# Kernel Heterogeneity Improves Sparseness of Natural Images Representations

Hugo J. Ladret[1,2], Christian Casanova[2], Laurent Udo Perrinet[1]

[1] Institut de Neurosciences de la Timone,
UMR 7289, CNRS and Aix-Marseille Université,
Marseille, 13005, France

[2] School of Optometry, Université de Montréal,
Montréal, QC H3C 3J7, Canada

December 25, 2023

## Abstract

Both biological and artificial neural networks inherently balance their performance with their operational cost, which balances their computational abilities. Typically, an efficient neuromorphic neural network is one that learns representations that reduce the redundancies and dimensionality of its input. This is for instance achieved in sparse coding, and sparse representations derived from natural images yield representations that are heterogeneous, both in their sampling of input features and in the variance of those features. Here, we investigated the connection between natural images' structure, particularly oriented features, and their corresponding sparse codes. We showed that representations of input features scattered across multiple levels of variance substantially improve the sparseness and resilience of sparse codes, at the cost of reconstruction performance. This echoes the structure of the model's input, allowing to account for the heterogeneously aleatoric structures of natural images. We demonstrate that learning kernel from natural images produces heterogeneity by balancing between approximate and dense representations, which improves all reconstruction metrics. Using a parametrized control of the kernels' heterogeneity used by a convolutional sparse coding algorithm, we show that heterogeneity emphasizes sparseness, while homogeneity improves representation granularity. In a broader context, these encoding strategy can serve as inputs to deep convolutional neural networks. We prove that such variance-encoded sparse image datasets enhance computational efficiency, emphasizing the benefits of kernel heterogeneity to leverage naturalistic and variant input structures and possible applications to improve the throughput of neuromorphic hardware.

**Keywords:** Sparseness; Vision; Heterogeneity; Efficiency; Coding; Representation; Deep Learning
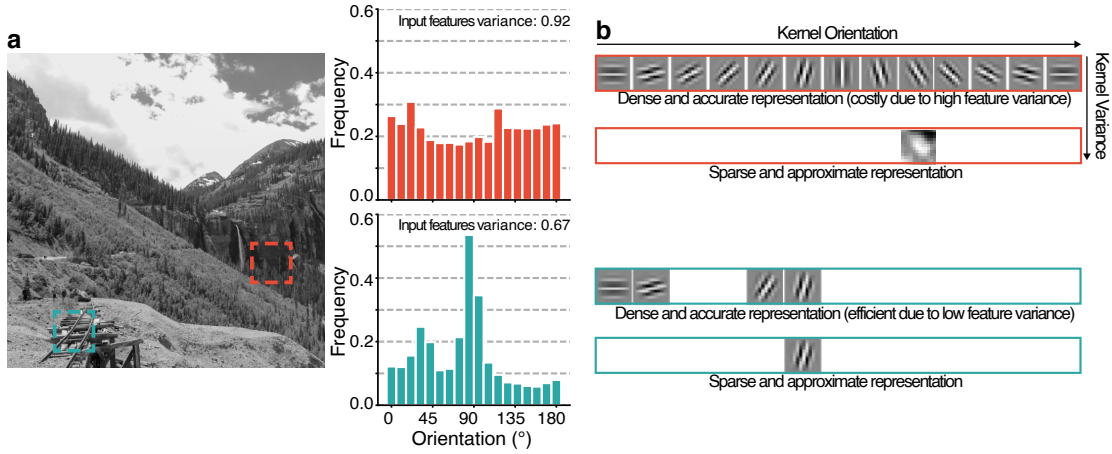
# 1 Introduction



Figure 1: Efficient coding of sensory inputs. (**a**) Orientation distributions with high (red) and low (blue) variance, in two $256^2$ pixel patches from a sample natural image. (**b**) Representation of these distributions and their efficiency depends on the structure of the input. The high-variance patch can be accurately represented with multiple oriented kernels, or approximated using one single kernel with high representational variance. Similarly, the low-variance patch can be encoded as a two-peaked orientation for an accurate representation, or using one kernel of low representation variance for a higher sparseness.

Neuromorphic neural networks are fundamentally designed to process inputs based on their statistical characteristics. This is particularly evident in vision-related tasks related to natural images, which exhibit a set of common statistical properties at multiple levels of complexity [1]. These statistical characteristics guide sensory processing, and are implicitly learned through efficient coding models [2, 3]. For example, natural images typically show a local redundancy in luminance patterns that biological neural network remove at early processing stages, enhancing computational efficiency [4]. In general, these images can be conceptualized as distributions of features (Figure 1), which are, at a low descriptive level, oriented edges that form the foundation of hierarchical representations of natural images [5]. The first moment of these distributions informs on the mean orientation in a given image patch, while the second central moment represents the heterogeneity of these features.

Modeling of such heterogeneity is crucial for sensory processing, both through input and representation bound variances [6]. Input variance, also referred to as *aleatoric* variance, stems from the intrinsic stochasticity in the processes that generate natural sensory inputs, such as sounds [7], textures [8] or images [9]. As its sources escape modeller control, it is challenging to predict, especially in computer vision models [10] or neuromorphic hardware [11], and mandates a robust approach to accurately represent and process naturalistic inputs.

Evidences from neurobiological networks support the notion that neural systems account for this variance in decision-making processes [12], following Bayesian-derived

rules [13]. In practice, this is supported through the variability of neuronal sparse activations [14], which depends directly on the variance of the input [15, 16]. This relationship ties input variance to representational variance : in feature space, the basis function of a neuron is intrinsically linked to its capacity to encode particular levels of aleatoric variance [17]. Neurons with broad kernels will more effectively encode broadly represented elements in orientation space, such as textures (see Figure 1). This neurobiological evidence can notably serve to "explain away" irrelevant input to neural networks, thereby optimizing neuromorphic designs at the hardware level.

Indeed, neuromorphic machine learning models which emulate the visual system, such as sparse coding, exhibit a dictionary of kernels which possess a wide range of tuning heterogeneity [3]. This heterogeneity is particularly notable in their convolutional forms, where feature activations, being both position- and scale-invariant, effectively mirror the aleatoric structure of natural images. This process is akin to maximum likelihood estimation, wherein modeling visual inputs involves capturing the variance of visual features through parametrized surrogate distributions. Thus, sparse coding, with its minimalistic yet effective neuromorphic approximation of the early visual system, provides a valuable theoretical framework for understanding how input variance is tied to representational variance.

Here, we aim to provide an empirical account of this relationship, namely by showcasing the advantages of incorporating kernels with heterogeneous feature representations in sparse coding models of natural images. We use a convolutional sparse coding model, trained to reconstruct a novel dataset of high-definition natural images, and manipulate the heterogeneity of its kernels to study its reconstruction performances. We show that optimal learning relies on balancing the heterogeneity of features, which reflects the aleatoric variance in natural images. In a general context, we provide a full PyTorch implementation of our convolutional sparse coding algorithms, and use these codes as inputs of a deep convolutional network, boosting resilience to adversarial input degradation. This underscores our finding that inherent heterogeneity of kernels in machine learning, akin to that of receptive fields in biology, enhances computational efficiency by effectively mirroring the statistical properties of inputs.

## 2 Methods

### 2.1 Convolutional Sparse Coding

Sparse coding (SC) is an unsupervised method for learning the inverse representation of an input signal [18]. Given the assumption that a signal can be represented as a linear mixture of kernels (or basis functions), SC aims to minimize the activation of kernels used to represent the input signal, yielding an efficient representation [19] that can be inverted for reconstruction. Here, SC was used to reconstruct an image $s$ from sparse representations $x$, while minimizing the $\ell_1$-norm of the representation:

$$\underset{x}{\mathrm{argmin}} \frac{1}{2}||s - Dx||_2^2 + \lambda||x||_1 \tag{1}$$

where $D$ is the set of kernels used to represent $s$ (called a dictionary) and $\lambda$ a regularization parameter that controls the trade-off between fidelity and sparsity. Conveniently,

this problem can be efficiently approached with a Basis Pursuit DeNoising (BPDN) algorithm [20]. As there is *a priori* no topology among elements of the dictionary, SC does not preserve the spatial structure of the input signal, which can be problematic in the context of the representation of natural images. Moreover, the overall decomposition is applied globally and handles poorly the overlap between redundant statistical properties of patches in the image [1], yielding a suboptimal representation of the input signal [21].

