

## COGNITIVE NEUROSCIENCE

# Incorporating intrinsic suppression in deep neural networks captures dynamics of adaptation in neurophysiology and perception

K. Vinken<sup>1,2,3\*</sup>, X. Boix<sup>1,2,4</sup>, G. Kreiman<sup>1,2</sup>

Adaptation is a fundamental property of sensory systems that can change subjective experiences in the context of recent information. Adaptation has been postulated to arise from recurrent circuit mechanisms or as a consequence of neuronally intrinsic suppression. However, it is unclear whether intrinsic suppression by itself can account for effects beyond reduced responses. Here, we test the hypothesis that complex adaptation phenomena can emerge from intrinsic suppression cascading through a feedforward model of visual processing. A deep convolutional neural network with intrinsic suppression captured neural signatures of adaptation including novelty detection, enhancement, and tuning curve shifts, while producing aftereffects consistent with human perception. When adaptation was trained in a task where repeated input affects recognition performance, an intrinsic mechanism generalized better than a recurrent neural network. Our results demonstrate that feedforward propagation of intrinsic suppression changes the functional state of the network, reproducing key neurophysiological and perceptual properties of adaptation.

## INTRODUCTION

The way we process and perceive the environment around us is not static but is continuously modulated by the incoming sensory information itself. This property of sensory systems is called adaptation and can markedly alter our perceptual experience, such as the illusory perception of upward motion after watching a waterfall for some time (1). In the brain, neural responses adapt to the recent stimulus history in a remarkably similar way across sensory modalities and across species, suggesting that neural adaptation is governed by fundamental and conserved underlying mechanisms (2). The effects of adaptation on both the neural and perceptual levels have been most extensively studied in the visual system, where they appear to play a central role in the integration of temporal context (3–5). Therefore, to understand vision under natural, dynamic conditions, we must consider the neural processes that contribute to visual adaptation and how these processes generate emergent functional states in neural networks. Yet, we do not have a comprehensive understanding of what the underlying mechanisms of adaptation are and how they give rise to changes in perception.

A fundamental question is whether the dynamics of adaptation are implemented by recurrent interactions in the neural network (6) or whether they can arise from established intrinsic biophysical mechanisms operating within each individual neuron (2). An important argument for the role of intrinsic cellular mechanisms in adaptation is that contrast adaptation in cat visual cortex leads to a strong afterhyperpolarization of the membrane potential (7). In other words, the more a neuron fires, the more its excitability is reduced, which is why the phenomenon is sometimes called neuronal fatigue (8). In this scenario, adaptation is caused by intrinsic properties of individual neurons that reduce their responsiveness proportional to

their previous activation. Throughout the paper, we use the term intrinsic suppression to refer to such neuronally intrinsic mechanisms, which suppress responses on the basis of recent activation.

However, adaptation phenomena in the brain go well beyond firing rate–based suppression, and it is not always clear whether those phenomena can be accounted for by intrinsic neuronal properties. First, the amount of suppression does not just depend on the preceding firing rate but can be stimulus specific; i.e., suppression depends on whether the current stimulus is a repetition or not (9). Second, adaptation can also lead to response enhancement of single neurons (5, 8, 10, 11), sometimes even at the population level (12). Last, adaptation can cause a shift in the neuron’s tuning function for a particular stimulus dimension such as orientation (13, 14), direction (15), or spatial and temporal frequency (16, 17). Tuning shifts include both response suppression and enhancement (13) and have been linked to perceptual aftereffects where adaptation produces a shift in the perception of a stimulus property (15). Complex adaptation phenomena such as tuning shifts have fueled the argument that recurrent network mechanisms should be involved (13, 15, 16, 18). The putative involvement of recurrent signals is supported by computational models, which implemented adaptation by changing recurrent interactions between orientation tuned channels to successfully produce peak shifts (18–20).

Adaptation effects cascade through the visual system and can alter the network interactions in unexpected ways (2, 21). For example, adaptation-induced shifts in spatial tuning of primary visual cortex (V1) neurons can be explained by a two-layer model where changes in response gain in lateral geniculate nucleus cascade to V1 through a fixed weighting (22). These findings highlight the need for deeper, multilayer models to capture the effects of adaptation, because previous models that lack the characteristic hierarchical depth and complexity of the visual cortex may not be sufficient to demonstrate the feedforward potential of intrinsic neuronal mechanisms. Moreover, the units in previous models are only designed to encode a particular stimulus dimension, such as orientation, and thus cannot provide a comprehensive framework of visual adaptation. In

Copyright © 2020  
The Authors, some  
rights reserved;  
exclusive licensee  
American Association  
for the Advancement  
of Science. No claim to  
original U.S. Government  
Works. Distributed  
under a Creative  
Commons Attribution  
NonCommercial  
License 4.0 (CC BY-NC).

<sup>1</sup>Boston Children’s Hospital, Harvard Medical School, Boston, MA 02115, USA. <sup>2</sup>Center for Brains, Minds and Machines, Cambridge, MA 02139, USA. <sup>3</sup>Laboratory for Neuro- and Psychophysiology, Department of Neurosciences, KU Leuven, 3000, Leuven, Belgium. <sup>4</sup>Department of Brain and Cognitive Sciences, MIT, Cambridge, MA 02139, USA. \*Corresponding author. Email: kasper.vinken@childrens.harvard.edu

contrast, deep convolutional neural networks have recently come forward as a powerful new tool to model biological vision (23–26) [see, however, discussion in (27)]. When trained to classify natural images, these models describe the stages of ventral visual stream processing of brief stimulus presentations with unprecedented accuracy (28–33), while capturing essential aspects of object recognition behavior and perceived shape similarity (29, 31, 34). In this study, we exploit another advantage of deep neural networks, i.e., their ability to demonstrate how complex properties can emerge from the introduction of biophysically inspired neural mechanisms. We implemented activation-based intrinsic suppression in a feedforward convolutional neural network (35) and tested the hypothesis that complex adaptation phenomena readily emerge without dedicated recurrent mechanisms.

A comprehensive model of visual adaptation should not only capture the neurophysiological dynamics of adaptation but also produce the perceptual consequences. Therefore, we evaluated the proposed computational model implementing intrinsic suppression with critical neurophysiological and psychophysical experiments. We first show that the model captures the fundamental neurophysiological hallmarks of repetition suppression, including stimulus-specific suppression, not only from one image to the next but also across several image presentations (5). Second, we show that the model readily produces the two fundamental perceptual after-effects of adaptation, namely, a perceptual bias in the estimate of a stimulus parameter and an enhanced discriminability between parameter levels (3). In contrast with previous models that were constrained to one low-level property such as orientation, we demonstrate these effects using face gender (36) as a stimulus parameter, to highlight the general applicability of the model. Third, we show that perceptual aftereffects coincided with response enhancements as well as tuning peak shifts, phenomena that are often considered to need the involvement of recurrent network mechanisms (13, 15, 16, 18). Response magnitude changes contributed mostly to the perceptual bias, but tuning changes were required to explain enhanced discriminability. Last, we show that a trained intrinsic neural mechanism is less likely to over-fit and thus provided a less complex solution than a recurrent network mechanism. Overall, while not ruling out any role of recurrent processes in the brain, these results demonstrate that the hallmark neural and perceptual effects of adaptation can be accounted for by activation-based suppression cascading through a complex feedforward sensory system.

## RESULTS

We investigate whether complex adaptation phenomena readily emerge from the propagation of activation-based intrinsic suppression, in a feedforward neural network model of ventral stream processing. We used a pretrained convolutional neural network (Fig. 1A) (35) as a bottom-up computational model of vision and introduced an exponentially decaying intrinsic adaptation state into each unit of each layer of the network, with its parameters set to impose suppression (Fig. 1B; Materials and Methods). The two neural adaptation parameters  $\alpha$  and  $\beta$  (Eqs. 1 and 2) were not trained to fit the neuronal responses or behavioral results but were the same for each unit and were chosen to lead to a gradual buildup and recovery of the adapted state over several time steps (Fig. 1B). Throughout the paper, we use  $\alpha = 0.96$  and  $\beta = 0.7$ , unless indicated otherwise. Because of the intrinsic suppression mechanism, the model units

display temporally evolving responses (Fig. 1C), and their activations can be directly compared to the neurophysiological dynamics.

### A neural network incorporating intrinsic suppression captures temporal dynamics of adaptation at the neurophysiological level

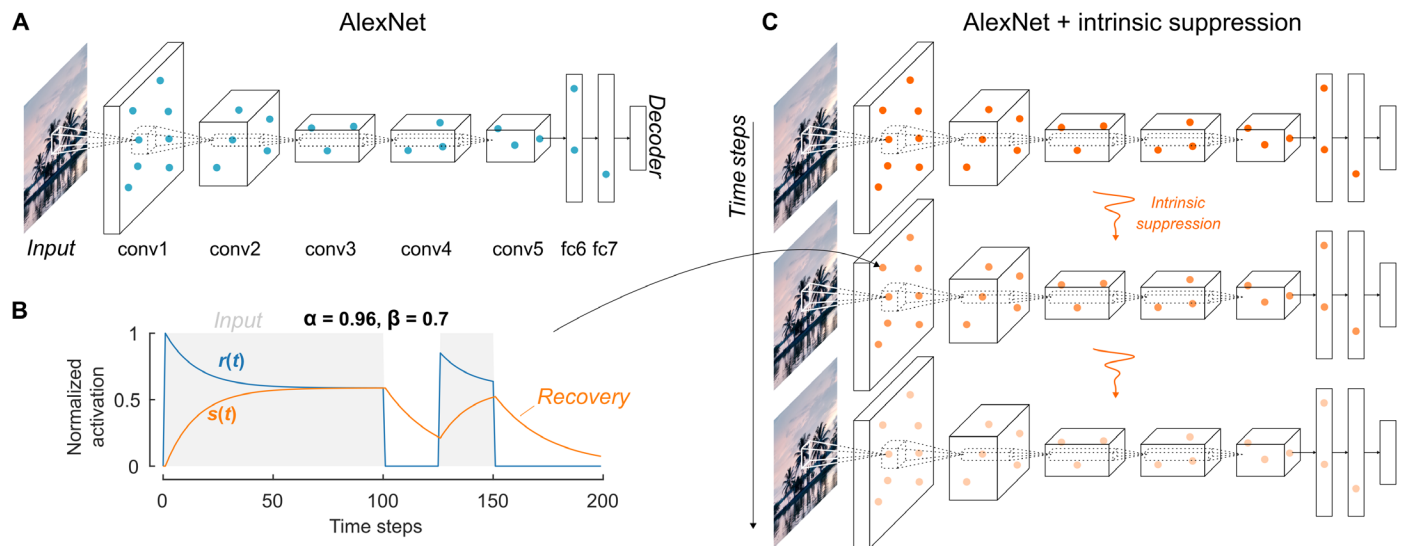
We start with the most prominent characteristic of neural adaptation: repetition suppression, which refers to a reduction in the neuronal responses when a stimulus is repeated. We illustrate this phenomenon using an experiment where face stimuli were presented to a macaque monkey in pairs of two: an adapter followed by a test stimulus (Fig. 2A) (37). In repetition trials, the test stimulus was identical to the adapter whereas, in alternation trials, the adapter and test stimuli were different. Neurons recorded in the middle lateral face patch of inferior temporal (IT) cortex showed a decrease in the response during stimulus presentation and after stimulus offset. In addition, the neurons showed a lower response to a face stimulus when it was a repetition trial (blue) compared to an alternation trial (orange; Fig. 2B).

