

---

# Stretching Beyond the Obvious: A Gradient-Free Framework to Unveil the Hidden Landscape of Visual Invariance

---

Lorenzo Tausani<sup>1</sup>, Paolo Muratore<sup>1</sup>, Morgan B. Talbot<sup>2,3,4</sup>, Giacomo Amerio<sup>1</sup>, Gabriel Kreiman<sup>2,3</sup>, and Davide Zoccolan<sup>1</sup>

<sup>1</sup>Neuroscience Area, International School for Advanced Studies (SISSA), Trieste (Italy)

<sup>2</sup>Boston Children’s Hospital, Harvard Medical School, Cambridge (USA)

<sup>3</sup>Center for Brains, Minds, and Machines, MIT, Cambridge (USA)

<sup>4</sup>Harvard-MIT Program in Health Sciences and Technology, MIT, Cambridge (USA)

## Abstract

Uncovering which features’ combinations high-level visual units encode is critical to understand how images are transformed into representations that support recognition. While existing feature visualization approaches typically infer a unit’s most exciting images, this is insufficient to reveal the manifold of transformations under which responses remain invariant, which is key to generalization in vision. Here we introduce Stretch-and-Squeeze (SnS), an unbiased, model-agnostic, and gradient-free framework to systematically characterize a unit’s invariance landscape and its vulnerability to adversarial perturbations in both biological and artificial visual systems. SnS frames these transformations as bi-objective optimization problems. To probe invariance, SnS seeks image perturbations that maximally alter the representation of a reference stimulus in a given processing stage while preserving unit activation. To probe adversarial sensitivity, SnS seeks perturbations that minimally alter the stimulus while suppressing unit activation. Applied to convolutional neural networks (CNNs), SnS revealed image variations that were further from a reference image in pixel-space than those produced by affine transformations, while more strongly preserving the target unit’s response. The discovered invariant images differed dramatically depending on the choice of image representation used for optimization: pixel-level changes primarily affected luminance and contrast, while stretching mid- and late-layer CNN representations altered texture and pose respectively. Notably, the invariant images from robust networks were more recognizable by human subjects than those from standard networks, supporting the higher fidelity of robust CNNs as models of the visual system. Overall, SnS offers a powerful tool to uncover the invariant manifold of a unit, moving beyond tests with predefined transformations and advancing our understanding of generalization and robustness in visual systems.

## 1 Introduction

Both visual neuroscience and deep learning seek to understand image processing systems composed of millions of interacting functional units (biological or artificial neurons), whose activity patterns are shaped by their experience with natural image statistics [1, 2]. This common goal raises a fundamental question in both fields: which combination of image features do visual neurons become tuned for? Traditionally, this question has been addressed by developing feature visualization approaches that discover the “preferred” stimuli that maximally activate a given unit within the network - often

referred to as the unit’s most exciting images (MEIs) [3, 4, 5, 6, 7]. MEIs however only reveal a few instances within the vast set of images that strongly activate a given unit [8, 9], offering poor insight into the manifold of transformations under which the unit’s activity remains invariant.

To overcome this limitation, we developed Stretch-and-Squeeze (SnS), an unbiased, model-agnostic method to probe visual invariance in both artificial and biological neurons. Based on previous work [10], SnS optimization exploits evolutionary algorithms, but with a different objective: to find invariant or adversarial images, rather than MEIs. SnS integrates the search for invariant images and adversarial examples within a unified optimization objective, generating image perturbations starting from a reference image (e.g., a MEI of the unit). For invariance, SnS explicitly looks for images that are maximally distinct from this reference stimulus in the representation of a chosen visual processing stage, while preserving the response of a target unit upstream. This allows SnS to characterize the invariance manifolds of the selected unit when the representation of an effective stimulus (either a MEI or a natural image) is stretched at different processing stages. As a result, our approach reveals the actual image variation axes the unit is able to tolerate, providing a richer, more veridical description of its invariance landscape, as compared to traditional tests based on predefined (e.g., affine) transformations [11, 12, 13].

In our study, we carried out comprehensive tests of the effectiveness of SnS to achieve the goals listed above, using a popular CNN (ResNet50) as a benchmark. We found dramatic differences among the invariance landscapes of ResNet50 readout units when stretching the representations at different depths of the processing hierarchy. Stretching early, middle, and late representations yielded invariant images that differed from the reference (and among themselves) in terms of luminance/contrast, texture, and pose, respectively. Such qualitative differences were confirmed by measuring how the three classes of invariant images were accurately discriminated by Support Vector Classifiers (SVCs) in the pixel space. We also discovered important differences in these hierarchical invariances between the standard and an adversarially trained version of the network, with invariant images from the robust network being more recognizable by human subjects and other *observer* networks. Finally, we checked the potential applicability of SnS to visual neuroscience experiments, by showing that the method works even if the experimenter can record the activity of just a small fraction of the units in the processing stage where the stretching is applied.

## 2 Related works

**Probing CNN representations: feature visualization, invariance, and adversarial examples.** The question of functional interpretability and feature visualization in deep learning has predominantly been approached by employing gradient-based optimization for image synthesis, taking advantage of the complete analytical description and differentiability of the network [3, 4, 5]. While most studies have focused on finding the MEIs of CNNs’ units, recent efforts have started to also explore their invariance landscape. A prominent approach relies on discovering *model metamers* - i.e., synthesized stimuli that match the internal representation of a reference (natural) image in a specific layer of a network [14]. Using metamers, Feather et al. [15] were able to demonstrate that standard CNNs display highly idiosyncratic invariances at the top layers of processing, with metamers being unintelligible to human observers or other neural networks. In contrast, CNNs trained to be robust to adversarial images (i.e., imperceptibly modified inputs capable of altering CNN object classification [16]) yielded metamers that were substantially more interpretable by human observers. This shows how invariance and adversarial vulnerability are conceptually related [17]. For instance, adversarially trained ("robust") networks can craft subtle image perturbations that not only fool the network but can also impair [18] or enhance [19] human object recognition, thereby reinforcing the perceptual parallels between robust network representations and invariance in human vision.

**Applications to visual neuroscience.** Advances in CNN interpretability have also impacted neuroscience, which has long suffered from a lack of effective methods to investigate the functional tuning of visual neurons in higher cortical areas. By leveraging the strong functional analogy between CNNs trained for image classification and the primate object recognition pathway (known as the *ventral stream* [20]), it is possible to create *digital twins* of the ventral stream using CNNs [21, 22, 23]. This allows employing gradient-based optimization to synthesize images that modulate the activity of biological neurons [24] or investigate properties like adversarial robustness in visual neurons [25]. The digital twin paradigm has been also extended beyond primate vision, informing models of the auditory cortex [26] and visual processing in rodents [27, 7, 28]. Importantly, related gradient-based

techniques [9] have also been used to map invariance in individual neurons, for instance, by looking for images that maximally activate a mouse primary visual neuron while being maximally distinct in pixel space [29].

Obviously, these approaches have an intrinsic limitation: they are only as good as the fidelity of the digital twin is in capturing the selectivity of visual neurons. To overcome this constraint, gradientless feature visualization methods like XDREAM have been developed, which successfully synthesize effective MEIs for units in both the primate ventral stream [30] and artificial neural networks [6] by relying on evolutionary algorithms. While recent work has also begun to explore the feature landscape around the MEIs obtained with XDREAM [31], current gradientless approaches have not been systematically applied to characterize the invariance of visual tuning in artificial and biological architectures. Hence the novelty of SnS, which is, to the best of our knowledge, the first gradientless approach to systematically infer the invariance manifolds of visual units.

### 3 Methods

#### 3.1 The Stretch-and-Squeeze algorithm

SnS consists of three key components: a generative model  $\psi$ , a test (or *subject*) network  $\phi$  and a gradient-free optimizer ( Fig. 1a). The generative model is a pretrained deep neural network [32] that maps  $n$ -dimensional vectors  $\xi^t \in \mathbb{R}^n$  (referred to as *codes*) to RGB images:  $\mathbf{x} = \psi(\xi) \in \mathbb{R}^{C \times H \times W}$ . Crucially, the generative model  $\psi$  was trained on naturalistic stimuli and embodies a powerful prior over the distribution of possible images. We use the Covariance Matrix Adaptation Evolutionary Strategy (CMA-ES) optimizer [33] to adjust the codes and iteratively improve on our objective. At each iteration  $t$  the optimizer yields a new set of codes  $\xi^{t+1} \in \mathbb{R}^n$ , which in turn is used to generate a new batch of images.