These problems are leveraged by Convolutional sparse coding (CSC), an extension of the SC method to a convolutional representation, which is closer to a rough neurally-inspired design [22] as used in deep convolutional network (CNNs) [23]. These CNNs use localized kernels similar to the receptive fields of biological neurons in the primary visual cortical areas. A convolutional architecture uses convolutional kernels (dictionary elements) that are spatially localized and replicated on the full input space (or possibly with a stride which subsamples that space). The number of kernels in the dictionary defines the number of features, or *channels*. In CSC, the total number of kernels with respect to standard SC is multiplied by the number of positions. As a result, a convolution allows to explicitly represent the spatial structure of the signal to be reconstructed. This further reduces the number of kernels required to achieve an efficient representation of an image, while providing shift-invariant representations. CSC extends equation (1) to:

$$\underset{\{x_k\}}{\mathrm{argmin}} \frac{1}{2} ||s - \sum_{k=1}^{K} \mathrm{d}_k * x_k||_2^2 + \lambda \sum_{k=1}^{K} ||x_k||_1 \qquad (2)$$

where $x_k$ is a $N^2$ dimensional coefficient map (given a $N^2$ sized image), $\mathrm{d}_k$ is one kernel (among $K$ channels) and $*$ is the convolution operator. As the convolution is a linear operator, CSC problems can be solved with convolutional BPDN algorithms [24]. Here, we used the Python SPORCO package [25] to implement CSC methods, using an Alternating Direction Method of Multipliers (ADMM) algorithm [26] which splits Convolutional Sparse Coding problems into two alternating sub-problems, as described in Appendix A. Additionally, CSC proves advantageous over other reconstruction techniques in its ability to learn interpretable and visualizable kernels from input data.

## 2.2 Dictionaries

Optimal dictionaries to reconstruct natural images are known to be localized, oriented elements [27, 2]. Here, we utilized log-Gabor filters, which have been shown to accurately model the receptive fields of neurons in the visual cortex. These filters have several advantages compared to Gabor filters, notably that they do not have a DC component and that they optimally capture the log-frequency structure of natural images to ensure its optimal reconstruction [28]. The log-Gabor filter [29] is defined in the frequency domain by polar coordinates $(f, \theta)$ as:

$$G(f, \theta) = \exp\left(-\frac{1}{2} \cdot \frac{\log(f/f_0)^2}{\log(1 + \sigma_f/f_0)^2}\right) \cdot \exp\left(\frac{\cos(2 \cdot (\theta - \theta_0))}{4 \cdot \sigma_\theta^2}\right) \qquad (3)$$

where $f_0$ is the center frequency, $\sigma_f$ the bandwidth parameter for the frequency, $\theta_0$ the center orientation and $\sigma_\theta$ the standard deviation for the orientation. This provides with a

parametrization of the dictionary, which is useful to compare the efficiency of different sparse coding models [30]. We kept $f_0 = \sigma_f = 0.4$ cpd, varying only the orientation-related parameters to build the dictionaries. The angular bandwidth $B_\theta$ of the log-Gabor filter, expressed in degrees, was defined as $B_\theta = \sigma_\theta \sqrt{2 \log 2}$ [31].

To titrate the impact of including heterogeneity in the dictionary, we created two log-Gabor dictionaries with the same number of channels, one with homogeneous (a single $\sigma_\theta$) the other with heterogeneous (multiple $\sigma_\theta$) variance of representations. We compared these dictionaries before and after fine-tuning on the dataset, using a dictionary learned from scratch over the dataset as a fifth reference. Such learning was done by performing convolutional sparse coding in a multi-image setting:

$$\underset{\{x_{k,j}\}}{\text{argmin}} \frac{1}{2} \sum_{j=1}^{J} || \sum_{k=1}^{K} d_k * x_{k,j} - s_j ||_2^2 + \lambda \sum_{k}^{K} \sum_{j}^{J} ||x_{k,j}||_1 \text{ s.t. } \forall k, ||d_k||_2 = 1 \quad (4)$$

where $s_j$ is the $j$-th image in the dataset and $x_{k,j}$ is the coefficient map for the $k$-th filter and the $j$-th image. This was alternated with an optimization step of the dictionary:

$$\min_{D} \sum_{i=1}^{N} \frac{1}{2} ||x_i - D * z_i||_2^2 \quad (5)$$

subject to the constraint $|d_k|_2 \leq 1$ for $k = 1, \ldots, K$.

Performance of these dictionaries was measured with two metrics. The peak signal-to-noise ratio (PSNR), a common metric to evaluate reconstruction quality of grayscale images, is defined as:

$$\text{PSNR}(I_1, I_2) = 20 \cdot \log_{10}(\max(I_1)) - 10 \cdot \log_{10} \left( \frac{1}{m \cdot n} \sum_{i=1}^{m} \sum_{j=1}^{n} (I_1 - I_2)^2 \right) \quad (6)$$

where $\max(I_1)$ is the maximum pixel intensity of the source image. The right hand-side term of the PSNR is the $\log_{10}$ of the mean squared error, where $I_1$ and $I_2$ represent the pixel intensity in the source and reconstructed images, respectively. Given that the natural images used here are encoded on 8 bits, common values of PSNR range between 20 (worse) to 50 (best) dB. We also measured the sparseness of the algorithm, which was defined as the fraction of basis coefficients used in a reconstruction which are equal to zero. This value is between 0 (no nonzero coefficient) and 1 (all coefficients are zero). Parametrization of the algorithm was chosen to balance sparseness and PSNR (Appendix A), i.e. $\lambda = 0.05$, with 750 iterations of the learning phase, a residual ratio of 1.05 with relaxation at 1.8, and dictionaries with $K = 144$ total elements of $12^2$ pixels each.

## 2.3 Histogram of oriented gradients

The distributions of oriented features in Figure 1 were computed using a histogram of gradient orientations. Using the 'scikit-image' library [32], given an input image $I$ of dimension $M \times N$, two gradients were computed at each pixel using Sobel filters

$G_h(x, y)$ and $G_v(x, y)$, respectively, for vertical and horizontal gradients. The maps of the magnitude $G_m$ and direction $\theta$ were then given as:

$$G_m(x, y) = \sqrt{G_h(x, y)^2 + G_v(x, y)^2}$$
$$\theta(x, y) = \arctan 2(G_v(x, y), G_h(x, y)) \tag{7}$$

The range of possible gradient directions over $[0, \pi]$ was divided into 18 bins. The orientation histogram $H$ for each bin $b$ was computed as:

$$H(b) = \sum_{(x,y)} I_b(\theta(x, y)) \tag{8}$$

where $I_b$ is an indicator function, ranging from 1 if $\theta(x, y)$ falls within the range of the bin $b$ and 0 otherwise. In that context, one can quantify the orientation content in natural images, then estimate the distribution of oriented features within the input: aleatoric variance can then be approximated as the inverse of the squared variance of this distribution in orientation space and is computed as $\text{Var}_{\text{circ}} = 1 - \sqrt{\bar{X}^2 + \bar{Y}^2}$, where $\bar{X}$ and $\bar{Y}$ are the average cosine and sine values respectively, yielding a scalar value between 0 (lowest orientation variance) and 1 (highest).

## 2.4 Dataset

Images for the CSC sections were captured using either a Canon EOS 650D or Canon EOS 6D camera, fitted with 28mm lenses. A total of $1145$ images was collected at a resolution of at least $5184 \times 3456$ pixels. For CSC, we extracted and used the central $256 \times 256$ pixel segment of each image. These images represent a variety of dynamic scenarios, and were carefully shot to ensure that the subjects of interest were in focus and entirely within the frame. We have made this dataset publicly available on Figshare [33].

## 2.5 Image classification using deep learning

To evaluate the role of sparse codes obtained, we decided to go further than only measuring representation performance by applying these codes on a common machine learning task: image classification. To perform such classification in a neuromorphic-inspired setting, we utilized a modified version of the CIFAR-10 dataset. This dataset, which is commonly used for image classification, originally contains $60,000$ color images of $32 \times 32$ pixel resolution across 10 balanced classes. We processed these images by first upscaling them to $128 \times 128$ resolution via bilinear interpolation. Subsequently, they were converted to grayscale and sparse-coded, as described above.

The dataset was divided into a training set containing $50,000$ sparse codes and a test set comprising $10,000$ sparse codes. The network was trained from scratch through a standard PyTorch implementation, with backpropagation of the gradient using the Adam optimizer [34]. The training objective was to minimize the categorical cross-entropy loss, defined as:

$$J(\theta) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{C} y_{ij} \log(\hat{y}_{ij}) \tag{9}$$

where $N$ is the number of samples, $C$ is the number of classes, $y_{ij}$ is the true label, and $\hat{y}_{ij}$ is the predicted label. The Adam update rule for each parameter $\theta$ is based on moment estimates given by:

$$\theta_{t+1} = \theta_t - \eta \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \qquad (10)$$

where $\eta$ is the learning rate, $\hat{m}_t$ and $\hat{v}_t$ are estimates of the mean and variance of the gradients, and $\epsilon$ is a small constant to prevent division by zero.

The sparse codes representing these images were then used as inputs for an adapted ResNet-18 architecture [35] which is a classically used CNN architecture. This deep residual neural network, typically composed of 18 layers and used for various vision tasks, was adapted to process the 144 dimensions of the sparse-coded inputs instead of the standard 3-channel (RGB) format. This dimensionality corresponds to the number of channels in our sparse coding dictionary. No other modifications were implemented in the network architecture design.

