

Reviewed Preprint

v1 • January 16, 2026

Not revised

Reviewed Preprint

v2 • June 15, 2026

Revised by authors

✉ For correspondence:

nikoskar@stanford.edu

kafranke@stanford.edu

* Equal contributions

Competing interests: No competing interests declared

Funding: See [page 36](#)

Reviewing editor: Tirin Moore, Stanford University, Howard Hughes Medical Institute, United States

© 2026, Karantzas et al. This article is distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use and redistribution provided that the original author and source are credited.

Dual-feature selectivity enables bidirectional coding in visual cortical neurons

Nikos Karantzas^{1,2,3,*} ✉, Katrin Franke^{1,2,3,4,*} ✉, Konstantin Willeke^{1,2,3}, Maria Diamantaki^{5,6}, Kandan Ramakrishnan⁷, Hasan Atakan Bedel^{1,2,3}, Pavithra Elumalai⁸, Kelli Restivo⁷, Paul Fahey^{1,2,3}, Cate Nealley^{1,2,3}, Tori Shinn¹¹, Gabrielle Garcia^{1,2,3}, Saumil Patel^{1,2,3}, Alexander Ecker^{8,9,10}, Edgar Y Walker^{12,13}, Emmanouil Froudarakis^{5,6}, Sophia Sanborn^{1,2,3}, Fabian H Sinz^{7,8,9}, Andreas Tolias^{1,2,3}

¹Department of Ophthalmology, Byers Eye Institute, Stanford University School of Medicine, Stanford, United States • ²Stanford Bio-X, Stanford University, Stanford, United States • ³Wu Tsai Neurosciences Institute, Stanford University, Stanford, United States • ⁴Institute for Ophthalmic Research, Tübingen University, Tübingen, Germany • ⁵Institute of Molecular Biology & Biotechnology, Foundation of Research & Technology - Hellas, Heraklion, Greece • ⁶Department of Basic Sciences, Faculty of Medicine, University of Crete, Heraklion, Greece • ⁷Department of Neuroscience & Center for Neuroscience and Artificial Intelligence, Baylor College of Medicine, Houston, United States • ⁸Institute of Computer Science and Campus Institute Data Science, University of Göttingen, Göttingen, Germany • ⁹Lower Saxony Center for AI and Causal Methods in Medicine, Hanover, Germany • ¹⁰Max Planck Institute for Dynamics and Self-Organization, Göttingen, Germany • ¹¹Department of Pediatrics; Allergy & Immunology, Baylor College of Medicine, Houston, United States • ¹²Department of Neurobiology & Biophysics, University of Washington School of Medicine, Seattle, United States • ¹³Computational Neuroscience Center, University of Washington, Seattle, United States

eLife Assessment

The authors combine a modeling approach, using a digital twin, with electrophysiological evidence in two species to assess the role of inhibition in shaping selectivity in the visual cortex. The results provide a **fundamental** advance beyond the classic view of sensory coding by proving **compelling** evidence that many neurons in visual areas exhibit dual-feature selectivity. Overall, the work compellingly showcases how in silico experiments can generate concrete hypotheses about neuronal coding that are difficult to discover experimentally.

<https://doi.org/10.7554/eLife.109861.2.sa2>

Abstract

Sensory neurons are traditionally viewed as feature detectors that respond with an increase in firing rate to preferred stimuli while remaining unresponsive to others. Here, we identify a dual-feature encoding strategy in macaque visual cortex, wherein many neurons in areas V1 and V4 are selectively tuned to two distinct visual features—one that enhances and one that suppresses activity—around an elevated baseline firing rate. By combining neuronal recordings with functional digital twin models—deep learning-based predictive models of biological neurons—we were able to systematically identify each neuron’s preferred and non-preferred features. These feature pairs served as anchors for a continuous, low-dimensional axis in natural image similarity space, along which neuronal activity varied approximately linearly. Within a single visual area, visual features that strongly or weakly activated individual neurons also had a high probability of modulating the activity of other neurons, suggesting a shared feature selectivity across the population that structures stimulus encoding. We show that this encoding strategy is conserved across species, present in both primary and lateral visual areas of mouse cortex. Dual-feature

selectivity is consistent with recent anatomical evidence for feature-specific inhibitory connectivity, complementing the feature-detector principle through circuit mechanisms in which selective excitation and inhibition may together enhance the representational capacity of the neuronal population.

Introduction

Early sensory physiology established the picture of the single neuron as a feature detector that remains mostly silent until a stimulus in its receptive field resembles a preferred pattern (e.g. [Lettvin et al., 1959](#)). This view has intuitive appeal because it suggests that each neuron signals the presence of a specific feature in the sensory world, resulting in a sparse ([Barlow, 1961](#); [Field, 1987](#); [Olshausen and Field, 1996](#)) and metabolically efficient code ([Levy and Baxter, 1996](#); [Attwell and Laughlin, 2001](#)). For example, in cat and monkey primary visual cortex (V1), a simple cell's preferred pattern can be described as a specific Gabor function ([Hubel and Wiesel, 1962](#), [1968](#); [Bishop et al., 1973](#))—a precise combination of orientation, spatial frequency, phase, size, and retinal location. As the similarity of a stimulus to this Gabor increases, the neuron's firing rate rises, while all dissimilar stimuli evoke little or no activity. Complex cells extend this principle by pooling across spatial phase, yielding phase-invariant orientation tuning ([Hubel and Wiesel, 1968](#); [Schiller et al., 1976](#)). This feature detector view has been extrapolated to higher visual areas, where, in the idealized model, a neuron might respond exclusively to a specific face, as in the so-called “Jennifer Aniston neuron” ([Quiroga et al., 2005](#)).

Under this framework, lifetime sparsity—the proportion of stimuli that elicit strong firing across a neuron's entire experience (see also [Willmore and Tolhurst, 2001](#); [Willmore et al., 2011](#))—is determined by the statistics of the preferred feature: neurons tuned to low spatial frequencies tend to fire more frequently because such content is abundant in natural scenes due to their $1/f$ spectral properties ([Field, 1987](#)), whereas neurons selective for high frequency edges or specific faces may fire only rarely. While traditionally described in terms of increased firing for preferred features, cortical responses also depend on the suppressive influences that modulate those activations.

Neuronal responses of real cortical neurons, however, reflect a balance of excitatory (i.e. preferred stimulus) and inhibitory inputs (i.e. non-preferred stimulus) within their receptive field. This inhibition is usually considered to be broadly tuned, providing a form of blanket gain control that normalizes neuronal activity across the population ([Heeger, 1992](#); [Carandini and Heeger, 2011](#)). A classic example of feature-specific inhibition constitutes crossorientation inhibition in V1, where responses to a neuron's preferred orientation are suppressed by the simultaneous presentation of an orthogonal grating ([Morrone et al., 1982](#); [Allison et al., 1995](#); [Ferster, 1986](#)). Such examples illustrate that selectivity emerges from the joint action of excitation and inhibition, whose relative balance shapes the final tuning curve. Yet, how this interplay operates in the space of natural images remains poorly understood.

Experimental constraints limit the number of stimuli that can be presented during neuronal recordings, making it challenging to comprehensively map the excitatory and inhibitory influences that determine a neuron's response across this high-dimensional stimulus space. Resolving this challenge is central to deciphering how neuronal populations transform complex sensory input into meaningful representations.

To overcome this challenge, we leveraged functional digital twin models trained on neuronal recordings from macaque primary (V1) and mid-level (V4) visual cortex, as well as mouse visual cortex. A *digital twin* of a brain area ([Walker et al., 2019](#); [Bashivan et al., 2019](#); [Franke et al., 2022](#); [Willeke et al., 2023](#)) can be constructed by training a deep neural network to learn the relationship between stimuli and the resulting neuronal responses. If trained with enough high-quality data, the model can be used to simulate the neuronal response to data never seen by the animal—allowing researchers to scale experiments in silico that would be difficult or impossible in vivo. This approach enabled us to perform a systematic characterization of neuronal selectivity across the full dynamic range of responses to naturalistic images.

We found that many neurons in visual cortex maintain non-zero baseline firing rates, enabling them to exhibit dual-feature selectivity: they respond strongly to preferred features while being systematically suppressed by distinct non-preferred features. We found this bidirectional selectivity not only in primate V1 and V4 but also in mouse V1 as well as lateral visual areas. In addition, we showed that the response function of these neurons reflects a continuum between the preferred and non-preferred features. That is, dual-feature selective neurons exhibit graded responses that vary continuously with perceptual stimulus similarity to their preferred and non-preferred features. By leveraging both excitatory and suppressive selectivity in single cells, this strategy may increase representational capacity by reducing the number of neurons that would be required to encode the equivalent number of features with unidirectional neurons (see section “Balancing sparsity and capacity? A hypothesized role for dual-feature selectivity” in Discussion).

Our results, particularly the observation that non-sparse neurons in visual cortex exhibit feature-selective suppression, align with recent connectomic studies in mice and flies showing specific inhibitory connectivity, where distinct interneuron types target defined excitatory populations (Schneider-Mizell et al., 2025 [↗](#); Matsliah et al., 2024 [↗](#)). This contrasts with blanket inhibition, in which inhibitory neurons provide non-selective input that broadly suppresses local excitatory neurons regardless of their feature selectivity (e.g. Heeger, 1992 [↗](#)). Such structured inhibitory connectivity suggests that inhibition contributes actively to normalization and selectivity within cortical circuits (see also Sebastian Seung, 2024 [↗](#)), thereby complementing the feature-detector principle through mechanisms operating under different circuit constraints to achieve efficient and interpretable neuronal codes.

Results

A continuum of response sparseness in macaque V1 and V4

We recorded spiking activity from macaque areas V1 and V4 using linear silicon probes while animals viewed a diverse set of naturalistic images. To align the visual stimulus with the recorded neurons’ receptive fields, we first mapped the population receptive field using sparse noise stimuli while the monkey was fixating on the center of the screen. Probe insertions were targeted orthogonal to the cortical surface, such that neurons sampled along the probe depth share overlapping receptive fields, allowing a single stimulus configuration to adequately drive the entire recorded population. For V1 recordings, the monkey maintained central fixation and we positioned grayscale stimuli (6.7×6.7 degrees of visual angle) such that we obtained full receptive field coverage (Fig. 1a [↗](#) left). For V4, we located the fixation dot that centered the population receptive field on the display (Fig. 1a [↗](#) right) and showed full-field RGB images (30×16.8 degrees of visual angle). Stimulus sizes were chosen to broadly cover receptive fields across the sampled neural population, including neurons with some spatial scatter around the probe axis (see section “Technical challenges and limitations” in Discussion).

Monkeys were trained to maintain fixation during 2.1-second trials while we presented 15 images per trial. Each session included between 7, 500 and 15, 000 unique naturalistic images, as well as a smaller set of repeated images to assess response reliability. The V1 dataset was obtained from a previously published study (Cadena et al., 2023 [↗](#)), while we collected new data from V4 for this work. After spike sorting, we analyzed data from 453 V1 neurons ($n = 2$ animals) and 394 V4 neurons ($n = 2$ animals).

To characterize each neuron’s stimulus-response function, we trained functional *digital twin* models (Walker et al., 2019 [↗](#); Bashivan et al., 2019 [↗](#); Franke et al., 2022 [↗](#); Willeke et al., 2023 [↗](#)). Each model combined a shared convolutional neural network (CNN) core pretrained on image classification (LeCun et al., 2015 [↗](#)) with neuron-specific readout layers (Fig. 1b [↗](#)). Our modeling approach leveraged prior findings (e.g. Cadena et al., 2023 [↗](#)) that different stages in the visual hierarchy align with different deep neural network layers—early layers with primary visual cortex and deeper layers with intermediate and higher-order areas. Accordingly, we used the first layer of a ConvNeXt architecture for V1 neurons, fine-tuning the convolutional core on the

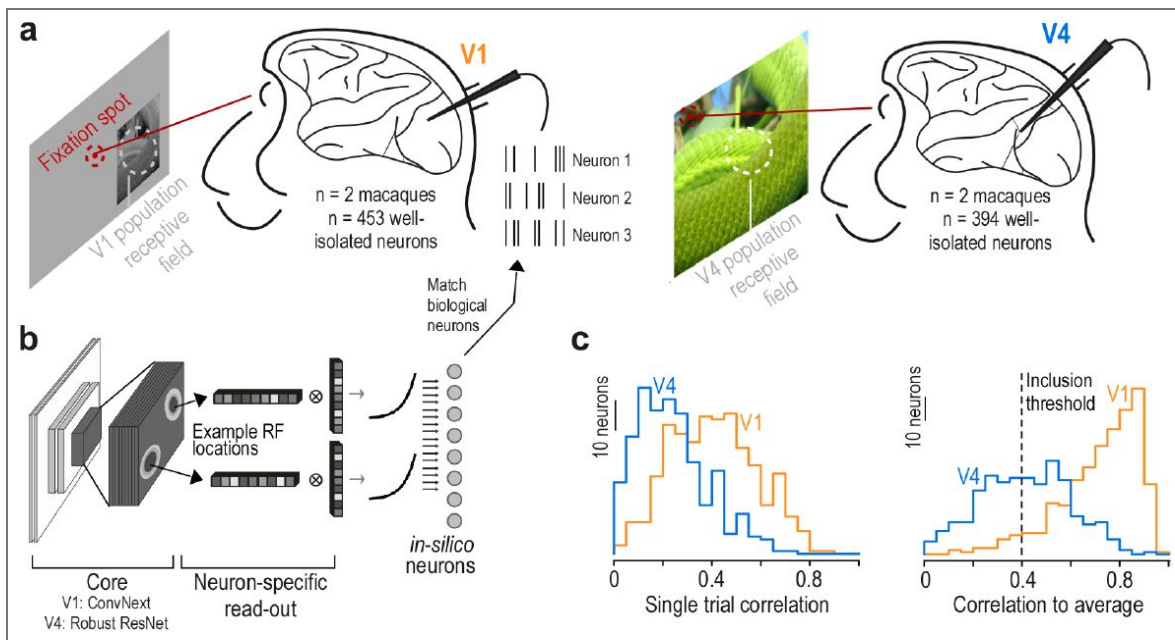


Fig. 1. Experimental approach.

a: Experimental recordings from macaque visual areas V1 (453 neurons from 2 monkeys) and V4 (394 neurons from 2 monkeys) during presentation of natural images. Monkeys fixated on the screen (fixation spot shown in red). The stimulus was centered on the population receptive field of neurons recorded that day, indicated by the white circle. **b:** Architecture of functional "digital twin" models, consisting of a core shared across neurons (V1: first layer of Conv Next; V4: layer 3 of Robust ResNet50) and a neuron-specific readout. This design creates *in-silico* neurons that model the response properties of individual biological neurons, with example receptive field locations of two neurons illustrated. **c:** Prediction accuracy of the model on test images not used during training, for both V1 (orange, $n = 75$ test images) and V4 (blue, $n = 150$ test images). Left panel shows distribution of single trial correlations, and right panel displays correlation to the average response across stimulus repeats. The dotted line indicates an inclusion threshold of 0.4. For further analysis, we only included neurons with correlation to average above that threshold ($n = 443$ (98% of total) neurons for V1, $n = 205$ (52% of total) neurons for V4).

recorded responses and learning neuron-specific readouts (see also Fu et al., 2024). For V4 neurons, we employed a pretrained ResNet50 model as a fixed core and trained a linear readout for each neuron on top of layer 3 of this representation (Willeke et al., 2023).

Model performance was evaluated on test images not used during training ($n = 75$ test images for V1, $n = 150$ test images for V4), using the correlation between predicted and observed responses averaged across stimulus repetitions (Fig. 1c). For further analysis, we included only neurons with correlation-to-average above 0.4, yielding high-confidence *digital twins* for 98% of V1 neurons ($n = 443$) and 52% of V4 neurons ($n = 205$). The lower proportion of high-confidence in-silico neurons in V4 likely reflects the greater complexity of V4 tuning compared to V1, as well as missing contextual information such as image surrounds and sequential image context—factors we discuss in detail below (see section “Technical challenges and limitations” in Discussion). The *digital twins* models enabled us to probe neuronal selectivity across a stimulus space orders of magnitude larger than possible in vivo, where recording time limits the number of stimuli that can be presented per neuron.

To assess response selectivity across natural images, we quantified each neuron’s lifetime sparsity (Willmore and Tolhurst, 2001), which measures how selectively a neuron responds—indicating whether it is activated by many stimuli or only by a few. Using the *digital twin* models, we predicted responses to 1.2 million natural images, generated activation curves by sorting predicted responses in ascending order, and computed their skewness (Fig. 2a,b). Higher skewness reflects one-sided asymmetry in the data and thus indicates sparser response profiles with few stimuli eliciting strong activity. Predicted activity was expressed relative to each neuron’s baseline, defined as the mean response during the 300 ms fixation window immediately preceding stimulus onset, when a uniform gray screen was presented.

Neurons in both V1 and V4 exhibited activation curves spanning a continuum from highly sparse to non-sparse (Fig. 2d). Skewness values derived from model predictions correlated strongly with those obtained from *in vivo* responses to repeated test images (Fig. 2c), confirming that our *digital twins* accurately capture neuronal stimulus selectivity. Baseline firing rates correlated negatively with response skewness, indicating that neurons with higher baseline activity tended to be less sparse (Fig. 2e). Consistent with this relationship, baseline activity to a gray screen correlated positively with the median response across natural images (Fig. 2f), showing that neurons with elevated baseline firing exhibit response distributions centered near baseline—fluctuating both above and below it. Together, these relationships demonstrate that non-sparse neurons operate in a more continuous, bidirectional modulation regime, while sparse neurons act as highly selective feature detectors. Lifetime sparsity distributions were similar across V1 and V4 (Fig. 2d), with both areas containing substantial proportions of non-sparse neurons (see also Willmore et al., 2011; Rust and DiCarlo, 2012). Model performance correlated only weakly with response skewness, indicating that our inclusion criterion did not systematically exclude sparse neurons from further analysis (Suppl. Fig. 1).

To facilitate subsequent analyses, we used a model-derived skewness of 2.0 to separate neurons with graded, non-sparse responses from those responding strongly to only a small subset of stimuli (Fig. 2d). We note that the underlying distribution of sparsity is continuous, consistent with recent findings (Gondur et al., 2025), and this threshold is adopted purely for analytical convenience to focus subsequent analyses on neurons with sufficiently graded response distributions; the key findings reported below are not dependent on the exact threshold chosen.

Identification of most and least activating stimuli of non-sparse macaque V1 and V4 neurons

Traditional approaches to visualizing neuronal selectivity have often focused on the most activating stimuli, providing valuable insights into preferred features. Examining weak or suppressive responses, however, can reveal additional aspects of selectivity, particularly in non-sparse neurons that respond in a graded fashion to most stimuli (see section “The role of suppression in visual cortical tuning: relating our findings to existing work” in Discussion). To test

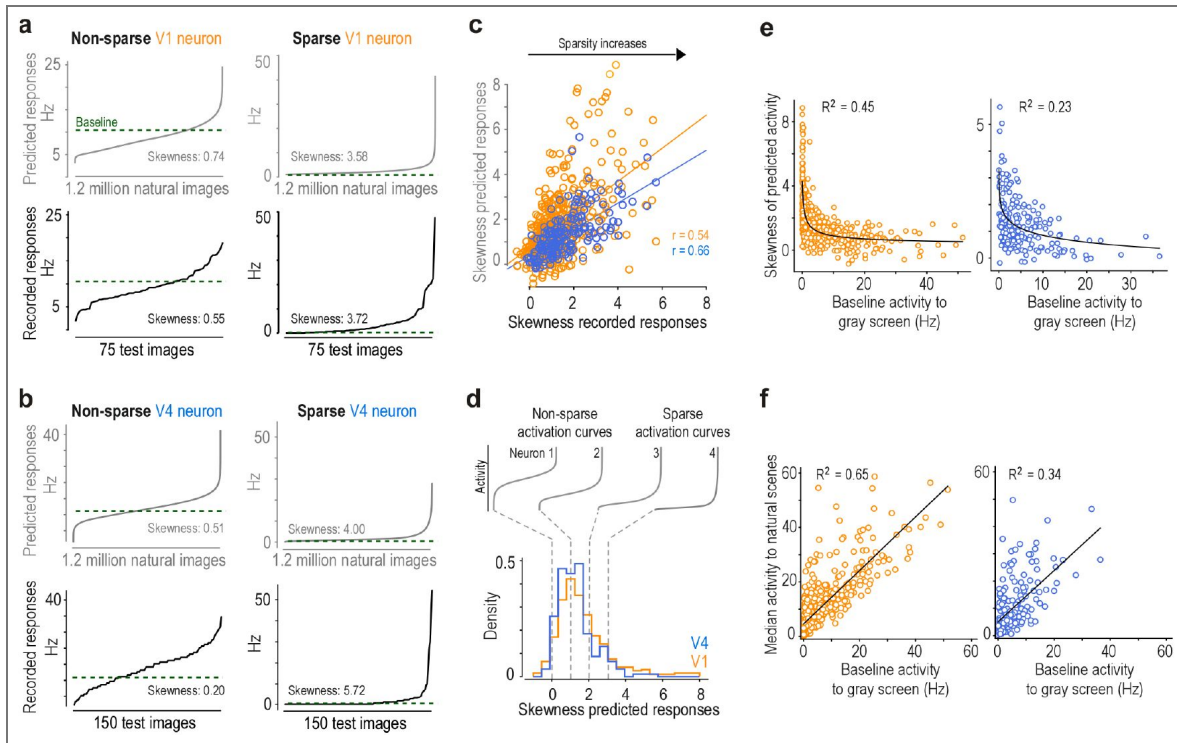


Fig. 2. Continuum of single neuron response sparseness in early and mid-level macaque visual cortex.

a, Response profiles of non-sparse (left) and sparse (right) V1 neurons. Curves display neuronal activity sorted from lowest to highest response, derived from model predictions across 1.2 million ImageNet images (gray, top row) and recorded responses to 75 test images (black, bottom row), averaged over stimulus repeats. We used responses to test images to obtain mean stimulus driven activity and average out stimulus-unrelated signals. Skewness values quantify lifetime sparsity, with higher values indicating neurons that respond selectively to fewer stimuli while remaining silent to most others. Green dotted lines indicate recorded baseline firing rate (Hz) during grey screen presentation prior to stimulus onset. **b**, Comparable response profiles for representative V4 neurons, demonstrating similar variations in lifetime sparsity in this higher-order visual area, with some neurons showing broadly tuned responses and others exhibiting highly selective activation patterns. **c**, Correlation analysis between prediction-based and recording-based skewness values for qualifying V1 (orange, $n = 443$) and V4 (blue, $n = 205$) neurons. The strong correlation ($r = 0.54$ for V1 and $r = 0.66$ for V4, both $p < 0.001$) validates that our model accurately captures intrinsic sparsity characteristics across natural scenes, supporting the use of model-predicted responses for systematic in silico analyses across large image datasets. **d**, Population-level distribution of lifetime sparsity across V1 (orange) and V4 (blue) neuronal populations (V1: $n = 443$ neurons, V4: $n = 205$ neurons), revealing a continuous spectrum rather than discrete categories. Representative activation curves above illustrate how response profiles change along this continuum. Neurons with skewness below 2.0 are defined as non-sparse, though this threshold represents a point along a gradual transition. This distribution highlights the functional diversity within each visual area. **e**, Baseline firing rate extracted from a 300 ms fixation window before stimulus onset plotted versus skewness of predicted responses. V1 (orange): $n = 443$ neurons, V4 (blue): $n = 205$ neurons. R^2 from exponential fit. Neurons with low baseline firing rates exhibit variable skewness that likely reflects the prevalence of their preferred features in natural scenes—rare features produce highly skewed responses while common features yield more symmetric distributions. **f**, Median predicted activity to natural scenes plotted versus baseline firing rate extracted from a 300 ms fixation window before stimulus onset. V1 (orange): $n = 443$ neurons, V4 (blue): $n = 205$ neurons. R^2 from linear regression.

whether such weak responses exhibit systematic structure, we characterized both ends of the activation spectrum by identifying stimuli that maximally and minimally activated each neuron using two complementary methods: gradient-based image synthesis and large-scale image screening. Because weak responses are most informative in neurons with graded activity, the following analyses focused specifically on non-sparse neurons, unless noted otherwise.