A core tenet of the SnS algorithm is the relational construction of its fitness function. We start by introducing two layer indices  $\kappa$  and  $\ell$  for our test network  $\phi$ , where for convention we include the input stage as  $\kappa = 0$ , and the measuring function  $\mathbf{a}^\ell = \Gamma(\mathbf{x}, \phi^\ell) \in \mathbb{R}^d$ , which returns the activations of all the  $d$  units in layer  $\ell$ , when the network  $\phi$  is presented with stimulus  $\mathbf{x}$ . We then identify a reference stimulus  $\mathbf{x}_{\text{ref}}$  from which we construct the pair of reference states  $(\mathbf{a}^\kappa, \mathbf{a}^\ell)$  as  $\mathbf{a}^{\kappa, \ell} = \Gamma(\mathbf{x}_{\text{ref}}, \phi^{\kappa, \ell})$ . Lastly, we introduce two optimization objectives as either the minimization  $\mathcal{L}_{\text{squeeze}}$  or maximization (i.e. minimization of the negative)  $\mathcal{L}_{\text{stretch}}$  of the euclidean distance of a given state  $\mathbf{a}^\kappa = \Gamma(\mathbf{x}, \phi^\kappa)$  from the corresponding reference state:

$$\mathcal{L}_{\text{stretch}}^\kappa(\mathbf{a}^\kappa, \mathbf{a}_{\text{ref}}^\kappa) = - \|\mathbf{a}^\kappa - \mathbf{a}_{\text{ref}}^\kappa\|_2, \quad \mathcal{L}_{\text{squeeze}}^\kappa(\mathbf{a}^\kappa, \mathbf{a}_{\text{ref}}^\kappa) = + \|\mathbf{a}^\kappa - \mathbf{a}_{\text{ref}}^\kappa\|_2. \quad (1)$$

The final bi-objective optimization problem is formulated as the simultaneous optimization of  $\mathcal{L}_{\text{stretch}}$  and  $\mathcal{L}_{\text{squeeze}}$  for a given choice of layer indices  $\kappa, \ell$  and reference states.

$$\Xi_{\text{SnS}} \equiv \arg \min_{\mathbf{x}} \left[ \mathcal{L}_{\text{stretch}}^\kappa(\Gamma(\mathbf{x}, \phi^\kappa), \mathbf{a}_{\text{ref}}^\kappa), \mathcal{L}_{\text{squeeze}}^\ell(\Gamma(\mathbf{x}, \phi^\ell), \mathbf{a}_{\text{ref}}^\ell) \right]. \quad (2)$$

Using this formalism, we can express the solutions for both the classical adversarial attack and stimulus invariance as the following special cases. First, we single out a target unit  $u_{\text{targ}}^\ell$  and identify its scalar activation  $a_{\text{ref}}^\ell = \Gamma_u(\mathbf{x}_*, \phi^\ell) \in \mathbb{R}$  as one of our anchor reference states, where we have introduced  $\mathbf{x}_{\text{ref}} = \mathbf{x}_*$  as the maximally exciting stimulus (i.e., the MEI) for our target unit, computed via 500 iterations of the XDREAM algorithm [30]. We can then express both the search for adversarial attacks or invariant stimulus for our target unit  $u_{\text{targ}}^\ell$  as the following two optimization problems:

$$\Xi_{\text{inv}} \equiv \arg \min_{\mathbf{x}} \left[ \mathcal{L}_{\text{stretch}}^{\kappa=0}(\mathbf{x}, \mathbf{x}_*), \mathcal{L}_{\text{squeeze}}^\ell(\Gamma_u(\mathbf{x}, \phi^\ell), a_{\text{ref}}^\ell) \right] \quad (3)$$

$$\Xi_{\text{adv}} \equiv \arg \min_{\mathbf{x}} \left[ \mathcal{L}_{\text{stretch}}^\ell(\Gamma_u(\mathbf{x}, \phi^\ell), a_{\text{ref}}^\ell), \mathcal{L}_{\text{squeeze}}^{\kappa=0}(\mathbf{x}, \mathbf{x}_*) \right] \quad (4)$$

This (dual) formalization reflects the fact that both invariance and robustness relate changes of high order representations to changes at the input level (Fig. 1b). However, we remark that a key flexibility

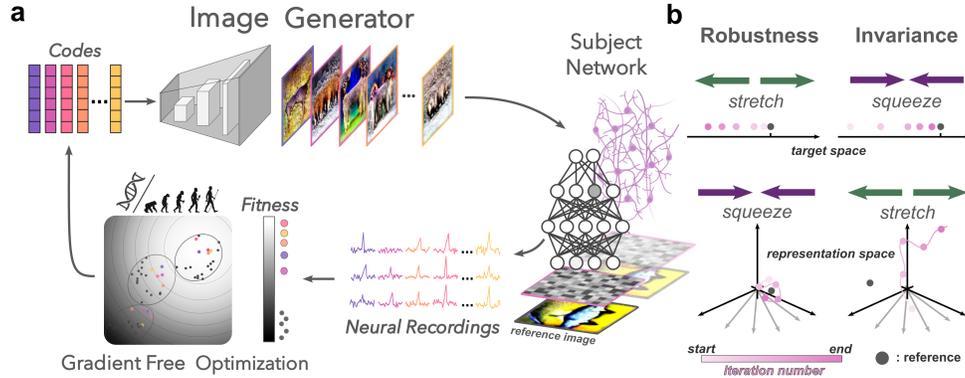


Figure 1: **The Stretch-and-Squeeze (SnS) algorithm.** (a) Overview of the SnS algorithm: candidate stimuli are synthesized from latent codes via the image generator, and the activation response of target units is recorded and used as fitness score by the optimizer, that adjusts a new set of codes. (b) Dual fitness objectives in SnS. To probe *invariance* (right), SnS maximizes stimulus distance from a reference in the representation space (stretch) and minimizes the variation in the activation of a target unit (squeeze). Conversely, to synthesize *adversarial examples* (left), SnS minimizes changes in the representation space, while maximizing the variation in the activation of the target unit.

of SnS is that its general formalization for the distance metric in the *stretch* and *squeeze* objectives allows computation not only in the input pixel space, but in general intermediate representation stages  $\kappa$  and  $\ell$  of the network. This allows probing invariance (or adversarial attacks) relative to different levels of feature abstraction.

For the characterization of invariance, we set this representation space at three distinct hierarchical levels within ResNet-50: (i)  $\kappa = 0$ , i.e., the input pixel space (denoted `low_level`), (ii) a mid-level convolutional layer (1<sup>st</sup> convolutional stage in layer 3, `mid_level`), and (iii) a deep convolutional layer (7<sup>th</sup> convolutional stage of layer 4, `high_level`). Our primary target units  $u_{\text{targ}}^\ell$  were chosen in the final readout layer (fully connected), as these units are typically selective for specific object categories and integrate information across the entire visual field, facilitating the interpretation of results. To demonstrate SnS’s generalizability beyond category-selective units, we also performed optimizations targeting units within `mid_level` and `high_level`, which are potentially closer to biological neurons in the visual system. In these cases, the input-side distance was computed in pixel space (i.e.  $k = 0$ ). For adversarial attack generation, we restricted our analysis to the configuration targeting readout layer units and defining input distance in pixel space. All SnS optimizations were conducted on the units of two specific instances of ResNet-50, our subject network: the standard ImageNet-pretrained model available in PyTorch [34], and a robust counterpart generated via adversarial training with an  $L_2$  perturbation norm constraint of  $\epsilon = 3$  [35].

Given that SnS fitness is defined relationally between both the target’s evoked activity and the input-level changes, it poses a multi-objective optimization problem. A priori, no unbiased trade-off between the competing metrics (*stretch* and *squeeze*) can be defined. Therefore, for each round of optimization, we sorted the fitness scores by organizing them in Pareto fronts [36], a method that groups together equally suitable solutions to the multi-objective problem. Initialization strategies and selection from the Pareto front were adapted based on the optimization goal:

**Adversarial Attacks:** The search was initialized from the MEI  $x_*$ . Solutions on the same Pareto front were selected for the next generation with uniform probability (random ordering). This promotes *exploration* along the front, preventing premature convergence to a single type of adversarial strategy.

**Invariance Experiments:** Initialization began from random normally distributed vectors. In this scenario, solutions on the same Pareto front were prioritized based on their proximity to the target unit’s reference activation level ( $\min_x \|\Gamma(\mathbf{x}, \phi^\ell) - \mathbf{x}_{\text{ref}}^\ell\|_2$ ). This exploitative strategy was motivated by the empirical difficulty of converging towards the target activation level while at the same time maximizing input distance.

For additional details regarding the computational experiments we refer to the Supplementary Material, section A.

### 3.2 Separability of the invariant images in the pixel representation

To assess whether the representation spaces where the stretching was applied yielded sets of invariant images that were distinguishable at the pixel level, we performed an image classification analysis. We generated invariant images for  $n = 77$  readout layer units (i.e., 77 ImageNet classes) by stretching the representations in three different processing stages: (i) `low_level` (pixel space), (ii) `mid_level`, and (iii) `high_level`. For each unit and each representation space condition, we performed 10 independent optimization runs with different random seeds, yielding a total dataset of  $77 \times 3 \times 10 = 2310$  images. Principal Component Analysis (PCA) was applied to the raw pixel representations of all images. The first  $k$  principal components were then used as input features to train a Support Vector Classifier (SVC) with a radial basis function kernel to discriminate the class to which the invariant images belonged (i.e., `low_level`, `mid_level`, or `high_level`). 5-fold cross-validation was used.

### 3.3 Representational Distance Analysis

To quantify how much the invariant images generated by SnS (when targeting a given processing stage) diverged from the reference MEI ( $x_*$ ) across the depth of ResNet-50, we computed the Euclidean distance between their activation vectors at each layer of the network. Again, 10 different invariant images were produced for  $n = 77$  distinct readout units by stretching representations at layers `low_level`, `mid_level`, and `high_level`. We calculated the mean ( $\pm$  SEM) Euclidean distance between each of these  $77 \times 10$  invariant images and the corresponding reference images at every convolutional, pooling and linear layer of the network, separately for each stretching stage. These distances were compared against three control metrics:

**Reference Variability distance:** Mean distance between all pairs drawn from 10 independent XDREAM runs for each unit ( $n = \binom{10}{2} = 45$  pairs per unit,  $n = 77$ )

**Within-Category distance:** Mean distance between pairs of images randomly sampled from the same ImageNet category (mean over 1000 categories, 10 images/category).

**Between-Category distance:** Mean distance between pairs of images randomly sampled across all ImageNet categories (using the same total number of pairs as the within-category control).

### 3.4 Interpretability of the invariant images by humans and observer networks

To assess whether the invariant images produced by SnS retain sufficient information for correct classification by other visual systems, we tested the ability of humans or other *observer* neural networks to classify them, using a 12 alternative forced choice task (AFC). For observer networks, we used several distinct ImageNet-pretrained architectures, both standard (AlexNet, VGG16, ResNet18, ResNet101, ConvNext base, RegNetX-16GF), obtained from PyTorch [34], and robust (ResNet-18 [37], ResNet-50 [37], ConvNext base [38], all  $L_\infty$ ,  $\epsilon = \frac{4}{255}$ ), obtained from RobustBench [39].