Hyperparameters were tuned via grid search to maximize accuracy on heterogeneous variance codes, with the resulting values: $\eta = 2e - 4$, $\hat{m}_t = 0.9$, $\hat{v}_t = 0.99$, $\epsilon = 1e - 08$. When training the network, CSC methods using ADMM algorithms were ported from SPORCO to a custom PyTorch implementation (available at `https://github.com/hugoladret/epistemic_CSC`) to speed up computations.

## 3 Results

### 3.1 Heterogeneous kernels improve the sparseness of natural images representations

We explored how variance in sensory inputs and neuromorphic representations controls the encoding strategies of natural images. We compared five distinct convolutional sparse coding dictionaries of similar sizes. Two dictionaries using Log-Gabor filters were constructed : one with a homogeneous level of orientation variance ($B_\theta = 12.0°$) and 72 orientations $\theta_0$ ranging from $0°$ to $180°$ (Figure 2a, green) compared to another one with heterogeneous orientation variance, spanning 12 orientation values $\theta_0$ and six $B_\theta$ ranging from $3°$ to $30°$ (Figure 2b, blue). We then benchmarked these constructed dictionaries against their learned counterparts, which were fine-tuned on the dataset (Figure 2a, orange; b, purple). A final comparison was made against a randomly initialized dictionary learned *de novo* on the same dataset (Figure 2c, black). Performance evaluation across the $1,445$ high-definition natural images revealed that dictionaries initialized with Log-Gabor filters consistently displayed highly variant performance from image to image (Figure 2d). Prior to learning, the dictionary integrating heterogeneous orientation variance outperformed its homogeneous counterpart in sparsity (Mann-Whitney U-test, $U = 1310760.0$, $p < 0.001$), but had significantly lower PSNR ($U = 262261.0$, $p < 0.001$). Post-learning, all dictionaries had similar performances in terms of both sparsity ($U = 634605$, $p = 0.18$ for homogeneous vs random initialized dictionaries ; $U = 634605.0$, $p = 0.97$ for heterogeneous vs random initialized dictionaries) and PSNR ($U = 694175$, $p = 0.46$ ; $U = 653943.0$, $p = 0.99$). This
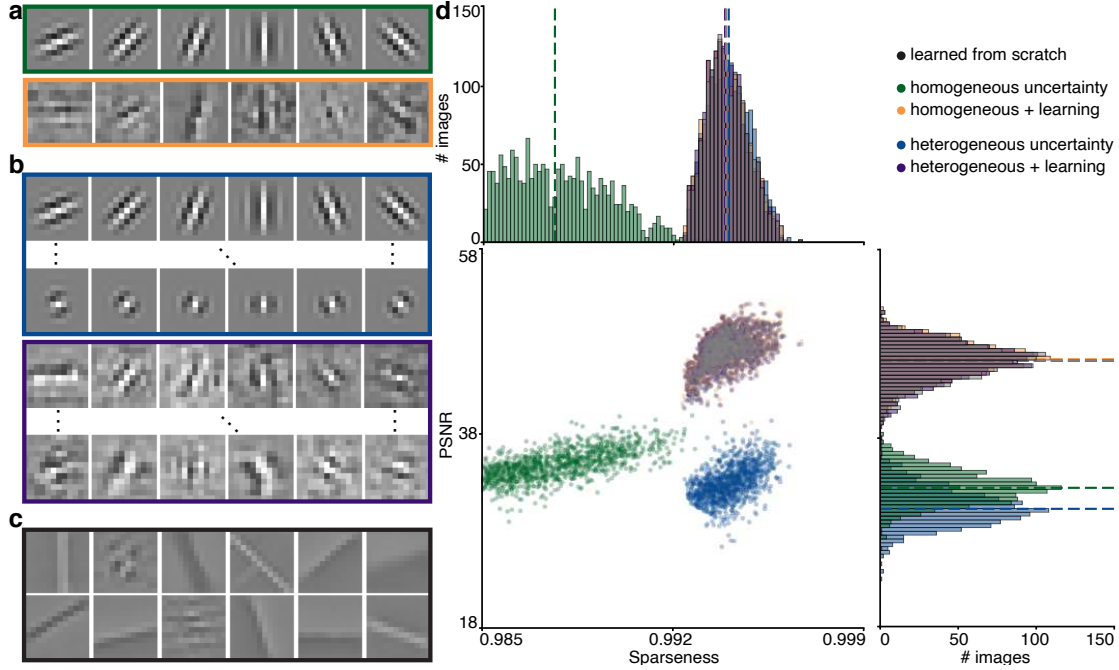
Figure 2: Kernel heterogeneity and reconstruction trade-off. **(a)** Elements from dictionaries with homogeneous kernel variance before (green) and after dictionary learning (orange). **(b)** Same, with heterogeneous kernel variance before (blue) and after learning (purple). **(c)** Elements from a dictionary learned from random initialization on the dataset. **(d)** Distribution of the sparseness (top) and Peak Signal-to-Noise Ratio (PSNR, right) of the five dictionaries. Median values are shown as dashed lines. All three post-learning dictionaries have overlapping (but not identical) distributions.

suggests that emphasis on heterogeneous variance modelling improves the sparsity, at the cost of reconstruction performance.

After learning from the dataset, whether from random initialization or from a pre-constructed log-Gabor dictionary, all dictionaries converge to qualitatively quite different filters, yet with a similar, superiorly sparse and performant form of encoding. The learning method indeed enhanced all Log-Gabor dictionaries, resulting in increased PSNR ($U = 0.0$, $p < 0.001$ ; $U = 181535.0$, $p < 0.001$, homogeneous and heterogeneous variance dictionaries, compared to their pre-learning version) and sparseness ($U = 23595.0$, $p < 0.001$ ; $U = 248667.0$, $p < 0.001$). Given the converging reconstruction and sparseness for all these dictionaries, we now focus on the heterogeneous variance dictionary, both pre- and post-learning, as well as the pre-learned homogeneous variance dictionary. Additional performance details for the homogeneous dictionary are provided in Appendix B.

What are then the kernel features changed through the learning process? While fine-tuned dictionaries do incur a significantly higher computational cost during the learning phase, they deliver substantial improvements in both PSNR and sparsity, compared to merely introducing heterogeneous variance into a pre-existing dictionary. These enhancements can be attributed to modifications in the dictionary coefficients following the learning phase, affecting both the feature orientations ($\theta_0$) and their associated lev-
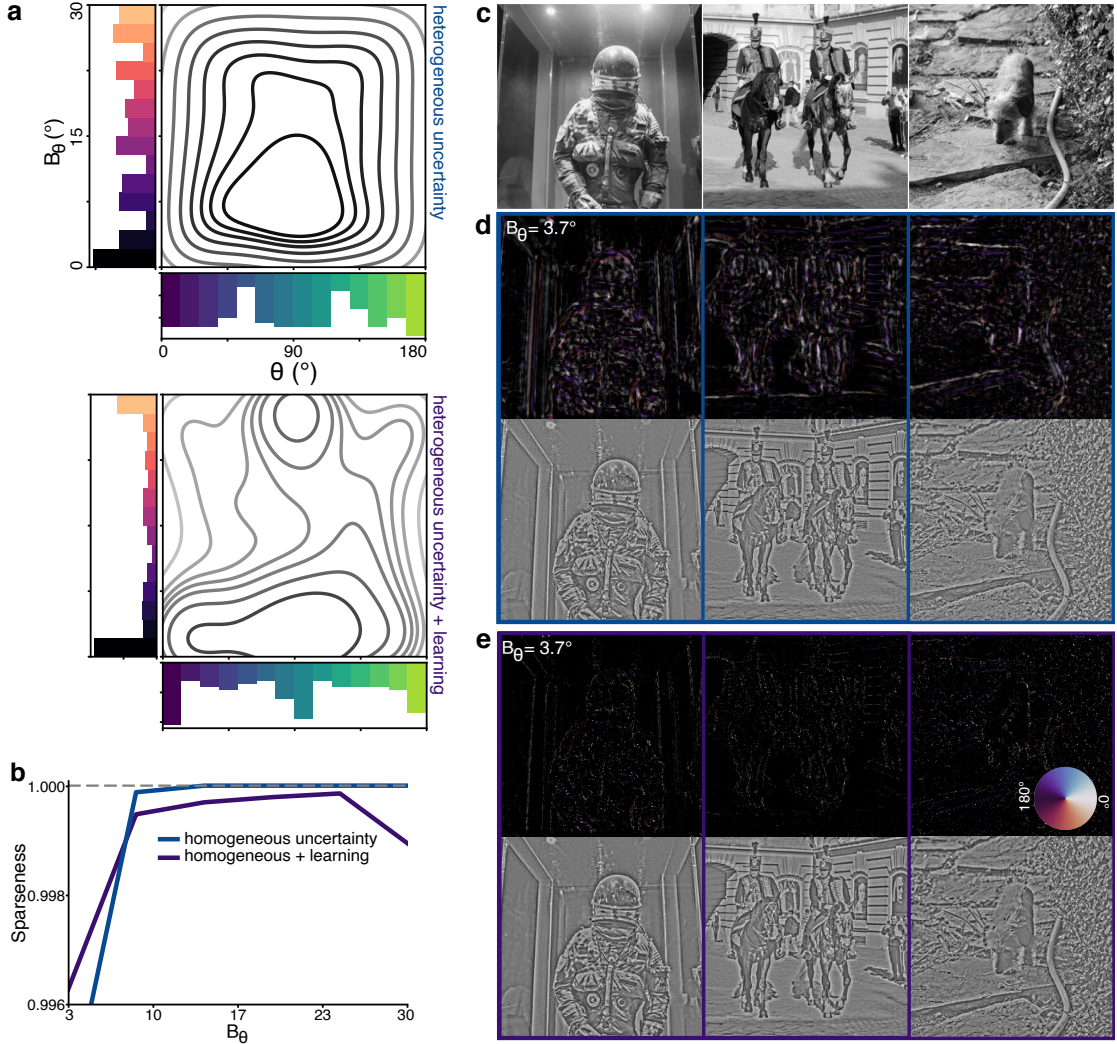
Figure 3: Learning balances coefficient distribution. **(a)** Kernel density estimation over $\theta$ and $B_\theta$ of the kernels before (top) and after (bottom) learning. **(b)** Sparseness of the dictionaries for kernel variance $B_\theta$. Sparseness $= 1$ (i.e. no activation, as in the case of the pre-learning encoding) is represented as a gray dashed line. **(c)** Example images from the dataset. **(d)** Sparse code for high $B_\theta$ values (color coded by each coefficient's $\theta$) and reconstructions for the pre-learned, heterogeneous variance dictionary. **(e)** Same as (d), for post-learned, heterogeneous variance dictionary. Orientation color code of the coefficients is shown on the rightmost coefficient map.