Based on previous work (Walker et al., 2019 [↗](#); Bashivan et al., 2019 [↗](#); Willeke et al., 2023 [↗](#)), the image synthesis approach used gradient ascent in the *digital twin* models to generate images that maximized model-predicted activity. Here, we extended this approach to also minimize neuronal responses. These synthetic stimuli—termed most exciting inputs (MEIs) and least exciting inputs (LEIs)—emerged through iterative modification of noise images to achieve the desired neuronal responses (Fig. 3a [↗](#)). The screening approach computed model activations over more than one million naturalistic images, ranking responses to identify the most activating images (MAIs) and least activating images (LAIs) for each neuron (Fig. 3b [↗](#)). Crucially, all images—both synthesized and screened—were normalized to identical ℓ_2 norms within each neuron's receptive field, reducing the influence of overall image energy on response differences. We note, however, that L2 normalization controls for root-mean-squared contrast but does not fully equate effective contrast in nonlinear cells, whose responses depend on the spatial structure of the stimulus beyond its total energy. Residual contrast-dependent effects, particularly in the suppressive regime, cannot be entirely excluded.

In V1, this approach revealed clear structure at both ends of the response spectrum. MEIs and MAIs consistently featured oriented edges and grating-like stimuli spanning various spatial frequencies (Fig. 4 [↗](#), Suppl. Fig. 2 [↗](#)), confirming well-established tuning properties of macaque V1 (e.g. Hubel and Wiesel, 1962 [↗](#); Schiller et al., 1976 [↗](#); Fu et al., 2024 [↗](#)). Remarkably, LEIs and LAIs exhibited systematic, feature-specific structure: suppression arose not only from orthogonal orientations but also from shifts in orientation, spatial frequency, phase, and texture structure, revealing inhibitory axes that extend beyond classic cross-orientation inhibition (e.g. Morrone et al., 1982 [↗](#), and see section “The role of suppression in visual cortical tuning: relating our findings to existing work” in Discussion). As a control analysis, we simulated idealized simple and complex cells and screened for their least and most activating images (Suppl. Fig. 4 [↗](#)). Simple cells showed phase-shifted versions of preferred stimuli as their least-activating images, while complex cells exhibited no coherent patterns in their low-activation regime, with diverse unrelated stimuli eliciting uniformly weak responses—reflecting their pooling mechanisms and nonlinear characteristics. Therefore, the empirical structure we observed—where least activating images differ from the most activating along specific features other than phase—cannot be accounted for by classical simple or complex cell models.

Non-sparse V4 neurons showed similarly structured selectivity at the activation extremes, but for more complex features. Aligning with previous work, their MEIs and MAIs revealed elaborate patterns—including curved contours, textured surfaces, and distinct color combinations (e.g. Willeke et al., 2023 [↗](#); Desimone and Schein, 1987 [↗](#); Pasupathy and Connor, 2002 [↗](#); Yamane et al., 2008 [↗](#))—as well as novel motifs such as eye-like configurations and branching structures (Fig. 5 [↗](#), Suppl. Fig. 3 [↗](#)). Their corresponding LEIs and LAIs depicted equally coherent feature configurations—alternative contour arrangements, different color combinations, or contrasting texture patterns that consistently suppressed neuronal responses.

In addition, visual inspection revealed that, in both V1 and V4, the images that most strongly activated a given neuron tended to be perceptually similar to one another, as did the images that elicited the weakest responses. For example, for many V1 neurons, the set of MAIs often shared the same edge orientation but differed in position within the receptive field—consistent with phase invariance in complex cells. In several V4 neurons, the MAIs typically preserved a common global texture or shape, while varying in attributes such as texture phase or color.

To quantify this, we examined the organization of MAIs and LAIs within a perceptual similarity space. We used DreamSim, a model of perceptual similarity fine-tuned to align with human visual judgments (Fu et al., 2023a [↗](#)), and embedded all naturalistic images in this high-dimensional space. For each neuron, we assessed the internal coherence of its preferred (MAIs) and non-

Fig. 3. Two complementary methods to study neuronal selectivity: optimization-based feature visualization and large-scale image screening.

a. Schematic of the feature visualization procedure, in which a starting noise image is iteratively optimized using gradient ascent on a neural predictive model to either maximize or minimize the activity of a single neuron. This yields the most exciting input (MEI) or the least exciting input (LEI), respectively. The process begins with a noise image and updates pixel values to achieve the target activation level. Example shows an LEI for a V4 neuron. **b.** Schematic of the image screening procedure, in which a large dataset of 1.2 million ImageNet images is used to probe neuronal responses. Each image elicits a predicted neuronal response from the model, allowing construction of a response profile across the dataset. Sorting these responses identifies the most activating (MAI) and least activating (LAI) images for each neuron.

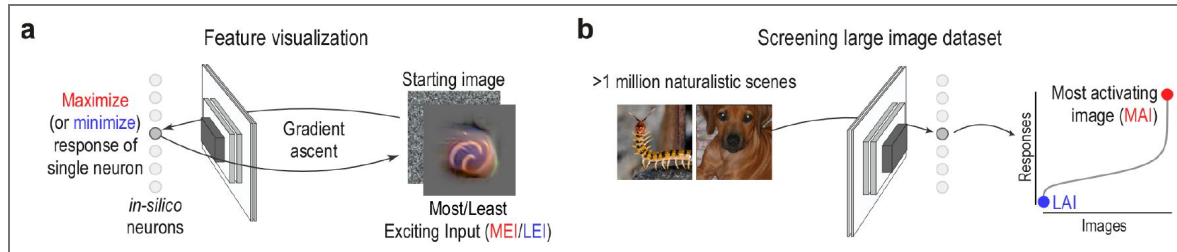


Fig. 4. Identification of most and least activating stimuli of macaque V1 neurons.

Least (left) and most (right) activating inputs for five example V1 neurons. For each neuron, top row show optimized images starting from different initialization (i.e. noise) seeds (LEIs on the left, MEIs on the right) and bottom row shows the most and least activating images identified through screening 1.2 million ImageNet images (LAIs on the left, MAIs on the right). Images are 2.3×2.3 degrees visual angle, with each neuron's receptive field located in the center of the image.

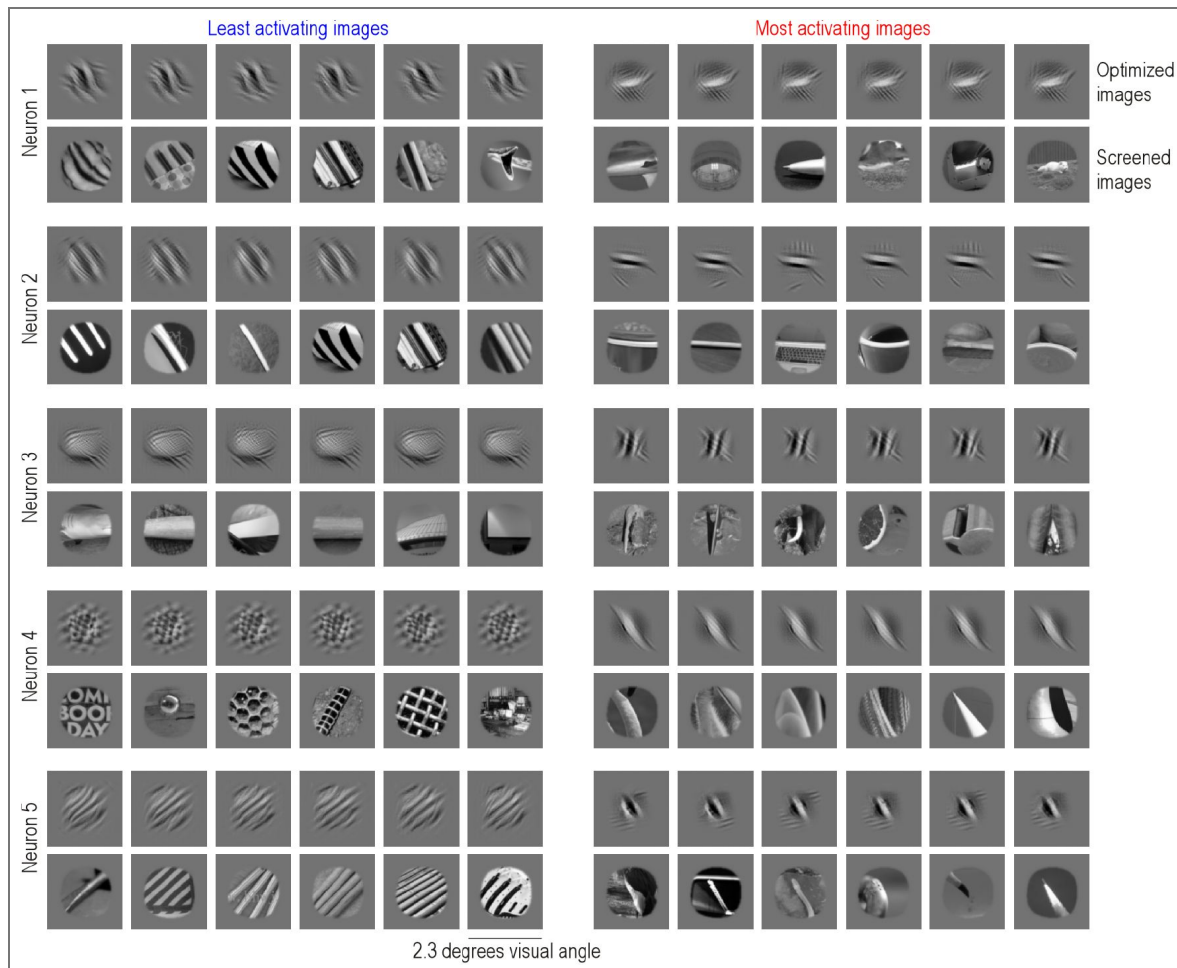




Fig. 5. Identification of most and least activating stimuli of macaque V4 neurons.

Least (left) and most (right) activating inputs for five example V4 neurons. For each neuron, top row show optimized images starting from different initialization (i.e. noise) seeds (LEIs on the left, MEIs on the right) and bottom row shows the most and least activating images identified through screening 1.2 million ImageNet images (LAIs on the left, MAIs on the right). Images are 14.81×14.81 degrees visual angle, with each neuron’s receptive field located in the center of the image.

preferred (LAIs) image sets by calculating pairwise cosine similarities within the MAIs and separately within the LAIs. We used representations from the penultimate layer of DreamSim, where distances are aligned with human similarity judgments (Fuet al., 2023a). As a baseline, we computed the similarity of the MAIs and the LAIs to randomly selected sets of naturalistic images (Fig. 6a). We found that, for each neuron, the MAIs were more similar to one another than to random images, and the same held for the LAIs (Fig. 6b), resulting in significantly positive d' -prime values (Fig. 6c)—a measure of how well each image set could be discriminated from random images. Notably, d' -prime values were similar for MAIs and LAIs, demonstrating that low-activating stimuli were just as structured and perceptually coherent as high-activating ones.

Verification of most and least activating images of V1 and V4 neurons

Having identified stimuli predicted to generate maximal and minimal neuronal responses, we next validated these model predictions through multiple approaches. Previous studies have confirmed that model-predicted MEIs reliably drive strong responses in both mouse and macaque visual cortex in vivo (Walker et al., 2019; Bashivan et al., 2019; Willeke et al., 2023; Fu et al., 2024). In contrast, whether model-predicted least-activating stimuli suppress neuronal activity remains largely unexplored experimentally (but see Unk, 2025).

For each neuron, we identified the single test-set image predicted by the model to produce the highest response, as well as the one predicted to produce the lowest response, and examined where these images ranked within the neuron's recorded response distribution over the test set (Fig. 7). For non-sparse V1 and V4 neurons, our models accurately identified both extremes: the image predicted to elicit the strongest response corresponded to a high recorded response, while the predicted least-activating image was associated with responses below baseline, indicating that activity in these neurons was modulated both above and below baseline (Fig. 7a–c, top panels). In contrast, for sparse V1 and V4 neurons, our models reliably predicted the strongest response but performed poorly in identifying weakly activating stimuli (Fig. 7a–c, bottom panels). Because sparse neurons exhibit near-baseline responses for most stimuli, the lower tail of their response distributions has minimal dynamic range. Consequently, the model cannot accurately predict a truly low-activating stimulus, and the recorded responses to predicted least-activating images are broadly distributed across the response range, yielding a nearly uniform distribution.

These results confirm that, for non-sparse neurons, the *digital twin* models capture meaningful stimulus selectivity across both high- and low-activity ranges relative to base-line (Fig. 7a). Although the distribution of response skewness is continuous, a threshold of $skewness = 2$ appears to provide a reasonable separation between neurons with graded, bidirectional modulation around an elevated base-line and those with sparse, predominantly unidirectional responses to a small subset of stimuli (Fig. 2d). For subsequent analyses of least-activating images, we therefore focused on non-sparse neurons with response skewness >2 .

Having established that the model predictions align with neuronal responses recorded in vivo, we next asked whether these identified most- and least-activating images reflect genuine neuronal tuning properties rather than artifacts of a particular model implementation. To test this, we evaluated both optimized (i.e. MEIs and LEIs) and screened stimuli (i.e. MAIs and LAIs) using independently trained models with different initializations or architectures. For V1 neurons, we trained an evaluator model with the same ConvNeXt architecture as the original generator model but initialized with independent weights (Fig. 8a). The core was fine-tuned on the same dataset, and neuron-specific readouts were trained from scratch. For V4 neurons, we used a completely different architecture: an attention-based convolutional network trained end-to-end to predict neuronal responses (Pierzchlewicz et al., 2023) (Fig. 8d).

Each evaluator model assessed the MEIs, LEIs, MAIs, and LAIs identified by the generator models. Responses were contextualized against a reference set of 200, 000 naturalistic images that were masked to each neuron's receptive field and contrast-matched to optimized and screened images. For each neuron, we computed response percentiles—reflecting the proportion of reference

Fig. 6. High and low activity reflect structured and perceptually coherent feature combinations.

a, Schematic illustrating the computation of image similarities. All naturalistic rendered images were embedded into DreamSim, a perceptual similarity space fine-tuned on human judgments. Within this high-dimensional space, we computed cosine similarity among the top 10 most activating images (MAIs) and among the least activating images (LAIs), as well as their similarity to random images. **b**, Distributions of cosine similarity among MAIs (top 10 images, red) and between MAIs and random images (gray). These distributions were used to compute discriminability using the d-prime metric. **c**, d-prime values for MAIs and LAIs across all non-sparse V1 and V4 neurons. Gray bars indicate a control condition comparing similarities between random image sets. Across the population, d-prime values for both MAIs and LAIs were significantly higher in V1 and V4 compared to random image sets (two-sample t-test, $p < 0.001$). After applying false discovery rate correction, all V4 neurons ($n = 168$) retained p -values below 0.05, indicating that the likelihood of observing such discriminability by chance was consistently low. In the larger population, $n = 293$ out of $n = 315$ neurons showed significant p -values for MAIs, and $n = 281$ neurons showed significant p -values for LAIs.

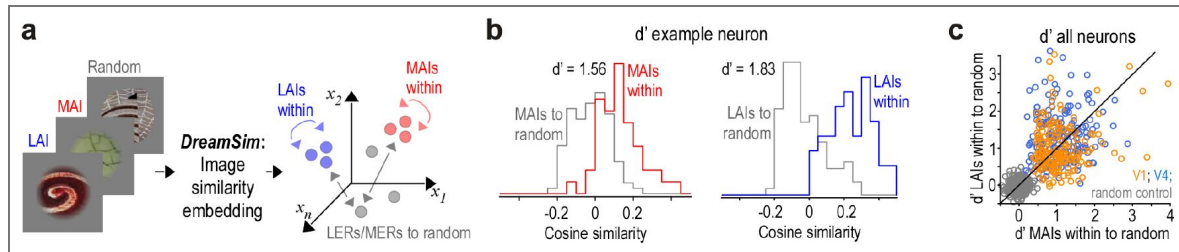
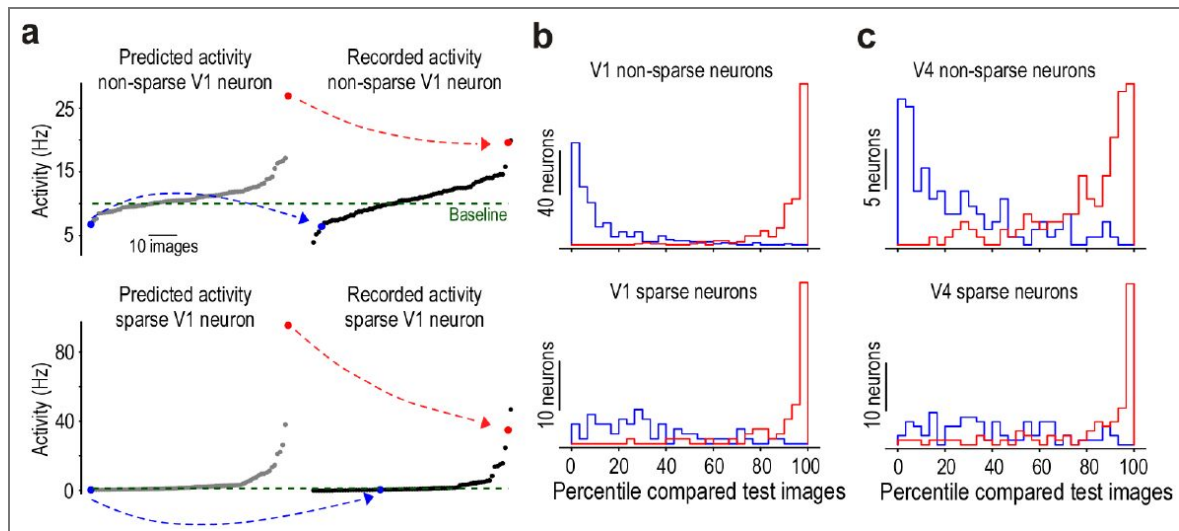


Fig. 7. Model predictions accurately identify extreme stimuli in recorded neuronal responses from V1 and V4.

a, Left: Model-predicted responses to 75 test ImageNet images for an example non-sparse (top) and sparse (bottom) V1 neuron. The predicted least and most activating test images are marked in blue and red, respectively. Right: Actual recorded responses of the same neurons to the same test images, averaged across stimulus repeats and sorted by recorded response magnitude. The blue and red dots indicate the images that the model predicted to elicit the lowest and highest responses. **b**, Distribution of response percentiles in recorded data for the predicted least (blue) and most (red) activating images, shown separately for non-sparse (top) and sparse (bottom) V1 neurons. Note that a non-selective ordering would be expected to yield a uniform distribution, similar to the blue distribution observed here. **c**, Same as (b), but for non-sparse and sparse neurons in area V4.