We chose 12 ImageNet categories (goldfish, chickadee, frilled lizard, Dungeness crab, Staffordshire bull terrier, cicada, candle, grand piano, minibus, reflex camera, soccer ball, cup). For each category, we selected 10 reference images  $x_{\text{ref}}$  and, for each of them, we generated one invariant image by stretching its representation in `low_level`, `mid_level`, and `high_level` in both standard and robust ResNet-50. This yielded 6 kinds of invariant images per reference, where invariance was always defined relative to the readout unit corresponding to the reference image’s category. For each category, we also included in the classification pool the 10 reference images and 20 MEIs (generated with XDREAM) of the readout unit corresponding to the chosen category (10 for the standard and 10 for the robust network). This yielded a total of 12 categories  $\times$  10 randomly sampled images  $\times$  9 stimulus types (i.e., 1 reference image, 2 MEIs and 6 invariant images). Observer networks had to classify all these 1080 images; humans had to classify a subset of 540 images (i.e., 5 rather than 10 randomly sampled images per category). Classification accuracy (12-AFC percent correct) was calculated for each of the 9 stimulus types.

To assess human recognizability of invariant images, we recruited 25 human participants using the online Prolific platform. After clicking on a fixation cross at the start of each trial, participants were presented with an image for 50 ms at  $\sim 10^\circ$  of visual angle. After a 20ms blank screen, a backward mask image (random RGB pixel noise) was then presented for 200ms. Images were presented in random order. Participants then clicked one of 12 category buttons that subsequently appeared in a circular arrangement, with a 10-second timeout (Fig.5a). To avoid directional bias, the sequence of

category buttons was randomly rotated for each trial. Further details on the task and the statistical analysis are reported in Supplementary Material, Section B.

## 4 Results

### 4.1 SnS: an effective dual formulation for invariance and adversarial attacks

**SnS generates effective adversarial and invariant images.** We first validated the SnS framework by generating invariant and adversarial images for a sample of 77 readout units of a robust ResNet50 using their MEIs as reference images  $x_*$ , and applying the stretching (to achieve invariance) or the squeezing (to achieve adversarial images) to the pixel representation (see Section 3.1). As shown in Fig. 2, SnS successfully generated effective adversarial examples (top left). These stimuli substantially suppressed the activation of the readout units relative to their MEIs (mean reduction of  $111\% \pm 7\%$ ) being displaced from the MEIs by a mean  $L_2$  distance of  $72 \pm 12$  pixels. Such pixel budget reflects both the network robustness and our stringent unit-silencing objective - a stricter criterion than mere misclassification. Consistent with [40, 18], perturbations were semantically relevant, not noise-like (Fig.2).

The invariant images generated by SnS (bottom right) were also very effective, yielding only a minor drop of the units' activation relative to the MEIs (mean reduction of  $34\% \pm 11\%$ ). At the same time, they substantially departed from the MEIs, with a mean  $L_2$  pixel distance ( $271 \pm 32$  pixels) that considerably exceeded the median distance between ImageNet images, as reported by [18]. More impressively, SnS discovered image transformations that were more extreme (in terms of pixel distance) than those achieved by applying to the MEI standard data augmentations (colored dots/lines), while altering the activation of the readout units substantially less, compared to the most extreme augmentations (darker dots). This indicates that SnS can discover the actual axes of image variation a unit learned to tolerate, exploring invariance manifolds way more precisely than traditional tests with predetermined (e.g., affine) transformations.

Finally, we note that SnS successfully generated invariant images also for units in intermediate convolutional layers (Supplementary Material, Section C). This demonstrates its applicability to the analysis of latent representations that are closer to those of biological neurons.

### 4.2 SnS reveals layer-specific invariant manifolds

Next, we compared the kinds of invariant images that SnS discovered by stretching the representations of the MEIs of the 77 readout units at different stages of processing along a *robust* ResNet-50 - i.e., in the pixel space (as already done in Fig. 2), as well as in an intermediate (layer 3, conv 1) and a deep (layer 4, conv 7) convolutional layer. Qualitatively, this yielded radically different sorts of invariances (Fig.3a). Stretching in the pixel-space produced luminance and contrast changes; stretching mid-level representations affected texture and color; stretching high-level representations produced abstract variations like viewpoint changes or multiple object instances. The same qualitative differences were observed when SnS was used to produce invariant images for 77 readout units of a *standard* ResNet-50 (Fig.3b), although, in this case, the images appeared to be less interpretable than those obtained for the robust counterpart. Also, the images lacked the high-frequency patterns that are

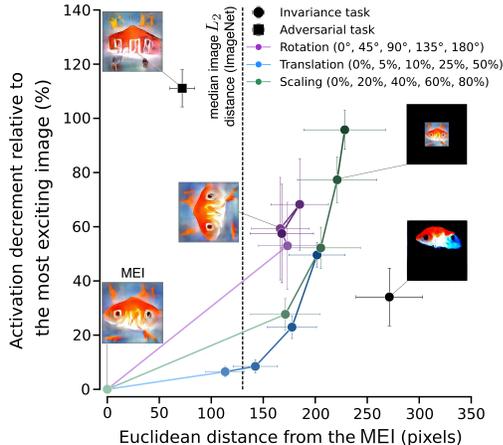


Figure 2: **SnS generates effective adversarial and invariant examples.** The average activation reduction ( $\pm$  SD) of  $n = 77$  readout units of a robust ResNet50 (relative to their MEIs) is plotted against the  $L_2$  pixel distance from the MEIs, when the latter were transformed with SnS to yield either adversarial or invariant images, or were subjected to affine transformations (rotation, translation and scaling). The pixel distance refers to  $224 \times 224$  RGB images with values between 0 and 1.

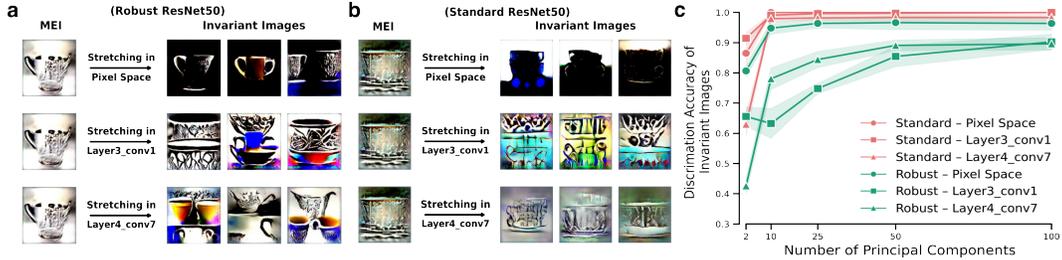


Figure 3: **SnS discovers layer-specific invariances.** (a) Example MEI and associated invariant images obtained for a readout neuron (“cup”) in robust ResNet-50 and computed for three different choices of the stretching representational stages (rows). Multiple results are shown for different random initialization seeds (columns) (b) Same as (a) but for a standard ResNet-50. (c) Accuracy of a SVC in discriminating the three classes of invariant images produced by stretching pixel-, mid-, and high-level representations as a function of the number of principal components fed to the classifier.

typically found in the MEIs of the units of standard networks [41], suggesting that the SnS generator strongly regularizes the image search towards natural statistics.

To quantify these observations, we applied PCA to the invariant images yielded by SnS. We then used the first  $k$  principal components as feature vectors to train an SVC to classify the invariant images according to the layer where the SnS stretching was applied. As shown in Fig.3c, a few components were enough to yield above chance performance (i.e.,  $> 33.3\%$  correct), and a few tens of components guaranteed a discrimination accuracy virtually perfect for the standard network (red) and above 80% correct for the robust one (green).

To characterize how the invariant images generated by stretching at a given processing stage were represented across the network, we measured their  $L_2$  distance from their reference MEIs in the representation provided by each layer of the network (Fig.4; this distance was normalized by the mean, within-category distance of Imagenet images; see Section 3.3). By construction, invariant images that were generated to be maximally different from the MEI in a specific layer had the largest normalized distance in that layer, as compared to invariant images obtained by applying the stretching in other layers. This distance consistently surpassed the one among multiple MEIs of the same unit, demonstrating SnS’s efficacy in exploring larger spans of a unit’s invariance manifold, as compared to searching multiple times for different MEIs. Interestingly, the representational distance between the invariant images and their MEIs was preserved in the layers that were adjacent to the one where the stretching was applied, revealing a clear hierarchical pattern. Stretching in the pixel-space yielded invariant images that remained dissimilar from the MEIs in the first layers of the network; stretching in mid and deep convolutional layers yielded images that were more different from the

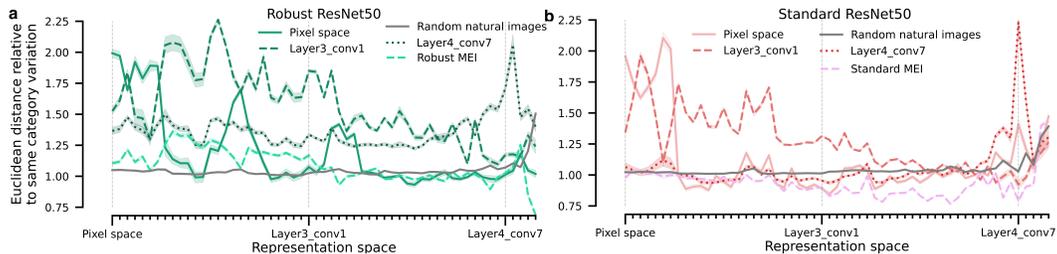


Figure 4: **Representation of the invariant images across the feedforward hierarchy.** (a) Normalized average distance between the invariant images generated by SnS and their reference MEIs across the different stages of a robust ResNet50. Different lines indicate the different stretching layers used in the optimization (also reported in the  $x$  axis). The average distance between randomly selected natural images from Imagenet (solid gray line) and the average distance between multiple MEIs generated via XDREAM are reported for comparison. (b) Same as (a) for standardly trained network.

