

---

# On Representation of Natural Image Patches

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 To optimize survival, organisms need to accurately and efficiently relay new in-  
2 formation throughout their systems for processing and responses. Furthermore,  
3 they benefit from predicting environmental occurrences, or in mathematical terms,  
4 understanding the probability distribution of their environment, based on both  
5 personal experiences and inherited evolutionary memory. These twin objectives  
6 of information transmission and learning environmental probabilistic distributions  
7 form the core of an organism's information processing system. While the early  
8 vision neuroscience field has primarily focused on the former, employing infor-  
9 mation theory as a guiding framework [3, 32, 19, 1, 9, 28], the latter is largely  
10 explored by the machine learning community via probabilistic generative models.  
11 However, the relationship between these two objectives has not been thoroughly  
12 investigated. In this paper, we study a biologically inspired information processing  
13 model and prove that these two objectives can be achieved independently. By  
14 evenly partitioning the input space to model input probability, our model bypasses  
15 the often intractable normalization factor computation. When applied to image  
16 patches, this model produces a sparse, nonlinear binary population code similar  
17 to early visual systems, with features like edge-detection and orientation-selective  
18 units. Our results not only offer potential new insights into the functioning of  
19 neurons in early vision systems, but also present a novel approach to represent  
20 natural image patches.

## 21 1 Introduction

22 Nature, through billions of years of evolution, has likely developed optimal methods for processing  
23 visual information within the constraints of biological feasibility. However, attempting to precisely  
24 emulate every detail of these biological systems [30, 21] in order to construct an optimal visual  
25 information processing model presents significant complexities, especially without a comprehensive  
26 understanding of the underlying principles.

27 In parallel, deep learning models, particularly Convolutional Neural Networks (CNNs), have demon-  
28 strated exceptional performance in various computer vision tasks, such as image classification, object  
29 detection, and semantic segmentation. Despite their success, these models still fall short of biological  
30 vision systems in several key areas such as the ability to generalize from limited data, robustness to  
31 variations, 3D understanding, and processing speed and efficiency.

32 Moreover, deep learning models pose several unresolved challenges. Firstly, their decision-making  
33 processes are often opaque, leading to the "black box" label. Secondly, deep learning architectures  
34 typically involve numerous layers without explicit functions associated with each layer [12, 8], unlike  
35 the biological brain where each stage of visual processing has a distinct role and purpose [7]. Lastly, he

36 sequence of information processing in deep learning models is largely dictated by their architectural  
37 structure. It is still a question how to ascertain when one processing stage is finished and when it's  
38 appropriate to pass information to the next layer.

39 Drawing inspiration from biological systems and prior studies [19, 1, 9, 28], and with a view to  
40 find an alternative approach to the current deep learning framework, this paper aims to explore the  
41 fundamentals of the first stage of an optimal visual information processing system. This exploration is  
42 undertaken incrementally, starting from a single pixel, and progressively advancing to image patches.

## 43 2 One Pixel

44 We begin with the simplest conceivable case, where the input to the model consists of a single pixel  
45 with one color channel. Although seemingly trivial, this model can represent various biological units.  
46 For example, it could model the eyespot of single-celled organisms like *Euglena*, the large monopolar  
47 cells found in an insect's compound eye or the bipolar cells in the retina. The single pixel case has  
48 been studied [14, 1]. We aim to review it to introduce the central concepts and the main theory, which  
49 will also be applicable to more complex scenarios.

50 Let us denote the light intensity of the pixel as  $x$ , and let  $p(x)$  represent its probability distribution.  
51 We define an information processing unit (IPU) as a model that receives inputs and passes processed  
52 information to subsequent stages. As the first stage in an organism's information processing system,  
53 the single pixel IPU carries the same dual objectives as later stages: transmitting information and  
54 learning environmental probabilistic distributions. Therefore, the two objectives for the single pixel  
55 IPU are to transmit information about  $x$  efficiently through its output and to learn  $p(x)$ .

