

## RESEARCH ARTICLE

# The attentive reconstruction of objects facilitates robust object recognition

Seoyoung Ahn<sup>1\*</sup>, Hossein Adeli<sup>2</sup>, Gregory J. Zelinsky<sup>3,4</sup>

**1** Department of Molecular and Cell Biology, University of California, Berkeley, California, United States of America, **2** Zuckerman Mind Brain Behavior Institute, Columbia University, New York City, New York, United States of America, **3** Department of Psychology, Stony Brook University, Stony Brook, New York, United States of America, **4** Department of Computer Science, Stony Brook University, Stony Brook, New York, United States of America

\* [ahnseoyoung@gmail.com](mailto:ahnseoyoung@gmail.com)

## Abstract

Humans are extremely robust in our ability to perceive and recognize objects—we see faces in tea stains and can recognize friends on dark streets. Yet, neurocomputational models of primate object recognition have focused on the initial feed-forward pass of processing through the ventral stream and less on the top-down feedback that likely underlies robust object perception and recognition. Aligned with the generative approach, we propose that the visual system actively facilitates recognition by reconstructing the object hypothesized to be in the image. Top-down attention then uses this reconstruction as a template to bias feedforward processing to align with the most plausible object hypothesis. Building on auto-encoder neural networks, our model makes detailed hypotheses about the appearance and location of the candidate objects in the image by reconstructing a complete object representation from potentially incomplete visual input due to noise and occlusion. The model then leverages the best object reconstruction, measured by reconstruction error, to direct the bottom-up process of selectively routing low-level features, a top-down biasing that captures a core function of attention. We evaluated our model using the MNIST-C (handwritten digits under corruptions) and ImageNet-C (real-world objects under corruptions) datasets. Not only did our model achieve superior performance on these challenging tasks designed to approximate real-world noise and occlusion viewing conditions, but also better accounted for human behavioral reaction times and error patterns than a standard feedforward Convolutional Neural Network. Our model suggests that a complete understanding of object perception and recognition requires integrating top-down and attention feedback, which we propose is an object reconstruction.

## OPEN ACCESS

**Citation:** Ahn S, Adeli H, Zelinsky GJ (2024) The attentive reconstruction of objects facilitates robust object recognition. *PLoS Comput Biol* 20(6): e1012159. <https://doi.org/10.1371/journal.pcbi.1012159>

**Editor:** Christoph Strauch, Utrecht University: Universiteit Utrecht, NETHERLANDS

**Received:** August 8, 2023

**Accepted:** May 11, 2024

**Published:** June 13, 2024

**Copyright:** © 2024 Ahn 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:** The code necessary to reproduce our results is available on the github repository at <https://github.com/ahnchive/ORA-recognition>.

**Funding:** This work was supported in part by NSF IIS awards 1763981 and 2123920 to G.Z., as well as American Psychological Association (APA) dissertation awards to S.A. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

## Author summary

Humans can dream and imagine things, and this means that the human brain can generate perceptions of things that are not there. We propose that humans evolved this generation capability, not solely to have more vivid dreams, but to help us better understand the world, especially when what we see is unclear or missing some details (due to occlusion,

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

changing perspective, etc.). Through a combination of computational modeling and behavioral experiments, we demonstrate how the process of generating objects—actively reconstructing the most plausible object representation from noisy visual input—guides attention towards specific features or locations within an image (known as functions of top-down attention), thereby enhancing the system’s robustness to various types of noise and corruption. We found that this generative attention mechanism could explain, not only the time that it took people to recognize challenging objects, but also the types of recognition errors made by people (seeing an object as one thing when it was really another). These findings contribute to a deeper understanding of the computational mechanisms of attention in the brain and their potential connection to the generative processes that facilitate robust object recognition.

## Introduction

One of the hallmarks of human visual recognition is its robustness against various types of noise and corruption [1, 2]. Convolutional Neural Networks (CNNs) trained on object classification have gained popularity as models of human and primate vision due to their ability to predict behavioral and neural activity during visual recognition [3–7]. However, CNNs can still be highly vulnerable to small changes in object appearances even when the changes are almost imperceptible to humans ([8, 9]; but see [7] for CNNs reaching human-level performance in recognizing distorted images). Recent studies also revealed that CNNs exploit fundamentally different feature representations than those used by humans (e.g., relying on texture rather than shape, [10, 11]), suggesting that current computational vision models still do not fully capture primate visual processing.

The question of how humans robustly recognize objects under challenging viewing conditions, such as noise and occlusion, has been a long-standing problem spanning various disciplines, from cognitive science [12, 13] and neuroscience [14, 15] to computer vision [16, 17]. One widely accepted hypothesis is that top-down feedback, broadly associated with attention control or bias, serves to refine object representations through the selective modulation of responses in early visual areas, thereby better aligning them with higher-level object hypotheses [18–22]. However, the computational mechanism underlying this top-down feedback remains elusive, particularly regarding how this modulation can selectively target specific locations and features of visual stimuli, as demonstrated in the experimental literature [23–26]. This challenge is particularly complex considering that object representations in higher cortical areas are largely invariant to changes in the lower-level representation of an object’s location or features [27]. In a recent perspective on the modeling of object-based attention [28], a failure to explain how coarse object representations from high-level cortex effectively modulate lower-level feature or spatial attention was identified as a critical shortcoming of current approaches.

Our approach is to apply generative modeling to this question. A generative model of perception suggests that the brain constantly attempts to generate percepts consistent with our hypotheses about the world [29–31], a process referred to as “prediction” in predictive coding theory [32] or “synthesis” in analysis-by-synthesis theory [30]. Although direct experimental evidence for these processes in the brain remains limited, various visual phenomena strongly suggest the involvement of top-down generative or reconstructive processes in shaping human visual perceptions (e.g., visual imagery, object completion, and pareidolia). We propose a new model of object recognition that generates an object reconstruction—a visualized prediction about the possible appearance of an object—and uses it as top-down attention feedback to bias

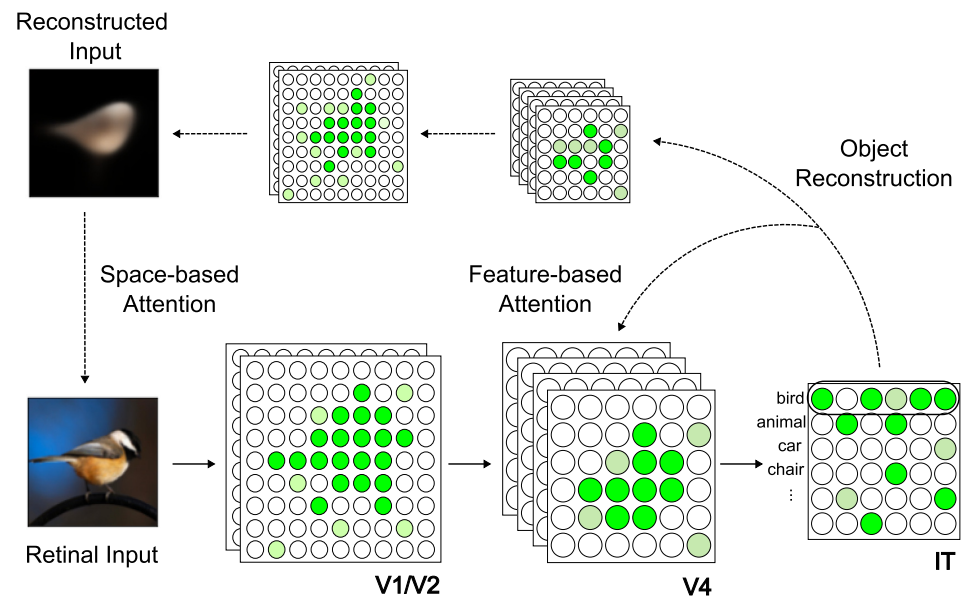
feed-forward processing. This model, called Object Reconstruction-guided Attention (ORA), is based on the fundamental premise that the process of object reconstruction serves as the underlying mechanism through which high-level object representations are used to recover specific location and feature information. This enables the model to exert precise top-down biasing control over both spatial and feature dimensions, a function typically associated with separate attention systems [28, 33, 34]. Under ORA, top-down biasing occurs by prioritizing the object features that offer the most accurate interpretation of the retinal input, that is, minimizing the reconstruction error (similar to “explain-away”, [32]). Although promising recent work has modeled human visual cortical processing using a neural network equipped with a decoder to reconstruct the input [35–42], no current models use top-down reconstruction feedback to improve the robustness of object perception and recognition.

ORA fills this gap by providing a new computational model that integrates object recognition and attention by leveraging object reconstructions to generate top-down attention feedback. As part of its novel neural network architecture, the model directs more attention to the areas and features of the input that align with the current best object reconstruction. This targeted attention effectively enhances otherwise weak bottom-up input signals, thereby reducing the interference from noise and occlusion and improving recognition performance. In extensive model evaluation, we found that ORA exhibits remarkable robustness as evidenced by its effective mitigation of noise and occlusion effects when applied to challenging recognition tasks such as MNIST-C (handwritten digits under corruptions, [43]) and ImageNet-C (real-world objects under corruptions, [44]). ORA also accounts for behavioral reaction time and error results that we observed from human participants performing the same tasks, better than a baseline feed-forward network (e.g., a ResNet; [45]). We hope that ORA’s use of object-centric representation and reconstruction will provide new insights into the neural mechanism endowing humans with robust object perception and recognition.

## Object Reconstruction-guided Attention

The concept of interactive top-down and bottom-up processing is a central theme in numerous influential cognitive neuroscience and vision theories [18–22, 29, 30, 46]. Among these theories, our model builds most closely upon Grossberg and Carpenter’s Adaptive Resonance Theory (ART). ART constructs expected category representations or *templates*, which it uses via top-down feedback to inhibit bottom-up features that do not align with these expectations (referred to as “top-down attentive matching”, [47]). Our model, ORA, advances the fixed categorical template approach of ART by incorporating a generative mechanism that dynamically selects the most probable object template based on the input, achieved through object reconstruction. This approach enables ORA to adaptively adjust its top-down expectations to accommodate variability in exemplars and noise in the visual input, thereby making it robust. ORA also demonstrates that top-down attentive biases discussed in the cognitive science literature, such as object-fitting attentional shrouds [24, 48] and the formation of robust object entities via attentive feature binding [25, 49], do not require explicit learning. Instead, these biases may emerge via the active reconstruction and use of hypothesized objects to guide bottom-up object recognition, as theorized by ORA.

Fig 1 shows how ORA might be implemented in the brain. The retinal image is first processed through the ventral stream, which consists of multiple layers of linear and non-linear filters (e.g., V1, V2, and V4) specialized in extracting the visual features of objects [27]. As information flows through these layers, neuronal receptive field sizes progressively increase and responses grow more selective and abstract (depicted by the greenish circles). For our implementation, we model the ventral stream using a Convolutional Neural Network (CNN)



**Fig 1. A potential implementation of ORA in the brain.** ORA uses object reconstruction as top-down attention feedback and iteratively routes the most relevant visual information in successive steps of a feed-forward object recognition process. This reconstruction-guided attention operates at two levels, 1) Spatial attention: a long-range projection that spatially constrains attention to the most likely object location and 2) Feature-based attention: more local feedback that biases feature binding strengths to favor the formation of a better object reconstruction. See main text for more details.

<https://doi.org/10.1371/journal.pcbi.1012159.g001>

incorporating skip connections between retinotopic feature maps [45]. At the top of this hierarchy is the inferotemporal cortex (IT), where the low-dimensional latent features are encapsulated into separate category-selective slots [50]. In computer vision, this encapsulation of information is sometimes referred to as a “capsule” [51, 52], although the term that we adopt is the more recently used “slot” [53, 54]. Note that in cognitive psychology, a slot corresponds to one “file” of information about an object as theorized by object file theory ([55]; see also the idea of category-consistent features in [56]). Unlike conventional CNNs and neural networks, which typically assign a single output neuron to represent each category, these object slots employ a vector representation in which a set of neurons represents each category and each neuron within this set encodes meaningful variations along specific feature dimensions of objects (e.g., scale, thickness, etc). This allows the model to construct robust yet flexible object representations [51, 57] that maintain selectivity for the distinctive features of objects while being tolerant to variability in these features, consistent with object file theory [55].

The object features that ORA encodes into its IT slots are then decoded to generate top-down reconstructions of objects in the visual input. This encoding-decoding process is similar to an autoencoder architecture [58], where a model first encodes a retinal input to a highly compressed representation and then expands this latent representation through a decoder to reconstruct the original input. ORA’s decoding process is different in that it iteratively generates reconstructions (rather than generating only one) and uses the early reconstructions to recover a representation of retinotopy and to facilitate the binding of object features to specific locations. ORA decodes an object reconstruction using dedicated feedback connections from an inversely mirrored CNN [59]. This CNN decoder takes the most activated object slot in ORA’s IT cortex (measured by the sum of squares of all neural activation within that slot, i.e., vector magnitude), and sequentially reconstructs the input features using a series of convolutional and unpooling layers. More details about ORA’s architecture can be found in the [Methods](#).