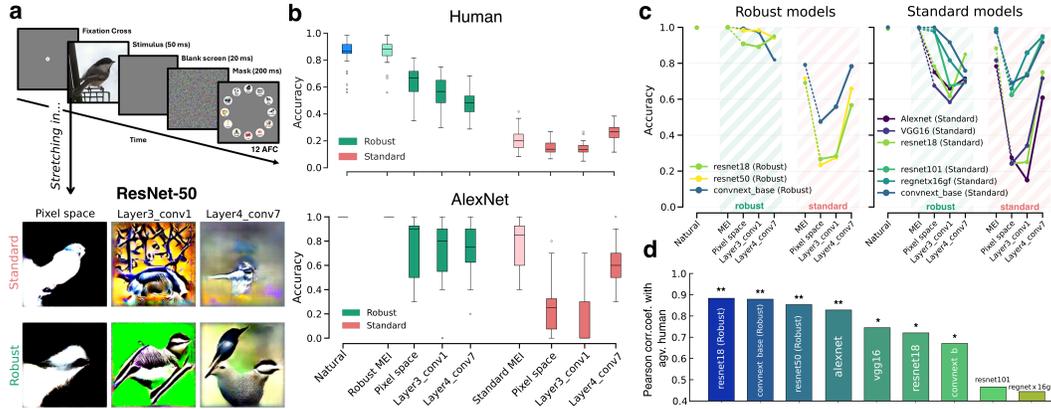


Figure 5: **Interpretability of the invariant images by humans and other networks.** (a) Illustration of the human classification task with the invariant images generated by SnS for standard and robust networks. (b) Classification accuracies of humans (top; averaged across subjects and categories) and AlexNet (bottom; averaged across categories) displayed for each experimental condition. (c) Classification accuracies across multiple robust (left) and standard (right) networks. (d) Correlation coefficient between average human performances and the performances of each observer model.

MEIs in the central and final portion of the network. These hierarchical trends were sharper for the standard network than the robust one, suggesting a smoother, more gradual processing of image transformations.

**Hierarchical invariances are robust to representation subsampling.** To assess SnS’s applicability to neuroscience, where only a small neuronal subpopulation is typically recorded from a visual area, we generated the invariant images using heavily subsampled hidden layer representations. Both qualitative inspection and representational distance analysis revealed that the resulting invariances closely resembled those obtained using the full layer representation, supporting SnS’s potential for neuroscience experiments (see Supplementary Material, Section D for further details).

### 4.3 Invariant images for robust and standard networks are perceptually different

We next evaluated how recognizable the invariant images generated by SnS were by humans and other CNNs. Based on [15], we hypothesized that the invariant images for the readout units of a robust network would be more semantically coherent and recognizable than those of a standard network. In our tests, invariant images for both kinds of networks were generated by stretching from the low- (pixel), mid- and high-level representations of the same set of natural images sampled from 12 different Imagenet categories (Fig.5a). Human subjects performed a 12 AFC classification task on these invariant images, their references, and MEI controls (Fig.5a, top; see Section 3.4 for a detailed description). Invariant images from the robust network (all generation stages) were significantly more recognizable by humans than their standard network counterparts (Fig.5b, top;  $p < 0.001$ , Supplementary Material Table 1 ). Invariant images from the robust network were more recognizable when stretching was applied to earlier-layer representations, with pixel-space stretching yielding the highest discrimination accuracy. Interestingly, an opposite trend was observed for the standard network: invariant images from a late layer were more recognizable than those from a middle layer or from pixel space ( $p < 0.001$  for each of the above comparisons: see Supplementary Material Table 1).

This classification task was replicated with various *observer* CNNs (both standard and robust). As exemplified by the performances obtained for AlexNet (Fig.5a, bottom), the similarity with the accuracy trends found in humans was remarkable (see Supplementary Material Table 1 ). More in general, we found that the patterns of accuracy of the robust networks were very similar (Fig.5c, left) and highly correlated (Fig.5d) with those of humans. Instead, in the case of the standard networks, smaller models (AlexNet, VGG16, ResNet18) mirrored human behavior, while larger, high-performing models (ResNet101, ConvNeXt-Base, RegNetX-16GF) showed no classification difference between the invariant images obtained for the robust and standard ResNet50 (Fig.5c, right

and Fig.5d). Nevertheless, all the networks displayed, as humans did, opposing accuracy trends based on the depth of ResNet50 where the stretching was applied: accuracy decreased across layers for the invariant images generated for robust ResNet50, while it increased in the case of the images generated for the standard ResNet50.

## 5 Discussion

We introduced Stretch-and-Squeeze, a novel, model-agnostic framework that reconceptualizes how we probe invariance of visual representations. Instead of testing pre-defined transformations (e.g., affine) or imposing strict representational identity (as with metamers [15]), SnS discovers the actual manifold of transformations a unit tolerates by maximizing stimulus dissimilarity in a chosen representational space while preserving as much as possible the unit’s activation. This exploration of functionally equivalent, yet representationally distinct inputs allows SnS to map a broader, potentially more “ecologically” relevant stretch of a unit’s response manifold and its boundaries. The hierarchical application of SnS shows that CNN units in readout layers are not necessarily that robust to the standard transformations (scaling, translation, etc) often applied to probe invariance [20, 11, 12, 13] (Fig. 2). Rather, they strongly tolerate image changes along radically different axes of abstraction (Fig. 3): from low-level properties (e.g., luminance) to mid-level features (e.g., texture) and, finally, to high-level semantic variations (e.g., object pose).

Our comparison of the invariant images generated for standard and robustly-trained ResNet50 extends prior work on the topic [15] with an innovative characterization. The superior semantic coherence and cross-system recognizability (by humans and other CNNs) of the invariances obtained for robust networks strongly corroborates that adversarial training sculpts representations towards human-aligned perceptual features [41, 18, 15]. More strikingly, standard networks revealed a complex relationship between model capacity and the nature of perceived invariance (Fig.5c). While small standard models often aligned with human judgments, larger, higher-performing standard models demonstrated a surprising ability to recognize idiosyncratic invariances generated from other standard networks, even when these were less understandable to humans or robust models. This is possibly a result of the previously described tradeoff between accuracy and robustness [40]: a model highly optimized for accuracy might learn more brittle invariances, leading to less robust representations.

Another interesting finding, consistent across all human and CNN observers, was the opposite trend for the recognizability of the invariant images generated for the robust and the standard network as a function of stretching depth (Fig.5c). This divergence may stem from fundamental differences in the way these networks process and "digest" image variations, whose representation appears to be diluted across wider stretches of adjacent layers in the robust architecture (see Fig.4).

Finally, the model-agnostic and gradientless nature of SnS positions it as a powerful new tool for neuroscience, particularly for systems where high-fidelity “digital twins” may be hard to develop. Our demonstration of SnS’s efficacy with simulated sparse neural recordings directly addresses a critical experimental constraint, paving the way for in vivo applications. This could lead to the discovery of new, hierarchical invariances in biological visual systems, extending our understanding of visual object representations in both primates [20, 30] and other species [42, 43, 44, 28, 45, 46, 47].

**Limitations and broader impacts.** While SnS has proven effective in uncovering novel, meaningful invariances, its current evolutionary, model-agnostic implementation presents opportunities for refinement and broader application. One avenue is to develop hybrid SnS variants that incorporate gradient-based optimization where network differentiability is available. Such an approach could potentially accelerate convergence or fine-tune the search, leveraging system-specific knowledge while retaining model-agnostic capabilities for biological or other black-box targets. The hierarchical invariance characterization we introduced could be significantly expanded. Future work should apply SnS across a more diverse array of subject networks used for generation and “observer” networks for evaluation, extending beyond CNNs to architectures with different inductive biases, such as Vision Transformers. This would help determine the generality of the observed hierarchical patterns and how architectural choices influence learned invariances. Furthermore, our finding that invariant images of standard ResNet50 are better understood by larger, high-performing standard observer networks warrants deeper investigation. Unraveling the mechanisms behind this phenomenon could shed light on the complex interplay between model scale, feature learning, and generalization. Finally, while SnS’ robustness to significant representational subsampling provides a strong foundation for its

neuroscientific applicability, the crucial next step is direct validation in biological systems. Future research should focus on adapting SnS for in vivo experiments, addressing challenges such as neural signal variability, and optimizing stimulus presentation within experimental constraints to truly bridge the gap between computational characterization and biological understanding of invariance.

# Supplementary Material

## A Additional details regarding computational experiments

### A.1 Pseudocode for the SnS algorithm

Here we present the pseudo-code (Algorithm 1) for the SnS algorithm. We used library implementations for the CMA-ES gradient-free optimizer [48] and the Pareto front computation [49]. Importantly, we introduce the following `should_stop` criteria for our experiments.

Optimizations were terminated by either reaching a maximum iteration limit (500) or satisfying an early stopping rule derived from the target unit’s response statistics to a large dataset of natural images (i.e., the full ImageNet train set,  $n \approx 1.2$  million). In particular, given the vastness of the dataset, we assumed that the activation elicited by the most effective image in the natural dataset ( $a_{\max\_nat}^\ell$ ) represented a conservative threshold for an image to be considered *invariant*. On the other hand, the activation from the least effective image ( $a_{\min\_nat}^\ell$ ) represents a non-response level, functioning as a threshold for an image to be considered *adversarial*. Therefore, when searching for adversarial images, termination was reached when a large fraction (i.e.,  $\geq 90\%$ ) of the current population of images elicited activations at or below  $a_{\min\_nat}^\ell$ , thus leading to a non-responsive status. Conversely, when looking for invariant images after initialization from random noise, it was terminated when a large fraction (i.e.,  $\geq 90\%$ ) of the population of images yielded activations equal or above  $a_{\max\_nat}^\ell$ , thus reaching the desired high-activation regime.