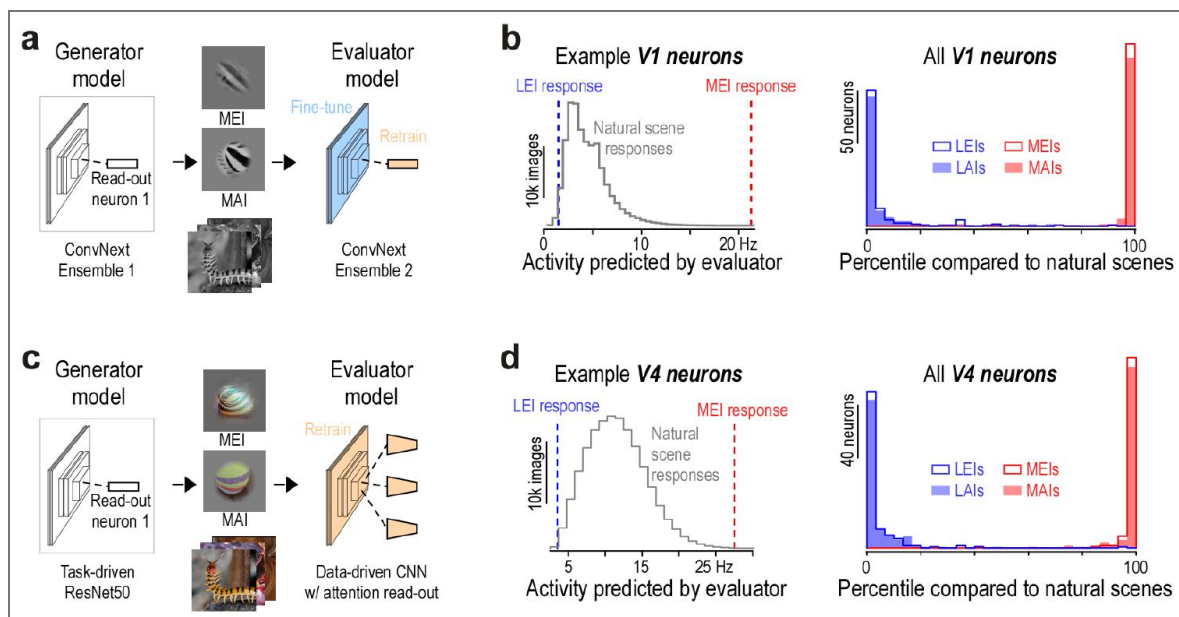


Fig. 8. Independent evaluator models confirm the identification of optimal stimuli for V1 and V4 neuronal responses.

a. Schematic of the verification pipeline. Least and most activating images for V1 neurons were identified using a generator model through both optimization and screening of images. These images were then passed to an independent evaluator model, comprising a separate model ensemble with the same core architecture but a readout trained from scratch. **b.** Example V1 neuron. Left: Distribution of predicted responses to size- and contrast-matched natural images, with the predicted responses to the MEI and LEI (identified by the generator model) highlighted by the evaluator model. Right: Distribution of response percentiles for MEIs, LEIs, MAIs, and LAIs as predicted by the evaluator model, relative to the distribution of natural image responses. MAIs and LAIs were identified based on 200k naturalistic rendered images. **c,d.** As in (a,b), but for V4 neurons. The evaluator model in this case had a distinct architecture and training objective and was trained from scratch on the neural data.

images eliciting weaker or stronger model-predicted responses than the identified least- or most-activating images (Fig. 8b,e, left panels). Across both visual areas, MEIs and MAIs consistently ranked in the top percentiles, while LEIs and LAIs fell reliably in the lowest percentiles (Fig. 8b,e, right panels).

These results confirm that for non-sparse neurons, *digital twin* models accurately predict both most- and least-activating images, with stimulus selectivity being robust across independently trained models and architectures, indicating they reflect true neuronal selectivity.

Feature similarity to most and least activating stimuli jointly shapes non-sparse neuronal responses

Having established that non-sparse neurons are both activated and suppressed around baseline by distinct features, we next asked how these opposing selectivity poles shape their tuning to natural images. Specifically, we sought to understand whether defining both a most- and a least-activating feature imposes structure on how neuronal activity varies across the natural image manifold—that is, whether responses to natural stimuli reflect graded encoding along a continuum between these two poles.

To investigate how neuronal activity varies within the high-dimensional perceptual similarity space, we first constructed a low-dimensional embedding defined by each neuron's most and least activating features. This allowed us to visualize and quantify how responses scale with perceptual similarity to these features, using a space defined independently of neuronal activity. Specifically, we used DreamSim (Fu et al., 2023a) to embed 200,000 naturalistic images into a 2-dimensional similarity space for each neuron (Fig. 9a), where each image was assigned coordinates based on its similarity to the MAI (x-axis) and LAI (y-axis; Fig. 9b). In the following, we focus on V4 neurons because DreamSim captures mid-level visual features like color and texture, which are better represented in V4 than in V1. This alignment is important for relating DreamSim's similarity space to neuronal activity, as both need to encode similar features. Results for macaque V1 neurons are shown in Suppl. Fig. 5.

Visualizing predicted V4 responses across this space revealed that many non-sparse neurons exhibited smooth response gradients along the diagonal, extending from high MAI similarity / low LAI similarity to the opposite corner (Fig. 9b). This gradient indicates that neuronal activity increased with similarity to the MAI and decreased with similarity to the LAI. In contrast, most sparse neurons exhibited less structure, with responses primarily varying with similarity to the MAI (Fig. 9c).

We performed regression to quantify how well the 2-dimensional space explained V4 neuronal activity. The resulting explained variance R^2 was significantly higher for non-sparse neurons (mean=0.23, std=0.12) than for sparse neurons (mean=0.13, std=0.08; Fig. 9d), indicating that the similarity space captures a substantial portion of response variance in the non-sparse population. For V1 neurons, we observed the same pattern of higher explained variance for non-sparse compared to sparse neurons (Suppl. Fig. 5), but overall the explained variance was lower than in V4, despite the *digital twin* model achieving higher prediction accuracy. This is likely because the Dream Sim space is not well aligned with low-level features represented in V1 neurons.

When we repeated the analysis for V4 neurons using randomly selected images instead of MAIs and LAIs (Fig. 9e), the R^2 values were significantly reduced to 0.04 and 0.02 for non-sparse and sparse neurons, respectively (Fig. 9f). This confirms that the observed response gradients depend specifically on the most and least activating images, rather than arising from properties of the similarity space per se. Additionally, for non-sparse neurons, replacing either the MAI or the LAI with random images significantly reduced predictive performance, indicating that both extremes contributed meaningfully to response variation. In contrast, for sparse neurons, only the most activating images carried predictive value; responses remained largely unchanged when the LAIs were replaced by random images.

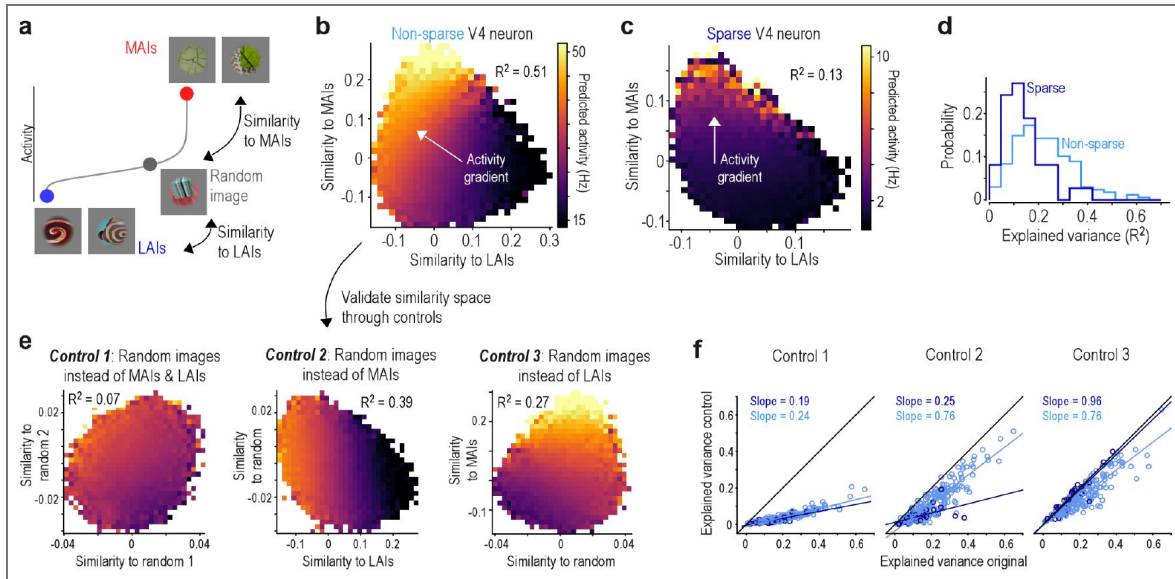


Fig. 9. Responses of V4 neurons vary continuously depending on preferred and non-preferred stimuli.

a, Schematic illustrating construction of the 2D image similarity space for each neuron. Each of ~200k naturalistic rendered images was assigned *x*- and *y*-coordinates based on its cosine similarity (computed using DreamSim) to the neuron’s least activating image (LAI, *x*-coordinate) and most activating image (MAI, *y*-coordinate), respectively. MAIs and LAIs were defined from the same image set. **b**, Example 2D similarity space for a non-sparse V4 neuron. Bins are color-coded by the mean predicted neuronal activity of the images within each bin. The arrow denotes the principal activity gradient, and R^2 indicates the variance explained by a linear fit. The color bar spans the 0.1 to 99.9th percentile of neuronal responses across all images. **c**, Same as (b), but for a sparse V4 neuron. Here, the activity gradient aligns primarily along the *y*-axis, indicating stronger modulation by similarity to the MAI than to the LAI. **d**, Variance explained (R^2) by linear regression predicting neuronal activity from the 2D similarity space (as in b,c), shown separately for sparse (dark blue) and non-sparse (light blue) V4 neurons. **e**, Results of three control analyses validating the similarity space for the V4 neuron shown in panel b. Control 1 replaces both MAIs and LAIs with random images. Control 2 replaces MAIs with random images but retains LAIs. Control 3 replaces LAIs with random images but retains MAIs. **f**, Variance explained (R^2) by linear regression for the original analysis (similarity to MAIs and LAIs on *x*- and *y*-axes) and for Control 1 (left), Control 2 (middle), and Control 3 (right). Sparse (dark blue) and non-sparse (light blue) neurons are shown separately. Lines indicate linear regression fits, with slopes reported in each panel.

Together, these findings suggest that the responses of non-sparse neurons are shaped by similarity to both preferred and distinct non-preferred features. Despite the complexity of natural stimuli, V4 responses could be well approximated within a low-dimensional subspace defined by these two feature types, with both contributing meaningfully to response variation.

A. Shared feature selectivity across the neuronal population

We found that, across neurons within a given visual area, the features that most strongly activate one neuron can resemble those that least activate another (Fig. 10a). Additionally, stimuli that most strongly activate one neuron can also strongly activate others, even when their least-activating images differ (Fig. 10c). These patterns suggest that neurons across the population exhibit shared feature selectivity—that is, they are tuned to a common set of features in perceptual space, which can excite some neurons while suppressing others.

To test this prediction, we examined how each neuron's MAIs and LAIs affected the entire population within the same cortical area (Fig. 10b,d). MAIs and LAIs of one neuron showed high probabilities of driving either strong or weak responses in other neurons, while intermediate responses were less common. MAIs elicited right-skewed distributions—with many neurons strongly activated, relatively few weakly activated. LAIs produced bimodal distributions with increased likelihood of both suppression and excitation across the population. This contrasted sharply with randomly selected control images that evoked uniform distributions. This organization extended across individual animals: MAIs and LAIs identified from neurons recorded in one monkey elicited similar activation patterns in neurons recorded from another monkey (Fig. 10e-g).

These results support the hypothesis that visual stimuli driving strong or weak responses in individual neurons also modulate responses across the population—exciting some neurons while suppressing others. This pattern reflects a population-level organization of dual-feature selectivity that generalizes across animals.

DUAL-feature selectivity is present in the mouse visual cortex

To assess whether dual-feature selectivity constitutes a general principle of visual coding across mammals, we extended our analyses to mouse visual cortex. While macaque and mouse visual cortex differ substantially in their functional organization and the complexity of neuronal selectivity (e.g. Fu et al., 2024), we asked whether the broader principle—that non-sparse neurons are jointly defined by distinct excitatory and suppressive feature sets—generalizes across mammalian visual systems. Using Neuropixels probes, we recorded spiking activity from neurons in V1 and two lateral visual areas (lateromedial area (LM) and laterointermediate area (LI)) while head-fixed mice viewed grayscale natural images (Fig. 11a). Prior work has shown that LM and LI share functional similarities with the ventral stream in primates (Wang et al., 2012) and are involved in object representation in mice (Froudarakis et al., 2021). Using the recorded data, we trained a *digital twin* model using the Sensorium competition model architecture (Willeke et al., 2022): specifically, a convolutional core shared across neurons combined with neuron-specific readouts, trained end-to-end to predict neuronal responses to natural images (Fig. 11b). For further analysis, we only used neurons with a correlation between predicted responses and mean test image responses larger than 0.4. Model performance was comparable to that obtained in macaque V1, yielding high-confidence digital twins for 92% of neurons across areas, supporting the validity of the subsequent analyses.

Consistent with the macaque data, neurons across all three mouse visual areas exhibited a continuum of life-time sparsity with substantial non-sparse populations, as measured by the skewness of their predicted responses to 200, 000 naturalistic images (Fig. 11c,d). The images were masked around the population receptive field and contrast-matched to ensure consistent contrast across images. The skewness of the predicted responses correlated strongly with the skewness of the recorded test responses (Fig. 11e-g). Importantly, skewness was negatively

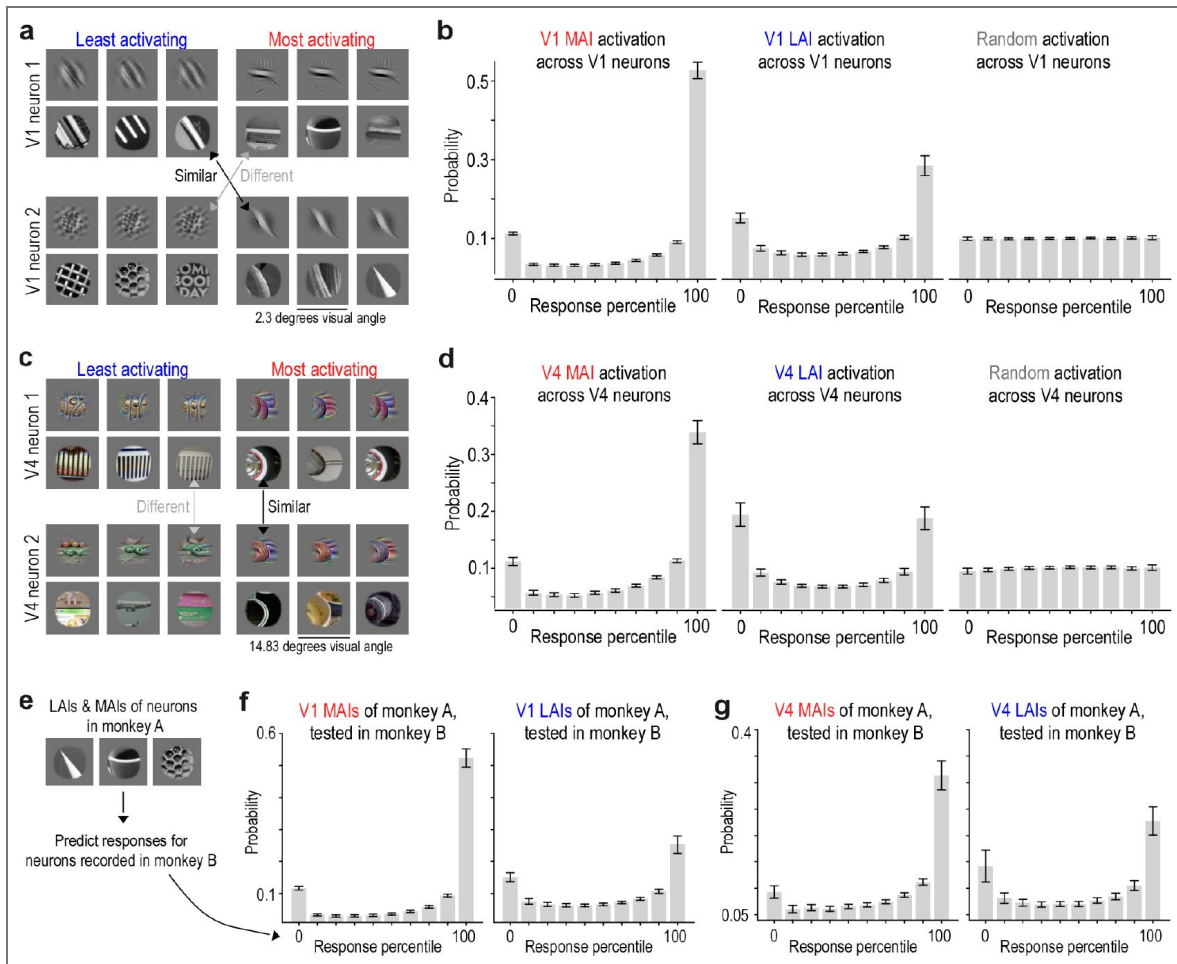


Fig. 10. Most and least activating stimuli reveal distributed feature tuning across neuronal populations.

a, Examples of least activating images (left) and most activating images (right) for two V1 neurons. Arrows indicate the perceptual similarity between MAIs and LAIs shared across both neurons. **b**, Distribution of response percentiles evoked by MAIs, LAIs, and randomly sampled natural images across the V1 neuronal population. The left histogram shows the probability that one neuron’s MAI will elicit a specific response percentile in other neurons. Response percentiles were calculated from each neuron’s responses to 1.2 million ImageNet images. Data represent population means with 99% confidence intervals. **c,d**, Same analyses as in (a,b) applied to neurons recorded in visual area V4. **e**, Schematic illustrating how we used the MAIs and LAIs identified from neurons recorded in one monkey to predict the responses of neurons recorded in another monkey. The resulting response percentiles indicate how these images ranked compared to the predicted responses to 1.2 million ImageNet images, as shown in panels (f,g). **f**, Cross-animal generalization: probability that MAIs (left) and LAIs (right) from V1 neurons in monkey A will evoke specific response percentiles in V1 neurons from monkey B. **g**, Same cross-animal analysis as in (e), applied to V4 neurons.

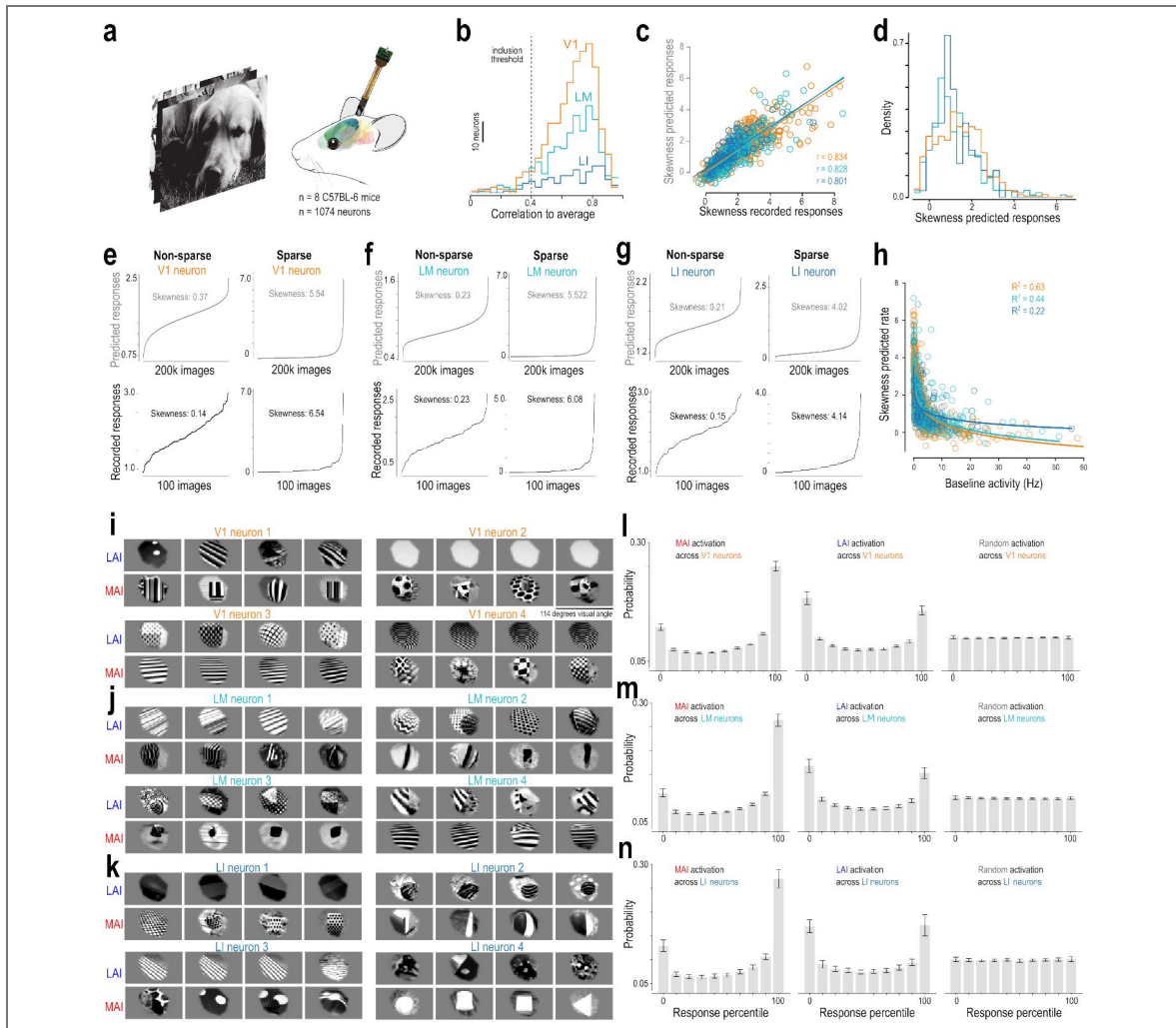


Fig. 11. Dual-feature selectivity in the mouse visual cortex.

a, Experimental recordings from mouse visual areas V1 (598 neurons), LM (350 neurons), and LI (126 neurons) from 8 C57BL/6 mice during presentation of natural images. Mice were head-fixed on a treadmill and passively viewing the stimuli presented on a screen in front of them. **b**, Prediction accuracy of the model on test images not used during training, for V1 (orange), LM (cyan), and LI (blue) as the correlation to the average response across stimulus repeats. The dotted line indicates an inclusion threshold of 0.4. For further analysis, we only included neurons with a correlation to average above that threshold ($n = 561$ neurons for V1, $n = 325$ neurons for LM, $n = 113$ neurons for LI). **c**, Correlation analysis between prediction-based and recording-based skewness values for qualifying V1 (orange, $n = 561$), LM (cyan, $n = 325$), and LI (blue, $n = 113$) neurons. Similar to the primates, there is strong correlation ($r = 0.834$ for V1, $r = 0.828$ for LM, and $r = 0.801$ for LI, all $p < 0.001$) indicating that the model accurately captures intrinsic sparsity characteristics across natural scenes. **d**, Population-level distribution of lifetime sparsity across V1 (orange), LM (cyan), and LI (blue) neuronal populations, revealing a continuous spectrum rather than discrete categories. Neurons with skewness below 2.0 are defined as non-sparse. **e**, Response profiles of non-sparse (left) and sparse (right) V1 neurons. Curves display neuronal activity sorted from lowest to highest response, derived from model predictions across 200,000 ImageNet images (gray, top row) and recorded responses to 100 test images (black, bottom row), averaged over stimulus repeats. Skewness values quantify lifetime sparsity, as in Fig. 2a,b. **f,g**, Same as (e) for LM neurons (f) and LI neurons (g). **h**, Baseline firing rate extracted from a 200 ms window before stimulus onset plotted versus skewness of predicted responses. V1 (orange, $n = 561$), LM (cyan, $n = 325$), LI (blue, $n = 113$) neurons. R^2 from exponential fit. **i**, Least (blue, LAI) and most (red, MAI) activating images for 4 example non-sparse V1 neurons. LAI and MAI are identified through screening of 200,000 ImageNet images. Images are 84×114 degrees visual angle. **j,k** Same as (i) for LM neurons (j) and LI neurons (k). **l**, Distribution of response percentiles evoked by MAIs, LAIs, and random images across the V1 neuronal population. The left histogram shows the probability that one neuron's MAI will elicit a specific response percentile in other neurons. Response percentiles were calculated from each neuron's responses to 200,000 ImageNet images. Data represent population means with 99% confidence intervals. **m,n** Same analysis as (l) for LM neurons (m) and LI neurons (n).