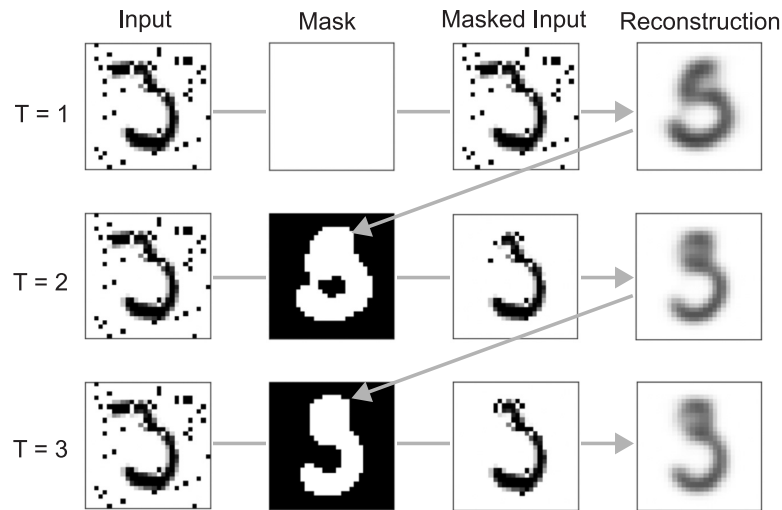
ORA relies on a neural network trained with two objectives: accurately classifying an input and generating its reconstruction. Margin loss [51] was used for classification training. By minimizing this loss, the activation of the target object slot (measured by vector magnitude) is pushed toward the desired value of 1, while the activation of non-target object slots is suppressed to near 0 (see [Methods](#)). We compute reconstruction loss as the mean squared differences in pixel intensities between the input image and the model's reconstruction. During model evaluation (i.e., inference), ORA leverages object reconstructions as top-down attentional templates, selectively routing the bottom-up processing of features that align with top-down expectations, similar to ART's iterative "matching" process [18, 47]. By reconstructing the details of object appearances, ORA can exert nuanced attention control over various levels of visual features, both spatial and non-spatial, indicated by the downward arrows in [Fig 1](#). The iterative process of top-down matching and bottom-up routing in ORA continues until the network stabilizes, determined by the model's classification confidence surpassing a predefined threshold. See Section Measuring ORA's reaction time of the [Methods](#) for details about how this confidence was estimated.

### Reconstruction-based spatial masking

Our goal in this study is to show that ORA's reconstruction-guided feedback can serve two significant roles with respect to attention. The first of these is the control of spatial attention. Spatial attention refers to the prioritization of information at a spatial location by filtering or masking out information from all other locations [23], with the neural expression of this being enhanced activity at the attended location and inhibition elsewhere. Several experimental studies have reported that spatial attention can dynamically configure itself to align with an object's shape [24, 60–62], and this has even been modeled as object-fitting attentional shrouds [48, 63]. ORA posits that the visual system constructs an attentional shroud around an object by iteratively comparing a model generated object reconstruction with incoming sensory signals. This iterative process involves generating a spatial priority mask, a function known to be mediated by the dorsal pathway encompassing posterior parietal cortex (PPC; [64, 65]), which we implement through simple binary thresholding. The resulting spatial mask is then utilized to modulate activity in early ventral areas, suppressing neural responses to features outside the expected object shape. ORA repeated this process until the model's classification confidence surpassed a predetermined threshold or until the limit of five forward processing steps was reached, whichever came first (see [Methods](#) for additional details). [Fig 2](#) demonstrates how the reconstruction-based spatial masking effectively filters out background noise, thereby enhancing classification accuracy. Initially, ORA hypothesized that the input was the digit 5, but after filtering out noise using its reconstruction mask the model correctly changed its hypothesis to a 3, matching the ground truth. On average, ORA's iterative spatial masking enabled it to successfully recover the object's correct identity on over 90% of the trials when the model started with an incorrect hypothesis (see [S1 Table](#)). Previous research has similarly argued for benefits from a reconstruction-based masking approach, but in the context of augmenting datasets for improving the learning of robust visual representations for use in downstream tasks [51, 66]. ORA is different from previous modeling work [35, 67, 68] in that it leverages object reconstructions learned during training to generate spatial attention masks during inference. These masks therefore bias model processing to the hypothesized shape of the object.

### Reconstruction-based feature binding

A second role of attention served by ORA's reconstruction-based top-down biasing is the binding of features into objects [53]. Feature binding refers to the process by which the varied features of an object are spatially integrated to form a coherent object representation (e.g.,



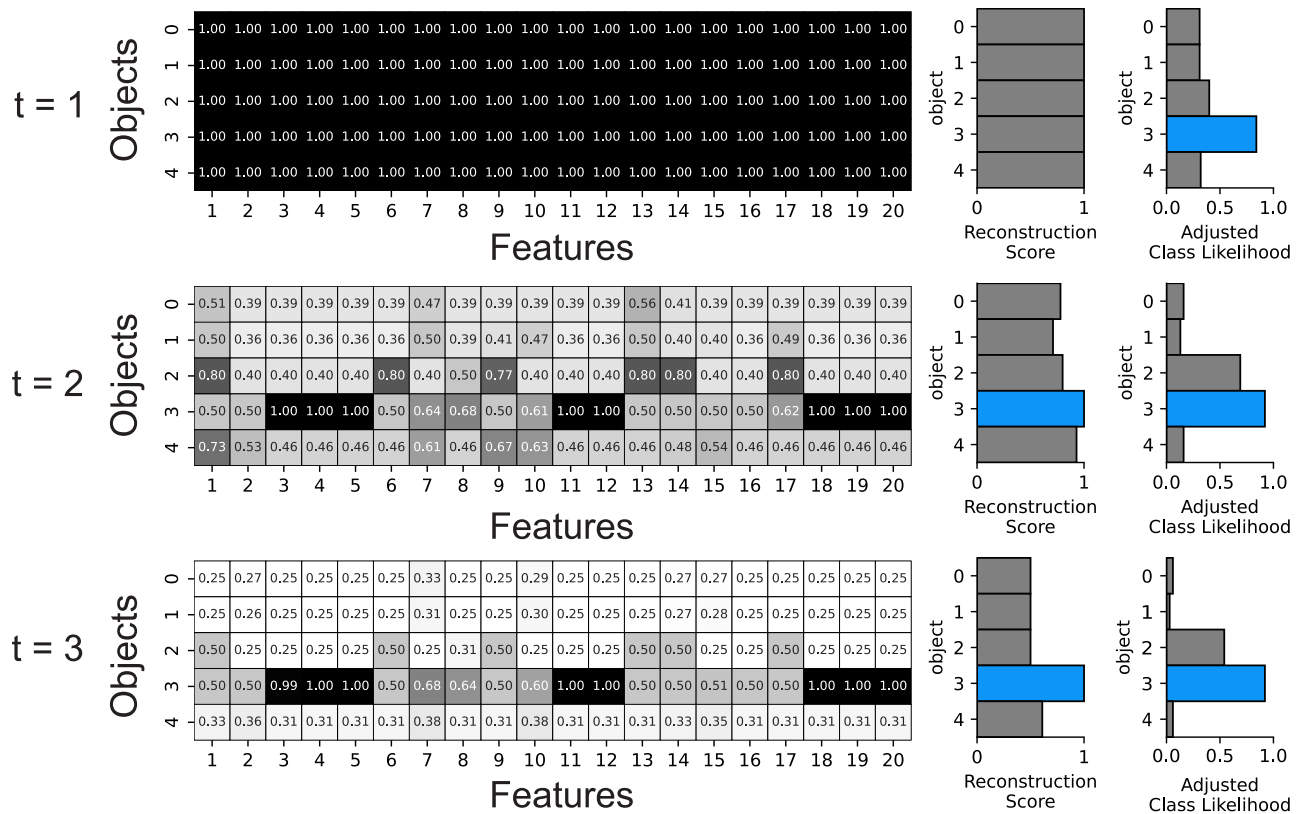
**Fig 2. Object reconstructions and generated spatial masks for three forward processing steps (out of a possible five in total).** The spatial mask serves to limit ORA's attention to the shape of the most likely object. In this example, the model's most likely digit hypothesis changed from a 5 to a 3. The groundtruth class is 3.

<https://doi.org/10.1371/journal.pcbi.1012159.g002>

horizontal, vertical, and curved lines are bound top-to-bottom to form the digit 5), and this process is believed to be fundamental to robust object perception and recognition [25, 48, 49]. Recent advances in computer vision demonstrated that training models to reconstruct objects resulted in robust feature binding enabling the grouping and tracking objects within a scene [51, 53, 54]. ORA accomplishes feature binding by leveraging an object reconstruction to determine the features that should be associated with specific objects. Specifically, ORA modulates the weights between the V4 and IT layers, suppressing the binding strength between a feature and an object if it leads to a poorer reconstruction match to the input (see [Methods](#) for pseudocode). By selectively suppressing information incongruent with ORA's top-down object reconstruction, this algorithm progressively refines neural selectivity to converge towards the most probable object hypothesis, thereby narrowing the population response. ORA's binding mechanism is consistent with the information tuning model of feature-based attention [26, 33], and it can be interpreted as a form of biased competition where the bottom-up control for neural resources (i.e., receptive fields) is biased by a higher-level representation of an object [69–71]. [Fig 3](#) shows ORA's binding coefficients (which modulate the original weights between the V4 and IT layers) sharpening over iterations and becoming more focused on the specific object features that best match the input. Binding coefficients were updated over three iterations during inference (testing); coefficients were not updated during model training. Available computing resources required us to limit ORA's reconstruction-based feedback to only three iterations, however we observed during piloting that three steps was usually sufficient for ORA to converge on a classification.

## Results

We evaluated ORA in challenging object recognition tasks. First, we compared ORA's recognition performance to that of a standard feed-forward CNN using a version of the MNIST dataset designed to test for effects of occlusions and specific forms of visual noise (MNIST-C; [43]). Both ORA and the CNN were trained on "clean" handwritten digits and then systematically tested on each form of corrupted digit. This evaluation revealed that ORA used its



**Fig 3. Step-wise visualizations of reconstruction-based feature binding.** The matrices show binding coefficients between object slots (rows) and feature slots (columns) for three time steps. The binding coefficients were initialized to one at  $t = 1$ , resulting in the dark matrix on the top, but in the subsequent time steps there is a rapid selective suppression of coefficients (matrix cells becoming lighter) as the model learns the object features leading to higher reconstruction accuracy. In the illustrated example, the coefficient matrix becomes sparser with each iteration as ORA focuses its attention on the features of the digit 3 (darker line forming along the fourth row). The middle column of bar graphs show the hypothesis yielding the highest reconstruction score (in blue), and the rightmost column shows the class likelihood adjusted by the reconstruction score. For clarity, only 5 object slots and 20 feature slots are illustrated, but the full figure can be found in [S1 Fig](#).

<https://doi.org/10.1371/journal.pcbi.1012159.g003>

reconstruction-based feedback to filter out irrelevant information and this led to greater recognition robustness to the local changes in texture introduced by the corruptions. Second, we ablated ORA's spatial masking and feature binding mechanisms and found the former to be most important for the accuracy of ORA's predictions and the latter to be most important for the speed (i.e., number of iterations) that ORA took to make its recognition judgment. Third, we collected reaction time (RT) data from humans performing the same recognition tasks and found that, on trials having the longest RTs, ORA also needed more reconstruction-based feedback before recognizing the corrupted digit. In a fourth analysis we found that ORA and humans also shared the same pattern of errors (i.e., the same digit confusions). Lastly, we conducted a limited evaluation of ORA's generalization to natural images by using the ImageNet-C dataset [44], which consists of various corruptions but now applied to ImageNet objects. Details about the datasets and experimental procedures can be found in the [Methods](#).

### ORA improves recognition performance under visual corruptions

We used the MNIST-C dataset [43] to evaluate the effectiveness of ORA's reconstruction-based feedback in improving recognition robustness under visual corruptions. This dataset

**Table 1. Model comparison results using MNIST-C and MNIST-C-shape datasets.** Recognition accuracy (means and standard deviations from 5 trained models, hereafter referred to as model “runs”) from ORA and two CNN baselines, both of which were trained using identical CNN encoders (one a 2-layer CNN and the other a Resnet-18), and a CapsNet model following the implementation in [51].