---

#### Algorithm 1 SnS gradient-free optimization algorithm

---

**Require:**  $\phi$  ▷ Target network  
**Require:**  $\psi$  ▷ Image generator  
**Require:**  $T \geq 0$  ▷ Maximum number of iterations  
**Require:**  $\kappa, \ell \in \{0, 1, \dots, \text{depth}(\phi)\}$   
**Require:**  $\mathbf{x}_{\text{ref}} \in \mathbb{R}^{C \times H \times W}$  ▷ Common choice is  $\mathbf{x}_{\text{ref}} = \mathbf{x}_*$

$t \leftarrow 0$   
 $\mathbf{a}_{\text{ref}}^{\kappa, \ell} \leftarrow \Gamma(\mathbf{x}_{\text{ref}}, \phi^{\kappa, \ell})$   
 $\xi_0 \leftarrow \mathcal{N}(0, \sigma_{\text{init}})$   
**while**  $t \leq T$  **and** `!early_stop` **do**  
     $\mathbf{x}_t = \psi(\xi_t)$  ▷ Compute new candidate stimuli  
     $\mathbf{a}^{\kappa, \ell} = \Gamma(\mathbf{x}_t, \phi^{\kappa, \ell})$  ▷ Measure activations in target layers  $\kappa, \ell$   
     $\mathcal{L}_{\text{stretch}}^\kappa \leftarrow - \|\mathbf{a}^\kappa - \mathbf{a}_{\text{ref}}^\kappa\|$   
     $\mathcal{L}_{\text{squeeze}}^\ell \leftarrow + \|\mathbf{a}^\ell - \mathbf{a}_{\text{ref}}^\ell\|$   
     $\mathcal{P} \leftarrow \text{compute\_pareto\_front}(\xi_t, \mathcal{L}_{\text{stretch}}^\kappa, \mathcal{L}_{\text{squeeze}}^\ell)$   
     $\xi_{t+1} \leftarrow \text{evolve}(\mathcal{P}, \xi_t)$  ▷ Evolution strategy CMA-ES optimization  
     $t \leftarrow t + 1$   
    `early_stop`  $\leftarrow \text{should\_stop}(\mathbf{a}^\kappa, \mathbf{a}^\ell)$   
**end while**

---

### A.2 Experimental Hyperparameters

The generative model  $\psi$  is instantiated as a pretrained deep neural network variant, specifically the `fc7` configuration from [32].

The CMA-ES optimizer is configured with the following hyperparameters:

- **Initial Step Size** ( $\sigma_0$ ): 1.0
- **Population Size**: 50

The initial step size  $\sigma_0$  determines the covariance of the sampling distribution at the onset of optimization, effectively controlling the exploration radius in the latent space.

The population size specifies the number of candidate solutions (i.e. codes) evaluated per iteration.

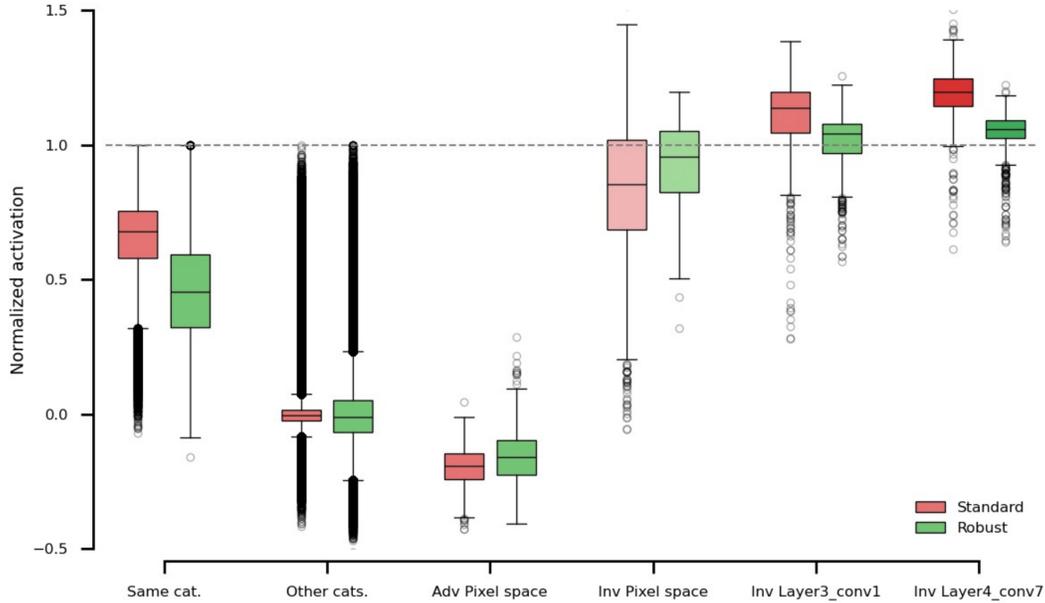


Figure S1: **SnS images fall into functionally relevant activation regimes.** Distributions of readout target unit activations (normalized by peak natural stimulus response) for SnS-generated images: invariant (stretched at pixel, mid-, high-levels) and adversarial (pixel-space). Each distribution pulls together all optimizations performed in that specific condition ( $n = 77$  target neurons  $\times$  10 random seeds = 770 experiments per condition). Comparison with natural image activations (same/different category) confirms that SnS invariant images elicit strong target responses, while adversarial examples achieve significant suppression.

Initial codes were randomly sampled from a Gaussian distribution  $\mathcal{N}(0, \sigma_{\text{init}})$ . The initial standard deviation,  $\sigma_{\text{init}}$ , was set to  $\sqrt{0.01 \times \mathbb{E}[|\xi_{\text{ref}}|]}$ , where  $\mathbb{E}[|\xi_{\text{ref}}|]$  is the mean absolute value of the reference code. For invariance experiments benchmarked against natural images, where direct reference codes were unavailable, initial codes were drawn from a standard normal distribution ( $\mathcal{N}(0, 1)$ ).

### A.3 Optimization Convergence

Although some optimizations terminated at maximum iterations (adversarial: 97.92% (robust), 93.12% (standard); invariant: 63.38%, 35.71%, 21.30% (robust: low\_level, mid\_level, high\_level); 86.75%, 32.60%, 5.84% (standard: low\_level, mid\_level, high\_level)), final populations consistently reached functionally relevant activation regimes. Most invariant images achieved activations within the top 0.1% of natural image responses, being comparable or more extreme than the readout activation of images of the same category encoded by the readout target unit (Fig. S1). Adversarial examples drove activations to the lowest 1<sup>st</sup> percentile of the natural images, thus significantly suppressing target unit activation w.r.t. average activation with a natural image.

### A.4 Computational Resources

All experiments reported in this work were executed on a Dell Precision 7960 Tower (Ubuntu 22.04.5 LTS) with the following configuration:

- **CPU:** Intel (R) Xeon(R) w5-3425 (24 cores, 48 threads);
- **GPU:** NVIDIA RTX A6000 (48 GB GDDR6);
- **System Memory:** 128 GiB;
- **Storage:** 8.2 TB.

The SnS source code is openly available in the GitHub repository at <https://github.com/zoccolan-lab/SnS>. Moreover the experimental data is accessible via [50].

## B Additional details regarding human subjects experiment

Participants for the human experiment were recruited using the online Prolific platform under a preexisting IRB protocol. Participants were compensated at an overall rate calibrated to \$15 USD per hour. To facilitate engagement with the task, compensation included a 1-cent bonus for each correct trial, and participants were shown a green check after correct responses or a black “X” otherwise. The typical bonus amount that participants would receive was overestimated based on pilot data. To adjust for this discrepancy, after recruitment for the study was complete, participants were uniformly provided with an additional bonus to bring compensation to the \$15/hour level.

The main task was preceded by a screening phase, in which participants classified 24 natural images (2 per category), and were allowed to proceed if they correctly classified  $\geq 1$  image per category and  $\geq 16$  in total (maximum 3 attempts allowed per participant). Data from the screening phase were not included in any analyses. Participants who did not pass the screening phase were compensated for the time spent during the screening process.

The experiment posed minimal risks to participants. It is unlikely but conceivable that the brief presentation of images and masks could be hazardous to a subset of individuals with photosensitive epilepsy: out of an abundance of caution, a warning was shown in the description of the task on the Prolific website to be viewed by participants before joining the study. The title and description of the task on the Prolific website are reproduced as follows:

**Title:** *Identify objects in photos, earn roughly \$4.00 bonus for accurate responses*

**Description:**

*View 540 photos and identify the animal or other object in each photo. Earn a bonus for each accurate response, typically around \$4.00 USD in total (maximum \$5.40).*

**PLEASE NOTE:**

- 1. The task begins with a screening phase you must pass with a certain accuracy level before starting the experiment proper. You may re-attempt the screening phase up to 2 times if you wish, but you will only be compensated for a maximum of 5 minutes spent in the screening phase.*
- 2. Please avoid resizing the task window during the experiment, as this can break the task. It's best to maximize your browser window and get comfortable before starting.*

**WARNING:** *this task contains bright, rapidly flashing images: it is not suitable for individuals with photosensitive epilepsy.*

Please see Fig. S2 for the additional in-task instructions provided to participants. The task was implemented using the JsPsych library [51] and the JsPsychPsychophysics plugin [52].

### B.1 Statistical Analysis for Classification Experiments

Analysis of data from classification experiments focused on identifying accuracy differences between 9 image types: natural images, MEI's from robust and vanilla networks, and invariant images from 3 different layers for both network types. We fitted a generalized linear mixed model (GLMM) with a logistic link function for binary correct vs. incorrect trial responses, including random intercepts for participants and semantic image categories. In addition to data from human participants (Fig.5b, top), we applied the same methodology to separately analyze AlexNet's responses to the full complement of 1080 images (see Section 3.4), but with GLMM random intercepts for semantic image categories only (Fig.5b, bottom).