els of variance ($B_\theta$) (Figure 2a). Specifically, learning from a dataset of natural images introduced a bias toward cardinal orientations (Figure 3a), mirroring inherent biases found in natural scenes [36], which is in contrast to the uniformly distributed initial dictionary. Furthermore, the learning process resulted in a non-uniform distribution of coefficients across multiple levels of orientation variance (Figure 3b). Notably, coefficients that were previously inactive (i.e., sparseness $= 1$) became activated at higher $B_\theta$ levels (Figure 3c-e). This led to consistent patterns in coefficient distribution across heterogeneous variance levels (Figure 3d,e). This uniformity is likely influenced by

the dataset's inherent variability. Consequently, the performance gains attributed to the learning process are contingent upon feature orientation biases ($\theta_0$) and a redistribution of the levels of variance ($B_\theta$), both of which should be reflective of the dataset's intrinsic structure.

## 3.2 Statistical properties of natural images reflect the variance of learned sparse code
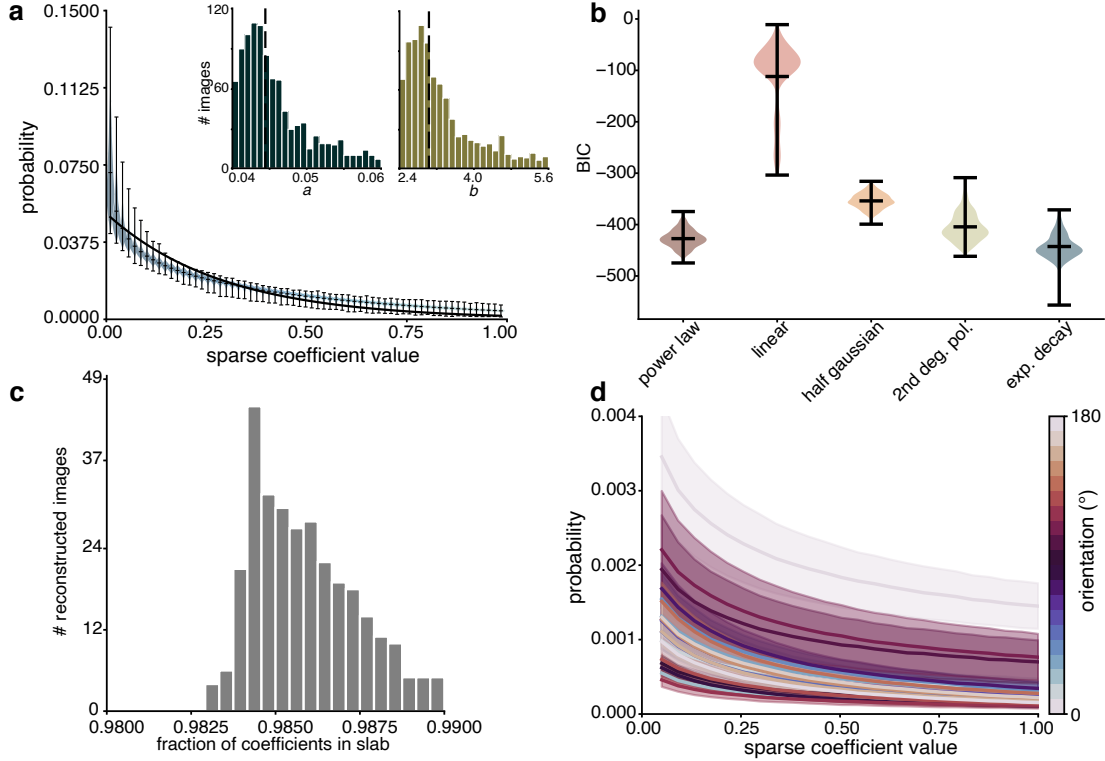


Figure 4: Spike-and-slab sparse representation of the natural images. **(a)** Distribution of the sparse coefficients values. Violin plots' central lines represent mean values, with top and bottom lines representing the extrema. For each image, this distribution was fitted with an exponential decay (black line) $y = a \cdot \exp(-b \cdot x)$, with the distributions for the parameters over the $1145$ images shown in inset **(b)** Bayesian Information Criterion (BIC) for the fitting of the distribution of spikes coefficients with different alternative functions. **(c)** Proportion of zero coefficients per image, i.e., belonging to the "spike" of the distribution. **(d)** Same as (a), with coefficients split by different encoded orientation.

The criteria for the relevance of features encoded in neural networks is dictated by the statistical properties of the environment itself [9, 1]. For instance, at a fundamental representational level, the neural code for light patterns in the retina is the cumulative sum of the Gaussian distribution of luminance found in natural images [4]. At higher levels, scale distributions of visual features, in the Fourier domain, obey a $1/f^2$ power law, which once again echoes the power-law behavior of cortical responses [37, 38]. At intermediate levels, the distribution of these oriented edges can be characterized along
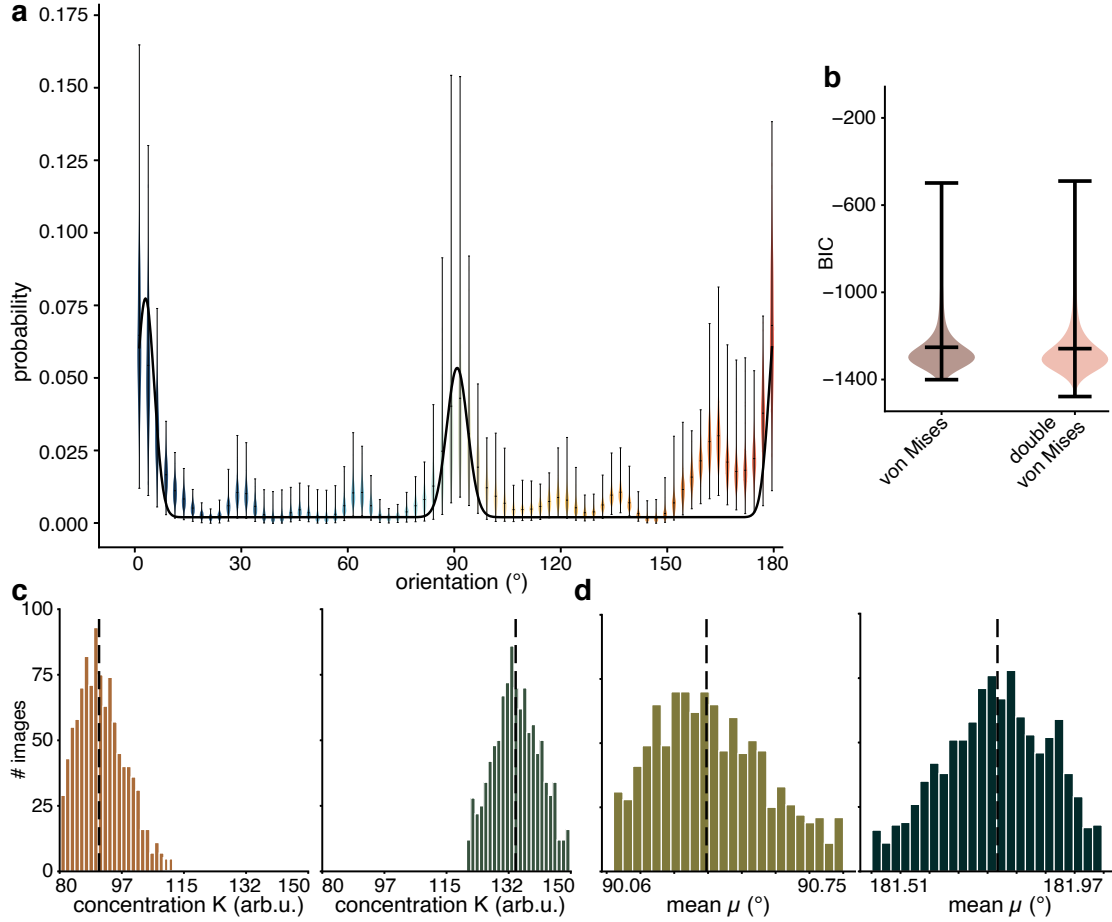
Figure 5: Orientations in natural images follow a double von Mises distribution. **(a)** Orientations of the sparse coefficients, fitted with a double von Mises distribution (black line). **(b)** Bayesian Information Criterion (BIC) for the fitting of the distribution of orientation coefficients. **(c)** Distribution of the concentration parameter $\kappa$ for the first (left) and second (right) peaks of the double von Mises distribution. **(d)** Same as (c), for the mean parameter $\mu$.