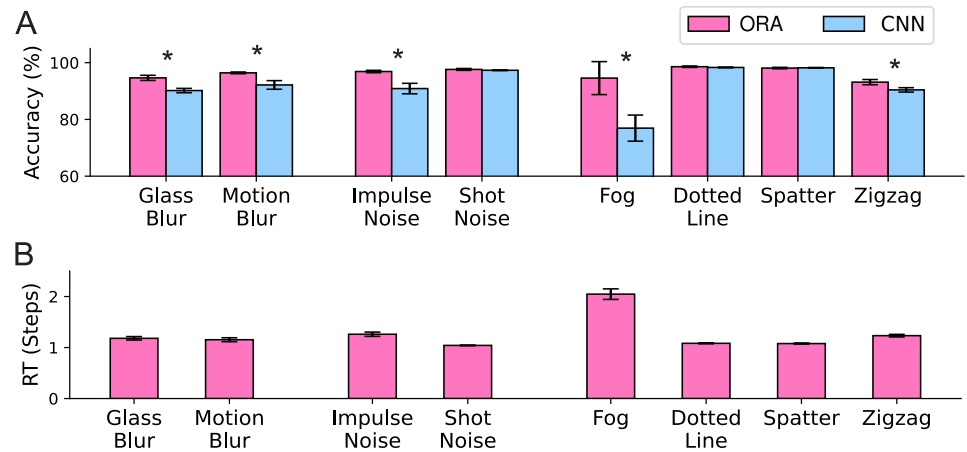
Dataset	ORA		CNN		CapsNet
	Resnet-18	2 Conv	Resnet-18	2 Conv	
MNIST-C	<b>91.84 (0.67)</b>	88.32 (0.34)	81.94 (0.43)	88.94 (0.64)	75.14 (1.00)
MNIST-C-shape	<b>96.24 (0.74)</b>	92.82 (0.90)	86.43 (0.44)	91.80 (0.88)	81.87 (0.55)

<https://doi.org/10.1371/journal.pcbi.1012159.t001>

includes 15 types of corruptions: affine transformation, noise, blur, occlusion, and others (See [S2 Fig](#) for example images of each). We also report results separately on a subset of the dataset referred to as **MNIST-C-shape**, which includes those corruptions relating mainly to object shape. Given ORA’s use of object-centric feedback, we hypothesized that robustness benefits would be greatest for these shape-specific corruptions. As shown in [Table 1](#), this hypothesis was largely confirmed. We report results for two types of CNN model backbones used by ORA to approximate the features extracted by ventral visual brain areas, one a shallow 2-layer CNN (2 Conv) and the other a deeper 18-layer CNN with skip connections (Resnet-18; [45]). We also implemented three baseline models for comparative evaluation. Two were CNNs having the same deep (Resnet-18) or shallow (2 conv) backbones as ORA, but with added fully-connected layers for digit classification. For statistical analyses, we only compared the better performing of the two CNN baselines to ORA, which from [Table 1](#) was the shallower 2 Conv CNN. We therefore refer to the better-performing 2 Conv baseline as “our CNN” for the remaining analyses. The third baseline was a CapsNet model [51]. Details about all three baselines can be found in the [Methods](#).

The best recognition accuracies on both MNIST-C and MNIST-C-shape were from ORA with a Resnet-18 encoder. Moreover, during piloting we observed that using only the low-spatial frequencies in the image to generate the object reconstruction was sufficient to produce recognition robustness comparable to what is achieved with full-resolution object reconstruction ([S2 Table](#)). This shows that effective attention biases can be mediated by low-frequency neural feedback [21]. ORA outperformed our CNN baseline by a significant margin ( $t(8) = 7.026$ ,  $p < 0.001$ , *Cohen’s d* = 4.444), particularly excelling on the MNIST-C-shape dataset where it showed an improvement of over 4.5% ( $t(8) = 8.617$ ,  $p < 0.001$ , *Cohen’s d* = 5.450). Despite incorporating object-centric representations for reconstruction (but lacking top-down attention feedback), the CapsNet model exhibited poor performance compared to ORA or our CNN baseline, indicating that the addition of reconstruction alone does not translate into improved accuracy and robustness. ORA was even more accurate than models specifically designed for generalization, such as CNNs trained using PGD/GAN adversarial noise ([S3 Table](#)). Lastly, we observed that increasing the capacity of the encoder from a 2-layer to an 18-layer CNN led to over-fitting by the CNN baseline and a consequent degradation of its accuracy from 88.94% to 81.94% on MNIST-C. ORA suffered no such degradation and achieved higher accuracy with the higher capacity encoder.

[Fig 4A](#) shows specific comparisons between ORA and our CNN baseline (the second best model) for each of the corruption types from MNIST-C-shape. For comprehensive results, refer to [S3 Fig](#). ORA consistently outperforms the CNN baseline and exhibits particularly strong performance under the fog corruption ( $t(8) = 5.318$ ,  $p < 0.001$ , *Cohen’s d* = 3.363). As shown in [Fig 4B](#), on average ORA confidently recognizes corrupted digits in 1.4 steps (measured as the number of forward processing steps required to reach a confidence threshold; see [Methods](#)). However, it takes significantly longer (more than 2 steps) under the fog corruption,



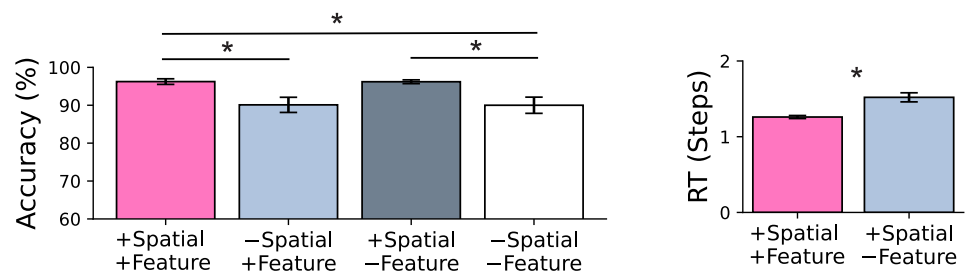
**Fig 4. Model accuracy and reaction time for specific corruption types from MNIST-C-shape.** **A:** ORA vs our CNN using each model’s best performing encoder (Resnet-18 and a 2-layer CNN, respectively). **B:** ORA’s reaction time (RT), estimated as the number of forward processing steps required to reach a confidence threshold. While many digits could be recognized in only one feed-forward pass, feedback mechanisms become useful for addressing specific types of noise, such as fog, or when the identity of the input is notably ambiguous, as demonstrated in S5 Fig. Error bars indicate standard deviations from 5 different model runs. Asterisks (\*) indicate statistical significance at a p-value of < 0.05.

<https://doi.org/10.1371/journal.pcbi.1012159.g004>

suggesting that top-down attention feedback is particularly helpful in cases of visual corruptions that remove high-spatial frequencies.

### Revealing distinct roles of spatial and feature-based attention

ORA uses two types of reconstruction-based feedback, roughly corresponding to spatial and feature-based attention, and we evaluated their unique contributions by ablating each from the best performing ORA model (i.e., Resnet-18 encoder) and observing the effect on performance. The left panel of Fig 5 shows that ORA’s recognition accuracy stemmed from its ability to use a generated reconstruction to perform spatial masking. Removing spatial masking led to a significant drop in performance of over 6% (-Spatial / +Feature;  $t(8) = 4.682$ ,  $p\text{-adjust} = 0.009$ ,  $Hedges' g = 2.674$ ). Ablating the feature binding component did not significantly affect prediction accuracy (+Spatial / -Feature;  $t(8) = 0.294$ ,  $p\text{-adjust} = 1.0$ ,  $Hedges' g = 0.168$ ), but



**Fig 5. Results from ablating (-) ORA’s reconstruction-based spatial masking and feature binding components.** Leftmost bars in pink indicate the complete (no ablation) ORA model. Left: Recognition accuracy. Right: reaction time, approximated by the number of forward processing steps taken to recognize digits with confidence. Error bars are standard deviations from 5 different model runs. Asterisks (\*) indicate statistical significance at a p-value of < 0.05 with Bonferroni corrections.

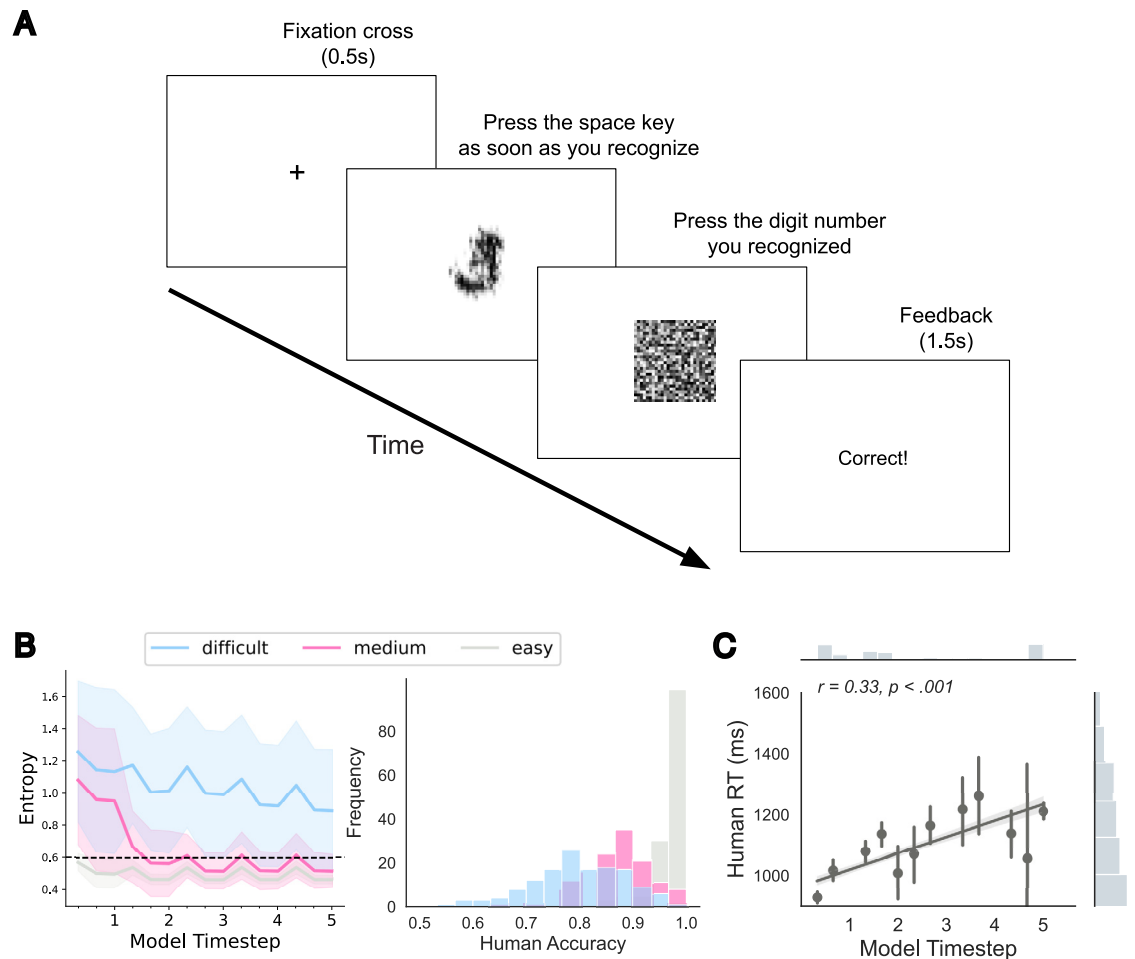
<https://doi.org/10.1371/journal.pcbi.1012159.g005>

this does not mean that it had no effect on performance. As shown in the right panel of Fig 5, feature-binding ablation instead lead to an increase in ORA's RT ( $t(8) = 10.207, p < 0.001$ , Cohen's  $d = 6.455$ ). Note that we cannot report the corresponding spatial ablation because ORA's RT is based on the number of spatial-mask iterations and would therefore always be zero. These findings suggest that ORA uses spatial and feature-based attention to perform distinct functional roles. Spatial attention improves recognition accuracy and enhances the neural signal-to-noise ratio by masking out noise from irrelevant spatial locations. On the other hand, feature-based attention enhances the overall discriminability of the encoder by suppressing features associated with competing object hypotheses. However, when the ventral encoder has a lower capacity for discriminability (i.e., the shallower 2-layer CNN), we observed that reconstruction-based feature binding also significantly influences both ORA's RT and accuracy (S4 Fig). The distinct roles played by feature and spatial attention in ORA's function are further supported by the differential impact that each had on the model's ability to recover from an initially incorrect hypothesis (S1 Table). On average, spatial attention led to an incorrect hypothesis changing to a correct prediction on approximately 6% of these trials, with a maximum of 40% incorrect-to-correct changes observed for the fog corruption. In contrast, iterations of feature-binding resulted in hypothesis changes on less than 1% of trials. This suggests that ORA uses feature-binding primarily to increase the efficiency of its recognition of MNIST-C digits whereas spatial masking is used to modify digit hypotheses and improve recognition accuracy.

### Comparison with human reaction time

As shown in Fig 4B, ORA's reaction time (RT) behavior changes depending on the difficulty of the recognition discrimination. The model iteratively generates the object reconstruction that it hypothesizes is the most plausible explanation of the input and uses it as a spatial mask for the subsequent feed-forward step. This process repeats until the model achieves a predetermined level of classification confidence measured by the entropy of the class likelihoods (softmax output). In the majority of cases, ORA reached the confidence threshold within a single feed-forward sweep (Fig 4B). However, on trials where the digit is inherently ambiguous or becomes ambiguous due to noise, ORA alternated between two object hypotheses, resulting in a longer time required for one to reach the confidence threshold (S5 Fig).