We evaluated the average time courses of the model unit activations for the same experiment (Fig. 2C; mean of all  $N = 43,264$  units in layer conv5). The model units showed a decrease in the response during the course of stimulus presentation. Consistent with repetition suppression in biological neurons, the response of model units to the test stimulus was lower for repetition than alternation trials. For this stimulus set, the largest difference between repetition and alternation trials was observed for layer conv5 (see other layers in fig. S1A).

The model units demonstrated several key features of adaptation at two time scales: (i) during presentation of any stimulus, including the first stimulus, there was a decrease in the response with time; (ii) the overall response to the second stimulus was smaller than the overall response to the first stimulus; and (iii) the response to the second stimulus was attenuated more when it was a repetition. However, the model did not capture more complex dynamics such as the second peak in neural responses. The model responses showed a smaller difference between repetitions and alternations than biological neurons: The average alternation-repetition difference was 0.07,  $SD = 0.12$  (model, five test time steps), and 0.11,  $SD = 0.15$  (IT neurons, 850 to 1000 ms window) in the normalized scale of Fig. 2 (B and C).

We hypothesized that the computer-generated faces were too similar for the model to display the full range of adaptation effects. Therefore, we ran the same experiment using natural images with more variability. Natural stimuli resulted in a considerably larger difference between repetition and alternation trials (fig. S1B), suggesting that the selectivity of adaptation at least partially reflects stimulus similarity in the model representations. Consistent with this idea, the stimulus similarity in preadaptation activation patterns for different adapter and test images was positively correlated with the amount of suppression for most layers (fig. S2).

An important property of repetition suppression in macaque IT is stimulus specificity: Even for two adapters that equally activate the same neuron, the suppression for an image repetition is still stronger than for an alternation (9). It is not straightforward to see how a neuronally intrinsic mechanism could account for this phenomenon, because an intrinsic firing rate-based mechanism is by itself not stimulus selective (5). However, fig. S3 demonstrates that when activation-based suppression propagates through the layers of the network, neural adaptation of single units becomes increasingly



**Fig. 1. Neural network architecture and incorporation of activation-based intrinsic suppression.** (A) Architecture of a static deep convolutional neural network, in this case AlexNet (35). AlexNet contains five convolutional layers (conv1 to conv5) and three fully connected layers (fc6, fc7, and the decoder fc8). The unit activations in each layer, and therefore the output of the network, are a fixed function of the input image. Photo credit: Kasper Vincken, Boston Children's Hospital, Harvard Medical School. (B) Intrinsic suppression was implemented for each unit using an intrinsic adaptation state  $s(t)$  (orange), which modulates the response  $r(t)$  (blue) and is updated at each time step based on the previous response  $r(t - 1)$  (Eqs. 1 and 2). The parameter values  $\alpha = 0.96$  and  $\beta = 0.7$  were chosen to impose a response suppression ( $\beta > 0$ ) that gradually builds up over time: For constant input (gray shaded areas), the value of the intrinsic state  $s(t)$  gradually increases, leading to a reduction in the response  $r(t)$ . The intrinsic adaptation state recovers in the absence of input (nonshaded areas). (C) Expansion over time of the network in (A), where the activation of each unit is a function of its inputs and its activation at the previous time step (Eqs. 1 and 2).

less dependent on their previous activation, until stimulus-specific suppression is present for most single units in fully connected layers.

In addition to the two temporal scales illustrated in Fig. 2 (A to C), adaptation also operates at longer time scales. For example, repetition suppression typically accumulates across multiple trials and can survive intervening stimuli (9). To illustrate this longer time scale, we present multi-unit data from rat visual cortex (12), recorded during an oddball paradigm where two stimuli, say *A* and *B*, were presented in a random sequence with different probabilities (Fig. 2D): A standard stimulus was shown with high probability ( $P = 0.9$ ; blue), and a deviant stimulus was shown with a low probability ( $P = 0.1$ ; purple). Stimulus (*A* or *B*) and condition (standard or deviant) were counterbalanced for each neural recording. The standard stimulus was far more likely to be repeated in the sequence, allowing adaptation to build up and therefore causing a decrease in the response for later trials in the sequence (Fig. 2, E and F, blue). Adaptation was evident both in V1 and in the extrastriate latero-intermediate visual cortex (LI).

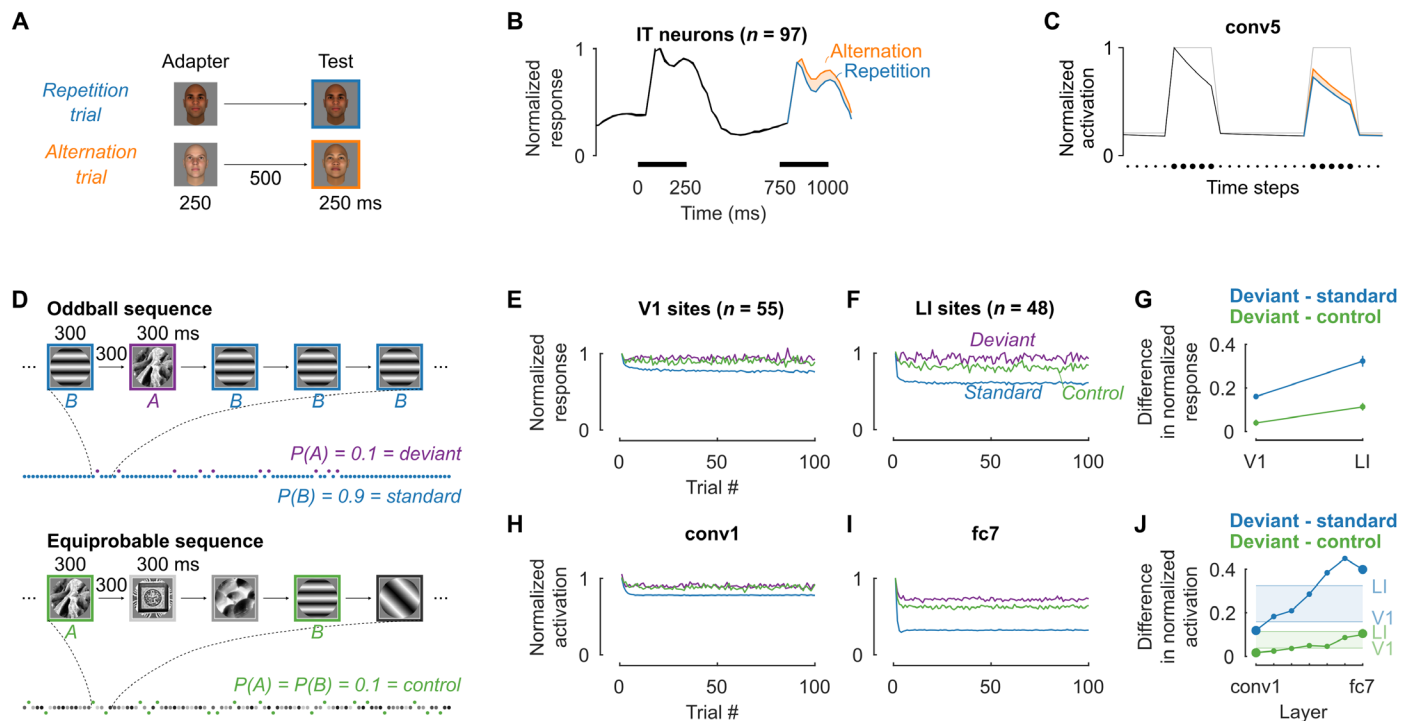
We evaluated the model in the oddball paradigm, without any tuning or parameter changes. The model qualitatively captured the response difference between standard and deviant stimuli (Fig. 2, H and I). Comparing Fig. 2E versus F, the effect of adaptation was stronger in LI compared to V1 (Fig. 2G). An increase in adaptation along the visual hierarchy is consistent with the idea of adaptation cascading through the visual system, with additional contributions at multiple stages. Like the neural data, the model showed increasing adaptation effects from one layer to the next (Fig. 2J), and this increase only occurred when intrinsic suppression was incorporated in multiple layers (fig. S7).

In the original experiment, images *A* and *B* were also presented in separate equiprobable control sequences, where each stimulus

was presented with an equally low probability ( $P = 0.1$ ) together with eight additional stimuli (Fig. 2D) (12). Equiprobable sequences are typically used to distinguish repetition from surprise effects, because the probability of a repetition in the control condition is the same as for the deviant, yet no image is more likely or unlikely than the others. Thus, if neural responses also signal the unexpectedness of the deviant, then the response to a deviant stimulus should be larger than the control condition, which was observed for recording sites in downstream visual area LI (Fig. 2F; purple > green). The model also showed a difference in response between deviant and equiprobable control conditions in higher layers (Fig. 2, I and J). Because the model only incorporates feedforward dynamics of intrinsic suppression, this response difference cannot be attributed to an explicit encoding of expectation. Instead, the lower response for the control condition results from higher cross-stimulus adaptation from the additional stimuli in the equiprobable sequences. This observation means that intrinsic suppression in a feedforward neural network captures not only response differences due to the repetition frequency of a stimulus itself (deviant versus standard) but also differences related to the occurrence probability of other stimuli (deviant surrounded by high-probability standard versus surrounded by several equiprobable stimuli).