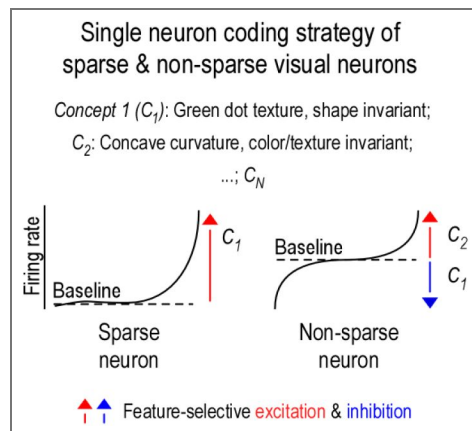


Fig. 12. Single neuron coding strategies in sparse and non-sparse visual neurons.

This schematic illustrates how sparse and non-sparse neurons exhibit selectivity for distinct *concepts*, where each concept (e.g., C_1 , C_2) represents a specific combination of latent visual features such as color, shape, and texture—dimensions known to be encoded in area V4. These concepts can be thought of as points in a high-dimensional perceptual space. **Left:** *Sparse neurons* respond selectively and strongly to a single concept (e.g., C_1 : green dot texture) and remain silent for most other stimuli. These neurons exhibit high lifetime sparseness and encode only a narrow portion of stimulus space. **Right:** *Non-sparse neurons* may exhibit *dual-feature selectivity*, characterized by excitation to one concept (e.g., C_2 : concave curvature) and suppression to another (e.g., C_1). These neurons tend to maintain non-zero baseline activity and modulate their firing rates bidirectionally, enabling graded responses to a broader range of stimuli. This bidirectional modulation likely reflects feature-selective excitation (red arrows) and inhibition (blue arrows), with firing rates encoding the similarity of each stimulus to both excitatory and suppressive features. Sparse and non-sparse neurons are intermingled within the same cortical area, forming a distributed population code in which different neurons anchor their selectivity to partially overlapping sets of concepts. This organizational motif—shared feature selectivity and bidirectional modulation—appears conserved across visual cortical areas, including earlier regions such as V1, though the underlying concepts differ with the feature space represented at each stage.

correlated with baseline firing rate during gray screen presentation (Fig.11h [↗](#)), suggesting that neurons that were non-sparse during image presentation tended to have higher spontaneous firing rates across areas.

For non-sparse neurons, we next identified LAIs and MAIs by screening a large-scale naturalistic image dataset (Fig.11i-k [↗](#)). As in the monkey, both MAIs and LAIs depicted coherent features. The MAIs of single neurons were perceptually coherent, often featuring specific orientations, textures, or particular visual patterns that drove strong excitation. Similarly, the LAIs showed coherent structure, but these patterns differed systematically from the MAIs along specific dimensions, suggesting that images eliciting minimal responses occupied distinct regions of feature space.

At the population level, the same principle as in the monkey emerged: LAIs and MAIs of a given neuron were more likely to strongly or weakly activate other neurons within the population, compared to randomly selected images which exhibited an approximately uniform likelihood of eliciting responses across the response range (Fig.11l-n [↗](#)).

Together, these findings demonstrate that dual-feature selectivity is present across mouse visual areas, with non-sparse neurons responding to distinct feature sets that drive excitation and suppression respectively. This suggests that the joint definition of neurons by both their preferred and suppressive stimulus features—rather than by excitatory tuning alone—may constitute a general principle of mammalian visual coding. We discuss the potential functional implications of this organization across species below.

Discussion

Visual cortex has long been understood to represent the visual world through neurons that increase their firing rates in response to specific visual features. Here, we identify an additional organizational principle of single neuron selectivity: many neurons exhibit dual-feature selectivity, responding not only to preferred features but also showing systematic suppression to distinct non-preferred features. This bidirectional selectivity enables individual neurons to continuously encode the contrast between two specific feature combinations in high-dimensional input space. Specifically, the neuron's firing rate is modulated in a graded fashion, with each rate reflecting the relative distance to the most and least activating stimuli, extending beyond simple feature detection. We find that this principle is present across species (macaque and mouse) and cortical areas (from primary visual cortex to higher-order regions), suggesting it may represent a common computational strategy for single neuron and population coding. These findings indicate that cortical representations integrate both specific excitation and suppression within individual neurons, potentially expanding representational capacity while maintaining interpretable single-neuron responses.

The role of suppression in visual cortical tuning: relating our findings to existing work

Inhibition plays a fundamental role in sensory processing, with extensive literature documenting its contributions to neuronal encoding (e.g. reviewed in [Isaacson and Scanziani, 2011 \[↗\]\(#\)](#)). Here, we focus specifically on visual cortex and suppression originating within the classical receptive field, excluding surround suppression. In primary visual cortex, a classic example is cross-orientation suppression, where a neuron's response to its preferred orientation is markedly reduced when an orthogonal (non-preferred) grating is presented simultaneously ([Bishop et al., 1973 \[↗\]\(#\)](#); [Morrone et al., 1982 \[↗\]\(#\)](#); [Allison et al., 1995 \[↗\]\(#\)](#); [Ferster, 1986 \[↗\]\(#\)](#); [Priebe, 2016 \[↗\]\(#\)](#); [Burr et al., 1981 \[↗\]\(#\)](#); [Hata et al., 1988 \[↗\]\(#\)](#)). Theoretical studies suggest that such inhibitory mechanisms sharpen orientation tuning by suppressing responses to non-optimal stimuli (e.g. [Ben-Yishai et al., 1995 \[↗\]\(#\)](#)). Our results reveal that suppression in visual cortex extends beyond well-characterized mechanisms such as cross-orientation inhibition, with single neurons being systematically suppressed by a diverse set of naturalistic features that are distinct from, and seemingly unrelated to (see also [Gondur et al., 2025 \[↗\]\(#\)](#)), their excitatory preferences. Using unbiased image synthesis and large-scale screening, we find that macaque V1 neurons' least activating stimuli are not limited to orthogonal

orientations, but comprise specific combinations of orientation, spatial frequency, phase, and size that systematically differ from the features driving maximal activation (Figs. 4 [↗](#), 5 [↗](#)). This aligns with previous work demonstrating that suppression extends beyond orthogonal orientations (Ringach et al., 2002 [↗](#); DeAngelis et al., 1992 [↗](#); Burg et al., 2021 [↗](#)) and to spatial frequencies outside a neuron's preferred range (Bauman and Bonds, 1991 [↗](#); De Valois and Tootell, 1983 [↗](#)), while going beyond what classical V1 models predict. For simple cells, phase-shifted stimuli are a well-established suppressive axis reflecting linear On-Off subfield structure, and for complex cells, phase pooling yields no coherent suppressive pattern—yet neither model class accounts for the multidimensional suppressive structure we observe (Suppl. Fig. 4 [↗](#)). Our unbiased approach reveals that this structure spans simultaneous changes across orientation, spatial frequency, phase, and texture, exceeding what any single known suppressive mechanism predicts.

Importantly, this structured suppression is conserved across species and visual cortical areas. In both mouse primary and lateral visual cortex and macaque V1 and V4, we observe similar feature-specific suppression. Notably, this does not imply that mouse and macaque visual cortex share similar functional organization or equivalent complexity of neuronal selectivity. Rather, within the representational regime of each area—whether mouse V1 or macaque V4—neurons are organized such that excitatory and suppressive feature sets are jointly structured and distinct, even as the specific features that drive or suppress the neurons differ substantially across species and areas.

Our results support previous evidence for diverse suppressive mechanisms in visual cortex, including cross-orientation inhibition in V2 (Rowekamp and Sharpee, 2017 [↗](#)), suppressive subfields in V4 receptive fields (Pollen et al., 2002 [↗](#)), suppression by non-optimal stimuli in inferior temporal cortex (Willmore et al., 2011 [↗](#); Rust and DiCarlo, 2012 [↗](#); Miller et al., 1993 [↗](#); Rolls and Tovee, 1995 [↗](#)), tuned-suppression to natural images (Tamura et al., 2004 [↗](#)), and biased competition among multiple stimuli within receptive fields, modulated by attention (Desimone and Duncan, 1995; Reynolds et al., 1999 [↗](#)). Our work, together with complementary evidence from macaque V4 recordings using multi-unit Utah arrays (Gondur et al., 2025 [↗](#))—which independently demonstrates that V4 neurons exhibit two-tailed response distributions with both preferred and anti-preferred stimuli that are seemingly unrelated in feature space—extends these prior studies by providing a more general framework for how suppression shapes neuronal tuning across the visual hierarchy and species. While prior work has established that inhibition can be structured and feature-selective (see above), our results suggest a broader organizing principle: within each visual area, there exists a set of feature combinations—appropriate to the area's level of abstraction—from which individual neurons draw both their excitatory and suppressive preferences. Rather than being idiosyncratic, these preferences are shared across neurons and animals within the same area, such that the features exciting one neuron frequently suppress another. This organization structures the single-neuron code along axes that are meaningful at the population level, giving rise to graded, bidirectional modulation around baseline that reflects a stimulus's position along feature dimensions collectively relevant to the area's computational role.

A directly analogous organizational principle has been described in face-selective cortex. Chang & Tsao demonstrated that face identity is encoded through a set of approximately 50 shared axes spanning a low-dimensional face space, with individual neurons responding as linear projections onto these axes (Chang and Tsao, 2017 [↗](#)). Rather than representing idiosyncratic faces, face-selective neurons collectively tile a shared manifold in which each axis captures a specific facial dimension. Our results suggest that a related principle may extend beyond the specialized domain of face perception to general visual cortex: each non-sparse neuron anchors a selectivity axis in natural image space—defined by its most and least activating stimuli—and responds in a graded, approximately linear fashion as images vary along this axis (Fig. 9 [↗](#)). Where Chang & Tsao characterized axes within the constrained geometry of face images, a central question for future work is whether the selectivity axes of V4 neurons—spanning arbitrary combinations of color, texture, curvature, and object structure—likewise decompose into a compact set of shared primitives reused across the neuronal population. A related geometric organization has been

observed in artificial vision systems, where self-supervised networks such as DINO exhibit population-level “concept axes” with antipodal poles encoding opposing semantics (Fel et al., 2025). Our results reveal a distinct principle operating at the single-neuron level: biological axes link structured, non-opponent features within a continuous representational manifold.

Balancing sparsity and capacity? A hypothesized role for dual-feature selectivity

Dual-feature selectivity arises when neurons exhibit bidirectional modulation to two distinct stimulus features, enabling a coding regime that may balance interpretability with representational capacity. Unlike classical feature detectors, which respond sparsely to a narrow set of stimuli (Barlow, 1961; Field, 1987; Olshausen and Field, 1996), dual-feature selective neurons respond to both excitatory and suppressive inputs, potentially spanning a broader dynamic range and likely supporting higher information-theoretic capacity.

Similar geometric organization has been observed in artificial vision systems, where self-supervised networks such as DINO exhibit population-level “concept axes” with antipodal poles encoding opposing semantics (Fel et al., 2025). Our results reveal a related but distinct principle operating at the single-neuron level: biological axes link structured, non-opponent features within a continuous representational manifold, extending the notion of bidirectional coding to the geometry of individual neurons.

By producing distinct firing rates for excitatory, neutral, and suppressive stimuli, such neurons increase the diversity of their response distributions. Mutual information between stimulus and response increases when the response distribution is both diverse—i.e., has high entropy—and reliable—i.e., exhibits low variability given the same stimulus. Formally, mutual information is defined as $I(x; r) = H(r) - H(r | x)$, where $H(r)$ is the entropy of the response distribution and $H(r | x)$ is the conditional entropy reflecting noise or ambiguity in the response. Neurons with bidirectional selectivity can increase $H(r)$ by utilizing a broader dynamic range—responding with different firing rates to excitatory, neutral, and suppressive stimuli—resulting in a more uniform and distributed response profile. If these response patterns are reliable across repeated presentations of the same stimulus (i.e., $H(r | x)$ remains low), then mutual information is increased. This strategy parallels mixed selectivity in higher cognitive areas, where neurons encode nonlinear combinations of features to support high-dimensional and flexible representations (Rigotti et al., 2013; Fusi et al., 2016).

The benefit of dual-feature coding, however, depends on the alignment between a neuron’s selectivity and the statistics of the input space. If excitatory and suppressive features occur frequently and independently, the neuron can fully exploit its dynamic range, distributing responses more uniformly and maximizing entropy. If the features are strongly correlated or rarely occur in natural scenes, responses become concentrated in a narrow range, limiting entropy and diminishing the coding benefit. Thus, the information-theoretic capacity of a dual-feature neuron is tightly constrained by how well its selectivity structure aligns with the distribution of stimuli it encounters.

Neuronal codes with high entropy—where individual neurons respond to many stimuli with variable firing rates—can enhance representational capacity but come with the trade-off of increased metabolic cost. This is because generating action potentials and driving postsynaptic glutamatergic currents are among the most energy-demanding processes in the brain (Levy and Baxter, 1996; Attwell and Laughlin, 2001). Yet, structured population responses may help offset this cost. We find that the same stimulus can increase firing in some neurons while suppressing others below baseline, which may reduce net activity and maintain population sparseness. The resulting response variance across the population accounts for a previous finding, where stimuli optimized to elicit high activity in one area simultaneously suppress certain neurons, thereby expanding the dynamic range (Tong et al., 2023).

Importantly, single neuron lifetime sparseness and population sparseness, while often correlated (e.g., Froudarakis et al., 2014 [↗](#)), capture distinct coding properties. As argued by Willmore and Tolhurst (2001) [↗](#); Willmore et al. (2011) [↗](#), population sparseness can remain high even when individual neurons are broadly tuned—so long as different neurons are active for different stimuli. Therefore, dual-feature selective neurons with high firing rates could be organized as a population to maintain balanced population sparseness, ensuring a broad dynamic range while minimizing metabolic costs.

Technical challenges and limitations

This study has several important limitations that merit careful consideration. First and foremost, our conclusions rely partially on in silico analyses employing *digital twin* models rather than direct experimental validation. Although digital twin models confer substantial advantages that direct analysis of experimentally recorded responses cannot match—enabling screening of more than one million images per neuron in silico, gradient-based synthesis of stimuli precisely optimized to drive or suppress individual neurons, and cross-model verification of identified selectivity patterns (Fig. 8 [↗](#))—the reliance on model-predicted rather than experimentally measured responses requires careful justification. We maintain confidence in our findings for several reasons. Prior research has consistently demonstrated that digital twin approaches can successfully predict neuronal responses to novel stimuli and identify maximally exciting inputs subsequently verified in vivo (Walker et al., 2019 [↗](#); Bashivan et al., 2019 [↗](#); Franke et al., 2022 [↗](#); Willeke et al., 2023 [↗](#); Fu et al., 2024 [↗](#); Tong et al., 2023 [↗](#)). We restricted analyses to neurons exhibiting high predictive accuracy on held-out test data, ensuring models faithfully captured neuronal response functions. We further confirmed that model accuracy remained consistent across both high- and low-activity regimes relative to baseline—though exclusively for non-sparse neurons (cf. Fig. 7 [↗](#)). Critically, cross-model verification confirms that identified stimuli reflect genuine neuronal tuning rather than model artifacts—a test that has no analog when working with fixed experimental image sets (cf. Fig. 8 [↗](#)).

A related consideration concerns receptive field coverage and its contribution to performance differences between V1 and V4. Although we targeted orthogonal probe insertions and centered stimuli on the population receptive field prior to recording, partial receptive field drive for individual neurons cannot be excluded. Given the larger and more variable receptive fields in V4, this may have contributed to lower model performance there. However, other factors are likely more consequential: we cropped images for computational tractability, potentially losing contextual information that may modulate V4 responses; our models did not account for sequential image context during training; and more recent backbone architectures such as DI-NOv2 (Oquab et al., 2023 [↗](#)) would likely improve predictivity. Future work moving towards continuous dynamic models and state-of-the-art architectures (e.g. Willeke et al., 2025 [↗](#)) will likely close this performance gap, particularly for higher visual areas where temporal context and largescale image statistics play a more prominent role. Crucially, the goal here was not to maximize predictive performance per se, but to identify response patterns—dual-feature selectivity—that are robust across neurons, areas, and species, and our restriction to high-confidence neurons provides resilience against these limitations.

A further consideration concerns generalizability across individual animals. Both V1 and V4 data were collected from 2 macaques. For V4, digital twin models were fit independently per neuron without sharing information across animals. For V1, while models shared a common backbone fine-tuned on recorded data, readout layers remained neuron-specific, ensuring individual response functions are learned from each neuron's own data. Critically, extreme image sets identified by the model elicited correspondingly extreme responses in neurons from another animal, confirming that identified selectivity patterns are not idiosyncratic to individual subjects. Equivalent results in mouse visual cortex. Together, these results suggest that dual-feature selectivity is a robust and general property of mammalian visual cortex, rather than an artifact of individual animals, recording sessions, or modeling choices.

Another important consideration concerns the scope of our characterization. While we demonstrate dual-feature selectivity in spiking responses, our analyses do not determine the underlying circuit mechanisms, for example whether the observed suppression arises from direct inhibitory synaptic input mediated by specific interneuron subtypes. Elucidating these pathways will prove essential for understanding how feature-selective inhibition integrates into the broader cortical processing hierarchy. In this context, emerging functional connectomics datasets (e.g. Ding et al., 2025 [DOI](#)) promise invaluable in-sights by enabling researchers to bridge functional response profiles—including the dual-feature selectivity described here—with precise anatomical connectivity maps and cell-type-specific wiring motifs.

A further methodological consideration is that DreamSim was trained on human perceptual similarity judgments, while our neuronal data are from macaques. This cross-species application is supported by the deep homology between primate ventral visual streams: the anatomical and functional organization of V4 and the ventral pathway is broadly conserved between macaques and humans, and natural-image similarity judgments have been found to be highly consistent across the two species (e.g. ?). Importantly, we deploy DreamSim not as a model of macaque perception but as an image-feature embedding to test whether stimuli that cluster in perceptual space evoke similar neuronal responses—a use that is robust to the precise calibration of the metric, as confirmed by control analyses (Fig. 9e,f [DOI](#)). Nevertheless, a macaque-specific embedding would provide stronger grounding. We plan to address this point directly in our ongoing work: we are developing custom embeddings trained to align with macaque response geometry, using contrastive learning on MAI/LAI pairs from our own recordings, which will allow us to replace Dream-Sim with a representation grounded in the neural data it-self.