To assess how this RT behavior of ORA compares to human RTs, we conducted a psychophysical experiment on humans who were asked to perform the same challenging MNIST-C recognition tasks (Fig 6A). Participants ( $N = 146$ ) viewed 480 images of digits and were instructed to press a keyboard spacebar as soon as they recognized the object in the image. Viewing time was response terminated and backward masked (see Methods for details). This RT judgment was followed by a second accuracy check in which participants used the keypad to enter the recognized digit. Only correct trials were used for the RT analysis. We selected MNIST-C digits for use as experimental stimuli based on the level of difficulty predicted by ORA, with our aim being to see whether human RTs will show a similar pattern. Specifically, we created easy, medium, and difficult conditions consisting of digits taking ORA either 1 step, 2–3 steps, or 4–5 steps to recognize, respectively (Fig 6B, left panel). Fig 6B (right panel) shows evidence for an alignment between human recognition accuracy and the three difficulty groups predicted by ORA. Human accuracy was by far the highest in the easy condition (i.e., only a single feed-forward step was needed by ORA for accurate model recognition), followed by the medium (2–3 steps) and difficult (4–5 steps) conditions, with the latter showing the most errors. The mean human RT was 1130ms ( $SD = 530ms$ ), with a significant effect of condition difficulty,  $F(2, 284) = 105.266, p < 0.001, \eta^2 = 0.426$ . Mean RTs in the difficult condition



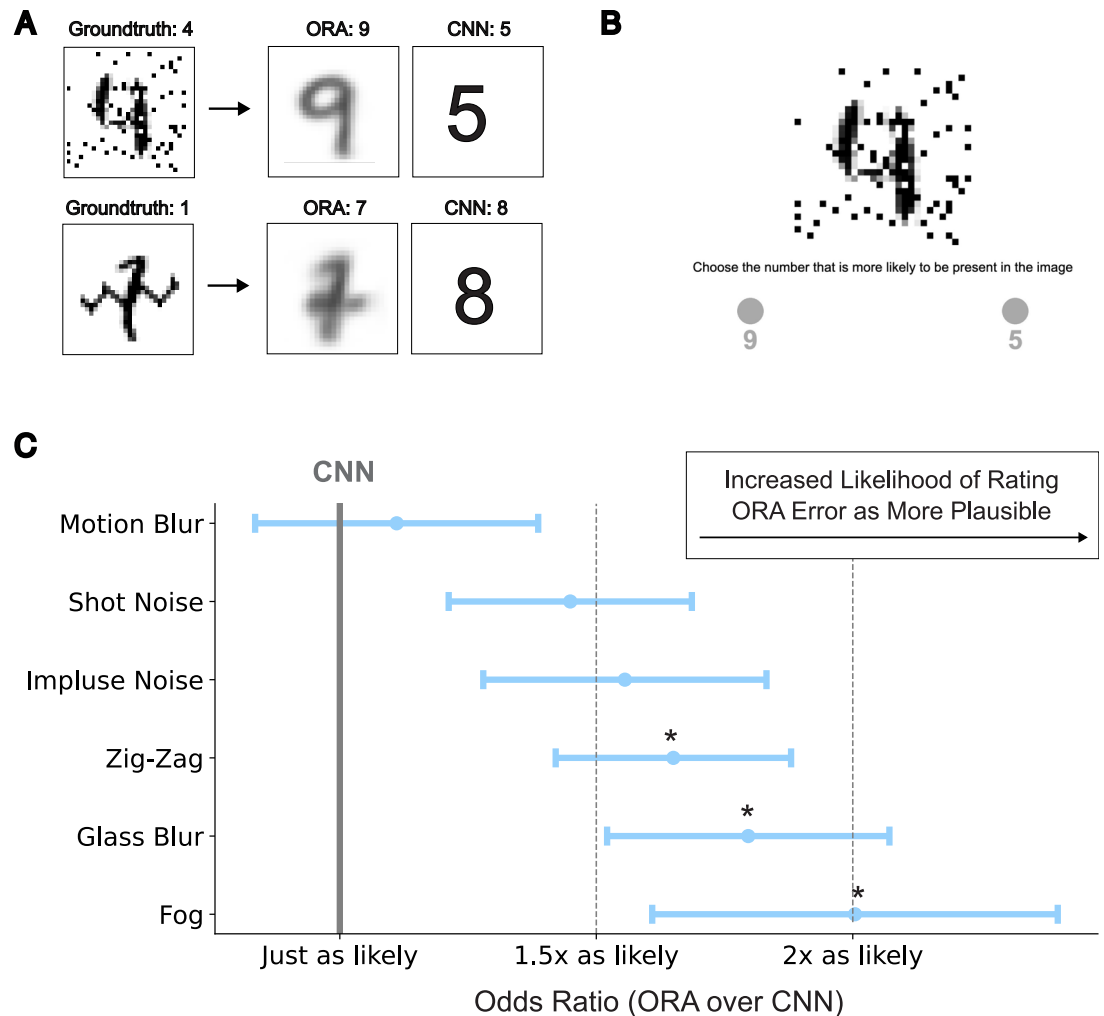
**Fig 6. Human recognition experiment and results.** **A:** Overview of the behavioral paradigm for measuring recognition reaction time (first button press) and accuracy (second button press). **B:** Stimuli were grouped into easy, medium, and difficult levels of recognition difficulty based on the number of forward processing steps needed for ORA to reach a predetermined recognition confidence threshold (dashed line; left panel). Colored regions indicate the corresponding standard deviations (SDs) across images. The unnumbered ticks on the x-axis indicate three local iterations of feature-based attention occurring within each forward processing (i.e., global spatial-masking iteration). The right panel shows distributions of human accuracy (averaged over participants) for the three difficulty conditions. **C:** Correlation between human RT and ORA's RT, with the normalized marginal distribution of each variable shown on the top and right panels. Error bars represent standard error estimates bootstrapped from 1000 samples. Correlation plots for all MNIST-C corruption types can be found in [S6 Fig](#).

<https://doi.org/10.1371/journal.pcbi.1012159.g006>

were 1265ms (SD = 649ms), but these decreased to 1170ms (SD = 517ms) in the medium condition and 985ms (SD = 473ms) in the easy condition. Lastly, [Fig 6C](#) shows a statistically significant positive Pearson correlation between human RTs and ORA's RTs ( $r = 0.33$ ,  $p < 0.001$ ). This relationship extends our previous analysis of recognition difficulty by suggesting that ORA and humans found individual trials to be similarly difficult.

### ORA's errors and explain-away behavior

To gain a clearer insight into underlying processing mechanisms, here we make a deeper comparison between the errors generated by ORA and those from the best-performing CNN from [Table 1](#). We found that ORA exhibits interesting systematic errors as it tries to explain-away random noise in the image and reconstruct it into a specific pattern that it has learned. For



**Fig 7. ORA often hallucinates a non-existing pattern out of noise, and these errors are perceived as more likely by humans compared to errors made by our CNN baseline.** A: Example of ORA's explain-way behavior and the resulting errors in classification. B: Two alternative forced choice experimental procedure requiring participants to choose which of two model predictions, one from ORA and the other from our CNN baseline (but both errors), is the more plausible human interpretation of the corrupted digit. C: The relative likelihood of ORA's errors being perceived as more like those from a human, quantified as an odds ratio (a higher value indicates a more plausible ORA error). Error bars indicate standard errors, and asterisks (\*) indicate odds ratios greater than 1 at a significance level of  $p < 0.05$ .

<https://doi.org/10.1371/journal.pcbi.1012159.g007>

example, the model often connects noise speckles to generate a line, as demonstrated in the first row of Fig 7A, thereby mistakenly reconstructing a digit 9 instead of the ground truth digit 4. Similarly, ORA occasionally misinterprets zig-zag noise as being part of the digit, thereby reconstructing a digit 7 instead of the ground truth digit 1 (second row). These types of interpretable errors were not observed in the CNN baseline.

To determine whether human recognition errors are more aligned with ORA's or our CNN baseline, we collected data from 184 participants in a two-alternative forced-choice task. On each trial, participants were shown a corrupted digit together with two mistaken model predictions, one from ORA and the other from the CNN, and they had to indicate which they thought would be more plausibly generated by a human (Fig 7B; see Methods for details). Note that we selected MNIST-C images as stimuli to satisfy two constraints: (1) that both ORA

and our CNN baseline made an error in predicting the ground truth digit, and (2) that these two model errors differed from each other. To evaluate our data we computed an odds ratio, which measures the relative perceived likelihood of ORA's errors compared to the CNN's errors. We found that this ratio was consistently greater than 1 (with 1 indicating equally likely) for every type of corruption (Fig 7C). This was particularly true for the fog corruption, where ORA's errors were rated twice as likely than those made by the CNN baseline ( $t(19) = 2.61$ ,  $p = 0.02$ , *cohen's d* = 0.58, testing against being > 1). Interestingly, the fog corruption type is also where ORA exhibited the highest performance advantage over the CNN (Fig 4). These results suggest that humans find ORA's errors to be plausible, at least more so than errors from our CNN baseline.

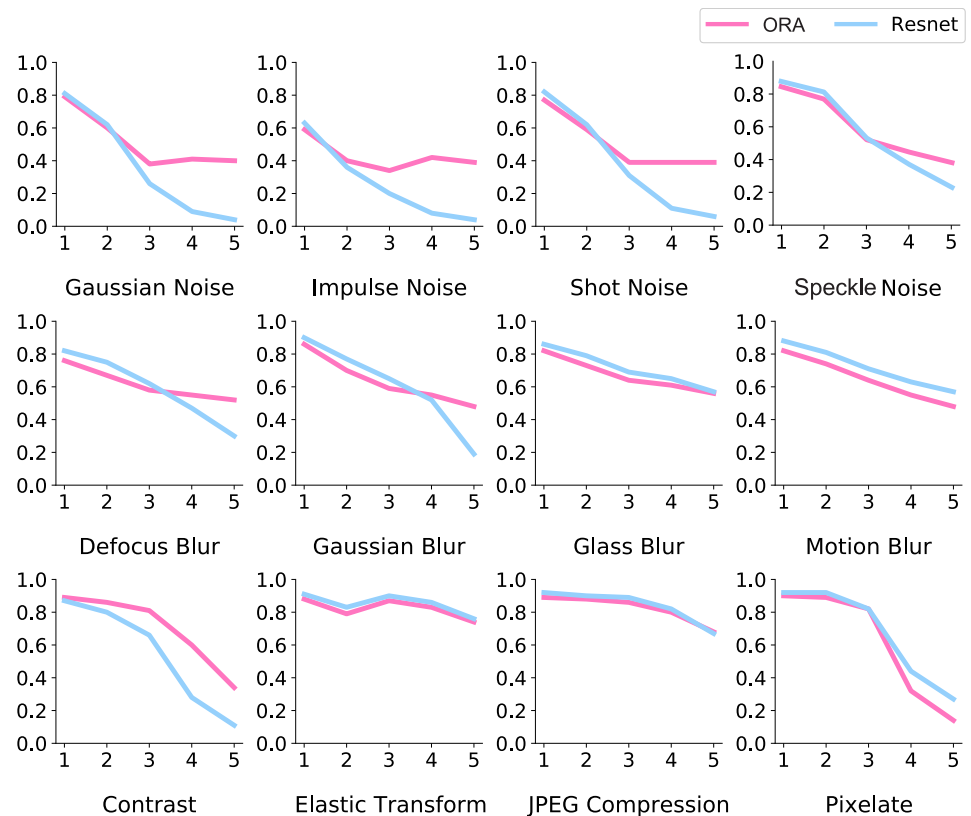
### Generalization to natural images

We evaluated whether ORA's reconstruction-based feedback mechanism, shown to be beneficial for digit recognition, would generalize to more complex images of natural objects. We did this using the ImageNet-C dataset [44], which is similar to MNIST-C but with visual corruptions applied to ImageNet objects [72]. For model training we used a dataset called 16-class ImageNet [73], which was developed to teach models a more human-like semantic structure by selecting from ImageNet a more comparable set of 16 basic-level categories (see Methods for the list of categories). To accommodate the complex backgrounds upon which objects appear in ImageNet images, we incorporated a recent segmentation model [74] to estimate the ground-truth object boundaries, then trained the model to reconstruct only the pixels within the object segments. The resulting object reconstructions capture coarse representations of object shape, missing texture detail (see S7 Fig).

Despite the coarseness of the generated reconstructions, their use as feedback provided sufficient spatial masking to significantly improve recognition performance. Fig 8 reports model recognition accuracy (y-axes) for five increasing levels of corruption (x-axes) from the ImageNet-C dataset. We compared ORA's performance to a CNN baseline model, now a ResNet-50 trained to recognize ImageNet objects. ORA's feed-forward ventral process was approximated by the identical ResNet-50 (see Methods). Due to limitations in computational resources, we were able to conduct only a single training for each model (i.e., no error bars). Still, a clear pattern emerged. Although both models exhibit similarly accurate recognition at low corruption levels, at higher levels ORA is more robust on several types (on average > 7% than the CNN baseline without reconstruction feedback; S8 Fig). The most notable benefits are observed with noise corruptions (top row), where accuracy increased for Gaussian noise by up to 35%. For some types of blur, the CNN baseline performed as well as ORA or slightly better, with the CNN working best for a pixelation corruption (bottom-right). These effects, however, are generally not as large as the improvements achieved by ORA (see S8 Fig for more results).

### Discussion

By whatever name is preferred, our common visual experience suggests that top-down, generative or reconstructive processes play a significant role in creating our visual perceptions. This is evident in phenomena such as visual imagery, object completion, and pareidolia. However, there are still open questions regarding *why* we have this ability to reconstruct objects and *what* functional role it plays in shaping our percepts [75–77]. Through a combination of computational modeling and behavioral experiments, we demonstrate that object reconstruction can serve as a mechanism for generating an effective top-down attentional template that can bias bottom-up processing to support challenging object recognition tasks. We



**Fig 8. Model performance on the ImageNet-C dataset.** Classification accuracy on the y-axes and level of corruption on the x-axes. ORA often demonstrates a clear advantage over the CNN baseline (ResNet-50), particularly with high levels of noise corruption.