### A neural network incorporating intrinsic suppression produces perceptual aftereffects

A comprehensive model of visual adaptation should not only capture the neural properties of repetition suppression but also explain perceptual aftereffects of adaptation. Aftereffects occur when recent exposure to an adapter stimulus biases or otherwise alters the perception of a subsequently presented test stimulus. For example, previous exposure to a male face will make another face appear more



**Fig. 2. Activation-based intrinsic suppression in a neural network captures the attenuation in neurophysiological responses during repetition suppression.** (A) Face stimuli (created with FaceGen: facegen.com) were presented in repetition trials (adapter = test) and alternation trials (adapter  $\neq$  test). (B) Responses in IT cortex ( $n = 97$ , shown normalized to average peak activity) are suppressed more for a repeated stimulus [orange, data from (37)]. Black bars indicate stimulus presentation. (C) The same experiment as in (A) and (B) produced similar repetition suppression in the model with intrinsic suppression (black, blue, and orange lines; gray: no adaptation mechanism; average activity after ReLU of all  $N = 43,264$  conv5 units). The x-axis units are time steps, mapping to bins of 50 ms in (B). (D) Example oddball sequence (top) with a high-probability standard (blue) and a low-probability deviant (purple) and example equiprobable sequence (bottom) as control (green, texture images from vismod.media.mit.edu/pub/VisTex/). (E and F) Average neural responses in rat V1 [ $n = 55$ , (E)] and LI [ $n = 48$ , (F)] (12) for the standard (blue), deviant (purple), and control (green) conditions (normalized by the response at trial one). (G) Deviant – standard (blue) and deviant – control (green) response differences increase from V1 to LI [error bars: 95% bootstrap confidence interval (CI), assuming no inter-animal difference]. (H to J) Running the experiment in the model captures response dynamics similar to rat visual cortex. (H) and (I) show the results for conv1 and fc7 [indicated by larger markers in (J)], respectively. Green and blue horizontal lines and shading in (J) indicate the neural data averages of (G).

female to an observer, and vice versa (Fig. 3A). In other words, adaptation biases the decision boundary for perceptual face-gender discrimination toward the adapter. A defining property of this type of aftereffect is that no perceptual bias should occur when the adapter corresponds to the original boundary stimulus (e.g., a gender-neutral face). Here, we focus on the face-gender dimension, but similar results for the tilt aftereffect (38) with gratings are shown in fig. S4.

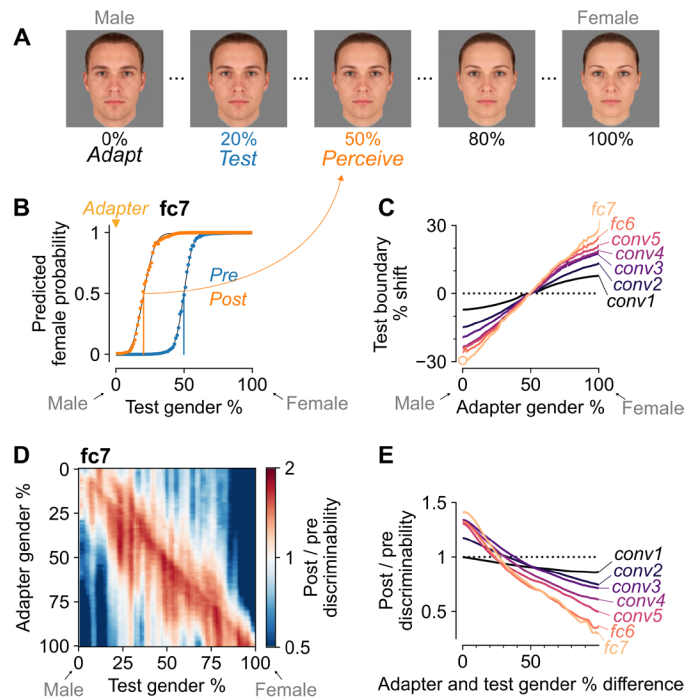
To evaluate whether the model can describe perceptual aftereffects, we created a set of face stimuli that morphed from average male to average female and measured the category decision boundary for each layer of the model before and after adaptation (Materials and Methods). Once again, we considered the same model from the previous section without any parameter changes. Exposing the model to an adapter face biased the decision boundary toward the adapter. Before adaptation, the predicted female probabilities for the model fc7 layer showed a typical sigmoidal curve centered around the gender-neutral face stimulus with morph level 50% (Fig. 3B, blue). After adapting to a male face with morph level 0%, the decision boundary shifted  $\sim 30$  percentage values toward the gender of the adapter (Fig. 3B, orange). Figure 3C shows that for all layers, adaptation to a face stimulus resulted in a boundary shift

toward the adapter. Consistent with perceptual aftereffects in human subjects, adapting to the original gender-neutral boundary stimulus with morph level 50% had no effect on the decision boundary (Fig. 3C). The perceptual bias did not suddenly emerge in later layers, but slowly built up with increasing layers (Fig. 3C, from black to purple to yellow colors), and already occurred within the first layer with intrinsic suppression (fig. S8A).

Although adapting to the boundary stimulus did not shift the decision boundary, it did increase the slope of the psychometric function for fc7 from 0.077 to 0.099 (29%; for layers conv1 to fc6, the slope changes were  $-3$ , 11, 9, 12, 16, and 31%, respectively). An increase in slope signifies a repulsion of more female and more male stimuli away from the adapter. This result is inconsistent with the perceptual renormalization hypothesis, which predicts that adaptation uniformly shifts the norm of the representational space toward the adapter and thus that adapting to the original norm (i.e., the boundary stimulus) should have no effect whatsoever [see figure 3 of (39)]. A series of previous experiments has shown that both tilt and face aftereffects involve repulsion rather than renormalization (40), which is consistent with the computational model proposed here.

Besides biasing the perception of a stimulus property, adaptation is also thought to increase sensitivity of the system for small differences





**Fig. 3. A neural network incorporating intrinsic suppression produces the perceptual bias and enhanced discriminability of aftereffects.** (A) Examples of the face-gender morph stimuli (created with webmorph.org) used in our simulated experiments. After exposure to a male adapter face, the gender decision boundary shifts toward the adapter and an observer perceives a subsequent test face as more female, and vice versa (36). The example adapt, test, and perceive morph levels were picked on the basis of the estimated boundary shift shown in (B). (B) Decision boundaries before (blue) versus after (orange) exposure to a male (0%) adapter based on the top layer (fc7) of the model with intrinsic suppression. Markers show class probabilities for each test stimulus, full lines indicate the corresponding psychometric functions, and vertical lines denote the classification boundaries. Adaptation to a 0% (male) face leads to a shift in the decision boundary toward male faces, hence perceiving the 20% test stimulus as gender-neutral (50%). (C) Decision boundary shifts for the test stimulus as a function of the adapter morph level per layer. The round marker indicates the boundary shift plotted in (B). (D) Relative face-gender discriminability (Materials and Methods, values  $>1$  signify increased discriminability and values  $<1$  denote decreased discriminability) for fc7 as a function of adapter and test morph level. See color scale on the right. The red diagonal indicates that face-gender discriminability is increased for morph levels close to the adapter. (E) Average changes in face-gender discriminability per layer as a function of the absolute difference in face-gender morph level between adapter and test stimulus.

from the current prevailing input characteristics, which could serve to maintain good stimulus discriminability (3, 4). In line with this hypothesis, Yang *et al.* (41) found that adapting to a female/male face improved gender discrimination around the face-gender morph level of the adapter. We evaluated whether intrinsic suppression in the model could account for such improved discrimination (Materials and Methods). Adaptation in the model indeed enhanced face-gender discriminability at morph levels close to the adapter (red diagonal in Fig. 3D) while decreasing discriminability at morph levels different from the adapter (blue). Like the perceptual bias (Fig. 3C), and consistent with the results shown in Fig. 2 (G and J), the discriminability effects built up monotonically across successive layers (Fig. 3E; see fig. S4, D and E, for similar results with oriented gratings).

Unlike boundary shifts, discriminability enhancements first occurred downstream of the first layer with intrinsic suppression (fig. S8B). Overall, the proposed model shows that activation-based suppression can account for discriminability improvements close to the adapter without any other specialized mechanisms and without introducing any model changes.

### Response enhancement and tuning curve shifts emerge from intrinsic suppression propagating to deeper layers

To better understand the mechanisms underlying perceptual after-effects, we investigated how adaptation affects the responses of individual units in the face-gender experiment (see fig. S5 for analyses of the tilt aftereffect). Figure 4A shows the preadaptation activation of each responsive fc7 unit across the female/male dimension (column 1) and how each unit's activation strength changed as a function of the adapter (columns 2 through 6). The rows in each heatmap are sorted according to the gender selectivity index ( $SI_g$ ; Materials and Methods), ranging from more responsive to male faces ( $SI_g < 0$ , units shown at the top) to more responsive to female faces ( $SI_g > 0$ , units shown at the bottom). After adaptation, most units showed an overall suppressed response (blue), regardless of the adapter gender morph level. However, units with a strong preference for male faces (top rows) showed an enhanced response (red) after neutral to female adapters (columns 3 to 5), whereas units with a strong preference for female faces (bottom rows) showed the opposite effect (columns 1 to 3). Thus, highly selective units showed response enhancement after adapting to the opposite of their preferred gender. This response enhancement can be explained by disinhibition (8), where adaptation reduces the inhibitory input for units that prefer morph levels further away from the adapter, much like response enhancements of middle temporal cells for their preferred direction, after adapting to the opposite direction (42).

To quantify and compare this response enhancement for all layers, we considered highly gender-selective units ( $|SI_g| > 0.6$ ) and calculated their response enhancement (averaged across all stimuli) after adapting to the opposite of their preferred gender. Figure 4B shows that the response enhancement for highly selective units (red) emerged in deeper layers (always downstream of the first layer with intrinsic suppression; fig. S9A), whereas less selective units mostly showed response suppression (blue) throughout all the layers.

Adaptation can lead to changes in response strength, but it can also cause a shift in the peak of a neuron's tuning curve. For example, in orientation-selective neurons, adapting to an oriented grating can produce a shift in the tuning curve's peak either toward the adapter [attractive shift (13, 14, 43)] or away from the adapter [repulsive shift (13, 18)]. Adaptation in the model produced both types of peak shifts in tuning curves (Fig. 4, D and E). For each unit, we calculated the proportion of adapters that produced an attractive shift or a repulsive shift (Fig. 4C). Adaptation-induced peak shifts emerged in deeper layers of the network, downstream from the first layer with intrinsic suppression (fig. S9B). Attractive shifts were more common overall, culminating to a proportion of  $\sim 0.5$  in the last layers.