Finally, we have not yet characterized the relationship between the most and least activating stimuli for individual neurons. Although our analyses identify images that elicit strong activation or suppression, we have not quantified the underlying tuning features—such as color, shape, texture, or their combinations—that give rise to these responses. Quantifying tuning along these feature dimensions will be essential for understanding how excitatory and suppressive stimuli are related. Qualitatively, our data suggest that there is no clear or consistent relationship between the two stimulus types; for example, we do not observe a systematic opponency in orientation or other visual dimensions. This observation aligns with Unk (2025), who similarly found that preferred and anti-preferred features appear to be independently distributed across neurons, with no apparent systematic relationship between the two feature types. More broadly, the concept of “opponency” itself requires clarification: does it reflect contrast along a single feature axis—such as high versus low spatial frequency—or coordinated shifts across multiple features, such as orientation, phase, and texture? Future work that integrates quantitative feature-tuning models with the functional selectivity profiles established here will be critical for determining whether excitatory and suppressive stimuli occupy distinct or systematically related regions within high-dimensional visual feature spaces.

Normalization revisited: A role for feature-selective inhibition

Inhibitory cells in the brain exhibit a diversity of cell types comparable to that of excitatory neurons (Tasic et al., 2018 [DOI](#); Gouwens et al., 2019 [DOI](#); Yao et al., 2021 [DOI](#)), despite their much lower density (Keller et al., 2018 [DOI](#)). This raises the question: why is such diversity necessary? In the retina, a classic test bed for central brain research, the vast variety of inhibitory amacrine cells provides selective drive to distinct retinal ganglion cell types (e.g. Diamond, 2017 [DOI](#); Matsumoto et al., 2025 [DOI](#)), shaping feature encoding and supporting parallel processing streams. Similarly, in the cortex, accumulating evidence indicates that inhibitory neurons target excitatory neurons with high specificity (Muñoz et al., 2017 [DOI](#); Lu et al., 2017 [DOI](#); Wu et al., 2023 [DOI](#)).

A recent millimetre-scale volumetric EM reconstruction of mouse visual cortex (Schneider-Mizell et al., 2025 [DOI](#)) mapped the connectivity of over 1,300 neurons, revealing that inhibitory neurons organize into motif groups with widespread target specificity. These motifs coordinate inhibition

onto precise combinations of perisomatic and dendritic compartments across excitatory cell types, enabling compartment- and cell-type-specific modulation of cortical circuits far beyond broad class connectivity.

Likewise, dense connectomic reconstructions in the fly revealed that inhibitory interneurons, though few in number, constitute the majority of cell types and form highly specific, feature-selective connections to excitatory neurons (Matsliah et al., 2024 [↗](#); Sebastian Seung, 2024 [↗](#)). For example, individual inhibitory cell types provide suppression to targeted single excitatory cell types at defined spatial scales, suggesting a division of labor across interneuron types. Sebastian Seung (2024) [↗](#) argues that such diversity is an inevitable consequence of implementing numerous highly specific normalization operations, paralleling the architectural principles seen in convolutional networks (see below).

Our results, together with the anatomical evidence discussed above, suggest that inhibition in sensory cortex is far more structured than traditionally assumed. Classical models treat cortical selectivity as driven by excitation to preferred features, with inhibition providing non-specific gain control via untuned normalization pools that sum across diverse feature preferences (Heeger, 1992 [↗](#); Carandini and Heeger, 2011 [↗](#)).

In contrast, the dual-feature selectivity we describe may represent a biological implementation of specific rather than untuned normalization, where neurons apply suppressive filters to target specific anti-preferred features, rather than uniformly inhibiting all non-preferred inputs, thus enabling functions that extend beyond simple gain control. Indeed, Schwartz and Simoncelli (2001) [↗](#) demonstrated that non-uniform, feature-dependent normalization naturally emerges when the normalization weights are optimized to reduce statistical dependencies in natural signals, supporting the idea that such structured inhibition serves an efficient coding purpose. Related modeling work further shows that single-neuron activation functions—including those that permit negative responses—can themselves shape the circuit mechanisms and population geometries that emerge in recurrent networks (Tolmachev and Engel, 2025 [↗](#)), suggesting that biological nonlinearities may likewise play a central role in organizing structured inhibition.

Structured suppression and normalization are also integral to modern artificial neural networks, albeit through different mechanisms. For example, batch normalization and attention mechanisms modulate activity based on the relationships between inputs (e.g. Ioffe and Szegedy, 2015 [↗](#); Ba et al., 2016 [↗](#); Vaswani et al., 2017 [↗](#)). In our case, suppression is a fixed property of individual neurons: each neuron encodes the input by responding to preferred features and being suppressed by non-preferred ones. In contrast, artificial networks apply input-dependent modulation that is computed dynamically across inputs, such as through attention between tokens. Despite this difference, both artificial and biological networks appear to benefit from selectively attenuating certain features to improve representational efficiency.

Overall, our findings highlight feature-selective suppression as a general and underappreciated characteristic of cortical neurons. By defining specific feature dimensions through structured excitation and suppression, inhibition may enhance the representational capacity of individual neurons while preserving structured and interpretable response profiles that support flexible downstream readout—a principle that may be shared across sensory modalities and cognitive functions.

Materials and Methods

Ethics and Animal Care

Data were collected from three healthy male rhesus macaques with approval from Baylor College of Medicine's Institutional Animal Care and Use Committee (permit AN-4367). The monkeys were housed individually in a room with approximately ten other monkeys, allowing for rich social interactions on a 12-hour light/dark cycle. Regular veterinary care, balanced nutrition, and environmental enrichment were provided. All surgical procedures were performed under general anesthesia using aseptic techniques, with post-operative analgesics administered for seven days.

Electrophysiological recordings

Non-chronic recordings were conducted using a 32-channel linear silicon probe (NeuroNexus V1×32-Edge-10mm-60-177). Custom titanium recording chambers and head posts were surgically implanted. Prior to recordings, small trephinations (2 mm) were made over either (i) lateral V4, with eccentricities ranging from 1.7° to 18.3° of visual angle; or Medial V1, with eccentricities ranging from 1.4° to 3.0° of visual angle. A Narishige Microdrive (MO-97) and guide tube were used to carefully position the probes through the dura, taking care to minimize tissue compression.

Mice: Eight mice (*Mus musculus*: 4 male, 4 female) aged from 14 to 27 weeks were selected for experiments, with 2 females and 1 male expressing GCaMP6s in excitatory neurons via Slc17a7-Cre and Ai162 transgenic lines (stock nos. 023527 and 031562, respectively; The Jackson Laboratory) and the rest being C57BL/6J wildtype (stock no. 000664; The Jackson Laboratory). We performed acute recordings using Neuropixels probes 1.0 in awake, head-fixed mice according to (Jun et al., 2017 [↗](#)). In brief, animals were implanted with a headpost and habituated to the experimental setup (head fixation on a treadmill) after recovery. On the recording day, the animals were briefly anaesthetized with isoflourane and a 1mm craniotomy was made above visual cortex (approximately 2.9 mm lateral to the midline sagittal suture and anterior to the lambda suture) (Froudarakis et al., 2014 [↗](#)). The animals were then transferred to the experimental setup and allowed to recover from anaesthesia. Location of probe insertion was chosen according to stereotaxic coordinates for targeting V1, LM, and LI using Pinpoint (Birman et al., 2023 [↗](#)), with all penetrations ranging from 600–1100 μm on the anteroposterior axis, 2900–3500 μm on the mediolateral axis, and at an angle of 55° or 60° with respect to the ventrodorsal axis. One probe was smoothly lowered through the craniotomy to the final depth according to the trajectory planning with Pinpoint (Birman et al., 2023 [↗](#)) to cover the whole cortex (covering 1800–2000 μm of the probe) and allowed to settle for approximately 20 minutes before any recording.

Data collection and processing

Electrophysiological data was recorded as a broadband signal (0.5Hz–16kHz) and digitized at 24 bits. For spike sorting, the 32-channel array was divided into 14 groups of six adjacent channels. Spikes were detected when signals exceeded five times the standard deviation of noise. Principal component analysis was used for feature extraction, and a Kalman filter mixture model tracked waveform drift. Single-unit isolation was manually verified by assessing stability, refractory periods, and principal component plots.

For mice, neuronal activity recordings were made with custom-written software in LabView and then automatically spike sorted with the Kilosort3 spike sorting software (Pachitariu et al., 2023 [↗](#)). Neurons automatically classified as “single units” and that showed reliable firing during visual stimuli presentation were used for the model, in total: 598 V1, 350 LM, and 126 LI neurons from 20 recording sessions.

Visual stimulation

Stimuli were displayed on a 23.8” LCD monitor (100 Hz refresh rate, 1920 × 1080 resolution) positioned 100 cm from the subjects (\approx 63 pixels/degree). A camera-based eye tracking system verified that monkeys maintained fixation within \approx 0.95° of a small red fixation target. After maintaining fixation for 300 ms, visual stimuli were presented. Successful fixation throughout the trial resulted in a juice reward. For mice, natural images were presented 15 cm away from the left eye with a 23.8” LCD monitor (100 Hz refresh rate, 1920 × 1080 resolution). We positioned the monitor so that it was centered on and perpendicular to the surface of the eye at the closest point, corresponding to a visual angle of 2.2°/cm on the monitor.

Receptive field mapping & stimulus placement

Receptive fields were mapped at the beginning of each session using a sparse random dot stimulus. A single dot (0.12–1° in size) was displayed on a gray background, changing position and color (black or white) every 30 ms during two-second fixation trials. Multi-unit receptive field

profiles were obtained through reverse correlation, and population receptive fields were estimated by fitting a 2D Gaussian to the spike-triggered average.

For V1 recordings, the fixation spot was kept at the center of the screen, and grayscale natural image stimuli (6.7° in size) were centered at the mean receptive field location. The remainder of the screen was kept gray. For V4 recordings, color natural image stimuli covered the entire screen, and the fixation spot was positioned to place the mean receptive field as close as possible to the screen's center. Due to recording site locations, this typically placed the fixation spot near the upper left border of the screen.

For mice: Visual area segmentation was performed by mapping the reversals of the retinotopy based on the RF progression along the probe as described previously (Tafa-zoli et al., 2017 [↗](#)).

Stimulus selection

We selected 24, 075 images from 964 ImageNet categories, cropped to 420 × 420 pixels with 8-bit intensity resolution. From this set, 75 images were designated as the test set, 20% of the remaining images as the validation set, and the rest (19, 200 images) as the training set. The same image sets were used for both V1 and V4 recordings, with some differences in preprocessing (see below).

For a subset of V4 experiments, the dataset was augmented with rendered scenes, resulting in an equal mix of ImageNet and rendered images in the training, validation, and test sets. This synthetic dataset of rendered 3D scenes was created using Kubric (Greff et al., 2022 [↗](#)) and Blender (Blender Foundation, 2024 [↗](#)). We first manually created 10 primitive 3D objects in Blender, including basic geometric shapes (spheres, cubes, cylinders, cones, pyramids, etc.) and simple composite forms. Each object was exported as a Wavefront OBJ file with accompanying material (MTL) and texture coordinate information to ensure consistent UV mapping across all renders. For surface textures, we utilized the Describable Textures Dataset (Cimpoi et al., 2014 [↗](#)), which contains 47 texture categories with diverse visual properties ranging from regular patterns (e.g., striped, checkered) to stochastic textures (e.g., marbled, bubbly). We developed an automated rendering pipeline that generated 200, 000 unique scenes at 420 × 236 pixel resolution, matching the aspect ratio used in our V4 experimental stimuli. Each rendered scene consisted of a single 3D object with UV-mapped texture randomly sampled from the DTD dataset, placed against a background with a different randomly selected DTD texture. To ensure comprehensive sampling of the visual parameter space, we systematically varied multiple scene attributes: (1) object identity (uniformly sampled from the 10 primitive shapes), (2) object position (x, y, z coordinates sampled within the camera frustum), (3) orientation (random rotation quaternions), (4) scale, and (5) lighting conditions (directional light with varying intensity and angle). Shadow rendering was enabled to introduce natural occlusion patterns and enhance depth cues. This systematic variation ensured that our synthetic dataset captured a broad range of feature combinations while maintaining precise control over individual visual attributes.

Stimulus presentation

Each trial involved 2.1 seconds of continuous fixation, including 300 ms of gray screen at the beginning and 15 consecutive images displayed for 120 ms each without gaps. Test images were repeated 20-50 times throughout the session, while training and validation images were shown only once. The animal completed up to 1400 trials per day, resulting in up to 20000 stimulus-response pairs per neuron.

For V1 recordings, grayscale images were displayed at their original resolution covering a 6.7° visual angle, with the fixation spot at the center of the screen and the images centered at the mean receptive field location. The rest of the screen remained gray. Neural responses were analyzed within a 40-160 ms window following stimulus on-set. In contrast, for V4 recordings, full-color images were upscaled to match the screen width while maintaining their aspect ratio, with upper and lower bands cropped to fill the entire screen. The fixation spot was positioned to place the mean receptive field as close as possible to the screen's center, typically near the upper left border due to recording site locations. Spike counts for V4 were collected within a 60-160 ms window after stimulus onset.

For mice, 5, 100 natural images from ImageNet (ILSVRC2012) were cropped to fit a 16 : 9 monitor aspect ratio and converted to gray scale. To collect data for training a predictive model of the brain, we showed 5, 000 unique images as well as 100 additional images repeated 10 times each. This set of 100 images were shown in every recordings for evaluating cell response reliability within and between recordings. Each image was presented on the monitor for 500 ms followed by a blank screen lasting between 100 and 200 ms, sampled uniformly.

Image pre-processing

We employed two distinct image processing pipelines to prepare stimuli for model training and evaluation. In the V4 pipeline, starting with original images of 420×420 pixels at a resolution of 14 px° , we cropped the upper and lower bands to fit the full screen for presentation, resulting in images of 420×236 pixels. We then extracted only the bottom center 200×200 pixel region because of the location of the receptive fields (RFs) and subsequently downsampled these images to 100×100 pixels (corresponding to either 5.8 px° or 7 px°) for model training. For the V1 approach, we cropped the central 2.65° (167 pixels) of the original 420×420 image at its resolution of 63 px° and applied bicubic interpolation to downsampled to 93×93 .

Model architectures & training

For macaque V1 data, following [Fu et al. \(2024\)](#), we used a ConvNext-v2-tiny ([Woo et al., 2023](#)) architecture as the core, with original weights from the huggingface transformers library. After hyperparameter search, we selected the *stages-1-layers-0* as the optimal output layer. We applied a Gaussian read-out approach ([Lurz et al., 2022](#)) to transform the core feature maps into neuronal responses. This readout learns the coordinates of each neuron's receptive field center on the feature maps and implements a 2D isotropic Gaussian distribution to extract features from this location. During training, the readout samples positions according to this distribution, gradually focusing on the optimal receptive field location, while at inference time it uses the learned fixed positions. The extracted features were then processed through a neuron-specific affine projection with ELU non-linearity to predict the scalar neuronal activity.

For V1 model training, following ([Fu et al., 2024](#)), we minimized the Poisson loss between recorded and predicted neuronal activity. We first trained the readout for 20 epochs with frozen core weights, then reduced the learning rate from 0.001 to 0.0001 and optimized both the ConvNext core and readout weights using the AdamW optimizer ([Loshchilov and Hutter, 2017](#)) for 200 epochs. We trained an ensemble of five models with different random seeds and used their averaged predictions for all analyses.

Our neural predictive model of primate V4 consisted of a pretrained *core* computing nonlinear features from input images, and a *Gaussian readout* ([Lurz et al., 2022](#)) mapping these features to single neuron responses. Following ([Willeke et al., 2023](#)), we used an adversarially trained ResNet50 ([Salman et al., 2019](#)) as the core, with the first residual block of layer 3 (layer3.0) providing the feature maps, as this configuration yielded the highest predictive performance. The parameters of this pretrained network remained fixed during training. After batch normalization and ReLU activation, we obtained a nonlinear feature space shared across all neurons. For each V4 neuron, the Gaussian readout learned the receptive field center position on the output tensor to extract a feature vector. The readout implemented a 2D isotropic Gaussian distribution, sampling

locations during training and using fixed positions during inference. The extracted features were then processed through a linear-nonlinear model with L_1 -regularized weights and an ELU+1 nonlinearity (Clevert et al., 2015) to ensure positive responses.

For the V4 model training, as in (Willeke et al., 2023), we minimized the summed Poisson loss across neurons between observed and predicted spike counts, with added L_1 regularization on the readout parameters. During training, we zeroed gradients for neurons not shown a particular image. We used a batch size of 64, the Adam optimizer (Kingma and Ba, 2014) with an initial learning rate of 3×10^{-4} and momentum of 0.1. We implemented early stopping based on validation loss, decaying the learning rate by 0.3 after five epochs without improvement, and stopping after four such decay steps. We trained an ensemble of five models with different random seeds and used their averaged predictions for all analyses.

For the mouse model, we trained a digital twin model using the Sensorium competition model architecture (Willeke et al., 2022) with neural data recorded across V1 ($n = 598$ neurons), LM ($n = 350$ neurons), and LI ($n = 126$ neurons). Specifically, a convolutional core shared across all neurons from all areas combined with neuron-specific readouts, trained end-to-end to predict neuronal responses to natural images.

Model evaluation

To evaluate both V1 and V4 models, we measured performance using correlation to average (Franke et al., 2022; Cadena et al., 2023; Willeke et al., 2022) on held-out test images. This metric computes the correlation between model predictions and the average neuronal responses across repeated presentations of the same stimuli. By comparing predictions to trial-averaged responses rather than single-trial responses, this approach focuses on how well the models capture the stimulus-driven component of neural activity while accounting for biological trial-to-trial variability. Following (Fu et al., 2024) for V1 and (Willeke et al., 2023) for V4, we applied this evaluation consistently across both areas to enable fair comparison of model performance. Models trained exclusively on ImageNet generalized well to rendered images, showing no significant differences in prediction performance (data not shown). Consequently, we refer to both datasets collectively as ‘naturalistic images’ in the Results section. Detailed information regarding which dataset was used for each analysis is provided in the figure legends and the Methods section below.

Identification of most and least activating images

To identify naturalistic images that elicited extreme responses from modeled neurons, we conducted a large-scale screening across two complementary sources: 200,000 synthetically rendered scenes with controlled shape and texture variations, and 1,281,167 natural images from the ImageNet-1K training set (Deng et al., 2009). Prior to neuronal response prediction, all images underwent standardized preprocessing to match the experimental conditions. Images were center-masked using the mean receptive field profile computed from the population of synthesized most and least exciting images for V1 and V4 neurons, respectively. This masking procedure ensured that visual features were evaluated within the approximated retinotopic locations while controlling the influence of peripheral image regions. Additionally, we normalized all images to fixed ℓ_2 norms (12.0 for V1, 40.0 for V4) to control for overall contrast differences and ensure fair comparison across diverse image content. These normalization values were empirically determined to match the typical contrast range of natural images while preventing saturation artifacts. For response prediction, we employed area-specific models: the ConvNeXt-based architecture for V1 neurons and the adversarially-trained ResNet50 for V4 neurons, as described in the model architecture section. Each model computed predicted firing rates for all neurons across both image datasets, resulting in response matrices of $200,000 \times N$ and $1,281,167 \times N$ for rendered and natural images respectively, where N represents the number of recorded neurons. From these combined predictions, ranked in ascending order per neuron, we identified the top and bottom images for each neuron, corresponding to the most and least activating images (MAIs and LAIs). Similarly for mice, screening was conducted across the 200,000 synthetically

rendered scenes and images were normalized to a fixed ℓ_2 norm (10.0 for V1, LM, and LI). For response prediction we employed the one model for all three areas, as described in the model architecture section. The model computed predicted firing rates for all neurons in the dataset, resulting in a response matrix of $200,000 \times N$, where N represents the total number of recorded neurons recorded across V1, LM, and LI ($N = 1074$).

Optimization of most and least exciting images

We employed gradient-based optimization in the pixel space (for V1 neurons) or frequency (for V4 neurons) domain to generate images that optimally drove or suppressed individual neuronal responses (Fel et al., 2023). Starting from random noise images, we iteratively modified pixel values to either maximize (for most exciting inputs, MEIs) or minimize (for least exciting inputs, LEIs) the predicted neuronal response. For V4, this optimization operated exclusively on the phase spectrum of the Fourier transform while constraining the amplitude spectrum to match a fixed mean amplitude computed over 10,000 randomly sampled ImageNet images, ensuring that synthesized images maintained realistic spatial frequency content. For V1, the image synthesis was achieved by directly modifying the image pixel values themselves. During optimization, we applied ℓ_2 norm constraints matching those used in the screening procedure (12.0 for V1, 40.0 for V4) to maintain consistent contrast levels across all analyses. To improve optimization robustness and avoid local minima, we implemented a multi-crop augmentation strategy. This operation generated 4 random crops per image at each optimization step, with crop centers sampled from a Gaussian distribution ($\mu = 0.5$, $\sigma = 0.15$) relative to image dimensions. Importantly, we maintained a fixed box size of 1.0 throughout optimization, meaning crops spanned the full image extent with only positional jittering. This approach effectively provided the optimizer with multiple gradient estimates per iteration by evaluating slightly shifted versions of the image, similar to translation data augmentation but applied during the optimization process itself. The crop sizes included additional Gaussian noise ($\sigma = 0.05$) clamped between 0.05 and 1.0, though with our box size fixed at 1.0, this primarily introduced minor scale variations around the full image size. Each crop was resized back to the original dimensions (93×93 pixels for V1, 100×100 pixels for V4) using bilinear interpolation before neural response prediction. By averaging gradients across these multiple crops, the optimization procedure became more robust to small spatial shifts, encouraging the emergence of features that drove consistent neuronal responses. We performed 256 optimization steps using the Adam optimizer (Kingma and Ba, 2017) with a learning rate of 0.05, generating both MEIs and LEIs for each recorded neuron.