56 Information is quantified by Shannon entropy. To compute the Shannon entropy of the input, we  
57 assume  $x$  is a discrete variable with  $M$  states, after all light intensity is quantized according to quantum  
58 mechanics, although  $M$  could be a very large number, making  $x$  practically indistinguishable from a  
59 continuous variable. The information obtained from the pixel when we know the intensity is  $x$  can be  
60 represented as  $I(x) = -\log p(x)$ . The average information a state of the pixel contains, the Shannon  
61 entropy, is given by:

$$H_p = -\sum_{i=1}^M p(x_i) \log p(x_i). \quad (1)$$

62 The IPU transforms the input  $x$  into the output  $y = f(x)$ , where the output space comprises  $N$   
63 distinct states. In contrast to previous studies that assumed one-to-one mapping between input and  
64 output, here we posit that  $N \ll M$ . This assumption is more congruent with biological constraints;  
65 for instance, the luminance resolution levels at synapse terminals in a zebrafish's retina are only about  
66 10 [25]. Moreover, this assures that after processing the information is significantly reduced, thereby  
67 simplifying the tasks for subsequent stages.

68 The function  $f(x)$  assigns  $x$  into  $N$  groups, each corresponding to a fixed  $y$ . We denote all  $x$  values  
69 in group  $j$  as  $G_j$ , and the size of this group as  $n_j$ . The entropy of the output is given by

$$H_Q = -\sum_{j=1}^N Q(y_j) \log Q(y_j), \quad (2)$$

70 where  $Q(y)$  represents the probability distribution of the output states.

71 Previous research [14, 1] has primarily emphasized the first objective of an IPU, which is maximizing  
72 the rate of transmission [19]. This goal is particularly relevant for early-stage IPUs, where the  
73 distinction between signal and noise is not yet clear. Maximizing the rate of transmission is equivalent  
74 to maximizing  $H_Q$  (see proof in Appendix A), leading to a constant  $Q(y)$ . Biological neurons have  
75 been observed to follow this coding scheme [14].

76 Simultaneously, an IPU should also strive to fulfill the second objective and model  $p(x)$  as accurately  
77 as possible. Mathematically, this involves minimizing the Kullback–Leibler divergence between

78  $p(x)$  and the distribution learned by the IPU. This raises an interesting question: Are these two  
 79 optimization objectives contradictory, or do they essentially represent the same task?

### 80 3 Even Code Principle

81 To determine how an IPU models  $p(x)$  we need to translate the output probability distribution  $Q(y)$   
 82 into the input space as  $q(x)$ .  $q(x)$  is a step function:

$$q(x) = q_j, \text{ for } x \in G_j, \quad (3)$$

83 and we have the following relations:

$$Q(y_j) = \sum_{x \in G_j} p(x) = \sum_{x \in G_j} q(x) = n_j q_j. \quad (4)$$

84 Minimizing the difference between  $p(x)$  and  $q(x)$  can be achieved by minimizing their Kull-  
 85 back–Leibler divergence:

$$D_{KL}(p||q) = H_{pq} - H_p, \quad (5)$$

86 where  $H_{pq}$  is the cross entropy. It can be proved that the cross entropy  $H_{pq}$  is equal to the entropy of  
 87 the learned distribution in the input space defined as (see proof in Appendix B):

$$H_q = - \sum_x q(x) \log q(x), \quad (6)$$

88 and we get

$$D_{KL}(p||q) = H_q - H_p. \quad (7)$$

89 Since  $H_p$  is fixed, minimizing the KL divergence requires minimizing  $H_q$ . The previous question now  
 90 transforms into understanding the relationship between maximizing the entropy of the distribution in  
 91 the output space ( $H_Q$ ) and minimizing the entropy of the learned distribution in the input space ( $H_q$ ).