<https://doi.org/10.1371/journal.pcbi.1012159.g008>

implemented a neural network model of Object Reconstruction-guided Attention (ORA) that uses an autoencoder-like architecture [51, 58]. ORA consists of an encoder that encodes the visual input into low-dimensional object-centric representations and a decoder that reconstructs objects from these representations. These object reconstructions are generated iteratively, enabling their use in routing the most relevant spatial and feature information to feed-forward object recognition processes. We found that ORA is highly robust in recognizing digits under a wide range of corruptions, outperforming other feed-forward model baselines. Theoretically, the functions performed by ORA's reconstruction-based attention mechanism are similar to functions performed by spatial and feature-based attention [25, 48]. Empirically, we demonstrated that ORA is more predictive of human RTs and error patterns in recognition tasks compared to our CNN baselines. Lastly, we showed that the benefits of ORA's reconstruction-based attention mechanisms generalize to recognition tasks using natural objects, with larger benefits observed in more challenging recognition tasks.

ORA shares many parallels with other influential cognitive neuroscience theories, including Mumford's pattern theory [78], Bar's top-down prediction model [21], and other modeling approaches relying on interacting top-down and bottom-up processes [18–20, 29, 30, 46]. Among these models, ORA builds most closely on adaptive resonance theory (ART; [18]). Both ORA and ART employ top-down templates to identify matches within bottom-up signals, selectively routing features that conform to these top-down expectations. However, a key difference lies in the formulation of the top-down templates. In ART, these templates are

categorical and similar to prototypical representations of a given category [79]. In contrast, ORA creates a top-down template from the visual input by using its generative attention mechanism to reconstruct an object category. It combines an autoencoder-like architecture with object-centric feature encapsulations to generate individualized top-down templates that are reconstructions from the individual image inputs, thus allowing for instance-specific features to be incorporated into categorical templates. This capability enables ORA to make detailed predictions about hypothesized objects and to facilitate precise attention control over low-level features—a critical missing link in the modeling of object-based attention [28]. ORA's iterative and reconstruction-guided attention control over low-level features enables it to adjust its confidence in hypotheses, correct its errors, and explore alternative hypotheses without needing explicit “resets” or relying on memory-based inhibition, as implemented in ART.

One question not answered by our current work is how completely an object must be reconstructed in order to serve as an effective attentional template. Our finding that even low-resolution object reconstructions are sufficient to generate attentional biases and performance benefits align with previous studies showing that the brain initially processes low spatial frequency information about objects for the purpose of generating top-down feedback signals for guiding the bottom-up processing of finer-grained details [21, 80, 81]. However, effects of feature-based attention are also found in a range of specific low-level features (e.g., color, orientation, etc, [82]), making further computational experiments needed to determine whether optimization of the information and level of object reconstruction might yield even greater performance benefits. Relatedly, although ORA's current high compression rate limited its ability to reconstruct fine-grained object detail, it may be possible to address this limitation by exploiting recent advancements in decoder architecture and composition-based methods in the computer vision literature to generate reconstructions of varying kind and quality [83, 84]. These questions, and broader questions regarding the biological plausibility of a reconstruction-based modeling approach, would benefit by additional study of the neural mechanisms underlying representation generation and object reconstruction in the brain. Studies on visual imagery have shown substantial overlap in neural processing between perception and mental imagery [77, 85], and more comparisons revealing similarities and differences between patterns of brain activity evoked by seen images and imagined images could provide valuable insights into the nature of the reconstruction process.

ORA has close theoretical ties to object-based attention mechanisms in that it uses object-centric representations (i.e., slots) to create top-down attentional biases [28, 86–88]. However, while theories and models of object-based attention [89, 90] have focused on the low-level grouping mechanisms responsible for combining visual elements into object parts and complete objects (e.g., gestalt rules; [91]), our generative approach highlights the role of top-down feedback in the creation of object entities. Specifically, we showed how top-down object-based feedback can mediate the selective binding of features and locations into objects, thereby implementing a form of object-based attention by integrating spatial and feature-based attentive processes [34, 92, 93]. ORA accomplishes this by using top-down generative feedback to produce the most plausible object appearance given a visual input; it selects image locations and features that align most closely with the reconstructed object, a process biasing low-level neural responses to accurately reconstruct the input object. This fills a gap in the modeling of object-based attention by providing a mechanistic explanation of how high-level top-down expectations, such as object representations in the inferotemporal cortex (IT), can modulate low-level bottom-up visual features [28].

Object-based attention has been studied using tasks ranging from object recognition to grouping to multiple-object tracking [88, 90, 94]. Our study focused on object recognition under conditions of occlusion and noise, making our work more aligned with the use of

object-based attention to create robust object entities (e.g., basic feature binding and grouping processes; [25, 49]) and less aligned with work emphasizing the operations performed on these entities and their limitations (e.g., multi-object tracking). Given the wide range of ways that researchers have conceptualized object-based attention, an important future direction will be to extend ORA and its evaluation to object-based attention tasks involving multiple objects. One potential area of investigation is perceptual grouping, which involves segregating an object of interest from other overlapping objects or a complex background. In our preliminary work on this question, shown in [S9 Fig](#), we found that ORA is showing promising results in segregating overlapping digits. This means that ORA has the potential to demonstrate classic effects of object-based attention, such as an ability to attend selectively to spatially overlapping faces or houses [94, 95]. Another interesting direction will be to apply ORA to the creation and maintenance of multiple objects in a scene over time and space. This task will require the model to maintain feature bindings for different objects and update these bindings as the objects' locations and appearances change over time, an essential step in extending ORA to multiple-object tracking. Such an extension would likely lead the use of dynamic slot representations and potentially reveal higher-level relationships between attention control and its modulation by short-term memory [71] or visual pointers [96].

Our study has broad implications for the fields of behavioral vision science and computer vision. While there is still only limited empirical evidence for the use of generative processes in vision [76, 97], there is growing interest in capturing the robustness of human perception by modeling it as a reconstructive process [35–42]. A reconstruction framework offers the potential to unify our understanding of top-down cognitive processes, including attention biases [28], memory retrieval [39], hypothesis testing [21], and mental imagery [77]. We also hope that our work will inspire the development of novel computational architectures in computer vision. Methods that train models to reconstruct objects have been shown to enhance generalization performance [54, 98], but this approach has primarily been explored in the context of pretraining or incorporating auxiliary loss [51, 66]. Our work is the first to show that a reconstruction approach can effectively model attentional biases during inference, which had been an open research question [99].

In summary, our study demonstrates that the attentive reconstruction of objects can serve as a plausible neural mechanism for achieving robust object perception and recognition. We model this mechanism as ORA, which uses object reconstruction as a top-down attention bias to focus bottom-up processing on information relevant to what it considers to be the most plausible (i.e., recognizable) object shape in the visual input. This object reconstruction-based mechanism unifies space-based and feature-based modes of attention into one computational attention system.

## Materials and methods

### Ethics statement

This study was approved by the Stony Brook University Institutional Review Board, and written consent was obtained from all participants.

### ORA implementation

For ORA's encoder, which models the ventral stream, we used different pytorch implementations of feedforward CNNs: one a shallow 2-layer CNN (<https://github.com/pytorch/examples/blob/main/mnist/main.py>; 2 Conv) and a deeper CNN (<https://github.com/pytorch/vision/blob/main/torchvision/models/resnet.py>) with skip connections (Resnet with 18 and 50 layers for the MNIST-C and ImageNet-C tasks, respectively; [45]). For ORA's decoder we used

Table 2. Model details.

	for MNIST-C task		for ImageNet-C task
Encoder type	2 Conv	Resnet-18	Resnet-50
Num of feature slots	1152	288	3136
Feature slot dimension	8	8	8
Num of object slots	10	10	16
Object slot dimension	16	16	49
Total number of parameters	2,904,720	4,581,296	67,116,160

<https://doi.org/10.1371/journal.pcbi.1012159.t002>

three fully-connected dense layers to reconstruct the MNIST image and an inversely mirrored Resnet to reconstruct the ImageNet images. As in the original CapsNet model [51], we used a vector implementation of object and feature slots, where the length of each slot vector represents the presence of an object class in the image. To make this value range from 0 to 1 (i.e., 0 for absence and 1 for presence), we used the following nonlinear squash function. We also provide specification details for object slot representations [51] and the total number of trainable parameters in Table 2.

$$\sigma_t^j = \frac{\|v_t^j\|^2}{1 + \|v_t^j\|^2} \frac{v_t^j}{\|v_t^j\|} \quad (1)$$

### Loss function

ORA outputs both object slot activations and image reconstructions, and the losses for each output are computed and combined using a weighting parameter of  $\lambda_{recon} = 0.392$ , Eq (2). For classification, we used margin loss [51] to match the object slot activations to the ground truth, Eq (3). The first term of the loss is active only when the target object is present ( $T_j > 0$ ), and minimizing this term pushes the object slot activation (measured by the vector magnitude  $\|o_j\|$ ) to be closer to 1, with a small margin ( $m$ ). The second term applies only when the target object is absent, and minimizing this term pushes the predicted object slot activation below the small margin ( $m$ ).

$$Total\ Loss = \sum_{j \in class} Class\ Loss_j + \lambda_{recon} Recon\ Loss \quad (2)$$

$$Class\ Loss_j = \begin{cases} T_j \cdot \max(0, (T_j - m) - \|o_j\|)^2, & \text{if } T_j > 0 \\ \max(0, \|o_j\| - m)^2, & \text{if } T_j = 0 \end{cases} \quad (3)$$

### Baseline implementation

The CNN baselines consist of the same encoders used for ORA (See Section ORA implementation) followed by two fully connected layers for classification readout. The readout layers are trained using the dropout technique for regularization [100]. Negative Log Likelihood loss is computed on the log softmax output for classification. The total number of model parameters are: 1,199,882 for the 2-layer CNN baseline, 2,803,882 for the ResNet-18 baseline in the MNIST-C experiment, and 23,540,816 for the ResNet-50 baseline in the ImageNet experiment.

Performance for the other baselines reported in [S3 Table](#) are taken from the original benchmarking paper [43].

### Training details

Both ORA and our CNN baselines were trained and validated on clean images (MNIST, [101] for a digit recognition task; ImageNet, [72] for a natural object recognition task) and tested using a separate dataset containing unseen images subjected to various types of corruptions (MNIST-C, [43] for a digit recognition task; ImageNet-C, [44] for a natural object recognition task). These separate training and testing datasets ensure that models are evaluated on their ability to generalize their learned object representations to previously unseen exemplars and types of corruption, a challenge commonly known as the out-of-distribution problem in computer vision [7]. All model parameters were updated using Adam optimization [102] with a mini-batch size of 128 for the MNIST-C experiment and 32 for the ImageNet-C experiment. For the MNIST-C experiment, we used the exponentially decaying learning rate by a factor of 0.96 at the end of each training epoch, starting from 0.1. Similarly, for the ImageNet-C experiment, the learning rate followed the same exponential decay pattern, but with a starting value of 0.0001. Model training was terminated based on the results from the validation dataset (10% of randomly selected training dataset) and was stopped early if there was no improvement in the validation accuracy after 20 epochs. Models were implemented in the PyTorch framework and trained on a single GPU using a NVIDIA GeForce RTX 3090 (24 GB) for the MNIST-C experiment and a RTX A6000 (48 GB) for the ImageNet-C experiment.

### Measuring ORA's reaction time

ORA's processing is limited to five forward processing steps, each followed by reconstruction-based spatial attention feedback (i.e., global spatial masking iteration). However, the model typically did not need all five of these forward processing steps in order to confidently recognize an object, and the number of steps actually taken by the model to reach a classification threshold is ORA's RT. This threshold is measured by entropy over the softmax distribution of class likelihoods, similar to the method used in [103], and we chose an entropy threshold of 0.6 based on a grid search. If this threshold is not reached within the five forward processing steps, we use the prediction from the final time step for classification.

### Spatial masking methods

In each forward processing step, the model generates a reconstruction-guided spatial mask based on the most likely object reconstruction. To create this mask, we simply apply a threshold to the reconstructed pixel values. We used thresholds of 0.1 for MNIST images and 0.05 for ImageNet images, where the pixel intensities in both range from 0 to 1. This thresholding process creates a boolean mask where pixels with values below the threshold are assigned 0 and pixels with values above the threshold are assigned 1. The masked input for the next forward processing step is obtained by element-wise multiplication of the reconstruction with the boolean mask. The non-zero areas of the masked input are normalized using MaxMin normalization, Eq (4), when  $p$  corresponds to the reconstructed pixel value at image coordinates  $i$  and  $j$ ). For ImageNet images, the RGB reconstructed images are first converted to grayscale before

applying the thresholding and normalization processes.

$$p^{ij} = lb + (ub - lb) \frac{p^{ij} - \min(p^{ij})}{\max(p^{ij}) - \min(p^{ij})} \tag{4}$$

### Feature binding algorithm

ORA’s reconstruction-guided feature binding process computes coefficients to dynamically modulate weights between features and objects. These coefficients, denoted as  $C_{ij}$ , are computed iteratively over three local iterations occurring within each forward processing step (i.e., global spatial masking iteration). We update the original weights by multiplying them by the final coefficients,  $C_{ij} \cdot W_{ij}$ . In the original CapsNet model [51], these binding coefficients are modulated based on representational similarity, which is calculated as the dot product between the features and objects (Line 6 in the Algorithm 1). In our model, we also adjust these binding coefficients to suppress the features that are connected to objects having high reconstruction error. By suppressing the features leading to inaccurate object reconstructions, those features that are most selective to the hypothesized object become bound. All coefficients are normalized over objects using MaxMin normalization like Eq (4).