Verification of LAIs & MAIs and LEIs & MEIs

To verify that model-identified least and most activating stimuli accurately represented neuronal preferences, we implemented two complementary validation approaches. To validate predictions against *in vivo* recordings, we selected the single images from held-out test sets predicted to elicit the strongest and weakest responses for each neuron. We then determined where these images ranked within each neuron's recorded response distribution, allowing us to assess whether model predictions aligned with actual neuronal behavior.

To test generalization across models, we trained independent evaluator models. For V1 neurons, we used the same ConvNeXt architecture as our generator model but with independent initialization, fine-tuning the core network while training neuron-specific readouts from scratch. For V4 neurons, we employed an entirely different architecture—an attention-based convolutional network trained end-to-end. A distinct architecture was necessary here because the V4 generator model uses a fixed, pretrained ResNet50 backbone whose weights are deterministic: any re-trained model sharing this backbone would not constitute a genuinely independent evaluation. By contrast, for V1, the ConvNeXt core is fine-tuned from different random initializations, producing architecturally equivalent but computationally independent models. To evaluate stimulus effectiveness, we presented each evaluator model with the four synthetic images (MEI, LEI, MAI, LAI) derived from the generator model for each neuron. To contextualize response strength, we compared these against a reference set of 200,000 contrast- and size-matched naturalistic images.

For each synthetic image, we calculated its response percentile, defined as the proportion of reference images evoking a weaker model-predicted response. This cross-model validation allowed us to determine whether identified stimuli represented genuine neuronal tuning properties rather than model-specific artifacts. The choice of distinct architectures for V1 and V4 evaluator models reflects a difference in how each generator is trained: the V1 generator fine-tunes its ConvNeXt core, yielding genuine independence across random seeds, whereas the V4 generator relies on a fixed pretrained ResNet50 backbone shared identically across all instances. A truly independent V4 evaluator therefore required a fundamentally different architecture.

Baseline firing rate

To extract baseline firing rates, we counted spikes during the 300 ms fixation window prior to stimulus onset. We converted these counts into firing rates (Hz) and computed the mean baseline activity across all valid trials. During this fixation period, the screen was set to mean gray level (127 in 8-bit). For mice the baseline firing rates (Hz) were calculated during the blank screen (200ms duration) before the stimulus onset.

DreamSim image similarity

To test whether least and most activating rendered images (LAIs and MAIs) contain robust feature combinations and exhibit similarity within categories, we embedded all rendered images into the image similarity model DreamSim (Fu et al., 2023b [↗](#)). Dream-Sim leverages deep neural network representations to quantify perceptual image similarity in a manner that correlates with human judgments. We used the penultimate layer as an embedding for each image. For V1, we preprocessed images by masking them with the mean receptive field mask of V1 neurons, normalizing them to have the same contrast, and converting them to grayscale. For V4, we masked images with the V1 mask, normalized them, and preserved their color information. To quantify similarity, we computed pairwise cosine similarity between DreamSim vectors for the top and bottom 10 images per neuron (within-MAI and within-LAI comparisons). Additionally, we calculated pairwise distances between MAIs and LAIs, as well as distances from both MAIs and LAIs to 10 randomly selected images. This analysis resulted in cosine similarity distributions per neuron and condition. We used d-prime as a measure of discriminability, with values below 0.5 indicating low discriminability between image categories.

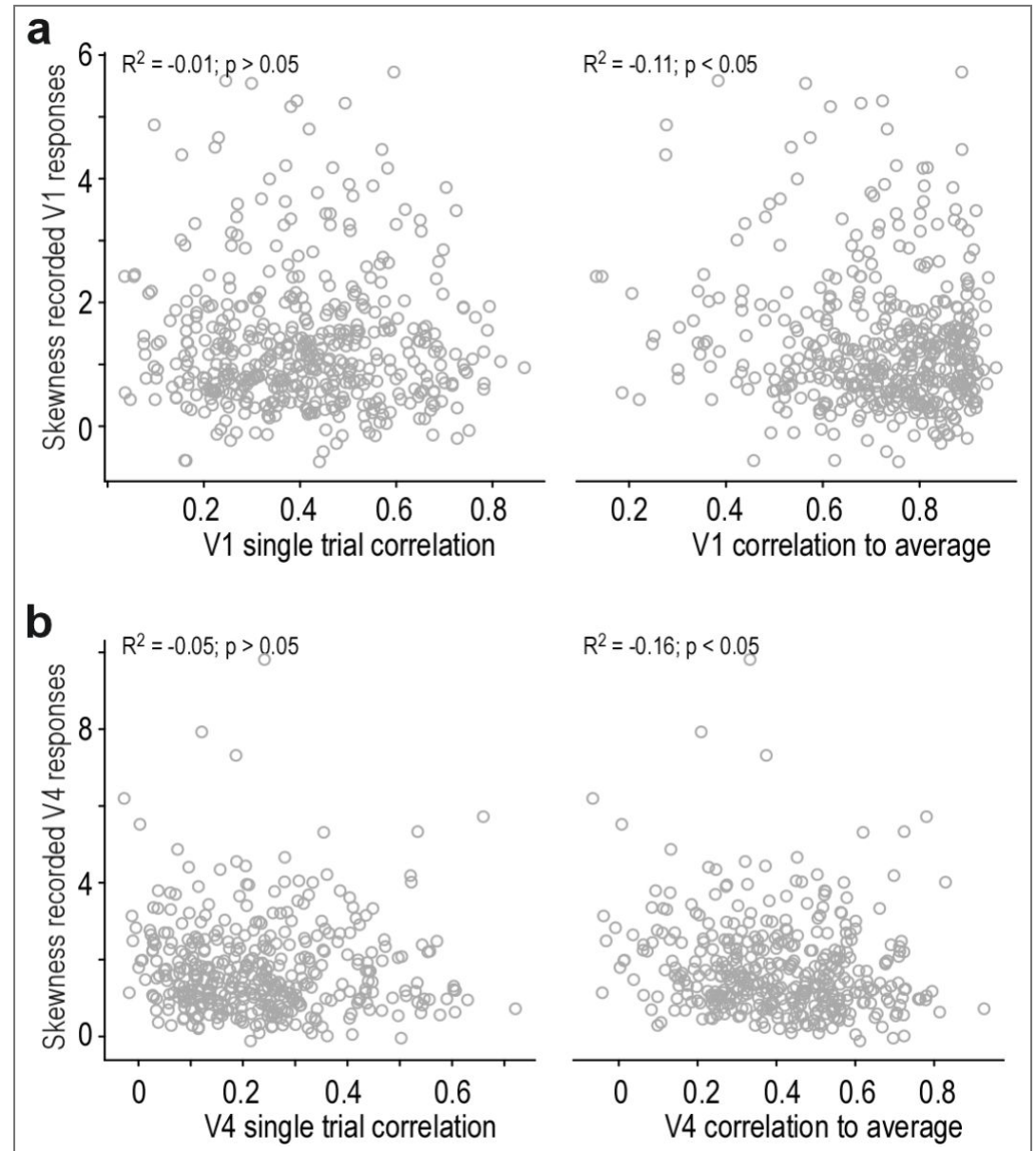
Population-level analysis

To assess the extent to which the most activating images (MAIs) and least activating images (LAIs) of a given neuron also activate other neurons within the same area, we performed the following analysis (for both species). For each neuron, we first identified its top 15 MAIs and bottom 15 LAIs. We then evaluated how these images ranked in terms of activation for every other neuron in the population. Specifically, for each other neuron, we computed the percentile rank of each MAI and LAI within its response distribution obtained from the entire set of size- and contrast-matched naturalistic images (> 1 million images). This yielded a distribution of percentile ranks for the MAIs (and LAIs) of each neuron across all other neurons. Percentile ranks were then averaged across the 15 MAIs and 15 LAIs for each neuron. As a control, we randomly sampled 15 images from the full set of naturalistic images (i.e. > 1 million images) and computed their percentile ranks in the same way. Additionally, to examine whether this effect generalizes across individuals, we conducted a cross-monkey analysis. For this, we identified the MAIs and LAIs of neurons recorded in one monkey and calculated the response percentiles of these images for neurons recorded in the other monkey.

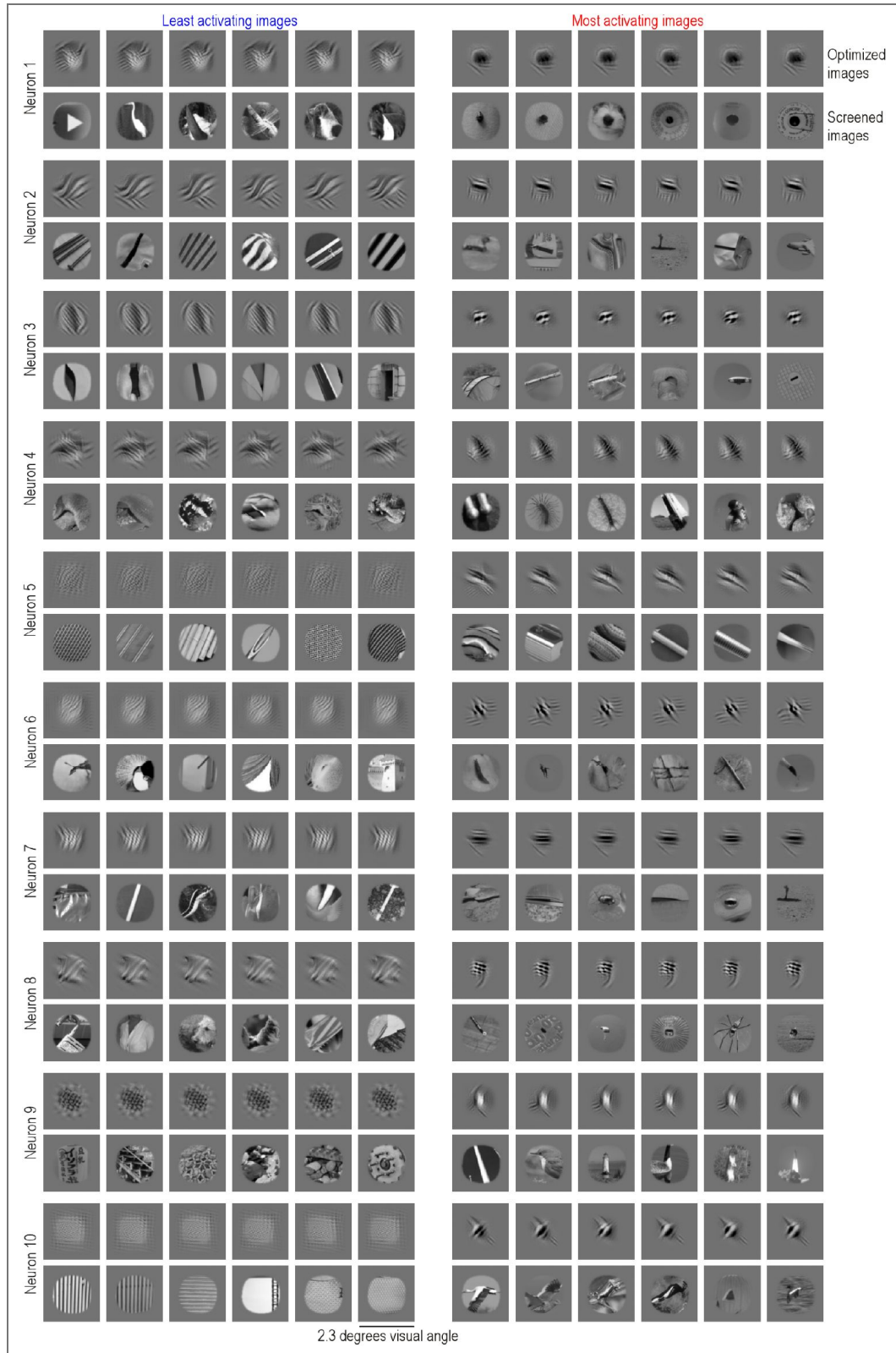
Data availability

Code for screening images in the digital twin, optimizing images, and using the digital twin is available at: <https://github.com/enigma-brain/dualneuron> [↗](#). Data, including predicted responses to ImageNet images, rendered images, as well as code to load ImageNet images, will be available at: <https://doi.org/10.5061/dryad.q573n5tx3> [↗](#).

Supplementary information

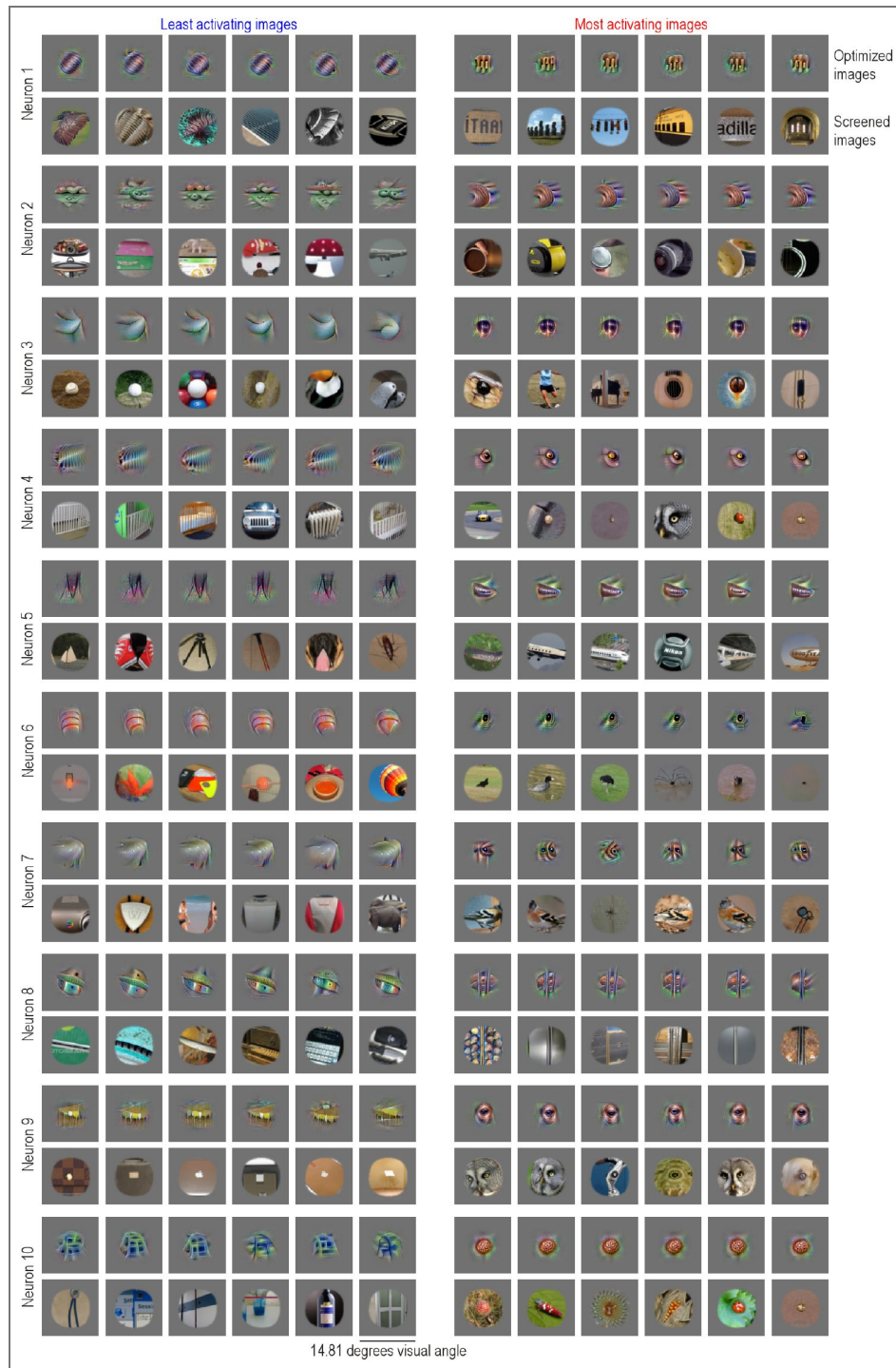


Supplemental Fig. 1. Relationship between model performance and response skewness **a**, Single trial correlation (left) and correlation to average (right) of V1 model plotted versus skewness of recorded V1 responses. **b**, Like (a), but for V4 model and responses.



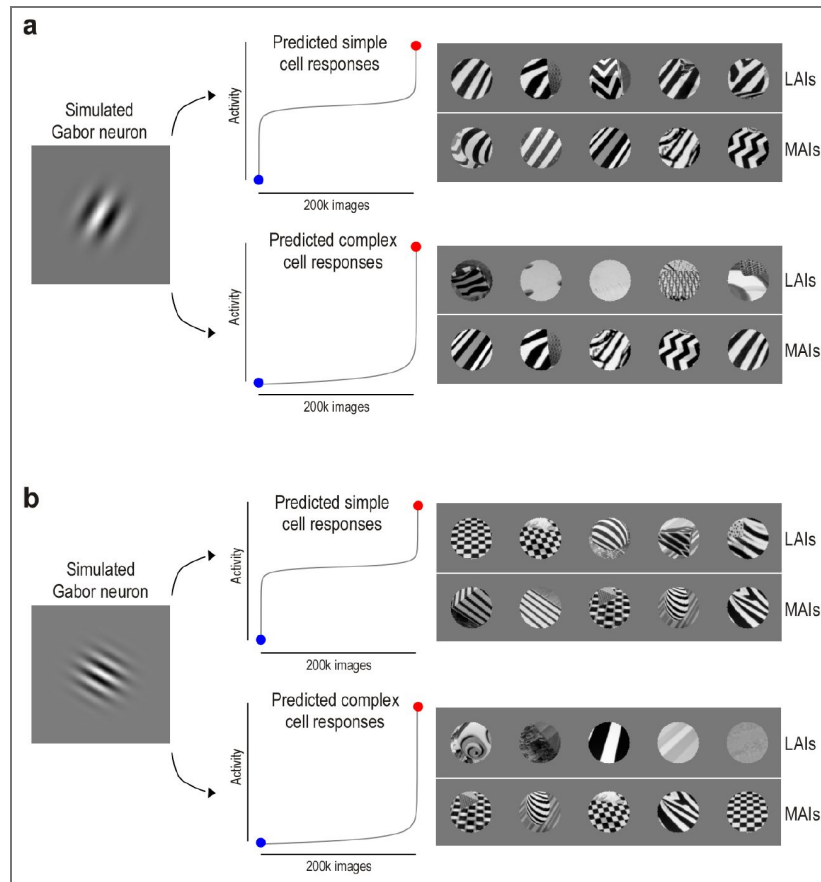
Supplemental Fig. 2. Most and least activating images of example V1 neurons

Least activating (left) and most activating (right) images of ten example V1 neurons. Per neuron, the top row shows optimized images (LEIs and MEIs) and the bottom row shows screened ImageNet images (LAIs and MAIs). Each image is 2.3×2.3 degrees visual angle, with the receptive field of the neuron in the center.



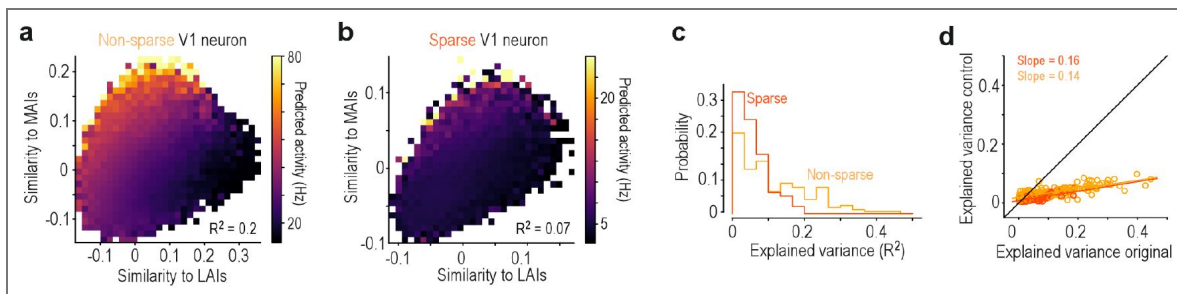
Supplemental Fig. 3. Most and Least Activating Images of Example V4 Neurons

Least activating (left) and most activating (right) images of ten example V4 neurons. Per neuron, the top row shows optimized images (LEIs and MEIs) and the bottom row shows screened ImageNet images (LAIs and MAIs). Each image is 14.83 × 14.83 degrees visual angle, with the receptive field of the neuron in the center.



Supplemental Fig. 4. Most and least activating images of simulated simple and complex cells

a, Example simulation of a V1 simple and complex cell based on the same Gabor receptive field (left). Activations to 200,000 naturalistic images were computed as the dot product between the Gabor filter and each image. For simple cells, this revealed least (blue) and most (red) activating images, showing a bimodal response distribution with LAIs sharing the same orientation but differing in phase. Complex cells were modeled by squaring and pooling responses across different phases, resulting in sparse activation profiles without coherent LAIs.



Supplemental Fig. 5. Responses of V1 neurons vary continuously depending on preferred and non-preferred stimuli

a, Example 2D DreamSim similarity space (cf. Fig. 9) for a non-sparse V1 neuron. Bins are color-coded by the mean predicted neuronal activity of the images within each bin. The R^2 indicates the variance explained by a linear fit. The color bar spans the 0.1 to 99.9th percentile of neuronal responses across all images. **b**, Same as (a), but for a sparse V1 neuron. Here, the activity gradient aligns primarily along the y-axis, indicating stronger modulation by similarity to the MAI than to the LAI. **c**, Variance explained (R^2) by linear regression predicting neuronal activity from the 2D similarity space (as in a,b), shown separately for sparse (dark orange) and non-sparse (light orange) V1 neurons. The explained variances values are generally lower for V1 compared to V4 neurons (cf. Fig. 9). This is likely due to the fact that the DreamSim representational space is more closely aligned with the representational space of V4 than V1. **d**, Variance explained (R^2) by linear regression for the original analysis (similarity to MAIs and LAIs on x- and y-axes) and for Control 1 (left), where both MAIs and LAIs were replaced with random images.