Tuning changes are thought to be necessary for producing perceptual aftereffects. For example, it has been argued that a repulsive perceptual bias, where the decision boundary shifts toward the adapter, requires tuning curves that shift toward the adapter (15, 19). The fact that intrinsic suppression in the model produces mostly attractive shifts (Fig. 4C) while also capturing boundary

shifts (Fig. 3C) seems consistent with this idea. To disentangle the separate contributions of tuning changes and response magnitude changes to the perceptual adaptation effects produced by the model, we manipulated the postadaptation layer activations to only contain either tuning changes or magnitude changes (Materials and Methods; Fig. 5). Changes restricted to response magnitude without tuning changes led to even larger boundary shifts than the original model, whereas changes restricted to tuning without any changes in response magnitude led to smaller boundary shifts (Fig. 5A). This observation suggests that while the perceptual bias of aftereffects might be the result of a complex interaction between changes in responsivity and tuning, the perceptual bias does not necessarily require attractive shifts as suggested by previous models (15, 19). On the other hand, an increased face-gender discriminability for morph levels close to the adapter did require changes in the tuning response patterns of single units. Magnitude changes only produced the opposite effect, with increased discriminability for morph levels furthest from the adapter (Fig. 5B).

### Intrinsic adaptation can be optimized by maximizing recognition performance

Thus far, we have considered a model with an intrinsic adaptation state for each unit, and the adaptation parameters  $\alpha$  and  $\beta$  (Eqs. 1 and 2) were chosen to impose response suppression. This leaves open the question of whether such adaptation mechanisms can be optimized or learned in a deep learning framework given a certain task goal. We considered two possible ways in which adaptation could be learned by artificial neural networks: (i) optimize  $\alpha$  and  $\beta$  by training a feedforward network with intrinsic adaptation state on a task where adaptation is useful for biological vision; and (ii) train a recurrent network without an intrinsic adaptation state on the same task.

To assess whether adaptation could be learned and to compare the two possible network mechanisms, we needed a task objective with a suitable goal where adaptation could affect visual performance. As mentioned earlier, one of the proposed computational roles of neural adaptation is to increase sensitivity to small changes in the sensory environment (3, 4). A system could increase sensitivity by decreasing the salience of recently seen stimuli or features (5, 21). Thus, we developed a task where the end goal was object classification, but the objects were hidden in a temporally repeated noise pattern. If adaptation serves to reduce the salience of a recent stimulus, then adapting to a noise pattern should increase the ability to recognize a subsequently presented target object embedded in the same noise pattern, and a network trained on this task could learn to reduce the salience of previously presented input. To keep the networks relatively lightweight, we chose a classification task with low-resolution hand-drawn doodles rather than natural images (Fig. 6A).

Before training any network, we evaluated human recognition performance in this task. For this experiment, adaptation to the noise pattern at early levels of processing is likely sufficient to enhance the object information of the doodle. We ran a psychophysics experiment where participants were exposed to an adapter image and then classified a test image (Fig. 6B; Materials and Methods). Recognizing the doodles in this task is not trivial: whereas subjects can readily recognize the doodles in isolation, when they are embedded in noise and in the absence of any adapter, categorization performance was 59.7% (SD = 8.1%) where chance is 20%. As

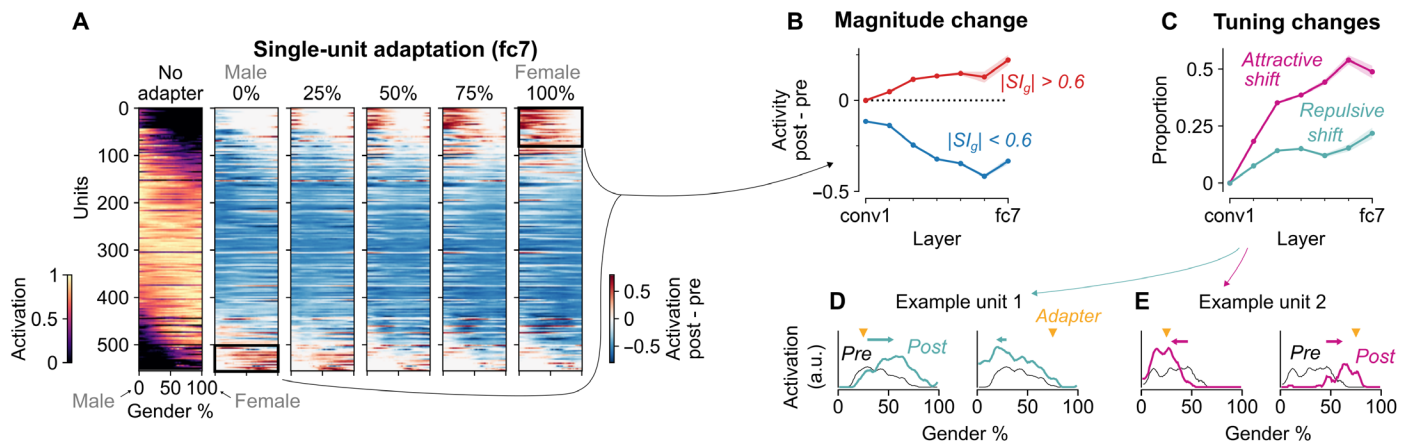
conjectured, adapting to the same noise pattern increased categorization performance by 9.3% (Fig. 6C;  $P = 0.0043$ , Wilcoxon signed-rank test,  $n = 15$  subjects). This increase in categorization performance was contingent upon the noise pattern presented during the test stimulus being the same as the noise pattern in the adapter: Performance in the same-noise condition was 9.6% higher than in the different-noise condition ( $P = 0.0015$ , Wilcoxon signed-rank test,  $n = 15$  subjects).

After establishing that adapting to the repeated noise pattern indeed improves the ability to recognize the target objects, we considered whether this behavior could be captured by the model. First, we considered the same model used in previous sections without any tuning. The same pattern of results was captured by the model with  $\alpha$  and  $\beta$  fixed to impose activation-based suppression (fig. S10). Next, we asked whether it is feasible to fit intrinsic adaptation parameters  $\alpha$  and  $\beta$  in the doodle experiment using recognition performance as the objective. We built a smaller network with an AlexNet-like architecture (Fig. 7A, without the recurrent connections shown in blue, which are discussed in the next paragraph; Materials and Methods). Each unit (excluding the decoder layer) had an exponentially decaying intrinsic adaptation state as defined by Eqs. 1 and 2. For simplicity, the trials were presented in three time steps: the adapter, a blank frame, and the test image (Fig. 7A). In addition to training the feedforward weights, we simultaneously optimized one  $\alpha$  and one  $\beta$  parameter per layer. The value of  $\alpha$  determines how fast the intrinsic adaptation state updates, ranging from no update ( $\alpha = 1$ ) to completely renewing at each time step ( $\alpha = 0$ ). The value of  $\beta$  determines whether the intrinsic adaptation state is used for suppression ( $\beta > 0$ ), enhancement ( $\beta < 0$ ), or nothing at all ( $\beta = 0$ ).

After training using 30 random initializations on same-noise trials, the resulting parameters revealed response suppression that was particularly strong for convolutional layers 1 and 2, as indicated by the positive high  $\beta$  and low  $\alpha$  values (Fig. 7B). The average categorization performance on the test set was 97.9% (blue), compared to 74.8% when no intrinsic adaptation state was included (black; Fig. 7C). Thus, when a network with intrinsic adaptation state was trained on an object recognition task with a temporally prevailing but irrelevant input pattern, the resulting adaptation parameters showed activation-based suppression.

A common way to model temporal dynamics in the visual system is by adding recurrent weights to a feedforward network (44–46). Recurrent neural networks can demonstrate phenomena similar to adaptation (47). Recurrent neural networks are the standard architectures used to process input sequences and should be able to perform well in the noisy doodle categorization task. To compare the intrinsic suppression mechanism with a recurrent circuit solution, we considered a network without intrinsic adaptation state and added lateral recurrent connections illustrated in blue in Fig. 7A (see the “Learning adaptation” section). After training on same-noise and different-noise trials, the recurrent architecture achieved the same categorization performance on the test set as the architecture with intrinsic adaptation (Fig. 7C). Thus, as expected, the recurrent network performed on par with the network with trained intrinsic adaptation.

Next, we asked whether there are any advantages of implementing adaptation via an intrinsic cellular mechanism versus lateral recurrent network mechanisms. We reasoned that a trained intrinsic suppression mechanism should generalize well across different input features or statistics, whereas the circuit-based solution learned



**Fig. 4. Response enhancements and tuning shifts emerge in deeper layers of a network incorporating intrinsic suppression.** (A) Effects of adapting to female/male faces on the activation strength of single units. Left: Heatmap showing the activation normalized to the maximum of all 556 responsive fc7 units (rows) for all face-gender morph images (columns). See the color scale on the bottom left. Rows are sorted according to the  $S_g$  (Eq. 3). The remaining five heatmaps show the difference (post – pre adaptation) in single-unit activations after adapting to five different adapters. See the color scale on the bottom right. (B) Mean response change (activity post – activity pre) across responsive units for each layer (shaded area = 95% bootstrap CI). For highly gender-selective units (red), the magnitude change (averaged across stimuli) was taken after adapting to a gender stimulus opposite to the unit’s preferred gender [0% adapter for  $S_g > 0.6$ , 100% adapter for  $S_g < -0.6$ ; black rectangles in (A)]. For less gender-selective units (blue), the magnitude change after both 0 and 100% adapters was used. (C) Proportion of adapters causing the preferred morph level to shift toward (attractive, magenta) or away (repulsive, green) from the adapter, averaged across units (shaded area = 95% bootstrap CI). (D) An example unit showing a repulsive shift in tuning curves for the 25% (left) and 75% (right) adapters [the y axes depict activation in arbitrary units (a.u.); black, preadaptation tuning curve; green, postadaptation tuning curve; yellow marker, adapter morph level]. (E) An example unit showing an attractive shift in tuning curves [magenta, postadaptation tuning curve; same conventions as (D)].

by a recurrent neural network might be less robust. Therefore, we considered situations where the distribution of noise patterns used during training and testing was different. The recurrent network failed to generalize well to higher standard deviations of Gaussian noise (Fig. 7D) and failed markedly when tested with uniformly distributed noise (Fig. 7E) or Gaussian noise with an offset (Fig. 7F). In stark contrast, the intrinsic mechanism generalized well across all of these different input noise changes (Fig. 7, D to F, magenta). This over-fitting cannot just be explained by a difference in the number of parameters and also occurs when the number of parameters is equalized between the two networks (fig. S11). Furthermore, depending on the number of parameters, the recurrent network did not necessarily demonstrate the hallmark property of repetition suppression (fig. S12). In sum, while a recurrent network implementation can learn to solve the same task, the solution is less robust than an intrinsic mechanism to deviations from the particular statistics of the adapter noise used for training the network. These results suggest that intrinsic neuronal mechanisms could provide sensory systems in the brain with a well-regularized solution to reduce salience of recent input, which is computationally simple and readily generalizes to novel sensory conditions.

