A large-scale hierarchical image database. *Computer Vision and Pattern Recognition, 2009 CVPR 2009 IEEE Conference on*; 2009: IEEE.
63. Liang M, Hu X, editors. Recurrent convolutional neural network for object recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2015.
64. Cichy RM, 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*. 2016; 6:27755. <https://doi.org/10.1038/srep27755> PMID: 27282108
65. Clarke A, Taylor KI, Tyler LK. The evolution of meaning: spatio-temporal dynamics of visual object recognition. *Journal of cognitive neuroscience*. 2011; 23(8):1887–99. <https://doi.org/10.1162/jocn.2010.21544> PMID: 20617883
66. Clarke A. Dynamic information processing states revealed through neurocognitive models of object semantics. *Language, cognition and neuroscience*. 2015; 30(4):409–19.
67. Tomita H, Ohbayashi M, Nakahara K, Hasegawa I, Miyashita Y. Top-down signal from prefrontal cortex in executive control of memory retrieval. *Nature*. 1999; 401(6754):699. <https://doi.org/10.1038/44372> PMID: 10537108
68. Garrido MI, Kilner JM, Kiebel SJ, Friston KJ. Evoked brain responses are generated by feedback loops. *Proceedings of the National Academy of Sciences*. 2007; 104(52):20961–6.
69. Goddard E, Carlson TA, Dermody N, Woolgar A. Representational dynamics of object recognition: Feedforward and feedback information flows. *NeuroImage*. 2016; 128:385–97. <https://doi.org/10.1016/j.neuroimage.2016.01.006> PMID: 26806290
70. Devereux BJ, Clarke AD, Tyler LK. Integrated deep visual and semantic attractor neural networks predict fMRI pattern-information along the ventral object processing pathway. *Scientific Reports*. 2018.
71. Clarke A, Tyler LK. Object-specific semantic coding in human perirhinal cortex. *Journal of Neuroscience*. 2014; 34(14):4766–75. <https://doi.org/10.1523/JNEUROSCI.2828-13.2014> PMID: 24695697
72. Güçlü U, van Gerven MA. Deep neural networks reveal a gradient in the complexity of neural representations across the ventral stream. *Journal of Neuroscience*. 2015; 35(27):10005–14. <https://doi.org/10.1523/JNEUROSCI.5023-14.2015> PMID: 26157000

73. Kubilius J, Bracci S, de Beeck HPO. Deep neural networks as a computational model for human shape sensitivity. *PLoS computational biology*. 2016; 12(4):e1004896. <https://doi.org/10.1371/journal.pcbi.1004896> PMID: 27124699
74. Kheradpisheh SR, Ghodrati M, Ganjtabesh M, Masquelier T. Deep networks can resemble human feed-forward vision in invariant object recognition. *Scientific reports*. 2016; 6:32672. <https://doi.org/10.1038/srep32672> PMID: 27601096
75. Kheradpisheh SR, Ghodrati M, Ganjtabesh M, Masquelier T. Humans and deep networks largely agree on which kinds of variation make object recognition harder. *Frontiers in computational neuroscience*. 2016;10.
76. Khaligh-Razavi S-M, Henriksson L, Kay K, Kriegeskorte N. Fixed versus mixed RSA: Explaining visual representations by fixed and mixed feature sets from shallow and deep computational models. *Journal of Mathematical Psychology*. 2017; 76:184–97. <https://doi.org/10.1016/j.jmp.2016.10.007> PMID: 28298702
77. Chen J, Liu B, Chen B, Fang F. Time course of amodal completion in face perception. *Vision research*. 2009; 49(7):752–8. <https://doi.org/10.1016/j.visres.2009.02.005> PMID: 19233227
78. Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:14091556. 2014.
79. Szegedy C, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, et al., editors. Going deeper with convolutions. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2015.
80. Taigman Y, Yang M, Ranzato MA, Wolf L, editors. Deepface: Closing the gap to human-level performance in face verification. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2014.
81. Bengio Y, LeCun Y. Scaling learning algorithms towards AI. *Large-scale kernel machines*. 2007; 34(5).
82. Dunbar RI. Neocortex size as a constraint on group size in primates. *Journal of human evolution*. 1992; 22(6):469–93.