For both human participants and AlexNet, We used an omnibus test (type II Wald chi-square) to assess overall differences among the 9 image types, followed by 8 planned contrasts via estimated marginal means: (i) robust vs. vanilla MEIs, (ii) robust vs. vanilla invariant images averaged across the 3 representation layers, and (iii-viii) within each network type, all 3 possible pairwise comparisons between invariant types at different layers. Our analysis employed the Benjamini-Hochberg procedure to control the false discovery rate under multiple comparisons [53]. The results

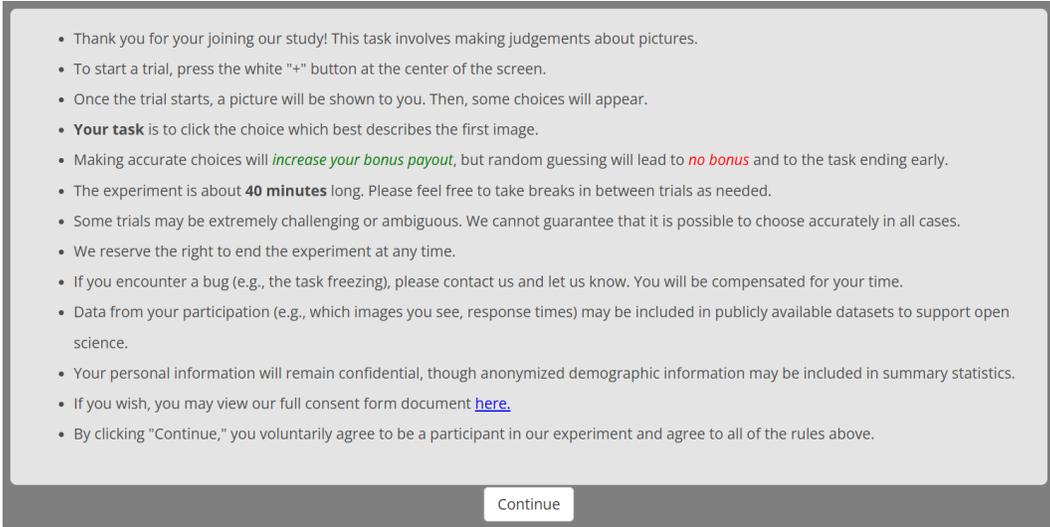


Figure S2: **Instructions shown to human participants.** These instructions were shown to the participants at the beginning of the task.

Table 1: **Statistical results from omnibus test (type II Wald chi-square) and post-hoc pairwise comparisons to test differences among stimulus types in the human/model image classification experiments.** All  $p$  values are adjusted using the Benjamini-Hochberg correction (separately for humans and AlexNet), with significant values in bold.

Comparison	Human subjects		AlexNet	
	z-value	p-value	z-value	p-value
Omnibus test (effect of stimulus type)	-	<b>&lt;0.0001</b>	-	<b>&lt;0.0001</b>
Average Robust Stretch vs Average Standard Stretch	36.63	<b>&lt;0.0001</b>	10.02	<b>&lt;0.0001</b>
Robust layer4_conv7 vs Robust layer3_conv1	-4.94	<b>&lt;0.0001</b>	1.06	0.39
Standard layer4_conv7 vs Standard layer3_conv1	-8.41	<b>&lt;0.0001</b>	7.28	<b>&lt;0.0001</b>
Robust layer4_conv7 vs Robust pixel space	-9.52	<b>&lt;0.0001</b>	-0.63	0.60
Standard layer4_conv7 vs Standard pixel space	7.49	<b>&lt;0.0001</b>	5.44	<b>&lt;0.0001</b>
Robust layer3_conv1 vs Robust pixel space	-4.66	<b>&lt;0.0001</b>	-1.68	0.15
Standard layer3_conv1 vs Standard pixel space	-0.99	0.32	-2.43	<b>0.030</b>
Robust MEI vs Standard MEI	32.96	<b>&lt;0.0001</b>	0.22	0.83

of the planned comparisons for both humans and AlexNet are provided in Table 1. Notably, the analysis of human data included a larger number of data points overall: while each human only viewed 540 images as opposed to the AlexNet’s 1080, there were 25 human participants responding independently unlike the single response per image from AlexNet.

## C SnS is effective in targeting units in hidden layers

Building upon our findings in Section 4.1, we extended the application of SnS to units within hidden convolutional layers. Unlike output layer units tuned to predefined classes, hidden units are hypothesized to function more analogously to biological visual neurons, responding to intermediate-level features rather than complete objects. This makes them compelling targets for deciphering learned internal representations.

To investigate this, we conducted SnS invariance experiments targeting 50 distinct units in both layer3\_conv1 and layer4\_conv7 of the ResNet50 network, using pixel space as the input representation. For each unit, we performed 10 optimization runs initiated with different random seeds. In the standard network, SnS terminated by early stopping in 67.20% of runs at layer3\_conv1 and

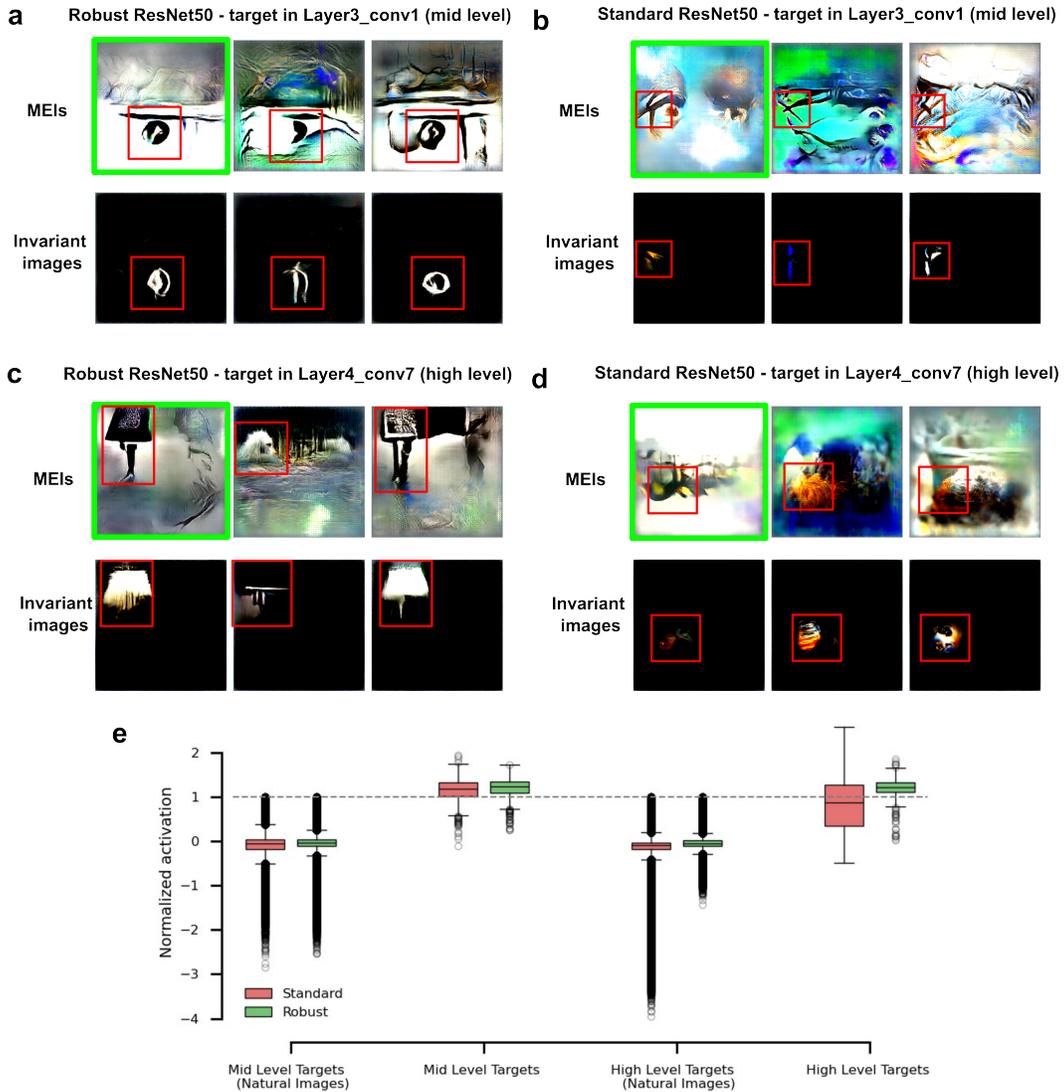


Figure S3: **SnS is effective with targets in hidden layers** (a) Comparison of MEIs (top row) with SnS-generated invariant images (bottom row) for a representative target unit in Layer3\_conv1 of robust ResNet50. SnS images, synthesized to be invariant in pixel space, were generated using the green-highlighted MEI as a reference. Red rectangles identify the unit’s key responsive regions. This comparison highlights SnS’s ability to distill the core visual features a unit detects, offering clearer insight than MEIs alone. (b-d) The same SnS-based feature visualization approach (as in (a)) is applied to different representative target units: (b) Layer3\_conv1 of a standard ResNet50; (c) Layer4\_conv7 of a robust ResNet50; and (d) Layer4\_conv7 of a standard ResNet50. (e) Final activation distributions for the target units from the SnS optimization in hidden layers. As in Fig. S1, convergence performance is compared to natural image statistics

27.40% at layer4\_conv7; in the robust network, convergence was 87.80% and 89.40%, respectively. Importantly, even in non-converging instances, the optimization process often guided the input towards stimuli that elicited high unit activations, consistent with functionally relevant regimes identified through natural image statistics on ImageNet (Fig. S3e).