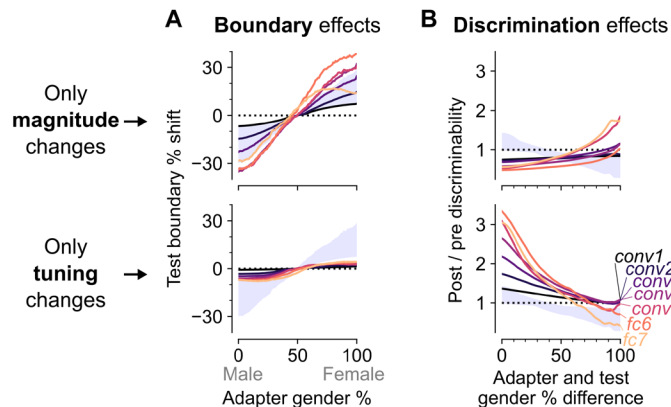
## DISCUSSION

We examined whether the paradigmatic neurophysiological and perceptual signatures of adaptation can be explained by a biologically inspired, activation-based, intrinsic suppression mechanism (7) in a feedforward deep network. The proposed computational model bridges the fundamental levels at which adaptation phenomena have been described: from intrinsic cellular mechanisms, to responses of neurons within a network, to perception. By implementing activation-based suppression (Fig. 1), our model exhibited stimulus-

specific repetition suppression (4, 5), which recovers over time but also builds up across repeats despite intervening stimuli (48) and increases over stages of processing (Fig. 2) (12, 49). Without any fine-tuning of parameters, the same model could explain classical perceptual aftereffects of adaptation (Fig. 3), such as the prototypical shift in perceptual bias toward the adapter (36, 38) and enhanced discriminability around the adapter (41, 50), thus suggesting that adaptation modulated the functional state of the network similarly to our visual system. In single units, perceptual aftereffects were associated with changes in overall responsivity (including response enhancements) as well as changes in neural tuning (Figs. 4 and 5). In addition, both intrinsic and recurrent circuit adaptation mechanisms can be trained in a task where reducing the salience of repeated but irrelevant input directly affects recognition performance (Fig. 6). However, the recurrent neural network converged on a circuit solution that was less robust to different noise conditions than the proposed model with intrinsic neuronal adaptation (Fig. 7). Together, these results show that a neuronally intrinsic suppression mechanism can robustly account for adaptation effects at the neurophysiological and perceptual levels.

The proposed computational model differs in fundamental ways from previous models of adaptation. Traditionally, adaptation has been modeled using multiple-channel models, where a fixed stimulus dimension such as orientation is encoded by a set of bell-shaped tuning functions (6, 19, 20). The core difference is that here we implemented adaptation in a deep, convolutional neural network model trained on object recognition (35). Even though current convolutional neural networks differ from biological vision in many ways (27), they constitute a reasonable first-order approximation for modeling ventral stream processing and provide an exciting opportunity for building general and comprehensive models of adaptation. First, in contrast with channel-based models, deep neural



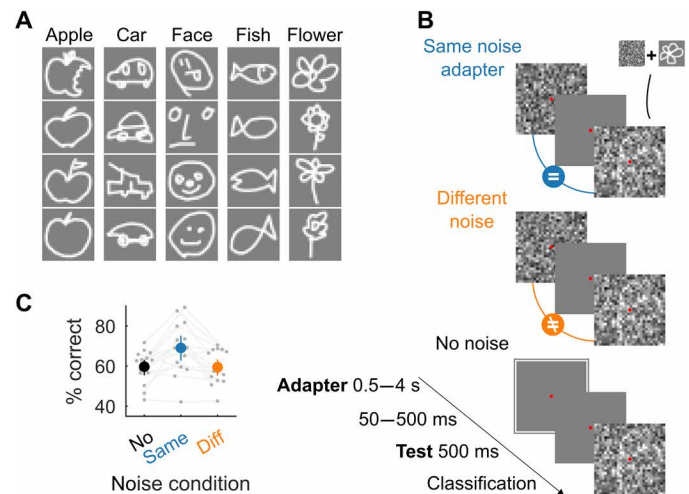


**Fig. 5. Response magnitude and tuning changes in the model differentially explain perceptual boundary shifts and discriminability changes.** (A) Face-gender boundary shifts toward the adapter were produced both by magnitude changes without tuning changes (top) and by tuning changes without magnitude changes (bottom). Gray shading indicates the range of original layer effects shown in Fig. 3C. (B) Face-gender discriminability enhancement for morph levels close to the adapter was produced by tuning changes without magnitude changes (bottom), but not by magnitude changes without tuning changes (top). Gray shading indicates the range of original layer effects shown in Fig. 3E.

networks can operate on any arbitrary image, from simple gratings to complex natural images. Second, the features encoded by the deep neural network model units are not hand-crafted tuning functions restricted to one particular stimulus dimension but consist of a rich set of increasingly complex features optimized for object recognition, which map reasonably well onto the features encoded by neurons along the primate ventral stream (28–32). A set of bell-shaped tuning curves might be a reasonable approximation of the encoding of oriented gratings in V1, but this scheme might not be appropriate for other visual areas or more complex natural images. Third, the realization that adaptation should be considered in the context of deep networks, where the effects can propagate from one stage of processing to the next (2, 21), calls for complex multilayer models that can capture the cascading of adaptation. Last, whereas several models implement adaptation by adjusting recurrent weights between channels (19, 20), we implemented an intrinsic suppression property for each unit and allowed adaptation effects to emerge from the feedforward interactions of differentially adapted units.

The goal was not to fit the model on specific datasets but to generally capture the phenomenology of adaptation in a model by giving its artificial neurons a biophysically plausible mechanism. The adaptation parameters  $\alpha$  and  $\beta$  were not fine-tuned for each simulated experiment and had the same value for each unit, showing that the ability of the model to produce adaptation phenomena did not hinge upon a carefully picked combination of parameters.

By using a feedforward deep neural network as the base for our computational model, we were able to empirically study the role of intrinsic suppression, without any contribution of recurrent interactions. These results should not be interpreted to imply that recurrent computations are irrelevant in adaptation. The results show that complex neural adaptation phenomena readily emerged in deeper layers, arguing that, in principle, they do not need to depend on recurrent mechanisms. Among the neural adaptation effects were enhanced responses of single units, as well as shifts in tuning curves, which are often thought to require recurrent network mech-



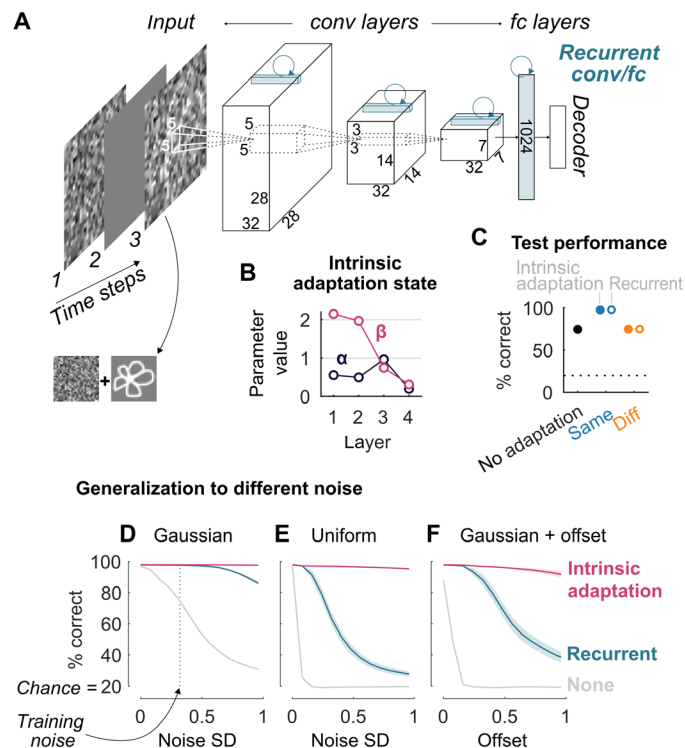
**Fig. 6. Adapting to prevailing but interfering input enhances object recognition performance.** (A) Representative examples for each of the five doodle categories from the total set of 540 selected images (63). (B) Schematic illustration of the conditions used in the doodle experiment. In each trial, participants or the model had to classify a hand-drawn doodle hidden in noise (test), after adapting to the same (middle), a different (right), or no (left) noise pattern. The trials with different or no noise adapters were control conditions where we expected to see no effect of adaptation. (C) Participants showed an increase in categorization performance after adapting to the same noise pattern. Gray circles and lines denote individual participants ( $n = 15$ ). The colored circles show average categorization performance; error bars indicate 95% bootstrap CIs. Chance = 20%.

anisms (13, 15, 16, 18). Any effect of intrinsic suppression could also be implemented by lateral inhibitory connections in the circuit, leaving open the question of why the brain would prefer one solution over the other. The generalization tests in Fig. 7 point to an intriguing possibility, which is that intrinsic suppression provides a simpler solution that is more constrained, yet sufficient to implement the goals of adaptation. In contrast, recurrent mechanisms require a complex combination of weights to achieve the same goals and tended to over-fit to the specific training conditions.