Acknowledgements

The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting PE. FHS is supported by the German Federal Ministry of Federal Ministry of Research, Technology and Space (BMFTR) via the Collaborative Research in Computational Neuroscience (CRCNS) (FKZ 01GQ2107), and FHS & KF are supported by the Collaborative Research Center (SFB 1233, Robust Vision, project number: 276693517). FHS and ASE acknowledges the support of the Lower Saxony Ministry of Science and Culture (MWK) with funds from the Volkswagen Foundation's zukunft.niedersachsen program (project name: CAIMed Lower Saxony Center for Artificial Intelligence and Causal Methods in Medicine; grant number: ZN4257). EF acknowledges support from a European Research Council (ERC) grant (ERC-2022-STG, NEURACT, Grant agreement No: 101076710), the Hellenic Foundation for Research and Innovation (HFRI) under the 2nd Call for HFRI Research Projects to Support Faculty Members and Researchers with Grant Agreement No. 4049, and the HFRI under the "Funding of Basic Research (Horizontal support of all Sciences)" of the National Recovery and Resilience Plan "Greece 2.0" with funding from the European Union - NextGenerationEU with Grant Agreement No. 016552. MD acknowledges support from the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie Actions with Grant Agreement No. 101025482. ASE acknowledges support from the European Research Council (ERC) under the European Union's Horizon Europe research and innovation programme (Grant agreement No. 101041669). KF acknowledges support from the European Research Council (ERC) under the European Union's Horizon Europe research and innovation programme (Grant agreement No. 101117156). AST acknowledges support from the NIH (R01EY026927), the NSF (Collaborative Research in Computational Neuroscience, IIS-2113173), and the Defense Advanced Research Projects Agency (DARPA), Contract No. N66001-19-C-4020 and Contract No. DARPA NESD N66001-17-C-4002. The views, opinions and/or findings expressed are those of the author and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government. AST and SS acknowledge support from the Amaranth Foundation (Enigma Project).

Additional information

Author contributions

NK: Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data Curation, Writing - Original Draft, Visualization, Project Administration **KF:** Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data Curation, Writing - Original Draft, Visualization, Project Administration, Funding Acquisition **KW:** Conceptualization, Methodology, Software, Data Curation **MD:** Investigation, Data Curation, Visualization, Writing - Review & Editing **KRa:** Conceptualization, Investigation, Methodology **PE:** Software **KRe:** Methodology, Investigation, Data Curation **PF:** Conceptualization, Data Curation **CN:** Methodology, Data Curation **TS:** Methodology, Data Curation **GG:** Methodology, Data Curation **SP:** Software, Methodology, Validation **AE:** Writing - Review & Editing **EYW:** Conceptualization, Data Curation **EF:** Writing - Review & Editing **SS:** Conceptualization, Funding Acquisition, Writing - Review & Editing **FHS:** Conceptualization, Supervision, Writing - Review & Editing, Funding Acquisition **AT:** Conceptualization, Supervision, Funding Acquisition, Writing - Review & Editing

Funding

Funder	Grant reference number	Author
Deutsche Forschungsgemeinschaft (DFG)	276693517	Katrin Franke Fabian H Sinz
Federal Ministry of Research, Technology and Space	01GQ2107	Fabian H Sinz

Lower Saxony Ministry of Science and Culture	ZN4257	Fabian H Sinz Alexander Ecker
EC European Research Council (ERC)	https://doi.org/10.3030/101076710	Emmanouil Froudarakis
European Research Council	https://doi.org/10.3030/101041669	Alexander Ecker
European Research Council	https://doi.org/10.3030/101117156	Katrin Franke
EC Horizon Europe Excellent Science HORIZON EUROPE Marie Skłodowska-Curie Actions (MSCA)	https://doi.org/10.3030/101025482	Maria Diamantaki
Hellenic Foundation for Research and Innovation (HFRI)	4049	Emmanouil Froudarakis
Hellenic Foundation for Research and Innovation (HFRI)	016552	Emmanouil Froudarakis
National Institute of Health Sciences (NIHS)	R01EY026927	Andreas Savas Tolias
National Science Foundation (NSF)	IIS-2113173	Andreas Savas Tolias
DARPA	N66001-19-C-4020	Andreas Savas Tolias
DARPA	N66001-17-C-4002	Andreas Savas Tolias

Author ORCID iDs

Katrin Franke:  <https://orcid.org/0000-0002-8649-4835>

Maria Diamantaki:  <https://orcid.org/0000-0002-8455-7628>

Alexander Ecker:  <https://orcid.org/0000-0003-2392-5105>

Emmanouil Froudarakis:  <https://orcid.org/0000-0002-3249-3845>

Fabian H Sinz:  <https://orcid.org/0000-0002-1348-9736>

Andreas Tolias:  <https://orcid.org/0000-0002-4305-6376>

References

Lettvin J. Y., Maturana H. R., McCulloch W. S., Pitts W. H. (1959) What the frog's eye tells the frog's brain. *Proceedings of the IRE* **47**:1940-1951 <https://doi.org/10.1109/jrproc.1959.287207>

Barlow H. B. (1961) Possible principles underlying the transformation of sensory messages. In: Rosenblith W. A. (Ed). *Sensory communication* Cambridge, MA: MIT Press. <https://doi.org/10.7551/mitpress/9780262518420.003.0013>

Field D. J. (1987) Relations between the statistics of natural images and the response properties of cortical cells. *J. Opt. Soc. Am. A* **4**:2379-2394 <https://doi.org/10.1364/josaa.4.002379> | [PubMed](#)

Olshausen B. A., Field D. J. (1996) Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature* **381**:607-609 <https://doi.org/10.1038/381607a0> | [PubMed](#)

Levy W. B., Baxter R. A. (1996) Energy efficient neural codes. *Neural Comput* **8**:531-543 <https://doi.org/10.1162/neco.1996.8.3.531> | [PubMed](#)

Attwell D., Laughlin S. B. (2001) An energy budget for signaling in the grey matter of the brain. *J. Cereb. Blood Flow Metab* **21**:1133-1145 <https://doi.org/10.1097/00004647-200110000-00001> | [PubMed](#)

Hubel D. H., Wiesel T. N. (1962) Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *J. Physiol* **160**:106-154 <https://doi.org/10.1113/jphysiol.1962.sp006837> | [PubMed](#)

- Hubel D. H., Wiesel T. N. (1968) Receptive fields and functional architecture of monkey striate cortex. *J. Physiol* **195**:215-243 <https://doi.org/10.1113/jphysiol.1968.sp008455> | PubMed
- Bishop P. O., Coombs J. S., Henry G. H. (1973) Receptive fields of simple cells in the cat striate cortex. *J. Physiol* **231**:31-60 <https://doi.org/10.1113/jphysiol.1973.sp010218> | PubMed
- Schiller P. H., Finlay B. L., Volman S. F. (1976) Quantitative studies of single-cell properties in monkey striate cortex. I. spatiotemporal organization of receptive fields. *J. Neurophysiol* **39**:1288-1319 <https://doi.org/10.1152/jn.1976.39.6.1288> | PubMed
- Quiroga R. Q., Reddy L., Kreiman G., Koch C., Fried I. (2005) Invariant visual representation by single neurons in the human brain. *Nature* **435**:1102-1107 <https://doi.org/10.1038/nature03687> | PubMed
- Willmore B., Tolhurst D. J. (2001) Characterizing the sparseness of neural codes. *Network* **12**:255-270 <https://doi.org/10.1080/713663277> | PubMed
- Willmore B. D. B., Mazer J. A., Gallant J. L. (2011) Sparse coding in striate and extrastriate visual cortex. *J. Neurophysiol* **105**:2907-2919 <https://doi.org/10.1152/jn.00594.2010> | PubMed
- Heeger D. J. (1992) Normalization of cell responses in cat striate cortex. *Vis. Neurosci* **9**:181-197 <https://doi.org/10.1017/s0952523800009640> | PubMed
- Carandini M., Heeger D. J. (2011) Normalization as a canonical neural computation. *Nat. Rev. Neurosci* **13**:51-62 <https://doi.org/10.1038/nrn3136> | PubMed
- Morrone M. C., Burr D. C., Maffei L. (1982) Functional implications of cross-orientation inhibition of cortical visual cells. I. neurophysiological evidence. *Proc. R. Soc. Lond. B Biol. Sci* **216**:335-354 <https://doi.org/10.1098/rspb.1982.0078> | PubMed
- Allison J. D., Casagrande V. A., Bonds A. B. (1995) The influence of input from the lower cortical layers on the orientation tuning of upper layer V1 cells in a primate. *Vis. Neurosci* **12**:309-320 <https://doi.org/10.1017/s0952523800007999> | PubMed
- Ferster D. (1986) Orientation selectivity of synaptic potentials in neurons of cat primary visual cortex. *J. Neurosci* **6**:1284-1301 <https://doi.org/10.1523/jneurosci.06-05-01284.1986> | PubMed
- Walker E. Y., Sinz F. H., Cobos E., Muhammad T., Froudarakis E., Fahey P. G., Ecker A. S., Reimer J., Pitkow X., Tolias A. S. (2019) Inception loops discover what excites neurons most using deep predictive models. *Nat. Neurosci* **22**:2060-2065 <https://doi.org/10.1038/s41593-019-0517-x> | PubMed
- Bashivan P., Kar K., DiCarlo J. J. (2019) Neural population control via deep image synthesis. *Science* **364**:eaav9436 <https://doi.org/10.1126/science.aav9436> | PubMed
- Franke K., Willeke K. F., Ponder K., Galdamez M., Zhou N., Muhammad T., Patel S., Froudarakis E., Reimer J., Sinz F. H., et al. (2022) State-dependent pupil dilation rapidly shifts visual feature selectivity. *Nature* <https://doi.org/10.1038/s41586-022-05270-3> | PubMed
- Willeke K. F., Restivo K., Franke K., Nix A. F., Cadena S. A., Shinn T., Nealley C., Rodriguez G., Patel S., Ecker A. S., et al. (2023) Deep learning-driven characterization of single cell tuning in primate visual area V4 unveils topological organization. *bioRxiv* <https://doi.org/10.1101/2023.05.12.540591> | PubMed
- Schneider-Mizell C. M., Bodor A. L., Brittain D., Buchanan J., Bumbarger D. J., Elabbady L., Gamlin C., Kapner D., Kinn S., Mahalingam G., et al. (2025) Inhibitory specificity from a connectomic census of mouse visual cortex. *Nature* **640**:448-458 <https://doi.org/10.1038/s41586-024-07780-8> | PubMed
- Matsliah A., Yu S.-C., Kruk K., Bland D., Burke A. T., Gager J., Hebditch J., Silverman B., Willie K. P., Willie R., et al. (2024) Neuronal parts list and wiring diagram for a visual system. *Nature* **634**:166-180 <https://doi.org/10.1038/s41586-024-07981-1> | PubMed
- Sebastian Seung H. (2024) Interneuron diversity and normalization specificity in a visual system. *bioRxiv* <https://doi.org/10.1101/2024.04.03.587837>
- Cadena S. A., Willeke K., Restivo K., Denfield G. H., Sinz F. H., Bethge M., Tolias A., Ecker A. S. (2023) Diverse task-driven modeling of macaque V4 reveals functional specialization towards semantic tasks. *PLoS Comput. Biol* **20** <https://doi.org/10.1371/journal.pcbi.1012056> | PubMed

- LeCun Y., Bengio Y., Hinton G. (2015) Deep learning. *nature* **521**:436-444 <https://doi.org/10.1038/nature14539> | PubMed
- Fu J., Shrinivasan S., Baroni L., Ding Z., Fahey P. G., Pierzchlewicz P., Ponder K., Froebe R., Ntanavara L., Muhammad T., et al. (2024) Pattern completion and disruption characterize contextual modulation in the visual cortex. *bioRxiv* 2023.03.13.532473 <https://doi.org/10.1101/2023.03.13.532473> | PubMed
- Rust N. C., DiCarlo J. J. (2012) Balanced increases in selectivity and tolerance produce constant sparseness along the ventral visual stream. *J. Neurosci* **32**:10170-10182 <https://doi.org/10.1523/jneurosci.6125-11.2012> | PubMed
- Gondur R., Stan P. L., Smith M. A., Cowley B. R. (2025) A tale of two tails: Preferred and antipreferred natural stimuli in visual cortex. In: The Fourteenth International Conference on Learning Representations.
- Desimone R., Schein S. J. (1987) Visual properties of neurons in area V4 of the macaque: sensitivity to stimulus form. *J. Neurophysiol* **57**:835-868 <https://doi.org/10.1152/jn.1987.57.3.835> | PubMed
- Pasupathy A., Connor C. E. (2002) Population coding of shape in area V4. *Nat. Neurosci* **5**:1332-1338 <https://doi.org/10.1038/nn972> | PubMed
- Yamane Y., Carlson E. T., Bowman K. C., Wang Z., Connor C. E. (2008) A neural code for three-dimensional object shape in macaque inferotemporal cortex. *Nat. Neurosci* **11**:1352-1360 <https://doi.org/10.1038/nn.2202> | PubMed
- Fu S., Tamir N., Sundaram S., Chai L., Zhang R., Dekel T., Isola P. (2023a) DreamSim: Learning new dimensions of human visual similarity using synthetic data. *arXiv* <https://doi.org/10.48550/arXiv.2306.09344>
- Pierzchlewicz P. A., Willeke K. F., Nix A., Elumalai P., Restivo K., Shinn T., Nealley C., Rodriguez G., Patel S., Franke K., et al. (2023) Energy guided diffusion for generating neurally exciting images. In: Thirty-seventh Conference on Neural Information Processing Systems.
- Wang Q., Sporns O., Burkhalter A. (2012) Network analysis of corticocortical connections reveals ventral and dorsal processing streams in mouse visual cortex. *Journal of Neuroscience* **32**:4386-4399 <https://doi.org/10.1523/JNEUROSCI.6063-11.2012> | PubMed
- Froudarakis E., Cohen U., Diamantaki M., Walker E. Y., Reimer J., Berens P., Sompolinsky H., Tolias A. S. (2021) Object manifold geometry across the mouse cortical visual hierarchy. *bioRxiv* <https://doi.org/10.1101/2020.08.20.258798>
- Willeke K. F., Fahey P. G., Bashiri M., Pede L., Burg M. F., Blessing C., Cadena S. A., Ding Z., Lurz K.-K., Ponder K., et al. (2022) The sensorium competition on predicting large-scale mouse primary visual cortex activity. *arXiv* <https://doi.org/10.48550/arxiv.2206.08666>
- Isaacson J. S., Scanziani M. (2011) How inhibition shapes cortical activity. *Neuron* **72**:231-243 <https://doi.org/10.1016/j.neuron.2011.09.027> | PubMed
- Priebe N. J. (2016) Mechanisms of orientation selectivity in the primary visual cortex. *Annu. Rev. Vis. Sci* **2**:85-107 <https://doi.org/10.1146/annurev-vision-111815-114456> | PubMed
- Burr D., Morrone C., Maffei L. (1981) Intra-cortical inhibition prevents simple cells from responding to textured visual patterns. *Exp. Brain Res* **43**:455-458 <https://doi.org/10.1007/bf00238391> | PubMed
- Hata Y., Tsumoto T., Sato H., Hagihara K., Tamura H. (1988) Inhibition contributes to orientation selectivity in visual cortex of cat. *Nature* **335**:815-817 <https://doi.org/10.1038/335815a0> | PubMed
- Ben-Yishai R., Bar-Or R. L., Sompolinsky H. (1995) Theory of orientation tuning in visual cortex. *Proc. Natl. Acad. Sci. U. S. A* **92**:3844-3848 <https://doi.org/10.1073/pnas.92.9.3844> | PubMed
- Ringach D. L., Bredfeldt C. E., Shapley R. M., Hawken M. J. (2002) Suppression of neural responses to nonoptimal stimuli correlates with tuning selectivity in macaque V1. *J. Neurophysiol* **87**:1018-1027 <https://doi.org/10.1152/jn.00614.2001> | PubMed
- DeAngelis G. C., Robson J. G., Ohzawa I., Freeman R. D. (1992) Organization of suppression in receptive fields of neurons in cat visual cortex. *J. Neurophysiol* **68**:144-163 <https://doi.org/10.1152/jn.1992.68.1.144> | PubMed

- Burg M. F., Cadena S. A., Denfield G. H., Walker E. Y., Tolia A. S., Bethge M., Ecker A. S. (2021) Learning divisive normalization in primary visual cortex. *PLoS Comput. Biol* **17**:e1009028 <https://doi.org/10.1371/journal.pcbi.1009028> | [PubMed](#)
- Bauman L. A., Bonds A. B. (1991) Inhibitory refinement of spatial frequency selectivity in single cells of the cat striate cortex. *Vision Res* **31**:933-944 [https://doi.org/10.1016/0042-6989\(91\)90201-f](https://doi.org/10.1016/0042-6989(91)90201-f) | [PubMed](#)
- De Valois K. K., Tootell R. B. (1983) Spatial-frequency-specific inhibition in cat striate cortex cells. *J. Physiol* **336**:359-376 <https://doi.org/10.1113/jphysiol.1983.sp014586> | [PubMed](#)
- Rowekamp R. J., Sharpee T. O. (2017) Cross-orientation suppression in visual area V2. *Nat. Commun* **8**:15739 <https://doi.org/10.1038/ncomms15739> | [PubMed](#)
- Pollen D. A., Przybyszewski A. W., Rubin M. A., Foote W. (2002) Spatial receptive field organization of macaque V4 neurons. *Cereb. Cortex* **12**:601-616 <https://doi.org/10.1093/cercor/12.6.601> | [PubMed](#)
- Miller E. K., Gochin P. M., Gross C. G. (1993) Suppression of visual responses of neurons in inferior temporal cortex of the awake macaque by addition of a second stimulus. *Brain Res* **616**:25-29 [https://doi.org/10.1016/0006-8993\(93\)90187-r](https://doi.org/10.1016/0006-8993(93)90187-r) | [PubMed](#)
- Rolls E. T., Tovee M. J. (1995) Sparseness of the neuronal representation of stimuli in the primate temporal visual cortex. *J. Neurophysiol* **73**:713-726 <https://doi.org/10.1152/jn.1995.73.2.713> | [PubMed](#)
- Tamura H., Kaneko H., Kawasaki K., Fujita I. (2004) Presumed inhibitory neurons in the macaque inferior temporal cortex: visual response properties and functional interactions with adjacent neurons. *J. Neurophysiol* **91**:2782-2796 <https://doi.org/10.1152/jn.01267.2003> | [PubMed](#)
- Desimone R., Duncan J. (1995) Neural mechanisms of selective visual attention. *Annu. Rev. Neurosci* **18**:193-222 <https://doi.org/10.1146/annurev.ne.18.030195.001205> | [PubMed](#)
- Reynolds J. H., Chelazzi L., Desimone R. (1999) Competitive mechanisms subserve attention in macaque areas V2 and V4. *J. Neurosci* **19**:1736-1753 <https://doi.org/10.1523/jneurosci.19-05-01736.1999> | [PubMed](#)
- Chang L., Tsao D. Y. (2017) The code for facial identity in the primate brain. *Cell* **169**:1013-1028.e14, <https://doi.org/10.1016/j.cell.2017.05.011> | [PubMed](#)
- Fel T., Wang B., Lepori M. A., Kowal M., Lee A., Balestrieri R., Joseph S., Lubana E. S., Konkle T., Ba D., et al. (2025) Into the rabbit hull: From task-relevant concepts in DINO to minkowski geometry. *arXiv* <https://doi.org/10.48550/arXiv.2510.08638>
- Rigotti M., Barak O., Warden M. R., Wang X.-J., Daw N. D., Miller E. K., Fusi S. (2013) The importance of mixed selectivity in complex cognitive tasks. *Nature* **497**:585-590 <https://doi.org/10.1038/nature12160> | [PubMed](#)
- Fusi S., Miller E. K., Rigotti M. (2016) Why neurons mix: high dimensionality for higher cognition. *Curr. Opin. Neurobiol* **37**:66-74 <https://doi.org/10.1016/j.conb.2016.01.010> | [PubMed](#)
- Tong R., da Silva R., Lin D., Ghosh A., Wilsenach J., Cianfarano E., Bashivan P., Richards B., Trenholm S. (2023) The feature landscape of visual cortex. *bioRxiv* <https://doi.org/10.1101/2023.11.03.565500>
- Froudarakis E., Berens P., Ecker A. S., Cotton R. J., Sinz F. H., Yatsenko D., Saggau P., Bethge M., Tolia A. S. (2014) Population code in mouse V1 facilitates readout of natural scenes through increased sparseness. *Nat. Neurosci* **17**:851-857 <https://doi.org/10.1038/nn.3707> | [PubMed](#)
- Oquab M., Darcet T., Moutakanni T., Vo H., Szafraniec M., Khalidov V., Fernandez P., Haziza D., Massa F., El-Nouby A., et al. (2023) DINOv2: Learning robust visual features without supervision. *arXiv* <https://doi.org/10.48550/arxiv.2304.07193>
- Willeke K. F., Turishcheva P., Gilbert A., Chakrabarty G., Bedel H. A., Fahey P. G., Qiu Y., Weis M. A., Vystrčilová M., Muhammad T., et al. (2025) Omni-Mouse: Scaling properties of multi-modal, multi-task brain models on 150B neural tokens. In: The Fourteenth International Conference on Learning Representations.