Qualitatively, the SnS-generated invariant images served as effective “feature visualizers.” This was particularly insightful as the precise selectivity of these hidden units is unknown a priori. The resulting visualizations clearly delineated the specific image features to which a target neuron was

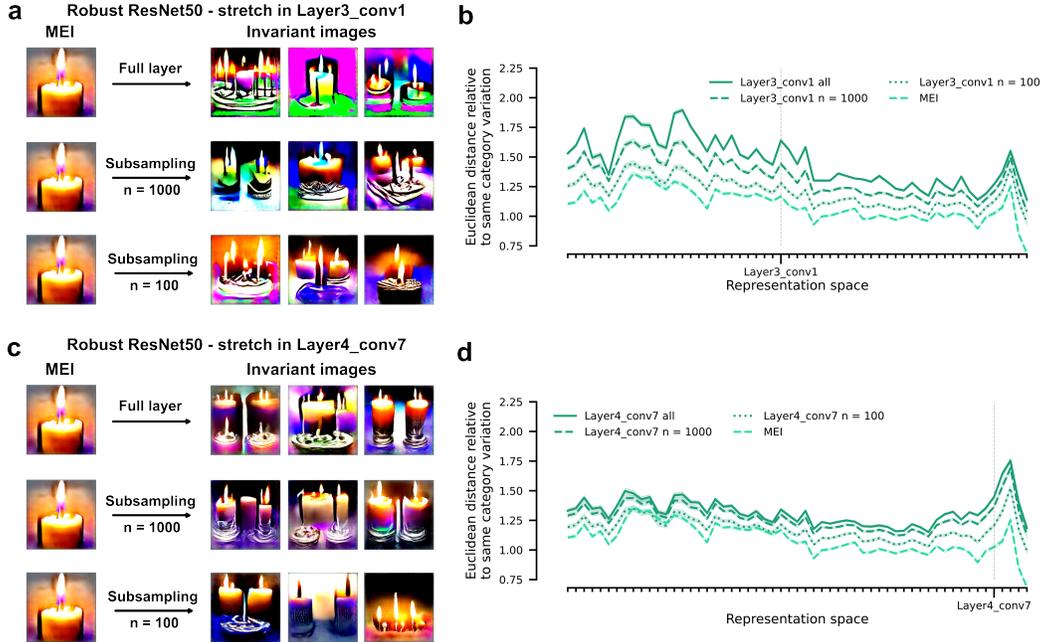


Figure S4: **SnS invariance is robust to subsampling of the representation space** (a) Example MEI and associated invariant images obtained for a readout neuron (“candle”) in robust ResNet-50 stretching the mid\_level layer (i.e. Layer3\_conv1). Each row represents examples for different levels of subsampling: using the full layer (1<sup>st</sup> row), subsampled spaces of size  $n = 1000$  and  $n = 100$  (2<sup>nd</sup> and 3<sup>rd</sup> rows respectively). Multiple results are shown for different random initialization seeds (columns). (b) Normalized average distance between the invariant images generated by SnS (stretching in Layer3\_conv1) and their reference MEIs across the different stages of a robust ResNet50. Different lines indicate different cardinalities of units in the representation space used for optimization (i.e. all,  $n = 1000$  and  $n = 100$ ). The average distance between multiple MEIs generated via XDREAM is reported for comparison. (c-d) Same as (a) and (b) respectively, but stretching was performed on robust ResNet-50 in the high\_level layer (i.e. Layer4\_conv7).

most sensitive (region of interest (ROI)), effectively isolating these core features from irrelevant contextual details (Fig. S3a-d (bottom rows)). This level of feature disentanglement and clarity surpassed that observed in images optimized by methods like XDREAM [10], which often retain more complex, less isolated visual elements (Fig. S3a-d (top rows)). Moreover, by comparing the reference image to multiple invariant instances (i.e. initialized with multiple random seeds), the images revealed specific transformations the identified ROI could tolerate while preserving the unit’s activation level, offering insights into the unit’s invariances (Fig. S3a-d (bottom rows)).

## D SnS is robust to subsampling of the representation space

SnS, as a gradientless and model-agnostic method, offers significant advantages for neuroscience research, particularly in scenarios where constructing detailed “digital twin” models (as described in Section 2) is impractical. However, a critical challenge for applying SnS to study hierarchical invariances in the brain is the inherent undersampling of representational spaces. Despite recent advances in electrophysiology and calcium imaging enabling simultaneous recording from hundreds or thousands of neurons [54, 55], these recordings invariably capture only a fraction of the total neural population within a given brain area.

To evaluate SnS’s performance under conditions analogous to experimental neuroscience, we investigated how subsampling the representational space impacts the qualitative nature and quantitative characteristics (assessed via representational distance analysis) of the invariances identified by SnS.

Specifically, we focused on the same 77 readout target units previously analyzed in Section 4.2. For these units, we identified invariances after subsampling their corresponding representational spaces (either `layer3_conv1` or `layer4_conv7` of the CNN) to either 1000 or 100 randomly selected units. To ensure robustness, this process was repeated 10 times for each target unit and subsampling condition, using different random selections of representational space units. As a reference, 1000 units represent 0.5% of the units in `layer3_conv1` ( $n = 200704$ ), while they constitute 4% of `layer4_conv7` ( $n = 25088$ ).

Our findings indicate that SnS is robust to such undersampling. Qualitatively, the invariances identified using subsampled spaces were comparable in their level of abstraction to those found using the full representational layer (Fig.S4a and c). This suggests that the SnS optimization process effectively identified relevant axes of image variation even with reduced unit information. This qualitative observation was corroborated by quantitative representation distance analysis (Fig. S4b and d), which demonstrated a strong alignment in the progression of invariance discovery between analyses using full and subsampled layers. These results highlight SnS’s potential for reliably characterizing neural invariances even when constrained by the partial observations typical of in vivo recordings.

A plausible explanation for this robustness to subsampling lies in the inherent representational redundancy present in both biological neural circuits and artificial neural networks. In biological systems, information is often encoded in a distributed manner across populations [56] [57], where multiple neurons may carry similar or overlapping information, contributing to robustness against noise and cell loss [58]. Similarly, deep neural networks, including CNNs, have been shown to learn redundant features or possess over-parameterization, where multiple units or weights contribute to similar computations [59]. Consequently, a sufficiently large and representative subsample of units can still capture the dominant computational characteristics and drive the overall representation into similar regimes as the full layer, explaining the observed consistency in identified invariances.

## References

- [1] Giulio Matteucci, Eugenio Piasini, and Davide Zoccolan. Unsupervised learning of mid-level visual representations. *Current Opinion in Neurobiology*, 84:102834, 2024.
- [2] David GT Barrett, Ari S Morcos, and Jakob H Macke. Analyzing biological and artificial neural networks: challenges with opportunities for synergy? *Current opinion in neurobiology*, 55: 55–64, 2019.
- [3] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2017. doi: 10.23915/distill.00007. <https://distill.pub/2017/feature-visualization>.
- [4] Anh Nguyen, Jason Yosinski, and Jeff Clune. Understanding neural networks via feature visualization: A survey. *Explainable AI: interpreting, explaining and visualizing deep learning*, pages 55–76, 2019.
- [5] Thomas Fel, Thibaut Boissin, Victor Boutin, Agustin Picard, Paul Novello, Julien Colin, Drew Linsley, Tom Rousseau, Rémi Cadène, Lore Goetschalckx, et al. Unlocking feature visualization for deep network with magnitude constrained optimization. *Advances in Neural Information Processing Systems*, 36:37813–37826, 2023.
- [6] Will Xiao and Gabriel Kreiman. Gradient-free activation maximization for identifying effective stimuli. *arXiv preprint arXiv:1905.00378*, 2019.
- [7] Edgar Y Walker, Fabian H Sinz, Erick Cobos, Taliyah Muhammad, Emmanouil Froudarakis, Paul G Fahey, Alexander S Ecker, Jacob Reimer, Xaq Pitkow, and Andreas S Tolias. Inception loops discover what excites neurons most using deep predictive models. *Nature neuroscience*, 22(12):2060–2065, 2019.
- [8] Anh Nguyen, Jason Yosinski, and Jeff Clune. Multifaceted feature visualization: Uncovering the different types of features learned by each neuron in deep neural networks. *arXiv preprint arXiv:1602.03616*, 2016.
- [9] Santiago A. Cadena, Marissa A. Weis, Leon A. Gatys, Matthias Bethge, and Alexander S. Ecker. Diverse feature visualizations reveal invariances in early layers of deep neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.