There are several functional goals that have been attributed to adaptation. Activation-based suppression could serve to decrease salience of recently seen stimuli or features (5, 21). We successfully exploited this principle to train adaptation in neural networks on a task with temporally repeated but irrelevant noise patterns. Reducing the salience of recently seen features has functional consequences beyond these artificial conditions. By selectively reducing the sensitivity of the system based on previous exposure, adaptation effectively changes the subjective experience of an observer, leading, for example, to a perceptual bias in the face-gender aftereffect. These changes in perception may more broadly reflect mechanisms that serve to maintain perceptual constancy by compensating for variations in the environment (51). The introduction of activation-based, intrinsic suppression to an artificial neural network subjected the network to the same perceptual biases characterizing perceptual aftereffects in humans (Fig. 3, B and C), suggesting that intrinsic suppression changed the model's functional state in a way that is similar to how exposure changes the functional state of our visual system.

Another proposed benefit of reducing sensitivity for recently seen stimuli may be to improve the detection of novel or less frequently





**Fig. 7. Intrinsic adaptation can be trained by maximizing recognition performance and is more robust to over-fitting than a recurrent neural network.** (A) A convolutional neural network with an AlexNet-like feedforward architecture. For the adaptation version, an exponentially decaying hidden state was added to each unit according to Eqs. 1 and 2 (except for the decoder). For the recurrent version, fully recurrent weights were added for the fully connected layer and convolutional recurrent kernels for the three convolutional layers (see drawings in blue; Materials and Methods). (B) Average fitted parameters  $\alpha$  and  $\beta$  for each layer after training 30 random initializations of the network with intrinsic adaptation state on same noise trials (SEM bars are smaller than the markers). (C) Test categorization performance on trials with the same Gaussian noise distribution as during training. Full markers: average categorization performance after training 30 random initializations on the same noise trials without intrinsic adaptation state (black), after training with intrinsic adaptation state on same noise trials (blue) or on different noise trials (orange). Empty markers: same as full markers but for the recurrent neural network. SEM bars are smaller than the markers. Chance = 20%, indicated by the horizontal dotted line. (D to F) Average generalization performance of the networks with an intrinsic adaptation state (magenta), recurrent weights (blue), or neither (gray) for same noise trials under noise conditions that differed from training. Performance is plotted as a function of increasing standard deviations (x axis) of Gaussian noise [(D), the vertical line indicates the SD = 0.32 used during training] and uniform noise (E) or as a function of increasing offset values added to Gaussian noise [(F), SD = 0.32, same as training]. Error bounds indicate SEM.

occurring stimuli (12, 48). For example, by selectively decreasing responses for more frequent stimuli, adaptation can account for the encoding of object occurrence probability, described in macaque IT (52, 53). Consistent with these observations, intrinsic suppression in the proposed computational model decreased the response strength for a given stimulus proportional to its probability of occurrence (Fig. 2, H to J). The model also produced stronger responses to a deviant stimulus compared to an equiprobable control condition. Thus, response strength in the model captured not only differences in occurrence probability (standard versus deviant) but

also relative differences in occurrence probability (control versus deviant): Compared to the control condition, the deviant is equally likely in terms of absolute occurrence probability, but it was unexpected merely by virtue of the higher occurrence probability of the standard stimulus.

Adaptation has also been suggested to increase coding efficiency of single neurons by normalizing their responses for the current sensory conditions (4). Neurons have a limited dynamic range with respect to the feature they encode and a limited number of response levels. Adaptation can maximize the information carried by a neuron by re-centering tuning around the prevailing conditions and thus increasing sensitivity and preventing response saturation (51). While AlexNet has ReLU activation functions, which do not suffer from the saturation problem, we did observe an abundance of attractive shifts of tuning curves (Fig. 4C). The collective result of these changes in tuning curves was an increased discriminability between stimuli similar to the adapter (Fig. 4D), consistent with reports for orientation, motion direction, and face-gender discrimination in humans (41, 50).

Besides direct functional benefits, adaptation may also serve an important role in optimizing the efficiency of the neural population code. Neurons use large amounts of energy to generate action potentials, which constrains neural representations (54). When a particular feature combination is common, the metabolic efficiency of the neural code can be improved by decorrelating responses of the activated cells and reducing their responsiveness. Adaptation has been shown to maintain existing response correlations and equality in time-averaged responses across the population (55), possibly resulting from intrinsic suppression at an earlier cortical stage, which we confirmed by running these experiments in the proposed computational model (fig. S13).

There are several possible extensions to the current model, including the incorporation of multiple time scales and recurrent circuit mechanisms. Adaptation operates over a range of time scales and thus may be best described by a scale-invariant power law, which could be approximated by extending the model using a sum of exponential processes (56). Our model also did not include any recurrent dynamics, because we focused on the feedforward propagation of intrinsic suppression. Yet, recurrent connections are abundant in sensory systems and most likely do contribute to adaptation. There is some evidence suggesting that recurrent mechanisms contribute to adaptation at very short time scales of up to 100 ms (57). During the first 50 to 100 ms after exposure, adaptation to an oriented grating produces a perceptual boundary shift in the opposite direction of the classical tilt aftereffect (58). This observation was predicted by a recurrent V1 model that only predicted repulsive tuning shifts (6). Repulsive shifts are indeed more common in V1 when each test stimulus is immediately preceded by an adapter (13, 18), whereas adaptation seems to produce mostly attractive shifts at longer gaps (14, 43, 59), consistent with the effects of intrinsic suppression in the proposed model (Fig. 4 and fig. S5; although repulsive shifts were more common in highly responsive units; fig. S6). These results seem to suggest that recurrent interactions contribute in the first (few) 100 ms, whereas qualitatively different longer adaptation effects might be best accounted for by intrinsic suppression.

The results of the noisy doodle experiment in humans (Fig. 6) could be explained by local light adaptation to the adapter noise patterns. It is unclear where in the visual system such local light adaptation would take place. In principle, it could take place

partly or totally at the level of photoreceptors in the retina. However, given that each noise pixel was only  $0.3 \times 0.3$  visual degrees and given that luminance was distributed independently across noise pixels, inherent variability in the gaze of a fixating subject poses a limit on the contribution of photoreceptor adaptation (60). Most likely, the increased performance observed in the behavioral data results from a combination of adaptation at different stages of processing, including the retina. The proposed computational model does not incorporate adaptation at the receptor level (i.e., pixels), but future models could incorporate adaptation in both the input layer and later processing layers.

Overall, the current framework connects systems to cellular neuroscience in one comprehensive multilevel model by including an activation-based, intrinsic suppression mechanism in a deep neural network. Response suppression cascading through a feedforward hierarchical network changed the functional state of the network similar to visual adaptation, producing complex downstream neural adaptation effects as well as perceptual aftereffects. These results demonstrate that intrinsic neural mechanisms may contribute substantially to the dynamics of sensory processing and perception in a temporal context.

## MATERIALS AND METHODS

### Computational models

#### Implementing intrinsic suppression

We used the AlexNet architecture (Fig. 1A) (35), with weights pretrained on the ImageNet dataset (61) as a model for the ventral visual stream. We implemented an exponentially decaying intrinsic adaptation state (62) to simulate neuronally intrinsic suppression. Specifically, in all layers (except the decoder), each unit had an intrinsic adaptation state  $s_t$ , which was updated at each time step  $t$  based on its previous state  $s_{t-1}$  and the previous response  $r_{t-1}$  (i.e., activation after the ReLU rectification and linearization operation)

$$s_t = \alpha s_{t-1} + (1 - \alpha) r_{t-1} \quad (1)$$

where  $\alpha$  is a constant in  $[0,1]$  determining the time scale of the decay (Fig. 1B). This intrinsic adaptation state is then subtracted from the unit's current input  $x_t$  (given weights  $W$  and bias  $b$ ) before applying the rectifier activation function  $\sigma$ , so that

$$r_t = \sigma(b + Wx_t - \beta s_t) \quad (2)$$

where  $\beta$  is a constant that scales the amount of suppression. Thus, strictly speaking, Eq. 2 modifies the bias and thus responsivity of the unit, before applying  $\sigma$ , to avoid negative activations. For  $\beta > 0$ , these model updating rules result in an exponentially decaying response for constant input that recovers in case of no input (Fig. 1B), simulating an activation-based suppression mechanism intrinsic to each individual neuron. Note that  $\beta < 0$  would lead to response enhancement and  $\beta = 0$  would leave the response unchanged. By implementing this mechanism across discrete time steps in AlexNet, we introduced a temporal dimension to the network (Fig. 1C). This model was implemented using TensorFlow v1.11 in Python. Throughout the paper, we use  $\alpha = 0.96$  and  $\beta = 0.7$  unless indicated otherwise (in Fig. 7, those parameters are tuned).

#### Decision boundaries

Perceptual aftereffects are typically measured by computing shifts in the decision boundary along a stimulus dimension. We evaluated

boundary shifts in the model using a set of face stimuli that morphed from average male to average female in 100 steps (created using webmorph.org) and measured category decision boundaries before and after adaptation using the 101 face-morph images (Fig. 3, A to C). The experiments were simulated by exposing the model to an adapter image for 100 time steps, followed by a gap of uniform gray input for 10 time steps before presenting the test image. The results were qualitatively similar when the number of time steps was changed.

To measure the pre- and postadaptation decision boundaries for a given layer, we trained a logistic regression classifier to discriminate between male and female faces using the preadaptation activations of responsive units for the full stimulus set. After training, the classifier can output female/male class probability estimates for any given activation pattern. Thus, we used the trained classifier to provide female/male probability estimates for each morph level, based on either the pre- or postadaptation activation patterns. The decision boundary is then given by the morph level associated with a female/male class probability of  $P = 0.5$ , which was estimated by fitting a psychometric function on the class probabilities (average  $R^2$  of at least 0.99 per layer).

#### Face-gender discriminability

To assess model changes in face-gender discriminability in Fig. 3J, we calculated the stimulus discriminability at each morph level of the stimulus dimension before and after adaptation. An increased discriminability between morph levels can be conceptualized as an increased perceived change in morph levels with respect to a certain physical change in morph level. Thus, to quantify discriminability, a linear mapping was fit to predict stimulus morph levels from preadaptation unit activations using partial least squares regression (using four components). We then used this linear mapping to predict morph levels from activation patterns before and after adaptation. If adaptation increases discriminability, then the change in model-estimated morph level  $y$  with respect to a physical change in morph level  $m$  should also increase. Thus, to quantify the change in discriminability at morph level  $m$ , we calculated the absolute derivative of the predicted postadaptation morph level ( $y_m^{post}$ ), normalized by the absolute derivative of the predicted preadaptation morph level ( $y_m^{pre}$ ):  $|\Delta y_m^{post}|/|\Delta y_m^{pre}|$ .