- Ding Z., Fahey P. G., Papadopoulos S., Wang E. Y., Celii B., Papadopoulos C., Chang A., Kunin A. B., Tran D., Fu J., *et al.* (2025) Functional connectomics reveals general wiring rule in mouse visual cortex. *Nature* **640**:459-469 <https://doi.org/10.1038/s41586-025-08840-3> | [PubMed](#)
- Tasic B., Yao Z., Graybiel L. T., Smith K. A., Nguyen T. N., Bertagnolli D., Goldy J., Garren E., Economo M. N., Viswanathan S., *et al.* (2018) Shared and distinct transcriptomic cell types across neocortical areas. *Nature* **563**:72-78 <https://doi.org/10.1038/s41586-018-0654-5> | [PubMed](#)
- Gouwens N. W., Sorensen S. A., Berg J., Lee C., Jarsky T., Ting J., Sunkin S. M., Feng D., Anastassiou C. A., Barkan E., *et al.* (2019) Classification of electrophysiological and morphological neuron types in the mouse visual cortex. *Nat. Neurosci* **22**:1182-1195 <https://doi.org/10.1038/s41593-019-0417-0> | [PubMed](#)
- Yao Z., Liu H., Xie F., Fischer S., Adkins R. S., Aldridge A. I., Ament S. A., Bartlett A., Behrens M. M., Van den Berge K., *et al.* (2021) A transcriptomic and epigenomic cell atlas of the mouse primary motor cortex. *Nature* **598**:103-110 <https://doi.org/10.1038/s41586-021-03500-8> | [PubMed](#)
- Keller D., Erö C., Markram H. (2018) Cell densities in the mouse brain: A systematic review. *Front. Neuroanat* **12** <https://doi.org/10.3389/fnana.2018.00083> | [PubMed](#)
- Diamond J. S. (2017) Inhibitory interneurons in the retina: Types, circuitry, and function. *Annu Rev Vis Sci* **3**:1-24 <https://doi.org/10.1146/annurev-vision-102016-061345> | [PubMed](#)
- Matsumoto A., Morris J., Looger L. L., Yonehara K. (2025) Functionally distinct GABAergic amacrine cell types regulate spatiotemporal encoding in the mouse retina. *Nat. Neurosci* **28**:1256-1267 <https://doi.org/10.1038/s41593-025-01935-0> | [PubMed](#)
- Muñoz W., Tremblay R., Levenstein D., Rudy B. (2017) Layer-specific modulation of neocortical dendritic inhibition during active wakefulness. *Science* **355**:954-959 <https://doi.org/10.1126/science.aag2599> | [PubMed](#)
- Lu J., Tucciarone J., Padilla-Coreano N., He M., Gordon J. A., Huang Z. J. (2017) Selective inhibitory control of pyramidal neuron ensembles and cortical subnetworks by chandelier cells. *Nat. Neurosci* **20**:1377-1383 <https://doi.org/10.1038/nn.4624> | [PubMed](#)
- Wu S. J., Sevier E., Dwivedi D., Saldi G.-A., Hairston A., Yu S., Abbott L., Choi D. H., Sherer M., Qiu Y., *et al.* (2023) Cortical somatostatin interneuron subtypes form cell-type-specific circuits. *Neuron* **111**:2675-2692.e9, <https://doi.org/10.1016/j.neuron.2023.05.032> | [PubMed](#)
- Schwartz O., Simoncelli E. P. (2001) Natural signal statistics and sensory gain control. *Nat. Neurosci* **4**:819-825 <https://doi.org/10.1038/90526> | [PubMed](#)
- Tolmachev P., Engel T. A. (2025) Single-unit activations confer inductive biases for emergent circuit solutions to cognitive tasks. *Nat. Mach. Intell* **7**:1742-1754 <https://doi.org/10.1038/s42256-025-01127-2> | [PubMed](#)
- Ioffe S., Szegedy C. (2015) Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv* <https://doi.org/10.48550/arxiv.1502.03167>
- Ba J. L., Kiros J. R., Hinton G. E. (2016) Layer normalization. *arXiv* <https://doi.org/10.48550/arxiv.1607.06450>
- Vaswani A., Shazeer N., Parmar N., Uszkoreit J., Jones L., Gomez A. N., Kaiser Ł., Polosukhin I. (2017) Attention is all you need. *arXiv* <https://doi.org/10.48550/arXiv.1706.03762>
- Jun J. J., Steinmetz N. A., Siegle J. H., Denman D. J., Bauza M., Barbarits B., Lee A. K., Anastassiou C. A., Aydin A. Andrei, Barbic M., *et al.* (2017) Fully integrated silicon probes for high-density recording of neural activity. *Nature* **551**:232-236 <https://doi.org/10.1038/nature24636> | [PubMed](#)
- Birman D., Yang K. J., West S. J., Karsh B., Browning Y., International Brain Laboratory, Siegle J. H., Steinmetz N. A. (2023) Pinpoint: trajectory planning for multi-probe electrophysiology and injections in an interactive web-based 3D environment. *bioRxiv* 2023.07.14.548952 <https://doi.org/10.1101/2023.07.14.548952> | [PubMed](#)
- Pachitariu M., Sridhar S., Stringer C. (2023) Solving the spike sorting problem with kilosort. *bioRxiv* <https://doi.org/10.1101/2023.01.07.523036>

- Tafazoli S., Safaai H., De Franceschi G., Rosselli F. B., Vanzella W., Riggi M., Buffolo F., Panzeri S., Zoccolan D. (2017) Emergence of transformation-tolerant representations of visual objects in rat lateral extrastriate cortex. *eLife* **6** <https://doi.org/10.7554/elife.22794> | PubMed
- Greff K., Belletti F., Beyer L., Doersch C., Du Y., Duckworth D., Fleet D. J., Gnanaprasam D., Golemo F., Herrmann C., *et al.* (2022) Kubric: a scalable dataset generator. *arXiv* <https://doi.org/10.48550/arXiv.2203.03570>
- Blender Online Community (2024) Blender - a 3D modelling and rendering package. <http://www.blender.org>
- Cimpoi M., Maji S., Kokkinos I., Mohamed S., Vedaldi A. (2014) Describing textures in the wild. In: Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). <https://doi.org/10.1109/cvpr.2014.461>
- Woo S., Debnath S., Hu R., Chen X., Liu Z., Kweon I. S., Xie S. (2023) ConvNeXt V2: Co-designing and scaling ConvNets with masked autoencoders. *arXiv* <https://doi.org/10.48550/arXiv.2301.00808>
- Lurz K.-K., Bashiri M., Willeke K., Jagadish A., Wang E., Walker E. Y., Cadena S. A., Muhammad T., Cobos E., Tolia A. S., *et al.* (2022) Generalization in data-driven models of primary visual cortex. *bioRxiv* <https://doi.org/10.1101/2020.10.05.326256>
- Loshchilov I., Hutter F. (2017) Decoupled weight decay regularization. *arXiv* <https://doi.org/10.48550/arXiv.1711.05101>
- Salman H., Yang G., Li J., Zhang P., Zhang H., Razenshteyn I., Bubeck S. (2019) Provably robust deep learning via adversarially trained smoothed classifiers. *arXiv* <https://doi.org/10.48550/arXiv.1906.04584>
- Clevert D.-A., Unterthiner T., Hochreiter S. (2015) Fast and accurate deep network learning by exponential linear units (ELUs). *arXiv* <https://doi.org/10.48550/arXiv.1511.07289>
- Kingma D. P., Ba J. (2014) Adam: A method for stochastic optimization. *arXiv* <https://doi.org/10.48550/arXiv.1412.6980>
- Deng J., Dong W., Socher R., Li L.-J., Li K., Fei-Fei L. (2009) Imagenet: A large-scale hierarchical image database. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition. pp. 248-255 <https://doi.org/10.1109/CVPR.2009.5206848>
- Fel T., Boissin T., Boutin V., Picard A., Novello P., Colin J., Linsley D., Rousseau T., Cadène R., Goetschalckx L., *et al.* (2023) Unlocking feature visualization for deeper networks with MAgnitude constrained optimization. *arXiv* <https://doi.org/10.48550/arXiv.2306.06805>
- Kingma D. P., Ba J. (2017) Adam: A method for stochastic optimization. *arXiv* <https://doi.org/10.48550/arXiv.1412.6980>
- Fu S., Tamir N. Y., Sundaram S., Chai L., Zhang R., Dekel T., Isola P. (2023b) DreamSim: Learning new dimensions of human visual similarity using synthetic data. In: Thirty-seventh Conference on Neural Information Processing Systems.
- Franke K., Karantzas N., Willeke K., Diamantaki M., Ramakrishnan K, Bedel HA, Elumalai P, Restivo K, Fahey P, N, *et al.* (2025) Dual-feature selectivity enables bidirectional coding in visual cortical neurons. Dryad. <https://doi.org/10.5061/dryad.q573n5tx3>

Peer reviews

Reviewer #1 (Public review):

The multi-species approach of testing the model in macaque and mouse is excellent, as it improves the chances that the observed findings are a general property of mammalian visual cortex. It would be useful to delineate however any notable differences between these species, which are to be expected given their lifestyle.

The overall performance of the model appears to be excellent in V1, with over 80% performance, but falls substantially in V4. It would be important to consider the implications of this finding; for example, in the context of studying temporal lobe structures that are central to recognizing objects. Would one expect that model performance decreases further here, and what measures could be taken to avoid this? Or is this type of model better restricted to V1 or even LGN?

While the manuscript delineates novel axes of inhibitory interactions, it remains unclear what exactly these axes are and how they arise. What are the steps that need to be taken to make progress along these lines?

Comments on revised version.

The authors have adequately addressed the points I raised in my review during the revision.

<https://doi.org/10.7554/eLife.109861.2.sa1>

Author response:

The following is the authors' response to the original reviews.

Public Reviews:

Reviewer #1 (Public review):

This manuscript used deep learning to highlight the role of inhibition in shaping selectivity in primary and higher visual cortex. The findings hint at hitherto unknown axes of structured inhibition operating in cortical networks with a potentially key role in object recognition.

The multi-species approach of testing the model in macaque and mouse is excellent, as it improves the chances that the observed findings are a general property of mammalian visual cortex. However, it would be useful to delineate any notable differences between these species, which are to be expected given their lifestyle.

The overall performance of the model appears to be excellent in V1, with over 80% performance, but it falls substantially in V4. It would be important to consider the implications of this finding; for example, in the context of studying temporal lobe structures that are central to recognizing objects. Would one expect that model performance decreases further here, and what measures could be taken to avoid this? Or is this type of model better restricted to V1 or even LGN?

While the manuscript delineates novel axes of inhibitory interactions, it remains unclear what exactly these axes are and how they arise. What are the steps that need to be taken to make progress along these lines?

Reviewer #2 (Public review):

The classic view of sensory coding states that (excitatory) neurons are active to some preferred stimuli and otherwise silent. In contrast, inhibitory neurons are considered broadly tuned. Due to the gigantic potential image space, it is hard to comprehensively map the tuning of individual neurons. In this tour de force study, Franke et al. combine electrophysiological recordings in macaque (V1, V4) and mouse (V1, LM, LI) visual cortex with large-scale screens based on digital twin models, as well as beautiful systems identification (most/least activating stimuli). Based on these digital twins, they discover dual-feature selectivity (which they validate both in macaques and mice). Dual-feature selectivity involves a bidirectional modulation of firing rates around an elevated baseline. Neurons are excited by specific preferred features and systematically suppressed by

distinct, non-preferred features. This tuning was identified by excellently combining advances in AI & high-throughput ephys.

The study is comprehensive and convincing. Overall, this work showcases how in silico experiments can generate concrete hypotheses about neuronal coding that are difficult to discover experimentally, but that can be experimentally validated! I think this work is of substantial interest to the neuroscience community. I'm sure it will motivate many future experimental and computational studies. In particular, it will be of great interest to understand when and how the brain leverages dual-feature selectivity. The discussion of the article is already an interesting starting point for these considerations.

Strengths:

(1) Using computational models to predict neuronal responses allowed them to go through millions of images, which may not be possible in vivo.

(2) The cross-species and cross-area consistency of the results is another major strength. Pointing out that the results may be a fundamental strategy of mammalian cortical processing.

(3) They show that the feature causing peak excitation in one neuron often drives suppression in another. This may be an efficient coding scheme where the population covers the visual manifold. I'd like to understand better why the authors believe that this shows that there are low-dimensional subspaces based on preferred and non-preferred stimulus features (vs. many more, but some axes are stronger).

We thank the reviewers for their constructive and helpful feedback on our manuscript. We are delighted that they found the study to be “comprehensive and convincing” and a “tour de force” in its combination of electrophysiological recordings with large-scale digital twin screening. We appreciate that the reviewers highlighted the strengths of our multi-species approach and the “cross-species and cross-area consistency” of the results, noting that the work showcases how *in silico* experiments can generate concrete, experimentally validatable hypotheses. Overall, we agree with the assessment of the reviewers. We have performed the following changes to the text to clarify and strengthen the manuscript, without introducing new analyses or altering the conclusions.

Recommendations for the authors:

Reviewer #1 (Recommendations for the authors):

(1) Page 3: The authors state that RFs were mapped using sparse noise, with the goal to ensure that the RFs align with the visual stimulus, but no data appear to be shown regarding this alignment. It would be important to provide a full analysis of the sparse noise-mapped RFs for both V1 and V4. Also, is it correct that the V4 data analyzed here came from a single animal? This could potentially be problematic and would need to be addressed, for example, by performing analyses also in V1 for participant animals separately. Please elaborate.

We have added a sentence to the Results section clarifying the sparse noise RF mapping procedure, noting that probe insertions were targeted orthogonal to the cortical surface so that neurons sampled along the probe depth share overlapping receptive fields, allowing a single stimulus configuration to adequately drive the entire recorded population. We have also corrected the text to clarify that V4 data were collected from 2 animals (not 3 as previously stated in an earlier draft), consistent with the Methods section.

(2) Page 4: Only half the neurons in V4 are "high confidence" in terms of test image performance, which seems a little low and probably significantly lower than the

corresponding value for V1 of 84%. It is unclear how to interpret this confidence, but it seems to suggest that half of the V4 neurons are not well captured by the model. If true, this fraction appears large enough to cast doubt on the validity of the V4 results. Please elaborate.

We have expanded the text to explicitly discuss the lower proportion of high-confidence *in-silico* neurons in V4 relative to V1. We attribute this to the greater complexity of V4 tuning compared to V1, as well as missing contextual information such as image surrounds and sequential image context—factors that likely limit model performance in higher visual areas. We note that our restriction of analyses to high-confidence neurons provides resilience against these limitations, and that the goal was not to maximize predictive performance per se but to identify response patterns—dual-feature selectivity—that are robust across neurons, areas, and species.

(3) Page 5: It seems that identical L2 norms are valid for discounting contrast variations, particularly if the neural responses are linear, since the L2 norm is computed on the entire RF. It might be judicious to attenuate the claim that contrast variation has no effect.

We have softened the claim that contrast variation has no effect. The revised text now states that L2 normalization controls for root-mean-squared contrast but does not fully equate effective contrast in nonlinear cells, whose responses depend on the spatial structure of the stimulus beyond its total energy. We note that residual contrast dependent effects, particularly in the suppressive regime, cannot be entirely excluded.

(4) Page 6: The authors acknowledge that, at least for simple cells, a phase shift in the grating and concomitant ON-OFF overlap is an inhibitory axis, which is correct. It does not really become clear what other axes were found, and whether any of these represent a novel discovery about V1.

We have clarified the description of inhibitory axes in V1, noting that while phase-shifted stimuli represent a well-established suppressive axis for simple cells reflecting linear On-Off subfield structure, and complex cells exhibit no coherent suppressive pattern due to phase pooling, neither model class accounts for the multidimensional suppressive structure we observe. We have made explicit that our unbiased approach reveals suppressive structure spanning simultaneous changes across orientation, spatial frequency, phase, and texture, exceeding what any single known suppressive mechanism predicts.

(5) Page 7: Dreamsim is based on human similarity judgements, whereas the data is from macaques. Is there any evidence suggesting that macaque similarity judgements might be similar to those of humans?

We have added a paragraph to the Discussion acknowledging that DreamSim was trained on human perceptual similarity judgments while our neuronal data are from macaques. We note that this cross-species application is supported by the deep homology between primate ventral visual streams, and that natural-image similarity judgments have been found to be highly consistent across macaques and humans. Importantly, we clarify that we deploy DreamSim not as a model of macaque perception but as an image feature embedding to test whether stimuli that cluster in perceptual space evoke similar neuronal responses—a use that is robust to the precise calibration of the metric. We also note that we are developing custom macaque-specific embeddings for future work.

(6) Page 7: How many images were in the test set?

We have added the number of test images to the relevant text ($n=75$ for V1, $n=150$ for V4) and to the Figure 1 caption.

(7) Page 8: As mentioned above, performing the analysis on V1 data of individual subjects and demonstrating similar digital twins might be an additional way to confirm the models' accuracy.

We have added text noting that for V4, 1digital twin models were fit independently per neuron without sharing information across animals, and that extreme image sets identified by the model elicited correspondingly extreme responses in neurons from the other animal, confirming that identified selectivity patterns are not idiosyncratic to individual subjects.

(8) Page 11: The mouse data is presented very briefly only, and the authors seem to imply that there is a high degree of coding similarity between this rodent species and macaques and, by extension, humans. Were there any notable differences between the mouse and macaque data?

We have added text explicitly noting that while macaque and mouse visual cortex differ substantially in their functional organization and the complexity of neuronal selectivity, the broader principle—that non-sparse neurons are jointly defined by distinct excitatory and suppressive feature sets—generalizes across mammalian visual systems. We clarify that this does not imply that mouse and macaque visual cortex share similar functional organization or equivalent complexity of neuronal selectivity; rather, within the representational regime of each area, neurons are organized such that excitatory and suppressive feature sets are jointly structured and distinct.

(9) Page 13: One main finding of the study is that inhibition appears to operate along additional dimensions that had not been previously recognized, but what is the nature of these dimensions, how do they arise and relate to known inhibitory effects in V1 such as centre-surround effects? The fact that suppression is tuned in response to natural images or other complex objects is not a new finding, and there is plenty of published work along these lines; the authors may want to cite Tamura et al 10.1152/jn.01267.2003. I am not sure introducing the term "dual feature selectivity" is really a major conceptual advance.

We have added a citation to Tamura et al. (2004) in the Discussion, alongside other prior work documenting suppression by non-optimal stimuli. We have also expanded the Discussion to more carefully position our findings relative to existing work on feature-selective suppression, noting that while prior work has established that inhibition can be structured and feature-selective, our results suggest a broader organizing principle: within each visual area, there exists a set of feature combinations from which individual neurons draw both their excitatory and suppressive preferences.

(10) Page 14: The authors enumerate a number of technical limitations, which is to be commended. It would be useful for them to comment on the particular advantages of the digital twin model, compared to a more traditional analysis of the responses to the thousands of natural images that were experimentally obtained. It seems likely that the main finding, i.e. tuned inhibition, is also evident directly in this population (?). While the digital twin is to some degree validated by the test images, its responses to the much larger set of images studied are not validated, and one must trust that the ResNet50 indeed captures V4 selectivity. It would be useful to discuss some of these points, and highlight a potential way that digital twins (maybe as a shared model between laboratories) can learn from a large number of animals and datasets, and maybe even be used to generate novel visual stimuli suitable to test emergent hypotheses.

We have added a paragraph to the Discussion explicitly contrasting the advantages of digital twin models with direct analysis of experimentally recorded responses, noting that digital twins enable screening of more than one million images per neuron *in silico*, gradient-based

synthesis of stimuli precisely optimized to drive or suppress individual neurons, and cross-model verification of identified selectivity patterns—a test that has no analog when working with fixed experimental image sets.

Reviewer #2 (Recommendations for the authors):

Minor comments:

(1) Call out Figure 1/b in the main text.

We have added a callout to Figure 1b in the main text

(2) Can you make a supplementary figure illustrating more examples with skewness around the middle (e.g. 1.5, 2, 2.5)? Namely, you state that 2 is a good threshold for deciding if it is non-sparse, but you only present clear-cut cases in Figure 2 (with <0.75 and >3.5). I am wondering if 2 is a good threshold?

We have revised the text to clarify that the skewness threshold of 2.0 is adopted purely for analytical convenience to focus subsequent analyses on neurons with sufficiently graded response distributions, and that the key findings are not dependent on the exact threshold chosen. We explicitly note that the underlying distribution of sparsity is continuous, consistent with recent findings (Gondur et al., 2025).

(3) The reference "A tale of two tails: Preferred and anti-preferred natural stimuli in visual cortex." Has no authors. I know it's anonymous, but maybe put that for now? I also congratulate including a paper that is anonymously under review at ICLR 2026. I don't find Unk, 2025 in the list of references. Perhaps related?

We have updated the reference "A tale of two tails" to include the authors (Gondur et al., 2025) and ensured it appears consistently in the reference list. We have also resolved the missing "Unk, 2025" citation, which now correctly refers to this same work.

(4) Why do you use a different model for the analysis in Figure 8?

We have added text to the Methods and Results clarifying why a distinct architecture was used for the V4 evaluator model in Figure 8. Specifically, the V4 generator model uses a fixed, pretrained ResNet50 backbone whose weights are deterministic; any re-trained model sharing this backbone would not constitute a genuinely independent evaluation. By contrast, for V1, the ConvNeXt core is fine-tuned from different random initializations, producing architecturally equivalent but computationally independent models. A truly independent V4 evaluator therefore required a fundamentally different architecture.

<https://doi.org/10.7554/eLife.109861.2.sa0>