its first- and second-order moments: a median orientation, and its corresponding variance. A proper model of natural images thus depends on a proper model of both these moments, which is reflected in the response properties of primary visual cortex neurons [27]. Which of these two parameters warrants greater emphasis? Previous studies suggested that heterogeneity on both orientation and variances arises from sparse learning processes, *in silico* [2] and *in vivo* [17].

Inherently, sparse coding enforces a prior on using a minimal number of coefficients to reconstruct an image, and is thus an encoding strategy that produces a "spike and slab" distribution of activations, characterized by a predominance of zero coefficients [37] (Figure 4a-c). This imposes a prior on the representation of images at the feature-level, with a decaying exponential variation of coefficients that unfolds heterogeneously across different types of orientations (Figure 4d). Lower BIC indicate less information lost in the fitting process, and thus a better fit. Such heterogeneity in feature space stems from the fact that orientations in natural images are biased to cardinal
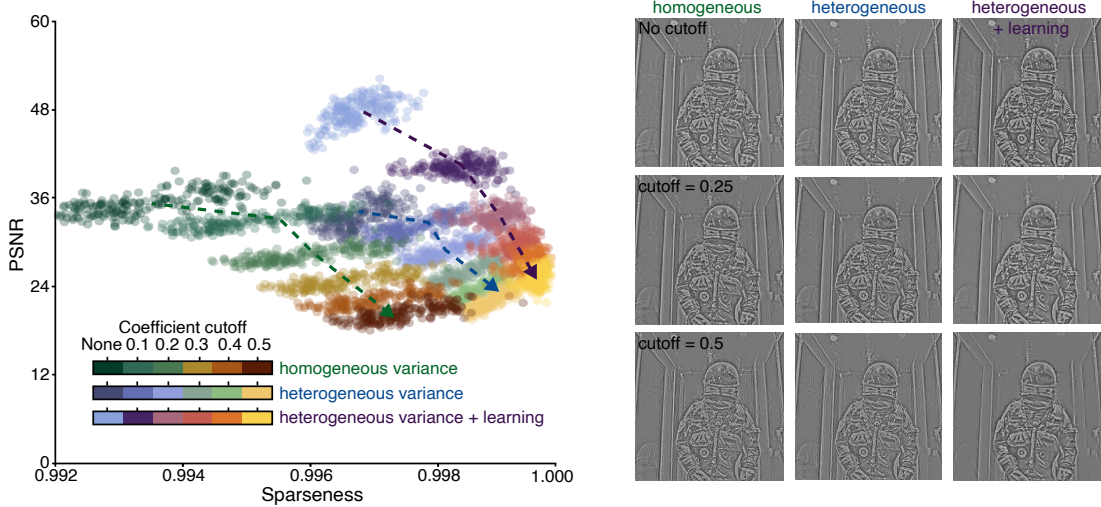
Figure 6: Sparse coefficients can be pruned for increased sparsity. **(a)** Pruning of the coefficients based on their values and resulting sparseness/PSNR for three dictionaries, with mean trajectory represented as a dashed arrow. **(b)** Reconstruction of an image with different cutoff levels.

(i.e., vertical and horizontal) orientations [39], which is echoed at the neuronal level by a cardinal bias in visual perception [40]. This biased distribution of orientation is well-captured by a double von Mises distribution in orientation space (Figure 5a,b):

$$f(x) = A_1 \exp\left(k_1\left(\cos\left(2\pi(x - \phi_1)\right) - 1\right)\right) + A_2 \exp\left(k_2\left(\cos\left(2\pi(x - \phi_2)\right) - 1\right)\right)$$
(11)

where $A_1, A_2$ are the amplitudes of the two von Mises distributions, $k_1, k_2$ are the concentration parameters for the two distributions, $\phi_1, \phi_2$ are the phase offsets for the two distributions.

This distribution is known for higher heterogeneity, and thus aleatoric variance, in natural images compared to synthetic ones [39]. At the cardinal orientations, this is also captured by the variation of the concentration parameters (Figure 5c,d) of the von Mises distributions, which underlies the notion that a proper description of natural images must be able to account for heterogeneous levels of aleatoric variance. This mandates a comparative evaluation of performance between dictionaries that emphasize a representation based on homogeneous or heterogeneous strategies, that is, emphasizing encoding mean features or their variances.

## 3.3 Heterogeneity improves resilience of the neural code

In addition to the previously described trade-off between performance and sparsity (Figure 2), the robustness of the representations can be further evaluated by modifying elements in the typical activation patterns. This then allows pruning less activated coefficients to further increase sparseness, testing the code's resilience to the adversarial degradation. We pruned coefficients with absolute values below a specific threshold, iterating from $0.001$ to $0.5$ in 6 steps. This pruning led to a construction-induced increase in sparseness, that correlated non-linearly with a decrease in PSNR for all dictionaries,
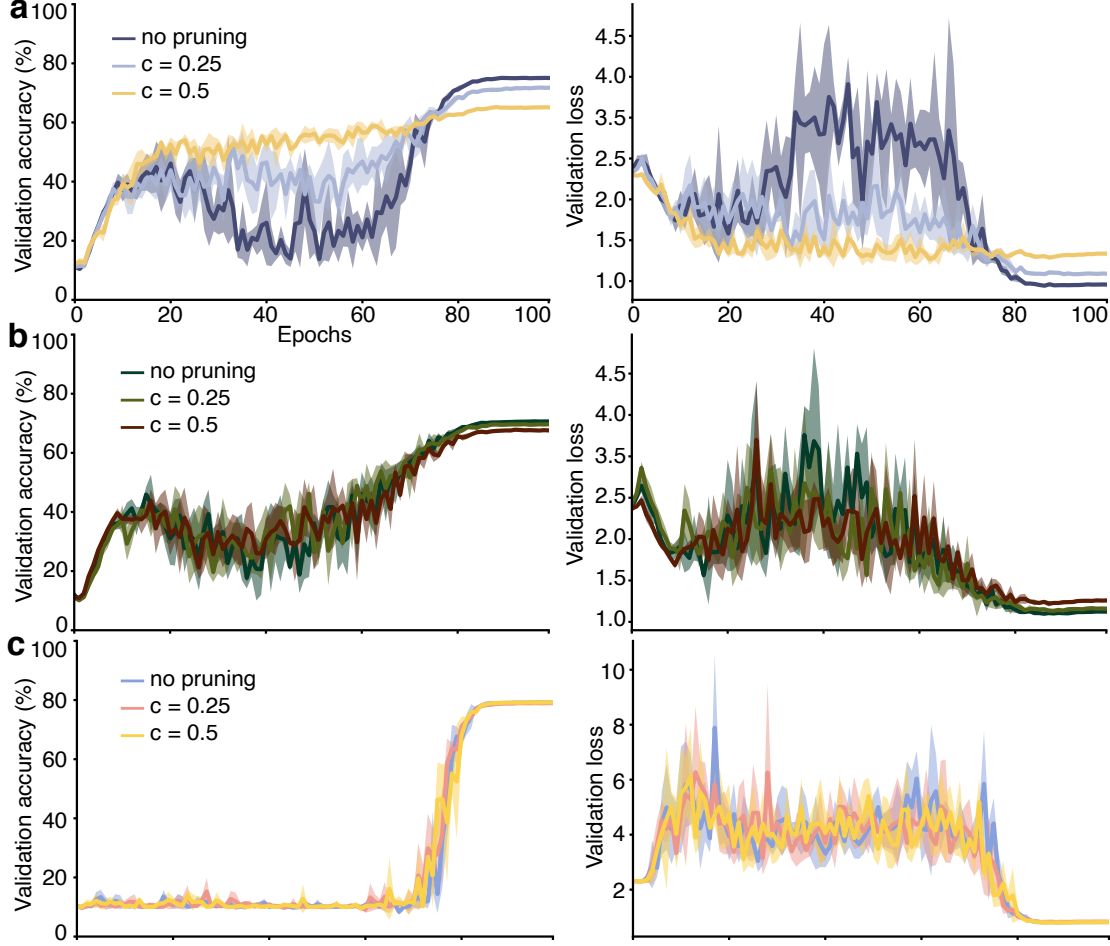
Figure 7: Deep Neural Networks (here, ResNet18), can be trained on sparse codes. **(a)** Validation accuracy (left) and losses (right) curves, for $3$ different pruning levels of coefficients for the heterogeneous variance dictionary. Each network is trained across $4$ random seeds, with the mean value shown as a solid line and the contour representing the standard deviation. **(b)** Same as (a), for the homogeneous variance dictionary. **(c)** Same as (a), for the heterogeneous variance dictionary, post-learning.