83. Kaas JH. Why is brain size so important: Design problems and solutions as neocortex gets bigger or smaller. *Brain and Mind*. 2000; 1(1):7–23.
84. Weaver AH. Reciprocal evolution of the cerebellum and neocortex in fossil humans. *Proceedings of the National Academy of Sciences of the United States of America*. 2005; 102(10):3576–80. <https://doi.org/10.1073/pnas.0500692102> PMID: 15731345
85. Isler K, van Schaik CP. The expensive brain: a framework for explaining evolutionary changes in brain size. *Journal of Human Evolution*. 2009; 57(4):392–400. <https://doi.org/10.1016/j.jhevol.2009.04.009> PMID: 19732937
86. Bosman CA, Aboitiz F. Functional constraints in the evolution of brain circuits. *Frontiers in neuroscience*. 2015; 9.
87. Breitmeyer BG, Ganz L. Implications of sustained and transient channels for theories of visual pattern masking, saccadic suppression, and information processing. *Psychological review*. 1976; 83(1):1. PMID: 766038
88. Portilla J, Simoncelli EP. A parametric texture model based on joint statistics of complex wavelet coefficients. *International journal of computer vision*. 2000; 40(1):49–70.
89. Taulu S, Simola J. Spatiotemporal signal space separation method for rejecting nearby interference in MEG measurements. *Physics in Medicine & Biology*. 2006; 51(7):1759.
90. Tadel F, Baillet S, Mosher JC, Pantazis D, Leahy RM. Brainstorm: a user-friendly application for MEG/EEG analysis. *Computational intelligence and neuroscience*. 2011; 2011:8.
91. Pelli DG. The VideoToolbox software for visual psychophysics: Transforming numbers into movies. *Spatial vision*. 1997; 10(4):437–42. PMID: 9176953
92. Kietzmann TC, Gert AL, Tong F, König P. Representational dynamics of facial viewpoint encoding. *Journal of cognitive neuroscience*. 2017; 29(4):637–51. https://doi.org/10.1162/jocn_a_01070 PMID: 27791433
93. Chang C-C, Lin C-J. LIBSVM: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*. 2011; 2(3):27.
94. Kriegeskorte N. Relating population-code representations between man, monkey, and computational models. *Frontiers in Neuroscience*. 2009; 3:35.
95. Kriegeskorte N, Kievit RA. Representational geometry: integrating cognition, computation, and the brain. *Trends in cognitive sciences*. 2013; 17(8):401–12. <https://doi.org/10.1016/j.tics.2013.06.007> PMID: 23876494

96. Khaligh-Razavi S-M, Bainbridge WA, Pantazis D, Oliva A. From what we perceive to what we remember: Characterizing representational dynamics of visual memorability. *bioRxiv*. 2016:049700.
97. Pedzahur E. Multiple regression in behavioral research: Explanation and prediction. London, UK: Wadsworth, Thompson Learning. 1997.
98. Gibbons JD, Chakraborti S. Nonparametric statistical inference. *International encyclopedia of statistical science*: Springer; 2011. p. 977–9.
99. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the royal statistical society Series B (Methodological)*. 1995:289–300.
100. Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N. A toolbox for representational similarity analysis. *PLoS Comput Biol*. 2014; 10(4):e1003553. <https://doi.org/10.1371/journal.pcbi.1003553> PMID: 24743308
101. Krizhevsky A, Sutskever I, Hinton GE, editors. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*; 2012.
102. Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*. 2015; 115(3):211–52.
103. He K, Zhang X, Ren S, Sun J, editors. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *Proceedings of the IEEE international conference on computer vision*; 2015.
104. Glorot X, Bengio Y, editors. Understanding the difficulty of training deep feedforward neural networks. *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*; 2010.
105. Mormann F, Kornblith S, Quiroga RQ, Kraskov A, Cerf M, Fried I, et al. Latency and selectivity of single neurons indicate hierarchical processing in the human medial temporal lobe. *Journal of Neuroscience*. 2008; 28(36):8865–72. <https://doi.org/10.1523/JNEUROSCI.1640-08.2008> PMID: 18768680