**Algorithm 1** Reconstruction-guided Feature Binding

- 1: **initialize** for all object  $j$  and feature  $i$ :
- 2: Binding coefficients  $c^{ij} \leftarrow 1$
- 3: Feature similarity score  $b^{ij} \leftarrow 0$
- 4: Reconstruction score coefficients  $r^{ij} \leftarrow 0$
- 5: **while** feature binding iteration  $t = 1, 2, 3$  **do**
- 6: For all object  $j$ :  $\hat{f}^{j|t} \leftarrow W^{tj} f^i$
- 7: For all object  $j$ :  $o_t^j \leftarrow \text{squash}(\sum_i c_t^{ij} \hat{f}^{j|t})$  ▷ Eq (1)
- 8: For all object  $j$  and feature  $i$ :  $b_t^{ij} \leftarrow b_t^{ij} + \hat{f}^{j|t} \cdot o_t^j$
- 9:  $b_t^{ij} \leftarrow \max(\maxmin(b_t^{ij}, \dim = j), 0.5)$  ▷ Eq (4)
- 10: For all object  $j$ :  $\text{reconstruction}_t^j \leftarrow \text{Decoder}(o_t^j)$
- 11: For all object  $j$ :  $r_t^j \leftarrow -\text{MSE}(\text{image}^j, \text{reconstruction}_t^j)$
- 12: For all feature  $i$ :  $r_t^{ij} \leftarrow \text{Repeat}(r_t^j)$
- 13:  $r_t^{ij} \leftarrow \max(\maxmin(r_t^{ij}, \dim = j), 0.5)$  ▷ Eq (4)
- 14:  $c_t^{ij} \leftarrow b_t^{ij} \cdot r_t^{ij}$
- 15: **end while**
- 16: For all object  $j$ :  $o^j \leftarrow \text{squash}(\sum_i c^{ij} \hat{f}^{j|i})$

### 16-Class-ImageNet training dataset

The original Imagenet dataset includes 1,000 subordinate-level labels. These include 120 different dog breeds, which results in a semantic structure that is highly atypical compared to that of an average person [104] that may influence model performance. To address this problem, we also trained ORA using 16-class-ImageNet, a subset of ImageNet that groups its 1000 fine-grained classes into 16 basic-level categories: airplane, bicycle, boat, car, chair, dog, keyboard, oven, bear, bird, bottle, cat, clock, elephant, knife, and truck [7, 73].

### Behavioral methods for the MNIST-C recognition experiment

**Participants.** A total of 143 Stony Brook University undergraduates, self-identified as 94 women, 47 men, and 2 non-binary individuals, participated in the experiment for course credit. This sample size was determined based on an a priori power analysis using G\*Power [105], which indicated that a minimum of 28 participants would be sufficient for a repeated-

measures ANOVA to detect a medium effect size with a power of .80 at a significance level of  $\alpha = .05$ . This estimate is based on each participant viewing a 96 image subset of the 480 total testing images (see Stimuli and Apparatus), a design choice motivated by a desire to avoid participant fatigue and to maintain their engagement in the task. Participants all had normal or corrected to normal vision, by self report.

**Stimuli and apparatus.** Stimuli consisted of 480 images from the MNIST-C testing dataset [43], 60 images for each of the 8 corruption types in MNIST-C-shape. Images were divided into 5 sets of 96 images, where each participant viewed a 96 image set. To experimentally test ORA's RTs we selected the images used for behavioral testing based on the RT predicted by the model. Specifically, for each corruption type we selected 20 images by randomly sampling from 3 different model RT conditions: short (1 step), medium (2–3 step), and long (4–5 step). Note that all images were correctly classified by the model, meaning that this manipulation is targeted to the number of forward processing steps needed by ORA to recognize an object. Selected images were presented on a 19-inch flat-screen CRT ViewSonic SVGA monitor at a screen resolution of  $1024 \times 768$  pixels and a refresh rate of 100 Hz. Participants were seated approximately 70 cm from the monitor and at this viewing distance the image height subtended approximately  $8^\circ$  of a visual angle. Space bar responses were made with the participant's dominant hand using a standard computer keyboard.

**Procedure.** Participants were instructed that they would see a series of centrally presented digits, ranging from 0 to 9, and that they should press the space bar as soon as they recognized the digit class. However, the task emphasized accuracy as well as speed. Participants were instructed to make their key press when they felt confident in their answer, although a reminder was provided if the response was too slow ( $> 5$  sec). Upon pressing the space bar, the original stimuli was masked by a randomly generated pixel image and participants were asked as part of a secondary judgment to report the digit that they recognized by pressing the digit on the numeric keypad. At the start of each trial a central cross appeared immediately prior to stimuli presentation (500 ms), and feedback was provided at the end of each trial on the accuracy of the digit recognition response (10 alternative forced choice). Each participant viewed 96 images, 4 images  $\times$  3 model RT conditions  $\times$  8 corruption types. Ten practice trials were provided before the experiment so that participants could familiarize themselves with the task. The order of presentation of the images was randomized over participants.

### Behavioral methods for the MNIST-C error likelihood rating experiment

**Participants.** A total of 184 undergraduate students from Stony Brook University, self-identified as 120 women, 61 men, and 3 non-binary individuals, participated in the experiment for course credit. Sample size was determined based on an a priori power analysis using G\*Power [105], which indicated that a minimum of 41 participants per stimulus set (4 stimulus sets in total) would be sufficient for a paired t-test to detect a medium effect size with a power of .80 at a  $\alpha = .05$  level of significance. All participants had normal or corrected to normal vision, by self report.

**Stimuli and apparatus.** Stimuli consisted of 120 images selected from the MNIST-C testing dataset [43], 20 images per corruption type, to satisfy the experimental constraint of having competing models make different errors to a stimulus. These ambiguous images were divided into 4 separate sets, with each set containing 30 images. Only six corruption types from the eight MNIST-C-shape categories were used because the number of images meeting the experimental constraint (i.e., different errors) were too few for analysis in the dotted-line and spatter corruption conditions. The stimuli image spanned 40% of their screen size in height. All responses were made with a computer mouse or trackpad.

**Procedure.** Participants accessed the rating experiment remotely using their desktop devices. On each trial, a single corrupted digit was presented at a size of 40% of the vertical screen dimension. Below this stimulus were the two model predictions, horizontally separated (the location was randomized), and participants were instructed to rate which of the two would be more plausibly made by a person (two alternative forced choice; see Fig 7B). A response was required on each trial and participants were not allowed to skip images, although they could change their answers anytime before proceeding to the next trial (which they initiated by clicking a “Next” button). Three practice trials were provided before the experiment so that participants would familiarize themselves with the task and interface. Each participant viewed 30 images. The order of presentation of the images was randomized over participants.

## Supporting information

### S1 Table. Percentage of changes in ORA’s class prediction during spatial masking and feature binding iterations.

(EPS)

### S2 Table. Results from manipulating the spatial frequencies used for image reconstruction.

To test how reconstruction efficacy is effected by spatial frequency, we manipulated the spatial frequencies of the visual input that ORA used to generate its reconstructions. we trained different models to reconstruct different spatial frequency components of the input (e.g., generating a low spatial frequency reconstruction of the original full spectrum image) and then tested them to evaluate the effect of spatial frequency on model performance. We extracted different spatial frequency information from the input using a broadband Gaussian filter with cutoff frequencies below 6 cycles per image and above 30 cycles per image for generating low and high spatial frequency images, respectively. As shown below, we found that the use of a low-resolution object reconstruction yields comparable performance to a model trained with the full spectrum. The figure on the left shows samples of reconstructions from full-spectrum images (top) and images filtered for low (middle) and high (bottom) spatial frequencies. The table on the right shows the average performance of models trained to reconstruct either full-spectrum (FS), low-spatial frequency (LSF), or high-spatial frequency (HSF) images. Values in parentheses indicate standard deviations from 5 model runs.

(EPS)

**S3 Table. Model performance comparison.** Extended model comparison. Average performance and standard deviations from 5 model runs are reported for **ORA**, **CNN**, and **CapsNet** models that we trained for this study (see text for details). Best model performance for the other baselines are taken from [43]. **ABS**: a generative model [106]. **PGD1**: a CNN trained against PGD adversarial noise [107]. **PGD2/GAN**: another CNN trained against PGD/GAN adversarial noise [108].

(EPS)

**S1 Fig. Step-wise visualization of reconstruction-guided feature binding process.** The left column of the figure below shows binding coefficient matrices between all 10 of the possible digit categories (rows) and the first 40 feature slots (columns) at each of the 3 feature binding iterations. Binding coefficients are initially unweighted (step 1, top) but those object features leading to greater reconstruction error become suppressed (become lighter) over iterations, resulting in the feature binding coefficients focusing on the object features that best match the input (middle and bottom matrices). The nine bar plots in the right of the figure illustrate the effect of reconstruction-guided feature binding on ORA’s digit classifications. The leftmost three plots show the raw likelihood of classification for each of the 10 object slots at the start of

each time step (i.e., before the application of reconstruction-guided feature binding). The middle three plots show reconstruction scores (inverse of reconstruction error) at each step, and the rightmost three plots show the updated class likelihood after reconstruction-guided feature binding is applied (i.e., the two are multiplied). Best viewed with zoom in PDF.  
(EPS)

**S2 Fig. MNIST-C dataset examples.**

(EPS)

**S3 Fig. Model comparison on all MNIST-C corruption types.** Average model accuracy and reaction time for all of the MNIST-C corruption types. A: ORA and CNN accuracy using each model's best performing encoder (Resnet-18 and a 2-layer CNN, respectively). B: ORA's reaction time (RT), estimated as the number of forward processing steps required to reach a confidence threshold. Error bars indicate the standard deviations from 5 different runs.

(EPS)

**S4 Fig. Ablation studies on ORA with 2 Convolutional encoder.** Results from ablating ORA's reconstruction-based spatial masking and feature binding components. Ablated components are indicated by a—symbol, none ablated components are indicated by a +. Leftmost bars in pink indicate the complete (no ablation) ORA model. Left: Recognition accuracy. Right: reaction time, approximated by the number of forward processing steps taken to recognize digits with confidence. Error bars are standard deviations from 5 different runs.

(EPS)

**S5 Fig. ORA's behavior on an ambiguous stimulus.** The figure below shows ORA taking longer to converge on a single digit hypothesis when the input is ambiguous. On the top is the input image, which could be plausibly interpreted as either the digit 3 or the digit 5. On the bottom are reconstructions hypothesized for all 10 digits, with the respective classification (C) and reconstruction (R) scores below each. ORA's most likely object hypothesis (with the highest class likelihood) at each step is indicated by the red square. Note that ORA's best guess about this object vacillates between the 3 and 5 hypotheses over time steps. The ground truth for this image is the digit 3.

(EPS)

**S6 Fig. Correlations between human and ORA reaction times for all MNIST-C corruption types.**

(EPS)

**S7 Fig. Examples of ORA reconstructions and spatial masks from Imagenet-C.**

(EPS)

**S8 Fig. Model comparison for all ImageNet-C corruption types.**

(EPS)

**S9 Fig. Applying ORA to an overlapping-object recognition task.** To apply ORA to an overlapping-digit recognition task we generated the MultiMNIST dataset as described in [51], where each image consists of two MNIST digits overlapped, randomly shifted 1–4 pixels horizontally and vertically in opposite directions. The resulting images are  $36 \times 36$  pixels, with an average 70% overlap between digit bounding boxes. A total of 10,000 images were randomly selected from the MNIST testing dataset for model evaluation. The figure shows a representative overlapping-digit stimulus and model outputs illustrating how ORA can solve the overlapping-digit task. The ORA model utilized in this demonstration is identical to the one presented in our paper, meaning it was trained for single-digit recognition and reconstruction,

and we performed no additional hyperparameter tuning. The only change that we made in our application of ORA to overlapping-digit recognition was algorithmic; once the model recognizes a first digit with sufficient certainty (using an entropy threshold of 1.0 as a measure of confidence), it proceeds to identify the next digit. It does this by expanding its spatial mask and focusing on the input signals that remain unexplained by first digit recognition. This approach therefore allows the model to sequentially attend to individual overlapping objects, effectively segregating and recognizing each one after the other, which parallels the seriality observed in human recognition [109, 110]. In the illustrated example, the ground truth consists of the digits 8 and 4 overlapped. ORA initially identifies the digit 8 before shifting its attention to 4, repeating this process until the confidence threshold is achieved. Using its sequential object attention mechanism, ORA outperforms the best-performing CNN baseline, a ResNet-18 (a shallower 2-layer CNN performed poorly on this task). Prediction accuracy was measured using two metrics: image-level accuracy, where the response is considered correct only if both objects in an image are identified accurately, and object-level accuracy, which reflects the average number of objects recognized (0.5 indicates one of two objects are recognized, while 1.0 indicates both objects are recognized). ORA's predictions are generated from the first two digits recognized with high confidence within a maximum of 10 forward iterations. The CNN baseline generated predictions using the top 2 most activated class outputs. Although there is room for improvement, this initial application of ORA to overlapping digits is already beating a strong baseline recognition model.