while maintaining interpretable representations (Figure 6), The pre-learning heterogeneous variance dictionary's PSNR demonstrated significantly greater resilience to coefficient degradation than the pre-learning homogeneous variance dictionary ($p < 0.05$ for pruning cutoff $c > 0.3$). Post-learning, both the homogeneous and heterogeneous variance dictionaries exhibited similar PSNR, reflective of their PSNR similarities before pruning (Figure 2). This emphasizes the advantage of heterogeneous variance in a dictionary, whether by construction or through learning, in bolstering resilience and efficiency for encoding natural images.

Overall, these findings show that sparse codes for natural images possess highly desirable properties when incorporating heterogeneous basis functions into a sparse model: enhanced sparseness (Figure 2d), more evenly distributed activation (Figure 3b), and increased resilience to code degradation (Figure 6a). Yet, the differences in PSNR may not necessarily translate to perceptible differences in image quality, depending on

Table 1: Mean top-1 accuracy (in %) ± standard deviation across 4 random initialization of ResNet-18 for varying sparse encoding schemes of CIFAR-10. $c = 0.25$ and $c = 0.5$ indicate the pruning level of the sparse coefficients, as done in Figure 6.

| Encoding scheme | No pruning | c=0.25 | c=0.5 |
|---|---|---|---|
| Homogeneous, pre-learning | $70.65 \pm 0.30$ | $69.70 \pm 0.26$ | $67.83 \pm 0.47$ |
| Homogeneous, post-learning | $67.31 \pm 0.20$ | $66.24 \pm 0.01$ | $67.40 \pm 0.12$ |
| Heterogeneous, pre-learning | $75.08 \pm 0.10$ | $71.81 \pm 0.41$ | $65.20 \pm 0.40$ |
| Heterogeneous, post-learning | $\mathbf{79.20 \pm 0.11}$ | $\mathbf{78.98 \pm 0.00}$ | $\mathbf{79.26 \pm 0.02}$ |

the context and application [41]. As such, it is necessary to investigate the potential of employing such codes in objective visual processing problems, for example, in image classification.

As a coarse analogy to a neuromorphic hierarchical sparse construction of visual processing [22, 42, 23], we trained a deep convolutional neural network to classify the sparse codes of natural images. The CIFAR-10 dataset, which was converted to grayscale in order to match the dimensionality of the dictionaries previously described, was sparse-coded and then classified using the Resnet-18 network, reaching a maximum top-1 accuracy of 79.20% in 100 epochs (Figure 7, Table 3.3). After sparse coding of the dataset, but without pruning of the coefficients, a learned dictionary initialized with a heterogeneous orientation variance basis achieved the highest classification accuracy (79.20%). This was followed by the pre-learned version of the network (75.08%), and was higher than homogeneous variance methods. Following degradation of the sparse code ($c = 0.5$), the post-learned heterogeneous variance kept similarly high performance, unlike all the other encoding scheme which showed loss of performance. The discrepancy between the deep learning performance and the previously noted similarities in PSNR and sparseness (Figure 2) underscores the significance of representing variance of low-level features in complex visual models.

## Discussion

Neural systems leverage heterogeneity for increased computational efficiency [43, 44]. Here, we have explored the effects of such heterogeneous encoding of orientation variance by integrating it into a convolutional sparse coding dictionary. Our findings show that this outperforms conventional feature-representing dictionaries with fixed variance, both in sparsity and robustness, at the cost of reconstruction performance. However, these representations can be effectively employed in subsequent visual processing stages, where they result in significantly improved performances of deep convolutional neural networks. Overall, these results imply that incorporating variance in sparse coding dictionaries can substantially improve the encoding and processing of natural images.

The connection between sparse models and neural codes, which underlies the mo-

tivation behind this approach, could be further showcased using biologically plausible algorithms, such as the Locally Competitive Algorithm (LCA) [45]. Rather than enforcing sparsity through convolution as done here, this model uses a mechanism of reciprocal inhibition between each of its elements, a process that mimics particular recurrent inhibition connectivity patterns observed in the cortex [46]. This method potentially mirrors a neural adaptation of winner-takes-all algorithms, reflecting innate competition and selective activation within neural networks, and highlights the potential role of feedback loops to improve sparse coding [47]. Under this analogy, LCA could reinforce the presented framework of heterogeneity by extending it from features space (i.e., receptive fields) to also include the connectivity matrix (i.e., synaptic weights). In terms of hardware, the use of variance weighting by such a lateral inhibition mechanism could provide dynamic computational allocation for significant, unpredictable fluctuations in the data, while reducing or bypassing routine, predictable data streams. This arguably reflects the response characteristics and dynamics of cortical neurons [15, 16]. Emphasizing these pronounced shifts could streamline the data transmitted across physical channels, addressing a primary source of thermal and computational efficiency bottlenecks in neuromorphic hardware [48, 49].

In the context of image classification, our approach employing sparse coding achieved a top-1 accuracy of 79.20% on the CIFAR-10 dataset. While this falls short of the state-of-the-art performance exceeding 99.0% accuracy using color images and transformer architectures [50], it is important to note that our primary objective centered on comparing model performance with heterogeneous degree of variance in the initial layer, rather than solely pursuing state-of-the-art results. Here, the high dimensionality of the sparse-coded CIFAR-10 dataset (144 input dimensions or sparse channels), in contrast to the standard 3 dimensions in RGB images, likely contributes to this difference of accuracy. Direct integration of sparse coding with deep neural networks is a promising avenue of research that aligns with recent developments in the fields of unsupervised learning, object recognition, and face recognition. Some approaches have emphasized the ability of sparse coding to generate succinct, high-level representations of inputs, especially when applied as a pre-processing step for unsupervised learning with unlabeled data using L1-regularized optimization algorithms [51]. In several instances, the mechanism of sparse coding has been seamlessly integrated into deep networks. For instance, the Deep Sparse Coding framework [52] maintains spatial continuity between adjacent image patches, boosting performance in object recognition. Likewise, a face recognition technique combining sparse coding neural networks with softmax classifiers effectively addresses aleatoric uncertainties, including changes in lighting, expression, posture, and low-resolution scenarios [53]. Classifiers relying on sparse codes, produced by lateral inhibition in an LCA, exhibit strong resistance to adversarial attacks [54]. This resilience, potentially enhanced by heterogeneous dictionaries as explored here, offers a promising avenue for research in safety-critical applications.

The empirical evidence presented here can be interpreted as an implicit Bayesian process, wherein initial beliefs about the coefficients are updated using input images to learn the variance of visual features to represent optimally (sparse) orientations. Models with explicit integration of both model and input variance have distinct advantages in that sense. Namely, this allows to maximize model performance and minimizing decision uncertainty. In contrast, we here focused on an implicit understanding of this

relationship, demonstrating through a simple approach that vision models can benefit from factoring-in feature variance without explicit learning rules.

# 4   Acknowledgments

# References

[1]   Eero P Simoncelli and Bruno A Olshausen. "Natural image statistics and neural representation". In: *Annual review of neuroscience* 24.1 (2001), pp. 1193–1216.

[2]   Bruno A Olshausen and David J Field. "Emergence of simple-cell receptive field properties by learning a sparse code for natural images". In: *Nature* 381.6583 (1996), pp. 607–609.

[3]   Bruno A Olshausen and David J Field. "Sparse coding with an overcomplete basis set: A strategy employed by V1?" In: *Vision research* 37.23 (1997), pp. 3311–3325.

[4]   Simon Laughlin. "A simple coding procedure enhances a neuron's information capacity". In: *Zeitschrift für Naturforschung c* 36.9-10 (1981), pp. 910–912.

[5]   Victor Boutin et al. "Sparse Deep Predictive Coding Captures Contour Integration Capabilities of the Early Visual System". In: *PLoS Computational Biology* (May 12, 2020).