#### Selectively retaining tuning or magnitude changes

For Fig. 4B, we manipulated the postadaptation layer activations to only contain either tuning changes or magnitude changes. To retain only tuning changes, we started with the postadaptation activation patterns and multiplied the activation of each unit by a constant so that the resulting mean activation matched the preadaptation mean value. On the other hand, to retain only magnitude changes, we started with the preadaptation activation patterns and multiplied the activation of each unit by a constant so that the resulting mean activation matched the postadaptation mean value.

#### Learning adaptation

In Fig. 7, we present two models where adaptation is learned for the noisy doodle classification task: a model with intrinsic adaptation state and a recurrent neural network model. The base feedforward part of the model was based on the AlexNet architecture (35) for the two networks, consisting of three convolutional layers and a fully connected layer followed by a fully connected decoder. The first convolutional layer filters a  $28 \times 28 \times 1$  input image with 32 kernels of size  $5 \times 5 \times 1$  with a stride of 1 pixel. The second convolutional layer filters the pooled (kernel =  $2 \times 2$ , stride = 2) output of the first convolutional layer with 32 kernels of size  $5 \times 5 \times 32$  (stride = 1).

The third convolutional layer filters the pooled (kernel =  $2 \times 2$ , stride = 2) output of the second convolutional layer with 32 kernels of size  $3 \times 3 \times 32$  (stride = 1). The fully connected layer has 1024 units that process the output of the third convolutional layer with 50% dropout during training.

The recurrent version was extended with lateral recurrent weights. For convolutional layers, lateral recurrence was implemented as 32 kernels of size  $1 \times 1 \times 32$  (stride = 1), which filtered the nonpooled outputs of the layer at time step  $t - 1$  (after ReLu) and were added to the feedforward-filtered inputs of the same layer at time step  $t$  (before ReLu). The fully connected layer was recurrent in an all-to-all fashion.

The intrinsic adaptation version was extended with adaptation states, as described in the “Implementing intrinsic suppression” section, of which the  $\alpha$  and  $\beta$  parameters were now also trained using back-propagation. The  $\beta$  parameters were initialized at 0 (i.e., no adaptation), and the  $\alpha$  parameters were initialized using a uniform distribution ranging from 0 to 1.

Both the recurrent and intrinsic adaptation models were trained on the doodle classification task using TensorFlow v1.11 in Python. We used a training set of 500,000 doodle images (<https://github.com/googlecreativelab/quickdraw-dataset>; 100,000 per category), with a separate set of 1000 images to select hyperparameters and evaluate the loss and accuracy during training. We used the Adam optimization algorithm (63) with a learning rate of 0.001, the sparse softmax cross entropy between logits and labels cost function, a batch size of 100, and 50% training dropout in fully connected layers. For the weights, we used Gaussian initialization, with the scale correction proposed by Glorot and Bengio (64). Each model was trained for five epochs on the training set, which was sufficient for the loss and accuracy to saturate. Generalization performance was then tested on a third independent set of 5000 images.

## Neurophysiology

We present neurophysiological data from two previously published studies to compare them with the neural adaptation effects of the proposed computational model: single-cell recordings from IT ( $n = 97$ ) cortex of one macaque monkey G (37) and multi-unit recordings from V1 ( $n = 55$ ) and latero-intermediate visual area (LI;  $n = 48$ ) of three rats (12). For methodological details about the recordings and the tasks, we refer to the original papers.

## Psychophysics

Before starting the data collection, we preregistered the study design and hypothesis on the Open Science Framework at <https://osf.io/tdb37/> where all the source code and data can be retrieved.

## Participants

A total of 17 volunteers (10 female, ages 19 to 50) participated in our doodle categorization experiments (Fig. 6). In accordance with our preregistered data exclusion rule, two male participants were excluded from analyses because we could not record eye tracking data. All subjects gave informed consent, and the studies were approved by the Institutional Review Board at Children’s Hospital, Harvard Medical School.

## Stimuli

The stimulus set consisted of hand-drawn doodles of apples, cars, faces, fish, and flowers from the Quick, Draw! dataset (<https://github.com/googlecreativelab/quickdraw-dataset>). We selected a total of 540 doodles (108 from each of the five categories) that were

judged complete and identifiable. We lowered the contrast of each doodle image ( $28 \times 28$  pixels) to either 22 or 29% of the original contrast, before adding a Gaussian noise pattern (SD = 0.165 in normalized pixel values) of the same resolution. The higher contrast level (29%) was chosen as a control so that the doodle was relatively visible in one-sixth of the trials and was not included in the analyses. The average categorization performance on these high-contrast trials was 74% (SD = 8.3%), versus 63% (SD = 8.9%) in the low-contrast trials.

## Experimental protocol

Participants had to fixate a cross at the center of the screen to start a trial. Next, an adapter image was presented (for 0.5, 2, or 4 s), followed by a blank interval (of 50, 250, or 500 ms), a test image (for 500 ms), and lastly a response prompt screen. The test images were noisy doodles described in the above paragraph. The adapter image could either be an empty frame (defined by a white square filled with the background color), the same mosaic noise pattern as the one of the subsequent test image, or a randomly generated different noise pattern (Fig. 6). Participants were asked to keep looking at the fixation cross, which remained visible throughout the entire trial, until they were prompted to classify the test image using keyboard keys 1 to 5. All images were presented at  $9^\circ \times 9^\circ$  from a viewing distance of approximately 52 cm on a 19-inch cathode ray tube monitor (Sony Multiscan G520;  $1024 \times 1280$  resolution), while we continuously tracked eye movements using a video-based eye tracker (EyeLink 1000, SR Research, Canada). Trials where the root mean square deviation of the eye movements exceeded  $1^\circ$  of visual angle during adapter presentation were excluded from further analyses. The experiment was controlled by custom code written in MATLAB using Psychophysics Toolbox Version 3.0 (65).

## Data analysis

### Selectivity index

For the face-gender experiments, we calculated a selectivity index based on the average activation of a unit to male (morph level < 50%) and female (morph level > 50%) faces

$$SI_g = (A_F - A_M) / (A_F + A_M) \quad (3)$$

A value >0 indicates stronger activation for female faces, and a value <0 indicates stronger activation for male faces.

## SUPPLEMENTARY MATERIALS

Supplementary material for this article is available at <http://advances.sciencemag.org/cgi/content/full/6/42/eabd4205/DC1>

[View/request a protocol for this paper from Bio-protocol.](#)

## REFERENCES AND NOTES

1. R. Addams, L.I. An account of a peculiar optical phenomenon seen after having looked at a moving body. *London Edinburgh Philos. Mag. J. Sci.* **5**, 373–374 (1834).
2. C. J. Whitmore, G. B. Stanley, Rapid sensory adaptation redux: A circuit perspective. *Neuron* **92**, 298–315 (2016).
3. C. W. G. Clifford, M. A. Webster, G. B. Stanley, A. A. Stocker, A. Kohn, T. O. Sharpee, O. Schwartz, Visual adaptation: Neural, psychological and computational aspects. *Vision Res.* **47**, 3125–3131 (2007).
4. A. Kohn, Visual adaptation: Physiology, mechanisms, and functional benefits. *J. Neurophysiol.* **97**, 3155–3164 (2007).
5. R. Vogels, Sources of adaptation of inferior temporal cortical responses. *Cortex* **80**, 185–195 (2016).
6. M. Del Mar Quiroga, A. P. Morris, B. Krekelberg, Adaptation without plasticity. *Cell Rep.* **17**, 58–68 (2016).