(EPS)

## Author Contributions

**Conceptualization:** Seoyoung Ahn, Hossein Adeli, Gregory J. Zelinsky.

**Data curation:** Seoyoung Ahn.

**Formal analysis:** Seoyoung Ahn.

**Funding acquisition:** Seoyoung Ahn, Gregory J. Zelinsky.

**Investigation:** Seoyoung Ahn.

**Methodology:** Seoyoung Ahn, Hossein Adeli.

**Project administration:** Gregory J. Zelinsky.

**Resources:** Gregory J. Zelinsky.

**Software:** Gregory J. Zelinsky.

**Supervision:** Gregory J. Zelinsky.

**Validation:** Seoyoung Ahn, Hossein Adeli, Gregory J. Zelinsky.

**Visualization:** Seoyoung Ahn, Hossein Adeli.

**Writing – original draft:** Seoyoung Ahn, Hossein Adeli, Gregory J. Zelinsky.

**Writing – review & editing:** Seoyoung Ahn, Hossein Adeli, Gregory J. Zelinsky.

## References

1. Vogels R, Biederman I. Effects of Illumination Intensity and Direction on Object Coding in Macaque Inferior Temporal Cortex. *Cerebral Cortex*. 2002; 12(7):756–766. <https://doi.org/10.1093/cercor/12.7.756> PMID: 12050087

2. Avidan G, Harel M, Hendler T, Ben-Bashat D, Zohary E, Malach R. Contrast Sensitivity in Human Visual Areas and Its Relationship to Object Recognition. *Journal of Neurophysiology*. 2002; 87(6):3102–3116. <https://doi.org/10.1152/jn.2002.87.6.3102> PMID: 12037211
3. 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 Computational Biology*. 2014; 10(12):e1003963. <https://doi.org/10.1371/journal.pcbi.1003963> PMID: 25521294
4. 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. <https://doi.org/10.1038/srep27755> PMID: 27282108
5. Kriegeskorte N. Deep Neural Networks: A New Framework for Modeling Biological Vision and Brain Information Processing. *Annual Review of Vision Science*. 2015; 1:417–446. <https://doi.org/10.1146/annurev-vision-082114-035447> PMID: 28532370
6. Schrimpf M, Kubilius J, Lee MJ, Ratan Murty NA, Ajemian R, DiCarlo JJ. Integrative Benchmarking to Advance Neurally Mechanistic Models of Human Intelligence. *Neuron*. 2020; 108(3):413–423. <https://doi.org/10.1016/j.neuron.2020.07.040> PMID: 32918861
7. Geirhos R, Narayanappa K, Mitzkus B, Thieringer T, Bethge M, Wichmann FA, et al. Partial Success in Closing the Gap between Human and Machine Vision. *Advances in Neural Information Processing Systems*. 2021; 34:23885–23899.
8. Szegedy C, Zaremba W, Sutskever I, Bruna J, Erhan D, Goodfellow I, et al. Intriguing Properties of Neural Networks. arXiv:13126199. 2013;.
9. Dodge S, Karam L. A Study and Comparison of Human and Deep Learning Recognition Performance under Visual Distortions. In: *International Conference on Computer Communication and Networks (ICCCN)*; 2017. p. 1–7.
10. Baker N, Lu H, Erlikhman G, Kellman PJ. Deep Convolutional Networks Do Not Classify Based on Global Object Shape. *PLOS Computational Biology*. 2018; 14(12):e1006613. <https://doi.org/10.1371/journal.pcbi.1006613> PMID: 30532273
11. Geirhos R, Rubisch P, Michaelis C, Bethge M, Wichmann FA, Brendel W. ImageNet-trained CNNs Are Biased towards Texture; Increasing Shape Bias Improves Accuracy and Robustness. In: *International Conference on Learning Representations*; 2018.
12. Biederman I. Recognition-by-Components: A Theory of Human Image Understanding. *Psychological Review*. 1987; 94(2):115. <https://doi.org/10.1037/0033-295X.94.2.115> PMID: 3575582
13. Tarr MJ, Pinker S. When Does Human Object Recognition Use a Viewer-Centered Reference Frame? *Psychological Science*. 1990; 1(4):253–256. <https://doi.org/10.1111/j.1467-9280.1990.tb00209.x>
14. Plaut DC, Farah MJ. Visual Object Representation: Interpreting Neurophysiological Data within a Computational Framework. *Journal of Cognitive Neuroscience*. 1990; 2(4):320–343. <https://doi.org/10.1162/jocn.1990.2.4.320> PMID: 23964758
15. Rolls ET. Brain Mechanisms for Invariant Visual Recognition and Learning. *Behavioural Processes*. 1994; 33(1-2):113–138. [https://doi.org/10.1016/0376-6357\(94\)90062-0](https://doi.org/10.1016/0376-6357(94)90062-0) PMID: 24925242
16. Marr D. *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. W. H. Freeman and Company; 1982.
17. Ullman S. Aligning Pictorial Descriptions: An Approach to Object Recognition. *Cognition*. 1989; 32(3):193–254. [https://doi.org/10.1016/0010-0277\(89\)90036-X](https://doi.org/10.1016/0010-0277(89)90036-X) PMID: 2752709
18. Carpenter GA, Grossberg S. A Massively Parallel Architecture for a Self-Organizing Neural Pattern Recognition Machine. *Computer vision, graphics, and image processing*. 1987; 37(1):54–115. [https://doi.org/10.1016/S0734-189X\(87\)80014-2](https://doi.org/10.1016/S0734-189X(87)80014-2)
19. Ullman S. Sequence Seeking and Counter Streams: A Computational Model for Bidirectional Information Flow in the Visual Cortex. *Cerebral cortex*. 1995; 5(1):1–11. <https://doi.org/10.1093/cercor/5.1.1> PMID: 7719126
20. Lee TS, Mumford D. Hierarchical Bayesian Inference in the Visual Cortex. *Journal of the Optical Society of America A*. 2003; 20(7):1434. <https://doi.org/10.1364/JOSAA.20.001434> PMID: 12868647
21. Bar M, Kassam KS, Ghuman AS, Boshyan J, Schmid AM, Dale AM, et al. Top-down Facilitation of Visual Recognition. *Proceedings of the National Academy of Sciences*. 2006; 103(2):449–454. <https://doi.org/10.1073/pnas.0507062103> PMID: 16407167
22. Gilbert CD, Li W. Top-down Influences on Visual Processing. *Nature Reviews Neuroscience*. 2013; 14(5):350–363. <https://doi.org/10.1038/nrn3476> PMID: 23595013
23. Posner MI. Orienting of Attention. *Quarterly Journal of Experimental Psychology*. 1980; 32(1):3–25. <https://doi.org/10.1080/00335558008248231> PMID: 7367577

24. Müller MM, Hübner R. Can the Spotlight of Attention Be Shaped like a Doughnut? Evidence from Steady-State Visual Evoked Potentials. *Psychological Science*. 2002; 13(2):119–124. <https://doi.org/10.1111/1467-9280.00422> PMID: 11933994
25. Treisman A. Features and Objects: The Fourteenth Bartlett Memorial Lecture. *The Quarterly Journal of Experimental Psychology Section A*. 1988; 40(2):201–237. <https://doi.org/10.1080/02724988843000104> PMID: 3406448
26. Martinez-Trujillo JC, Treue S. Feature-Based Attention Increases the Selectivity of Population Responses in Primate Visual Cortex. *Current Biology*. 2004; 14(9):744–751. <https://doi.org/10.1016/j.cub.2004.04.028> PMID: 15120065
27. DiCarlo JJ, Zoccolan D, Rust NC. How Does the Brain Solve Visual Object Recognition? *Neuron*. 2012; 73(3):415–434. <https://doi.org/10.1016/j.neuron.2012.01.010> PMID: 22325196
28. Cavanagh P, Caplovitz GP, Lytchenko TK, Maechler MR, Tse PU, Sheinberg DL. The Architecture of Object-Based Attention. *Psychonomic Bulletin & Review*. 2023. <https://doi.org/10.3758/s13423-023-02281-7> PMID: 37081283
29. Dayan P, Hinton GE, Neal RM, Zemel RS. The Helmholtz Machine. *Neural Computation*. 1995; 7(5):889–904. <https://doi.org/10.1162/neco.1995.7.5.889> PMID: 7584891
30. Yuille A, Kersten D. Vision as Bayesian Inference: Analysis by Synthesis? *Trends in Cognitive Sciences*. 2006; 10(7):301–308. <https://doi.org/10.1016/j.tics.2006.05.002> PMID: 16784882
31. de Lange FP, Heilbron M, Kok P. How Do Expectations Shape Perception? *Trends in Cognitive Sciences*. 2018; 22(9):764–779. <https://doi.org/10.1016/j.tics.2018.06.002> PMID: 30122170
32. Clark A. Whatever next? Predictive Brains, Situated Agents, and the Future of Cognitive Science. *Behavioral and Brain Sciences*. 2013; 36(3):181–204. <https://doi.org/10.1017/S0140525X12000477> PMID: 23663408
33. Carrasco M. Visual Attention: The Past 25 Years. *Vision Research*. 2011; 51(13):1484–1525. <https://doi.org/10.1016/j.visres.2011.04.012> PMID: 21549742
34. Kravitz DJ, Behrmann M. Space-, Object-, and Feature-Based Attention Interact to Organize Visual Scenes. *Attention, Perception, & Psychophysics*. 2011; 73(8):2434–2447. <https://doi.org/10.3758/s13414-011-0201-z> PMID: 22006523
35. Adeli H, Ahn S, Zelinsky GJ. A Brain-Inspired Object-Based Attention Network for Multiobject Recognition and Visual Reasoning. *Journal of Vision*. 2023; 23(5):16. <https://doi.org/10.1167/jov.23.5.16> PMID: 37212782
36. Fleming RW, Storrs KR. Learning to See Stuff. *Current Opinion in Behavioral Sciences*. 2019; 30:100–108. <https://doi.org/10.1016/j.cobeha.2019.07.004> PMID: 31886321
37. Xing J, Chrastil ER, Nitz DA, Krichmar JL. Linking Global Top-down Views to First-Person Views in the Brain. *Proceedings of the National Academy of Sciences*. 2022; 119(45):e2202024119. <https://doi.org/10.1073/pnas.2202024119> PMID: 36322732
38. Al-Tahan H, Mohsenzadeh Y. Reconstructing Feedback Representations in the Ventral Visual Pathway with a Generative Adversarial Autoencoder. *PLOS Computational Biology*. 2021; 17(3): e1008775. <https://doi.org/10.1371/journal.pcbi.1008775> PMID: 33760819
39. Hedayati S, O'Donnell RE, Wyble B. A Model of Working Memory for Latent Representations. *Nature Human Behaviour*. 2022; 6(5):709–719. <https://doi.org/10.1038/s41562-021-01264-9> PMID: 35115675
40. Yildirim I, Belledonne M, Freiwald W, Tenenbaum J. Efficient Inverse Graphics in Biological Face Processing. *Science Advances*. 2020; 6(10):eaax5979. <https://doi.org/10.1126/sciadv.aax5979> PMID: 32181338
41. Boutin V, Zerroug A, Jung M, Serre T. Iterative VAE as a Predictive Brain Model for Out-of-Distribution Generalization. In: *Advances in Neural Information Processing System Workshops*; 2020.
42. Csikor F, Meszéna B, Szabó B, Orbán G. Top-down Inference in an Early Visual Cortex Inspired Hierarchical Variational Autoencoder. *arXiv preprint arXiv: 220600436*. 2022;.
43. Mu N, Gilmer J. MNIST-C: A Robustness Benchmark for Computer Vision. *arXiv:190602337*. 2019;.
44. Hendrycks D, Dietterich TG. Benchmarking Neural Network Robustness to Common Corruptions and Perturbations. *International Conference On Learning Representations*. 2019;.
45. He K, Zhang X, Ren S, Sun J. Deep Residual Learning for Image Recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*; 2016. p. 770–778.
46. Friston K. The Free-Energy Principle: A Unified Brain Theory? *Nature Reviews Neuroscience*. 2010; 11(2):127–138. <https://doi.org/10.1038/nrn2787> PMID: 20068583
47. Grossberg S. Towards Solving the Hard Problem of Consciousness: The Varieties of Brain Resonances and the Conscious Experiences That They Support. *Neural Networks*. 2017; 87:38–95. <https://doi.org/10.1016/j.neunet.2016.11.003> PMID: 28088645