[6]   Eyke Hüllermeier and Willem Waegeman. "Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods". In: *Machine Learning* 110.3 (2021), pp. 457–506.

[7]   Keisuke Nakamura and Kazuhiro Nakadai. "Robot audition based acoustic event identification using a bayesian model considering spectral and temporal uncertainties". In: *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE. 2015, pp. 4840–4845.

[8]   Charles E Pettypiece, Melvyn A Goodale, and Jody C Culham. "Integration of haptic and visual size cues in perception and action revealed through cross-modal conflict". In: *Experimental brain research* 201.4 (2010), pp. 863–873.

[9]   Daniel L Ruderman. "The statistics of natural images". In: *Network: computation in neural systems* 5.4 (1994), p. 517.

[10]  Yann Gousseau and Jean-Michel Morel. "Are natural images of bounded variation?" In: *SIAM Journal on Mathematical Analysis* 33.3 (2001), pp. 634–648.

[11] Kaitlin L. Fair et al. "Sparse Coding Using the Locally Competitive Algorithm on the TrueNorth Neurosynaptic System". In: *Frontiers in Neuroscience* 13 (2019). ISSN: 1662-453X. URL: https://www.frontiersin.org/articles/10.3389/fnins.2019.00754 (visited on 12/21/2023).

[12] Hermann LF von Helmholtz. *Treatise on physiological optics*. 1867.

[13] Karl Friston. "A theory of cortical responses". In: *Philosophical transactions of the Royal Society B: Biological sciences* 360.1456 (2005), pp. 815–836.

[14] Gergő Orbán et al. "Neural variability and sampling-based probabilistic representations in the visual cortex". In: *Neuron* 92.2 (2016), pp. 530–543.

[15] Olivier J Hénaff et al. "Representation of visual uncertainty through neural gain variability". In: *Nature communications* 11.1 (2020), pp. 1–12.

[16] Hugo J Ladret et al. "Cortical recurrence supports resilience to sensory variance in the primary visual cortex". In: *Communications Biology* 6.1 (2023), p. 667.

[17] Robbe LT Goris, Eero P Simoncelli, and J Anthony Movshon. "Origin and function of tuning diversity in macaque visual cortex". In: *Neuron* 88.4 (2015), pp. 819–831.

[18] Honglak Lee et al. "Efficient sparse coding algorithms". In: *Advances in neural information processing systems* 19 (2006).

[19] Laurent U Perrinet. "Sparse Models for Computer Vision". In: *Biologically Inspired Computer Vision*. Ed. by Matthias Keil, Gabriel Cristóbal, and Laurent U Perrinet. Weinheim, Germany: Wiley-VCH Verlag GmbH & Co. KGaA, 2015, pp. 319–346.

[20] Scott Shaobing Chen, David L Donoho, and Michael A Saunders. "Atomic decomposition by basis pursuit". In: *SIAM review* 43.1 (2001), pp. 129–159.

[21] Michael Lewicki and Terrence J Sejnowski. "Coding time-varying signals using sparse, shift-invariant representations". In: *Advances in neural information processing systems* 11 (1998).

[22] Thomas Serre, Aude Oliva, and Tomaso Poggio. "A feedforward architecture accounts for rapid categorization". In: *Proceedings of the national academy of sciences* 104.15 (2007), pp. 6424–6429.

[23] Victor Boutin et al. "Pooling Strategies in V1 Can Account for the Functional and Structural Diversity across Species". In: *PLOS Computational Biology* 18.7 (2022), e1010270. ISSN: 1553-7358. DOI: 10.1371/journal.pcbi.1010270. URL: https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1010270 (visited on 09/14/2022).

[24] Brendt Wohlberg. "Efficient algorithms for convolutional sparse representations". In: *IEEE Transactions on Image Processing* 25.1 (2015), pp. 301–315.

[25] Brendt Wohlberg. "SPORCO: A Python package for standard and convolutional sparse representations". In: *Proceedings of the 15th Python in Science Conference, Austin, TX, USA*. 2017, pp. 1–8.
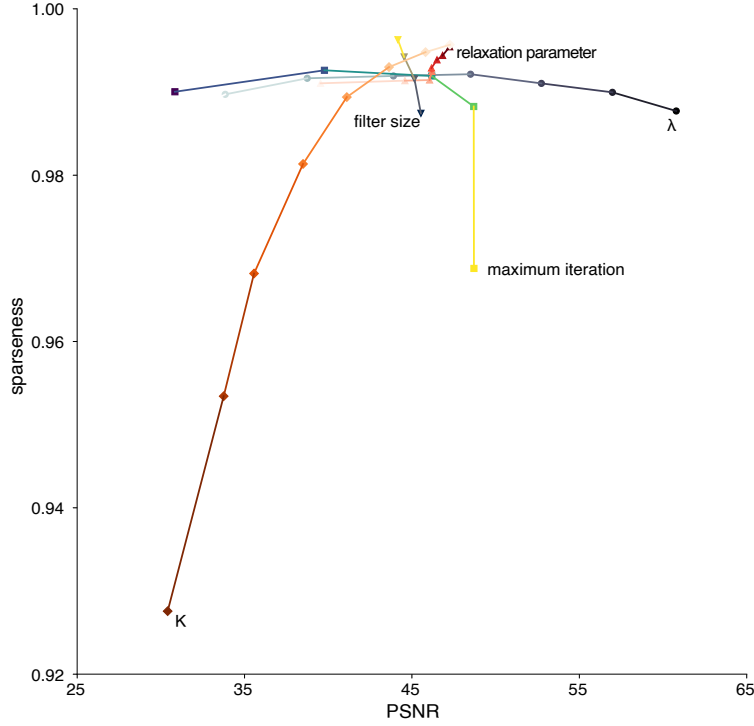
[26] Yu Wang, Wotao Yin, and Jinshan Zeng. "Global convergence of ADMM in nonconvex nonsmooth optimization". In: *Journal of Scientific Computing* 78.1 (2019), pp. 29–63.

[27] David H Hubel and Torsten N Wiesel. "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex". In: *The Journal of physiology* 160.1 (1962), p. 106.

[28] Sylvain Fischer et al. "Sparse Approximation of Images Inspired from the Functional Architecture of the Primary Visual Areas". In: *EURASIP Journal on Advances in Signal Processing* 2007.1 (2007), pp. 1–17. ISSN: 16876172.

[29] Sylvain Fischer et al. "Self-invertible 2D log-Gabor wavelets". In: *International Journal of Computer Vision* 75.2 (2007), pp. 231–246.

[30] Sylvain Fischer et al. "Sparse Approximation of Images Inspired from the Functional Architecture of the Primary Visual Areas". In: *EURASIP Journal on Advances in Signal Processing* 2007.1 (Dec. 2006), pp. 1–17. ISSN: 16876172. DOI: 10.1155/2007/90727. URL: http://dx.doi.org/10.1155/2007/90727.

[31] Nicholas V Swindale. "Orientation tuning curves: empirical description and estimation of parameters". In: *Biological cybernetics* 78.1 (1998), pp. 45–56.

[32] Stefan Van der Walt et al. "scikit-image: image processing in Python". In: *PeerJ* 2 (2014), e453.

[33] Hugo Ladret. "HD natural images database for sparse coding". In: *FigShare* (2023). DOI: "10.6084/m9.figshare.24167265.v1".

[34] Diederik P Kingma and Jimmy Ba. "Adam: A method for stochastic optimization". In: *arXiv preprint arXiv:1412.6980* (2014).

[35] Kaiming He et al. "Deep residual learning for image recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.

[36] Stuart Appelle. "Perception and discrimination as a function of stimulus orientation: the" oblique effect" in man and animals." In: *Psychological bulletin* 78.4 (1972), p. 266.

[37] David J Field. "Relations between the statistics of natural images and the response properties of cortical cells". In: *Josa a* 4.12 (1987), pp. 2379–2394.

[38] Carsen Stringer et al. "High-dimensional geometry of population responses in visual cortex". In: *Nature* 571.7765 (2019), pp. 361–365.

[39] David M Coppola et al. "The distribution of oriented contours in the real world". In: *Proceedings of the National Academy of Sciences* 95.7 (1998), pp. 4002–4006.

[40] Bruce C Hansen and Edward A Essock. "A horizontal bias in human visual processing of orientation and its correspondence to the structural components of natural scenes". In: *Journal of vision* 4.12 (2004), pp. 5–5.

[41] Anastasia Mozhaeva et al. "Full reference video quality assessment metric on base human visual system consistent with PSNR". In: *2021 28th Conference of Open Innovations Association (FRUCT)*. IEEE. 2021, pp. 309–315.

[42] Martin Schrimpf et al. "Brain-score: Which artificial neural network for object recognition is most brain-like?" In: *BioRxiv* (2020), p. 407007.

[43] Nicolas Perez-Nieves et al. "Neural heterogeneity promotes robust learning". In: *Nature communications* 12.1 (2021), p. 5791.