92 Suppose we have two adjacent zones in the transformed space where the corresponding  $Q(y_1)$  and  
 93  $Q(y_2)$  are not equal, let's assume  $Q(y_1) > Q(y_2)$ . One can reduce the inequality by shifting the  
 94 boundary between these two zones and moving one  $x$  value from  $G_1$  to  $G_2$ . This shift corresponds  
 95 to a small change of probability,  $\delta$ , for both zones. Note that  $\delta$  is comparable to  $q_1$  and  $q_2$ , as we  
 96 assume the distribution is smooth. We know that reducing the inequality of  $Q(y_1)$  and  $Q(y_2)$  always  
 97 increases  $H_Q$ . If the two optimization problems are the same, then  $H_q$  should increase; if they are  
 98 contradictory,  $H_q$  should decrease. The change of  $H_q$  can be calculated as:

$$\begin{aligned} \Delta H_q &= -[Q(y_1) - \delta] \log \frac{Q(y_1) - \delta}{n_1 - 1} - [Q(y_2) + \delta] \log \frac{Q(y_2) + \delta}{n_2 + 1} \\ &\quad + Q(y_1) \log \frac{Q(y_1)}{n_1} + Q(y_2) \log \frac{Q(y_2)}{n_2} \end{aligned} \quad (8)$$

$$= q_2 - q_1 + \delta(\log q_1 - \log q_2 + \frac{1}{n_1} + \frac{1}{n_2}) + O(\delta^2) + O(\frac{1}{n_1^2}) + O(\frac{1}{n_2^2}) \quad (9)$$

$$\approx q_2 - q_1 + \delta \log \frac{q_1}{q_2}. \quad (10)$$

99 The change can either be positive or negative depending on  $q_1$  and  $q_2$ . Since minimizing  $H_q$  and  
 100 maximizing  $H_Q$  are not contradictory, these objectives can be tackled independently. Given a fixed  
 101 number of output levels  $N$ , we first maximize  $H_Q$  to retain as much input information as possible. If  
 102 further refinement of  $p(x)$  modeling is required, we can increase the output resolution  $N$ .

103 The aforementioned reasoning extends naturally to multivariate scenarios, as no assumptions about  
 104 one-dimensionality of the input were made. We articulate the goal of a general information processing  
 105 unit as follows: An information processing unit (IPU) transforms input space with  $M$  states into  
 106 output space with  $N$  states, where  $N \ll M$ . Given a smooth input probability distribution as  
 107  $M \rightarrow \infty$  and a piecewise smooth transformation function, the *sole* goal of an IPU with a fixed output  
 108 resolution  $N$  is to yield an even output probability distribution, hence retaining maximum information  
 109 from the input. To attain better modeling precision, the output resolution  $N$  of the IPU should be  
 110 increased. This will be referred to as the principle of even code. In the next sections, we will apply  
 111 the even code principle to more complex inputs.

112 **4 Two Pixels**

113 For two pixels  $(x_1, x_2)$ , we can either use one IPU directly to model  $p(x_1, x_2)$  or use two IPUs to  
 114 model  $p(x_1)$  and  $p(x_2)$  separately, followed by another IPU to model the outputs  $p(y_1, y_2)$ . We will  
 115 use the second approach, as processing as much information locally reduces the cost of information  
 116 transfer. In fact, when images are stored on computers, gamma encoding is utilized to create an  
 117 approximately even distribution of pixel values. When these images are displayed, pixel values  
 118 undergo gamma correction to recover the original statistics for human eyes to process. In the  
 119 following sections, we will assume that all pixel values  $x$  have already been processed by dedicated  
 120 IPUs, resulting in a roughly even probability distribution.