48. Fazl A, Grossberg S, Mingolla E. View-Invariant Object Category Learning, Recognition, and Search: How Spatial and Object Attention Are Coordinated Using Surface-Based Attentional Shrouds. *Cognitive Psychology*. 2009; 58(1):1–48. <https://doi.org/10.1016/j.cogpsych.2008.05.001> PMID: 18653176
49. Hummel JE, Biederman I. Dynamic Binding in a Neural Network for Shape Recognition. *Psychological Review*. 1992; 99(3):480. <https://doi.org/10.1037/0033-295X.99.3.480> PMID: 1502274
50. Rajalingham R, DiCarlo JJ. Reversible Inactivation of Different Millimeter-Scale Regions of Primate IT Results in Different Patterns of Core Object Recognition Deficits. *Neuron*. 2019; 102(2):493–505.e5. <https://doi.org/10.1016/j.neuron.2019.02.001> PMID: 30878289
51. Sabour S, Frosst N, Hinton GE. Dynamic Routing between Capsules. In: *Advances in Neural Information Processing Systems*; 2017. p. 3856–3866.
52. Doerig A, Schmittwilken L, Sayim B, Manassi M, Herzog MH. Capsule Networks as Recurrent Models of Grouping and Segmentation. *PLOS Computational Biology*. 2020; 16(7):e1008017. <https://doi.org/10.1371/journal.pcbi.1008017> PMID: 32692780
53. Greff K, van Steenkiste S, Schmidhuber J. On the Binding Problem in Artificial Neural Networks. *arXiv:201205208*. 2020;.
54. Locatello F, Weissenborn D, Unterthiner T, Mahendran A, Heigold G, Uszkoreit J, et al. Object-Centric Learning with Slot Attention. *Advances in Neural Information Processing Systems*. 2020; 33:11525–11538.
55. Kahneman D, Treisman A, Gibbs BJ. The Reviewing of Object Files: Object-specific Integration of Information. *Cognitive Psychology*. 1992; 24(2):175–219. [https://doi.org/10.1016/0010-0285\(92\)90007-O](https://doi.org/10.1016/0010-0285(92)90007-O) PMID: 1582172
56. Yu CP, Maxfield JT, Zelinsky GJ. Searching for Category-Consistent Features: A Computational Approach to Understanding Visual Category Representation. *Psychological Science*. 2016; 27(6):870–884. <https://doi.org/10.1177/0956797616640237> PMID: 27142461
57. Peters B, Kriegeskorte N. Capturing the Objects of Vision with Neural Networks. *Nature Human Behaviour*. 2021; 5(9):1127–1144. <https://doi.org/10.1038/s41562-021-01194-6> PMID: 34545237
58. Hinton GE, Salakhutdinov RR. Reducing the Dimensionality of Data with Neural Networks. *Science*. 2006; 313(5786):504–507. <https://doi.org/10.1126/science.1127647> PMID: 16873662
59. Dumoulin V, Visin F. *A Guide to Convolution Arithmetic for Deep Learning*. 2018;.
60. Blaser E, Pylyshyn ZW, Holcombe AO. Tracking an Object through Feature Space. *Nature*. 2000; 408(6809):196–199. <https://doi.org/10.1038/35041567> PMID: 11089972
61. Moore CM, Fulton C. The Spread of Attention to Hidden Portions of Occluded Surfaces. *Psychonomic Bulletin & Review*. 2005; 12(2):301–306. <https://doi.org/10.3758/BF03196376> PMID: 16082810
62. Chen Y, Zelinsky GJ. Is There a Shape to the Attention Spotlight? Computing Saliency over Proto-Objects Predicts Fixations during Scene Viewing. *Journal of Experimental Psychology: Human Perception and Performance*. 2019; 45(1):139–154. PMID: 30596438
63. Tyler CW, Kontsevich LL. Mechanisms of Stereoscopic Processing: Stereoattention and Surface Perception in Depth Reconstruction. *Perception*. 1995; 24(2):127–153. <https://doi.org/10.1068/p240127> PMID: 7617422
64. Behrmann M, Geng JJ, Shomstein S. Parietal Cortex and Attention. *Current Opinion in Neurobiology*. 2004; 14(2):212–217. <https://doi.org/10.1016/j.conb.2004.03.012> PMID: 15082327
65. Xu Y. The Posterior Parietal Cortex in Adaptive Visual Processing. *Trends in Neurosciences*. 2018; 41(11):806–822. <https://doi.org/10.1016/j.tins.2018.07.012> PMID: 30115412
66. He K, Chen X, Xie S, Li Y, Dollár P, Girshick R. Masked Autoencoders Are Scalable Vision Learners. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*; 2022. p. 16000–16009.
67. Deco G, Rolls ET. A Neurodynamical Cortical Model of Visual Attention and Invariant Object Recognition. *Vision Research*. 2004; 44(6):621–642. <https://doi.org/10.1016/j.visres.2003.09.037> PMID: 14693189
68. Jeurissen D, Self MW, Roelfsema PR. Serial Grouping of 2D-image Regions with Object-Based Attention in Humans. *Elife*. 2016; 5:e14320. <https://doi.org/10.7554/eLife.14320> PMID: 27291188
69. Desimone R, Duncan J. Neural Mechanisms of Selective Visual Attention. *Annual Review of Neuroscience*. 1995; 18(1):193–222. <https://doi.org/10.1146/annurev.ne.18.030195.001205> PMID: 7605061
70. Beck DM, Kastner S. Top-down and Bottom-up Mechanisms in Biasing Competition in the Human Brain. *Vision Research*. 2009; 49(10):1154–1165. <https://doi.org/10.1016/j.visres.2008.07.012> PMID: 18694779
71. Bundesen C, Habekost T, Kyllingsbæk S. A Neural Theory of Visual Attention and Short-Term Memory (NTVA). *Neuropsychologia*. 2011; 49(6):1446–1457. <https://doi.org/10.1016/j.neuropsychologia.2010.12.006> PMID: 21146554

72. Deng J, Dong W, Socher R, Li LJ, Li K, Fei-Fei L. Imagenet: A Large-Scale Hierarchical Image Database. In: IEEE Conference on Computer Vision and Pattern Recognition. IEEE; 2009. p. 248–255.
73. Geirhos R, Temme CR, Rauber J, Schütt HH, Bethge M, Wichmann FA. Generalisation in Humans and Deep Neural Networks. In: Advances in Neural Information Processing System. vol. 31; 2018.
74. Cheng B, Schwing A, Kirillov A. Per-Pixel Classification Is Not All You Need for Semantic Segmentation. Advances in Neural Information Processing Systems. 2021; 34:17864–17875.
75. DiCarlo JJ, Haefner R, Isik L, Konkle T, Kriegeskorte N, Peters B, et al. How Does the Brain Combine Generative Models and Direct Discriminative Computations in High-Level Vision?; 2021.
76. Gershman SJ. The Generative Adversarial Brain. Frontiers in Artificial Intelligence. 2019; 2:18. <https://doi.org/10.3389/frai.2019.00018> PMID: 33733107
77. Breedlove JL, St-Yves G, Olman CA, Naselaris T. Generative Feedback Explains Distinct Brain Activity Codes for Seen and Mental Images. Current Biology. 2020; 30(12):2211–2224.e6. <https://doi.org/10.1016/j.cub.2020.04.014> PMID: 32359428
78. Mumford D. Pattern Theory: A Unifying Perspective. In: First European Congress of Mathematics: Paris, July 6-10, 1992 Volume I Invited Lectures (Part 1); 1994. p. 187–224.
79. Minda JP, Smith JD. Prototypes in Category Learning: The Effects of Category Size, Category Structure, and Stimulus Complexity. Journal of Experimental Psychology: Learning, Memory, and Cognition. 2001; 27(3):775–799. PMID: 11394680
80. Bar M. A Cortical Mechanism for Triggering Top-down Facilitation in Visual Object Recognition. Journal of Cognitive Neuroscience. 2003; 15(4):600–609. <https://doi.org/10.1162/089892903321662976> PMID: 12803970
81. Bi Z. Top-down Generation of Low-Precision Representations Improves the Perception and Imagination of Fine-Scale Visual Information. bioRxiv. 2021; p. 2021.05.07.443208.
82. Maunsell JH, Treue S. Feature-Based Attention in Visual Cortex. Trends in Neurosciences. 2006; 29(6):317–322. <https://doi.org/10.1016/j.tins.2006.04.001> PMID: 16697058
83. Ho J, Jain A, Abbeel P. Denoising Diffusion Probabilistic Models. Advances in Neural Information Processing Systems. 2020; 33:6840–6851.
84. Singh G, Deng F, Ahn S. Illiterate DALL-E Learns to Compose. In: International Conference on Learning Representations; 2022.
85. Dijkstra N, Bosch SE, van Gerven MA. Shared Neural Mechanisms of Visual Perception and Imagery. Trends in Cognitive Sciences. 2019. <https://doi.org/10.1016/j.tics.2019.02.004> PMID: 30876729
86. Egly R, Driver J, Rafal RD. Shifting Visual Attention between Objects and Locations: Evidence from Normal and Parietal Lesion Subjects. Journal of Experimental Psychology: General. 1994; 123(2):161. <https://doi.org/10.1037/0096-3445.123.2.161> PMID: 8014611
87. Vecera SP. Toward a Biased Competition Account of Object-Based Segregation and Attention. Brain and Mind. 2000; 1(3):353–384. <https://doi.org/10.1023/A:1011565623996>
88. Scholl BJ. Objects and Attention: The State of the Art. Cognition. 2001; 80(1-2):1–46. [https://doi.org/10.1016/S0010-0277\(00\)00152-9](https://doi.org/10.1016/S0010-0277(00)00152-9) PMID: 11245838
89. Logan GD. The CODE Theory of Visual Attention: An Integration of Space-Based and Object-Based Attention. Psychological Review. 1996; 103(4):603. <https://doi.org/10.1037/0033-295X.103.4.603> PMID: 8888649
90. Roelfsema PR, Houtkamp R. Incremental Grouping of Image Elements in Vision. Attention, Perception, & Psychophysics. 2011; 73(8):2542–2572. <https://doi.org/10.3758/s13414-011-0200-0> PMID: 21901573
91. Wagemans J, Elder JH, Kubovy M, Palmer SE, Peterson MA, Singh M, et al. A Century of Gestalt Psychology in Visual Perception: I. Perceptual Grouping and Figure–Ground Organization. Psychological Bulletin. 2012; 138(6):1172. <https://doi.org/10.1037/a0029333> PMID: 22845751
92. Scolaro M, Ester EF, Serences JT. Feature- and Object-Based Attentional Modulation in the Human Visual System. In: The Oxford Handbook of Attention. Oxford University Press; 2014.
93. Shomstein S, Behrmann M. Cortical Systems Mediating Visual Attention to Both Objects and Spatial Locations. Proceedings of the National Academy of Sciences. 2006; 103(30):11387–11392. <https://doi.org/10.1073/pnas.0601813103> PMID: 16840559
94. Cohen EH, Tong F. Neural Mechanisms of Object-Based Attention. Cerebral Cortex. 2015; 25(4):1080–1092. <https://doi.org/10.1093/cercor/bht303> PMID: 24217991
95. O’Craven KM, Downing PE, Kanwisher N. fMRI Evidence for Objects as the Units of Attentional Selection. Nature. 1999; 401(6753):584–587. <https://doi.org/10.1038/44134> PMID: 10524624
96. Pylyshyn Z. The Role of Location Indexes in Spatial Perception: A Sketch of the FINST Spatial-Index Model. Cognition. 1989; 32(1):65–97. [https://doi.org/10.1016/0010-0277\(89\)90014-0](https://doi.org/10.1016/0010-0277(89)90014-0) PMID: 2752706

97. Dijkstra N, Fleming SM. Subjective Signal Strength Distinguishes Reality from Imagination. *Nature Communications*. 2023; 14(1):1627. <https://doi.org/10.1038/s41467-023-37322-1> PMID: 36959279
98. Dittadi A, Papa SS, Vita MD, Schölkopf B, Winther O, Locatello F. Generalization and Robustness Implications in Object-Centric Learning. In: *International Conference on Machine Learning*. vol. 162 of *Proceedings of Machine Learning Research*. PMLR; 2022. p. 5221–5285.
99. Shi B, Darrell T, Wang X. Top-down Visual Attention from Analysis by Synthesis. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*; 2023. p. 2102–2112.
100. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *The Journal of Machine Learning Research*. 2014; 15(1):1929–1958.
101. Deng L. The MNIST Database of Handwritten Digit Images for Machine Learning Research. *IEEE Signal Processing Magazine*. 2012; 29(6):141–142. <https://doi.org/10.1109/MSP.2012.2211477>
102. Kingma DP, Ba J. Adam: A Method for Stochastic Optimization. arXiv:14126980. 2014;.
103. Spoerer CJ, Kietzmann TC, Mehrer J, Charest I, Kriegeskorte N. Recurrent Neural Networks Can Explain Flexible Trading of Speed and Accuracy in Biological Vision. *PLOS Computational Biology*. 2020; 16(10):e1008215. <https://doi.org/10.1371/journal.pcbi.1008215> PMID: 33006992
104. Ahn S, Zelinsky GJ, Lupyan G. Use of Superordinate Labels Yields More Robust and Human-like Visual Representations in Convolutional Neural Networks. *Journal of Vision*. 2021; 21(13):13. <https://doi.org/10.1167/jov.21.13.13> PMID: 34967860
105. Faul F, Erdfelder E, Lang AG, Buchner A. G\*Power 3: A Flexible Statistical Power Analysis Program for the Social, Behavioral, and Biomedical Sciences. *Behavior Research Methods*. 2007; 39(2):175–191. <https://doi.org/10.3758/BF03193146> PMID: 17695343
106. Schott L, Rauber J, Bethge M, Brendel W. Towards the First Adversarially Robust Neural Network Model on MNIST. In: *International Conference on Learning Representations*; 2018.
107. Madry A, Makelov A, Schmidt L, Tsipras D, Vladu A. Towards Deep Learning Models Resistant to Adversarial Attacks. In: *International Conference on Learning Representations*; 2018.
108. Wang X, Yu K, Wu S, Gu J, Liu Y, Dong C, et al. ESRGAN: Enhanced Super-Resolution Generative Adversarial Networks. In: *Proceedings of the European Conference on Computer Vision Workshops*; 2018.
109. Broadbent DE. A Mechanical Model for Human Attention and Immediate Memory. *Psychological Review*. 1957; 64(3):205. <https://doi.org/10.1037/h0047313> PMID: 13441856
110. Pashler H. *The Psychology of Attention*. The MIT Press; 1998.