[44] Matteo Di Volo and Alain Destexhe. "Optimal responsiveness and information flow in networks of heterogeneous neurons". In: *Scientific reports* 11.1 (2021), p. 17611.

[45] Christopher J Rozell et al. "Sparse coding via thresholding and local competition in neural circuits". In: *Neural computation* 20.10 (2008), pp. 2526–2563.

[46] Robert Coultrip, Richard Granger, and Gary Lynch. "A cortical model of winner-take-all competition via lateral inhibition". In: *Neural networks* 5.1 (1992), pp. 47–54.

[47] Victor Boutin et al. "Effect of Top-down Connections in Hierarchical Sparse Coding". In: *Neural Computation* 32.11 (2020-02-04, November 2020), pp. 2279–2309.

[48] Jason K Eshraghian, Xinxin Wang, and Wei D Lu. "Memristor-based binarized spiking neural networks: Challenges and applications". In: *IEEE Nanotechnology Magazine* 16.2 (2022), pp. 14–23.

[49] Mostafa Rahimi Azghadi et al. "Complementary metal-oxide semiconductor and memristive hardware for neuromorphic computing". In: *Advanced Intelligent Systems* 2.5 (2020), p. 1900189.

[50] Alexey Dosovitskiy et al. "An image is worth $16 \times 16$ words: Transformers for image recognition at scale". In: *arXiv preprint arXiv:2010.11929* (2020).

[51] Raghavendran Vidya, GM Nasira, and RP Jaia Priyankka. "Sparse coding: a deep learning using unlabeled data for high-level representation". In: *2014 World Congress on Computing and Communication Technologies*. IEEE. 2014, pp. 124–127.

[52] Yunlong He et al. "Unsupervised feature learning by deep sparse coding". In: *Proceedings of the 2014 SIAM international conference on data mining*. SIAM. 2014, pp. 902–910.

[53] Zhuomin Zhang, Jing Li, and Renbing Zhu. "Deep neural network for face recognition based on sparse autoencoder". In: *2015 8th International Congress on Image and Signal Processing (CISP)*. IEEE. 2015, pp. 594–598.

[54] Dylan M Paiton et al. "Selectivity and robustness of sparse coding networks". In: *Journal of vision* 20.12 (2020), pp. 10–10.

[55] Brendt Wohlberg. "Efficient convolutional sparse coding". In: *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 2014, pp. 7173–7177.

# Appendix A - Additional Convolutional Sparse Coding details



Appendix A Figure 1: Parametrization of the CSC learning algorithm. $\lambda$ was varied in 8 steps in a $[0.001 : 0.1]$ range, max iteration in 5 steps in a $[10 : 1000]$ range, relaxation parameter $\rho$ in 8 steps in a $[0.2 : 1.8]$ range, filter size in 8 steps in a $[5 : 21]$ pixels range and $K$ in 8 steps in a $[89 : 2351]$ range.

Convolutional Sparse Coding was implemented using an Alternating Direction Method of Multipliers (ADMM) algorithm, which decomposes the problem into a standard form:

$$\underset{x,y}{\operatorname{argmin}} f(x) + g(y) \tag{12}$$

with the constraint $x = y$. This is then solved iteratively by alternating between the two sub-problems:

$$x_{i+1} = \underset{x}{\operatorname{argmin}} f(x) + \frac{\rho}{2}||x + y_i + \mathbf{u}_i||_2^2 \tag{13}$$

$$y_{i+1} = \underset{y}{\operatorname{argmin}} g(y) + \frac{\rho}{2}||x_{i+1} + y + \mathbf{u}_i||_2^2 \tag{14}$$

where $\rho$ is a penalty parameter that controls the convergence rate of the iterations, also called the relaxation parameter. $x$ and $y$ are residuals whose equality is enforced by the prediction error:

$$\mathbf{u}_{i+1} = \mathbf{u}_i + x_{i+1} + y_{i+1} \tag{15}$$

ADMM can be readily applied to equation (2) by introducing an auxiliary variable $Y$ [55], such that the problem to solve becomes:

$$\operatorname*{argmin}_{\{x_k\},\{y_k\}} \frac{1}{2} || \sum_{k=1}^{K} \mathrm{d}_k * x_k - s ||_2^2 + \lambda \sum_{k=1}^{K} ||y_k||_1 \text{ s.t. } \mathrm{x}_k = \mathrm{y}_k \tag{16}$$

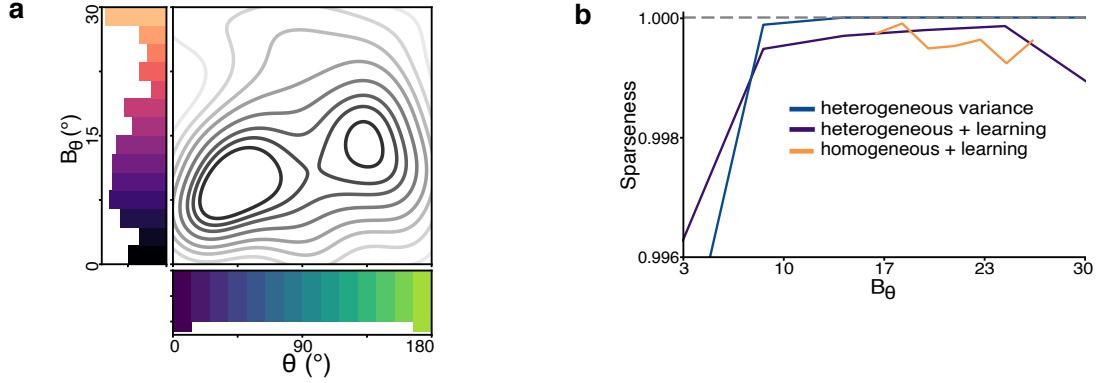which, following the ADMM alternation in equations (13)-(15), is solved by alternating:

$$\{x_k\}_{i+1} = \operatorname*{argmin}_{\{x_k\}} \frac{1}{2} || \sum_{k=1}^{K} \mathrm{d}_k * x_k - s ||_2^2 + \frac{\rho}{2} ||x_k - y_{k,i} + \mathrm{u}_{k,i}||_2^2 \tag{17}$$

$$\{y_k\}_{i+1} = \operatorname*{argmin}_{\{y_k\}} \lambda \sum_{k=1}^{K} ||y_k||_1 + \frac{\rho}{2} ||x_{k,i+1} - y_k + \mathrm{u}_{k,i}||_2^2 \tag{18}$$
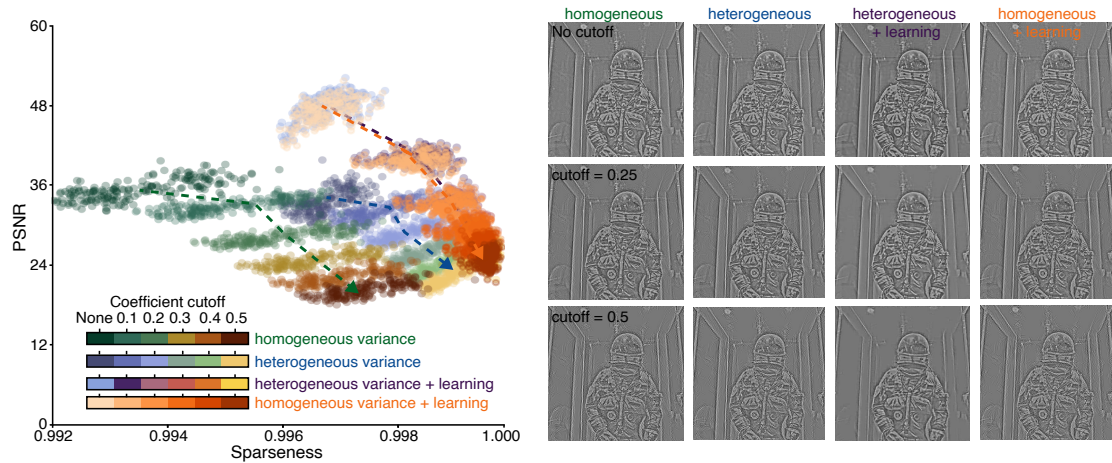
$$\mathrm{u}_{k,i+1} = \mathrm{u}_{k,i} + x_{k,i+1} - y_{k,i+1} \tag{19}$$

# Appendix B - Homogeneous variance dictionary

Results from the main text are shown here for the homogeneous variance dictionary, post-learning.



Appendix B Figure 1: Learning balances coefficient distribution. **(a)** Kernel density estimation of coefficients over $\theta_0$ and $B_\theta$ after learning from the homogeneous variance dictionary. **(b)** Sparseness of coefficients for each $B_\theta$. Sparseness $= 1$ is represented as a gray dashed line.



Appendix B Figure 2: Sparse coefficients can be pruned to boost sparsity. **(a)** Pruning of the coefficients based on their values and resulting sparseness/PSNR for both dictionaries. **(b)** Reconstruction of the image shown in Figure 1 with different cutoff levels.