121 The probability distribution  $p(x_1, x_2)$  of natural images is relatively simple. The majority of the  
 122 probability is concentrated around the diagonal line  $x_1 - x_2 = 0$ , with  $p(x_1, x_2)$  rapidly decaying  
 123 as  $|x_1 - x_2|$  increases (see Fig. 1 (a) for example). Intuitively, we can use lines parallel or/and  
 124 perpendicular to  $x_1 - x_2 = 0$  to divide the probability distribution into even partitions.

125 **4.1 One Basis**

126 To investigate how IPUs learn  $p(x_1, x_2)$ , we conduct numerical experiments using a multilayer  
 127 perceptron (MLP) as the IPU to approximate  $y = f(x)$  and model  $p(x)$  [23]. Other function  
 128 approximation methods may also be applicable. To partition the input probability distribution with  
 129 one set of parallel lines, only one IPU with  $N$  output nodes is needed. According to the even code  
 130 principle, for each input, only one of the  $N$  output nodes should be activated, and the probability  
 131 of activating any one of the  $N$  output nodes should be equal. We use the softmax function as the  
 132 last layer of the MLP to ensure each output value is within  $[0, 1]$ , and that if a node is activated  
 133 (output value equals 1), it is the only node being activated. We use stochastic gradient descent and  
 134 the following loss function to train the MLP:

$$E = \sum_i \langle y_{si} \rangle_s \log \langle y_{si} \rangle_s + k \langle - \sum_i y_{si} \log y_{si} \rangle_s. \quad (11)$$

135  $y_{si}$  represents the value of the  $i$ -th output node for the  $s$ -th input sample, while  $\langle \cdot \rangle_s$  denotes the average  
 136 over all samples in a training batch. The first term in the loss function ensures each output node has  
 137 an equal chance to be activated on average. The second term promotes activation of only one node  
 138 per input while suppressing the remaining nodes, mimicking lateral inhibition when combined with  
 139 the softmax function. The factor  $k$  balances the two terms to achieve the desired result. Fig. 1 (a)  
 140 show the results learned by MLPs with 16 output nodes.

141 **4.2 Multiple Bases**

142 To partition the input space with two sets of orthogonal lines we need two MLPs. The orthogonality  
 143 is achieved by enforcing

$$Q(y_1, y_2) = \frac{1}{N_1 N_2}, \quad (12)$$

144 where  $N_1$  and  $N_2$  represent the number of output nodes of the two MLPs (refer Appendix C for  
 145 proof). If more than two orthogonal bases are required for partitioning the space, we can enforce  
 146 Eq. (12) for each combination of two bases to ensure orthogonality between them. The loss function  
 147 for multiple orthogonal bases with independent states is

$$E = \frac{1}{\binom{B}{2}} \sum_{\langle b, b' \rangle} \sum_{ij} \langle y_{bsi} y_{b'sj} \rangle_s \log \langle y_{bsi} y_{b'sj} \rangle_s + \frac{k}{B} \langle - \sum_{b=1}^B \sum_i y_{bsi} \log y_{bsi} \rangle_s, \quad (13)$$

148 where  $b$  is the base index, and  $B$  is the number of bases.  $\sum_{\langle b, b' \rangle}$  denotes the sum over all  $\binom{B}{2}$   
 149 combinations of two distinct bases.  $y_{bsi} y_{b'sj}$  is the probability  $Q(y_b, y_{b'})$  for the sample  $s$  when  $y_b$   
 150 and  $y_{b'}$  take their  $i$ -th and  $j$ -th value respectively.

151 Fig. 1 (b) shows an example of two-pixel input space partitioning using two orthogonal bases. For a  
 152 more detailed discussion on the experiments and additional results, refer to Appendix D.