- [10] Will Xiao and Gabriel Kreiman. Xdream: Finding preferred stimuli for visual neurons using generative networks and gradient-free optimization. *PLoS computational biology*, 16(6): e1007973, 2020.
- [11] Ian Goodfellow, Honglak Lee, Quoc Le, Andrew Saxe, and Andrew Ng. Measuring invariances in deep networks. *Advances in neural information processing systems*, 22, 2009.
- [12] Saeed Reza Kheradpisheh, Masoud Ghodrati, Mohammad Ganjtabesh, and Timothée Masquelier. Deep networks can resemble human feed-forward vision in invariant object recognition. *Scientific reports*, 6(1):32672, 2016.
- [13] Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. Exploring the landscape of spatial robustness. In *International conference on machine learning*, pages 1802–1811. PMLR, 2019.
- [14] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, page 5188–5196. IEEE, June 2015. doi: 10.1109/cvpr.2015.7299155. URL <http://dx.doi.org/10.1109/CVPR.2015.7299155>.
- [15] Jenelle Feather, Guillaume Leclerc, Aleksander Madry, and Josh H McDermott. Model metamers reveal divergent invariances between biological and artificial neural networks. *Nat. Neurosci.*, 26(11):2017–2034, November 2023.
- [16] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [17] Jörn-Henrik Jacobsen, Jens Behrmann, Richard Zemel, and Matthias Bethge. Excessive invariance causes adversarial vulnerability, 2020. URL <https://arxiv.org/abs/1811.00401>.
- [18] Guy Gaziv, Michael Lee, and James J DiCarlo. Strong and precise modulation of human percepts via robustified ansns. *Advances in Neural Information Processing Systems*, 36:65936–65947, 2023.
- [19] Morgan B. Talbot, Gabriel Kreiman, James J. DiCarlo, and Guy Gaziv. L-wise: Boosting human visual category learning through model-based image selection and enhancement. *International Conference on Learning Representations (ICLR)*, 2025.
- [20] James J. DiCarlo, Davide Zoccolan, and Nicole C. Rust. How does the brain solve visual object recognition? *Neuron*, 73(3):415–434, February 2012. ISSN 0896-6273. doi: 10.1016/j.neuron.2012.01.010. URL <http://dx.doi.org/10.1016/j.neuron.2012.01.010>.
- [21] Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the national academy of sciences*, 111(23):8619–8624, 2014.
- [22] Daniel LK Yamins and James J DiCarlo. Using goal-driven deep learning models to understand sensory cortex. *Nature neuroscience*, 19(3):356–365, 2016.
- [23] Martin Schrimpf, Jonas Kubilius, Ha Hong, Najib J Majaj, Rishi Rajalingham, Elias B Issa, Kohitij Kar, Pouya Bashivan, Jonathan Prescott-Roy, Franziska Geiger, et al. Brain-score: Which artificial neural network for object recognition is most brain-like? *BioRxiv*, page 407007, 2018.
- [24] Pouya Bashivan, Kohitij Kar, and James J DiCarlo. Neural population control via deep image synthesis. *Science*, 364(6439):eaav9436, 2019.
- [25] Chong Guo, Michael Lee, Guillaume Leclerc, Joel Dapello, Yug Rao, Aleksander Madry, and James DiCarlo. Adversarially trained neural representations are already as robust as biological neural representations. In *International Conference on Machine Learning*, pages 8072–8081. PMLR, 2022.

- [26] Alexander JE Kell, Daniel LK Yamins, Erica N Shook, Sam V Norman-Haignere, and Josh H McDermott. A task-optimized neural network replicates human auditory behavior, predicts brain responses, and reveals a cortical processing hierarchy. *Neuron*, 98(3):630–644, 2018.
- [27] Aran Nayebi, Nathan CL Kong, Chengxu Zhuang, Justin L Gardner, Anthony M Norcia, and Daniel LK Yamins. Mouse visual cortex as a limited resource system that self-learns an ecologically-general representation. *PLOS Computational Biology*, 19(10):e1011506, 2023.
- [28] Rudi Tong, Ronan da Silva, Dongyan Lin, Arna Ghosh, James Wilsenach, Erica Cianfarano, Pouya Bashivan, Blake Richards, and Stuart Trenholm. The feature landscape of visual cortex. November 2023. doi: 10.1101/2023.11.03.565500. URL <http://dx.doi.org/10.1101/2023.11.03.565500>.
- [29] Zhiwei Ding, Dat T Tran, Kayla Ponder, Erick Cobos, Zhuokun Ding, Paul G Fahey, Eric Wang, Taliah Muhammad, Jiakun Fu, Santiago A Cadena, et al. Bipartite invariance in mouse primary visual cortex. *bioRxiv*, 2023.
- [30] Carlos R Ponce, Will Xiao, Peter F Schade, Till S Hartmann, Gabriel Kreiman, and Margaret S Livingstone. Evolving images for visual neurons using a deep generative network reveals coding principles and neuronal preferences. *Cell*, 177(4):999–1009, 2019.
- [31] Binxu Wang and Carlos R Ponce. Tuning landscapes of the ventral stream. *Cell Reports*, 41(6), 2022.
- [32] Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks. *Advances in neural information processing systems*, 29, 2016.
- [33] Nikolaus Hansen, Sibylle D Müller, and Petros Koumoutsakos. Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (cma-es). *Evolutionary computation*, 11(1):1–18, 2003.
- [34] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf).
- [35] Logan Engstrom, Andrew Ilyas, Hadi Salman, Shibani Santurkar, and Dimitris Tsipras. Robustness (python library), 2019. URL <https://github.com/MadryLab/robustness>.
- [36] Kalyanmoy Deb. *Multi-objective Optimisation Using Evolutionary Algorithms: An Introduction*, page 3–34. Springer London, 2011. ISBN 9780857296528. doi: 10.1007/978-0-85729-652-8\_1. URL [http://dx.doi.org/10.1007/978-0-85729-652-8\\_1](http://dx.doi.org/10.1007/978-0-85729-652-8_1).
- [37] Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversarially robust imagenet models transfer better?, 2020. URL <https://arxiv.org/abs/2007.08489>.
- [38] Chang Liu, Yinpeng Dong, Wenzhao Xiang, Xiao Yang, Hang Su, Jun Zhu, Yuefeng Chen, Yuan He, Hui Xue, and Shibao Zheng. A comprehensive study on robustness of image classification models: Benchmarking and rethinking, 2023. URL <https://arxiv.org/abs/2302.14301>.
- [39] Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo Debenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial robustness benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. URL <https://openreview.net/forum?id=SSKZPJct7B>.
- [40] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.

- [41] Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Brandon Tran, and Aleksander Madry. Adversarial robustness as a prior for learned representations, 2019. URL <https://arxiv.org/abs/1906.00945>.
- [42] Paolo Muratore, Alireza Alemi, and Davide Zoccolan. Unraveling the complexity of rat object vision requires a full convolutional network and beyond. *Patterns*, 6(2):101149, February 2025. ISSN 2666-3899. doi: 10.1016/j.patter.2024.101149. URL <http://dx.doi.org/10.1016/j.patter.2024.101149>.
- [43] Paolo Muratore, Sina Tafazoli, Eugenio Piasini, Alessandro Laio, and Davide Zoccolan. Prune and distill: similar reformatting of image information along rat visual cortex and deep neural networks, 2022. URL <https://arxiv.org/abs/2205.13816>.
- [44] Sina Tafazoli, Houman Safaai, Gioia De Franceschi, Federica Bianca Rosselli, Walter Vanzella, Margherita Riggi, Federica Buffolo, Stefano Panzeri, and Davide Zoccolan. Emergence of transformation-tolerant representations of visual objects in rat lateral extrastriate cortex. *eLife*, 6, April 2017. ISSN 2050-084X. doi: 10.7554/elife.22794. URL <http://dx.doi.org/10.7554/eLife.22794>.
- [45] Fabian A. Soto and Edward A. Wasserman. Promoting rotational-invariance in object recognition despite experience with only a single view. *Behavioural Processes*, 123:107–113, February 2016. ISSN 0376-6357. doi: 10.1016/j.beproc.2015.11.005. URL <http://dx.doi.org/10.1016/j.beproc.2015.11.005>.
- [46] Samantha M. W. Wood and Justin N. Wood. A chicken model for studying the emergence of invariant object recognition. *Frontiers in Neural Circuits*, 9, February 2015. ISSN 1662-5110. doi: 10.3389/fncir.2015.00007. URL <http://dx.doi.org/10.3389/fncir.2015.00007>.
- [47] V. Schluessel, G. Fricke, and H. Bleckmann. Visual discrimination and object categorization in the cichlid pseudotropheus sp. *Animal Cognition*, 15(4):525–537, March 2012. ISSN 1435-9456. doi: 10.1007/s10071-012-0480-3. URL <http://dx.doi.org/10.1007/s10071-012-0480-3>.
- [48] Nikolaus Hansen, Youhei Akimoto, and Petr Baudis. CMA-ES/pycma on Github. Zenodo, DOI:10.5281/zenodo.2559634, February 2019. URL <https://doi.org/10.5281/zenodo.2559634>.
- [49] Félix-Antoine Fortin, François-Michel De Rainville, Marc-André Gardner, Marc Parizeau, and Christian Gagné. DEAP: Evolutionary algorithms made easy. *Journal of Machine Learning Research*, 13:2171–2175, jul 2012.
- [50] Anonymous Author. Stretching beyond the obvious: A gradient-free framework to unveil the hidden landscape of visual invariance, 2025. URL <https://zenodo.org/doi/10.5281/zenodo.15491763>.
- [51] Joshua R De Leeuw. jspsych: A javascript library for creating behavioral experiments in a web browser. *Behavior research methods*, 47:1–12, 2015.
- [52] Daiichiro Kuroki. A new jspsych plugin for psychophysics, providing accurate display duration and stimulus onset asynchrony. *Behavior Research Methods*, 53:301–310, 2021.
- [53] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(1):289–300, 1995.
- [54] James J Jun, Nicholas A Steinmetz, Joshua H Siegle, Daniel J Denman, Marius Bauza, Brian Barbarits, Albert K Lee, Costas A Anastassiou, Alexandru Andrei, Çağatay Aydın, et al. Fully integrated silicon probes for high-density recording of neural activity. *Nature*, 551(7679): 232–236, 2017.
- [55] Nicholas James Sofroniew, Daniel Flickinger, Jonathan King, and Karel Svoboda. A large field of view two-photon mesoscope with subcellular resolution for in vivo imaging. *elife*, 5:e14472, 2016.

- [56] Ehud Zohary, Michael N Shadlen, and William T Newsome. Correlated neuronal discharge rate and its implications for psychophysical performance. *Nature*, 370(6485):140–143, 1994.
- [57] Carsen Stringer, Marius Pachitariu, Nicholas Steinmetz, Charu Bai Reddy, Matteo Carandini, and Kenneth D Harris. Spontaneous behaviors drive multidimensional, brainwide activity. *Science*, 364(6437):eaav7893, 2019.
- [58] Alexandre Pouget, Peter Dayan, and Richard Zemel. Information processing with population codes. *Nature Reviews Neuroscience*, 1(2):125–132, 2000.
- [59] Misha Denil, Babak Shakibi, Laurent Dinh, Marc’Aurelio Ranzato, and Nando De Freitas. Predicting parameters in deep learning. *Advances in neural information processing systems*, 26, 2013.