7. M. V. Sanchez-Vives, L. G. Nowak, D. A. McCormick, Membrane mechanisms underlying contrast adaptation in cat area 17 in vivo. *J. Neurosci.* **20**, 4267–4285 (2000).
8. B. Krekelberg, G. M. Boynton, R. J. A. van Wezel, Adaptation: From single cells to BOLD signals. *Trends Neurosci.* **29**, 250–256 (2006).
9. H. Sawamura, G. A. Orban, R. Vogels, Selectivity of neuronal adaptation does not match response selectivity: A single-cell study of the fMRI adaptation paradigm. *Neuron* **49**, 307–318 (2006).
10. D. A. Kaliukhovich, R. Vogels, Divisive normalization predicts adaptation-induced response changes in macaque inferior temporal cortex. *J. Neurosci.* **36**, 6116–6128 (2016).
11. S. C. Wissig, A. Kohn, The influence of surround suppression on adaptation effects in primary visual cortex. *J. Neurophysiol.* **107**, 3370–3384 (2012).
12. K. Vinken, R. Vogels, H. O. de Breeck, Recent visual experience shapes visual processing in rats through stimulus-specific adaptation and response enhancement. *Curr. Biol.* **27**, 914–919 (2017).
13. V. Dragoi, J. Sharma, M. Sur, Adaptation-induced plasticity of orientation tuning in adult visual cortex. *Neuron* **28**, 287–298 (2000).
14. J. Jeyabalaratnam, V. Bharamuria, L. Bachatene, S. Cattani, A. Angers, S. Molotchnikoff, Adaptation shifts preferred orientation of tuning curve in the mouse visual cortex. *PLOS ONE* **8**, e64294 (2013).
15. A. Kohn, J. A. Movshon, Adaptation changes the direction tuning of macaque MT neurons. *Nat. Neurosci.* **7**, 764–772 (2004).
16. A. B. Saul, M. S. Cynader, Adaptation in single units in visual cortex: The tuning of aftereffects in the spatial domain. *Vis. Neurosci.* **2**, 593–607 (1989).
17. A. B. Saul, M. S. Cynader, Adaptation in single units in visual cortex: The tuning of aftereffects in the temporal domain. *Vis. Neurosci.* **2**, 609–620 (1989).
18. G. Felsen, Y.-s. Shen, H. Yao, G. Spor, C. Li, Y. Dan, Dynamic modification of cortical orientation tuning mediated by recurrent connections. *Neuron* **36**, 945–954 (2002).
19. A. F. Teich, N. Qian, Learning and adaptation in a recurrent model of V1 orientation selectivity. *J. Neurophysiol.* **89**, 2086–2100 (2003).
20. Z. M. Westrick, D. J. Heeger, M. S. Landy, Pattern adaptation and normalization reweighting. *J. Neurosci.* **36**, 9805–9816 (2016).
21. S. G. Solomon, A. Kohn, Moving sensory adaptation beyond suppressive effects in single neurons. *Curr. Biol.* **24**, R1012–R1022 (2014).
22. N. T. Dhruv, M. Carandini, Cascaded effects of spatial adaptation in the early visual system. *Neuron* **81**, 529–535 (2014).
23. D. L. K. Yamins, J. J. DiCarlo, Using goal-driven deep learning models to understand sensory cortex. *Nat. Neurosci.* **19**, 356–365 (2016).
24. B. A. Richards, T. P. Lillicrap, P. Beaudoin, Y. Bengio, R. Bogacz, A. Christensen, C. Clopath, R. P. Costa, A. de Berker, S. Ganguli, C. J. Gillon, D. Hafner, A. Kepecs, N. Kriegeskorte, P. Latham, G. W. Lindsay, K. D. Miller, R. Naud, C. C. Pack, P. Poirazi, P. Roelfsema, J. Sacramento, A. Saxe, B. Scellier, A. C. Schapiro, W. Senn, G. Wayne, D. Yamins, F. Zenke, J. Zylberberg, D. Thierien, K. P. Kording, A deep learning framework for neuroscience. *Nat. Neurosci.* **22**, 1761–1770 (2019).
25. T. C. Kietzmann, P. M. Clure, N. Kriegeskorte, Deep neural networks in computational neuroscience, in *Oxford Research Encyclopedia of Neuroscience* (Oxford Univ. Press, 2019), pp. 1–28.
26. R. M. Cichy, D. Kaiser, Deep neural networks as scientific models. *Trends Cogn. Sci.* **23**, 305–317 (2019).
27. T. Serre, Deep learning: The good, the bad, and the ugly. *Annu. Rev. Vis. Sci.* **5**, 399–426 (2019).
28. C. F. Cadieu, H. Hong, D. L. K. Yamins, N. Pinto, D. Ardila, E. A. Solomon, N. J. Majaj, J. J. DiCarlo, Deep neural networks rival the representation of primate it cortex for core visual object recognition. *PLOS Comput. Biol.* **10**, e1003963 (2014).
29. D. L. K. Yamins, H. Hong, C. F. Cadieu, E. A. Solomon, D. Seibert, J. J. DiCarlo, Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proc. Natl. Acad. Sci. U.S.A.* **111**, 8619–8624 (2014).
30. I. Kalfas, S. Kumar, R. Vogels, Shape selectivity of middle superior temporal sulcus body patch neurons. *eNeuro* **4**, ENURO.0113-17 (2017).
31. I. Kalfas, K. Vinken, R. Vogels, Representations of regular and irregular shapes by deep Convolutional Neural Networks, monkey inferotemporal neurons and human judgments. *PLOS Comput. Biol.* **14**, e1006557 (2018).
32. D. A. Pospisil, A. Pasupathy, W. Bair, 'Artphysiology' reveals V4-like shape tuning in a deep network trained for image classification. *eLife* **7**, e38242 (2018).
33. S.-M. Khaligh-Razavi, N. Kriegeskorte, Deep supervised, but not unsupervised, models may explain it cortical representation. *PLOS Comput. Biol.* **10**, e1003915 (2014).
34. J. Kubilius, S. Bracci, H. P. O. de Breeck, Deep neural networks as a computational model for human shape sensitivity. *PLOS Comput. Biol.* **12**, e1004896 (2016).
35. A. Krizhevsky, I. Sutskever, G. E. Hinton, ImageNet classification with deep convolutional neural networks. *Adv. Neural Inform. Process. Syst.* **2012**, 1097–1105 (2012).
36. M. A. Webster, D. Kaping, Y. Mizokami, P. Duhamel, Adaptation to natural facial categories. *Nature* **428**, 557–561 (2004).
37. K. Vinken, H. P. O. de Breeck, R. Vogels, Face repetition probability does not affect repetition suppression in macaque inferotemporal cortex. *J. Neurosci.* **38**, 7492–7504 (2018).
38. J. J. Gibson, M. Radner, Adaptation, after-effect and contrast in the perception of tilted lines. I. Quantitative studies. *J. Exp. Psychol.* **20**, 453–467 (1937).
39. M. A. Webster, D. I. A. Macleod, Visual adaptation and face perception. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **366**, 1702–1725 (2011).
40. K. R. Storrs, D. H. Arnold, Face aftereffects involve local repulsion, not renormalization. *J. Vis.* **15**, 1–18 (2015).
41. H. Yang, J. Shen, J. Chen, F. Fang, Face adaptation improves gender discrimination. *Vision Res.* **51**, 105–110 (2011).
42. S. E. Petersen, J. F. Baker, J. M. Allman, Direction-specific adaptation in area MT of the owl monkey. *Brain Res.* **346**, 146–150 (1985).
43. N. Ghisovan, A. Nemri, S. Shumikhina, S. Molotchnikoff, Long adaptation reveals mostly attractive shifts of orientation tuning in cat primary visual cortex. *Neuroscience* **164**, 1274–1283 (2009).
44. H. Tang, M. Schrimpf, W. Lotter, C. Moerman, A. Paredes, J. O. Caro, W. Hardesty, D. Cox, G. Kreiman, Recurrent computations for visual pattern completion. *Proc. Natl. Acad. Sci. U.S.A.* **115**, 8835–8840 (2018).
45. K. Kar, J. Kubilius, K. Schmidt, E. B. Issa, J. J. DiCarlo, Evidence that recurrent circuits are critical to the ventral stream's execution of core object recognition behavior. *Nat. Neurosci.* **22**, 974–983 (2019).
46. T. C. Kietzmann, C. J. Spoeer, L. K. A. Sørensen, R. M. Cichy, O. Hauk, N. Kriegeskorte, Recurrence is required to capture the representational dynamics of the human visual system. *Proc. Natl. Acad. Sci. U.S.A.* **116**, 21854–21863 (2019).
47. W. Lotter, G. Kreiman, D. Cox, A neural network trained for prediction mimics diverse features of biological neurons and perception. *Nat. Mach. Intell.* **2**, 210–219 (2020).
48. N. Ulanovsky, L. Las, I. Nelken, Processing of low-probability sounds by cortical neurons. *Nat. Neurosci.* **6**, 391–398 (2003).
49. D. A. Kaliukhovich, H. O. de Breeck, Hierarchical stimulus processing in rodent primary and lateral visual cortex as assessed through neuronal selectivity and repetition suppression. *J. Neurophysiol.* **120**, 926–941 (2018).
50. C. W. G. Clifford, Perceptual adaptation: Motion parallels orientation. *Trends Cogn. Sci.* **6**, 136–143 (2002).
51. M. A. Webster, J. S. Werner, D. J. Field, Adaptation and the phenomenology of perception, in *Fitting the Mind to the World Adaptation and After-Effects in High-Level Vision* (Oxford Univ. Press, 2005), chap. 10, pp. 241–278.
52. A. H. Bell, C. Summerfield, E. L. Morin, N. J. Malecek, L. G. Ungerleider, Encoding of stimulus probability in macaque inferior temporal cortex. *Curr. Biol.* **26**, 2280–2290 (2016).
53. K. Vinken, R. Vogels, Adaptation can explain evidence for encoding of probabilistic information in macaque inferior temporal cortex. *Curr. Biol.* **27**, R1210–R1212 (2017).
54. P. Lennie, The cost of cortical computation. *Curr. Biol.* **13**, 493–497 (2003).
55. A. Benucci, A. B. Saleem, M. Carandini, Adaptation maintains population homeostasis in primary visual cortex. *Nat. Neurosci.* **16**, 724–729 (2013).
56. P. J. Drew, L. F. Abbott, Models and properties of power-law adaptation in neural systems. *J. Neurophysiol.* **96**, 826–833 (2006).
57. M. Wehr, A. M. Zador, Synaptic mechanisms of forward suppression in rat auditory cortex. *Neuron* **47**, 437–445 (2005).
58. M. Del Mar Quiroga, A. P. Morris, B. Krekelberg, Short-term attractive tilt aftereffects predicted by a recurrent network model of primary visual cortex. *Front. Syst. Neurosci.* **13**, 67 (2019).
59. V. Bharamuria, L. Bachatene, S. Molotchnikoff, The speed of neuronal adaptation: A perspective through the visual cortex. *Eur. J. Neurosci.* **49**, 1215–1219 (2019).
60. M. Rucci, M. Poletti, Control and functions of fixational eye movements. *Annu. Rev. Vis. Sci.* **1**, 499–518 (2015).
61. O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, L. Fei-Fei, ImageNet large scale visual recognition challenge. *Int. J. Comput. Vis.* **115**, 211–252 (2015).
62. G. Bellec, D. Salaj, A. Subramoney, R. Legenstein, W. Maass, Long short-term memory and learning-to-learn in networks of spiking neurons, in *Proceedings of the 32nd International Conference on Neural Information Processing Systems* (December 2018), pp. 795–805.
63. D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, in *Proceedings of the 3rd International Conference on Learning Representations (ICLR)* (San Diego, 2015), pp. 1–15.
64. X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, in *Proceedings of the 13th International Conference On Artificial Intelligence and Statistics* (2010), vol. 9, pp. 249–256.
65. D. H. Brainard, The psychophysics toolbox. *Spat. Vis.* **10**, 433–436 (1997).

#### Acknowledgments

**Funding:** This work was supported by Research Foundation Flanders, Belgium (fellowship of K.V.), by NIH grant R01EY026025, and by the Center for Brains, Minds and Machines, funded by



NSF Science and Technology Centers Award CCF-1231216. **Author contributions:** K.V. conceived the model and experiment; K.V., X.B., and G.K. designed the model and experiment; K.V. collected the data, implemented the model, and carried out analyses; K.V. and G.K. wrote the manuscript, with contributions from X.B. **Competing interests:** The authors declare that they have no competing interests. **Data and materials availability:** All data needed to evaluate the conclusions in the paper are present in the paper and/or the Supplementary Materials. Additional data related to this paper may be requested from the authors. All the psychophysics data and source code are available at <https://osf.io/tdb37/>.

Submitted 21 June 2020  
Accepted 26 August 2020  
Published 14 October 2020  
10.1126/sciadv.abd4205

**Citation:** K. Vinken, X. Boix, G. Kreiman, Incorporating intrinsic suppression in deep neural networks captures dynamics of adaptation in neurophysiology and perception. *Sci. Adv.* **6**, eabd4205 (2020).