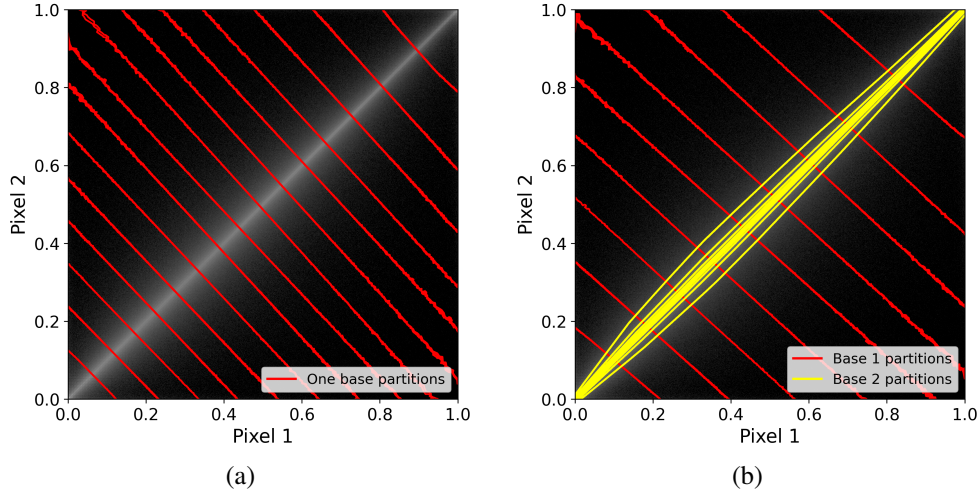


Figure 1: Evenly partitioning the two-pixel probability distribution learned by multilayer perceptrons (MLPs). The X and Y axes represent the rescaled intensities  $x_1$  and  $x_2$  of the two pixels in the range  $[0, 1]$ . The quantity  $n(x_1, x_2) + 1$  is plotted in gray on a log scale, where  $n(x_1, x_2)$  denotes the number of occurrences of the two-pixel values among the sampled data. Color lines indicate the boundaries of states for each basis learned by an MLP, with one color representing one basis. (a) One basis with 16 independent states, which partitions the space based on the total intensity  $x_1 + x_2$ . (b) Two orthogonal bases, each with 10 independent states, dividing the space based on the total intensity  $x_1 + x_2$  and the contrast  $x_1 - x_2$  approximately.

153 Additionally, it's worth noting that orthogonal bases with independent states might model grid cells  
 154 [11] in the entorhinal cortex, though this topic is beyond the scope of the current paper.

## 155 5 Image Patches

156 Next we move on to study gray and color image patches. We use  $\mathbf{x}$  to represent the vector of input  
 157 pixel values of the image patch. The multivariate input probability distribution  $p(\mathbf{x})$  is considerably  
 158 more complex compared to the previous examples. If we use only one basis to discretize the input  
 159 probability space (e.g. Fig. 1 (a) in Appendix D), the required number of independent states for a good  
 160 approximation would be very large, making the evaluation of the softmax function computationally  
 161 expensive. On the other hand, using multiple orthogonal bases would also significantly increase the  
 162 computational cost to ensure orthogonality if the number of bases is more than just a few. Additionally,  
 163 determining the optimal number of bases and the number of independent states for each basis are  
 164 challenging.

165 Aside from the computational cost, another issue arises when working with image patches: we want  
 166 the representation to capture the similarity between inputs. However, using orthogonal bases with  
 167 independent states makes it difficult to gauge input similarity through methods such as calculating  
 168 the difference between representations, even if we can establish an order for the states of each basis.  
 169 Therefore, we need a more suitable coding scheme for complex inputs like image patches.

170 Real-valued vectors are a natural choice, given their extensive use in representing a variety of entities  
 171 such as images, texts, and categorical variables [13, 22, 10]. The norm of the difference between  
 172 two vectors can function as a measure of similarity. Nevertheless, if we want the representation  $\mathbf{y}$  to  
 173 mirror input similarity, each value of  $\mathbf{y}$  should encapsulate all samples perceived as identical within  
 174 the same group  $G$ . Under this constraint,  $Q(\mathbf{y})$  cannot remain constant, thereby conflicting with the  
 175 even code principle.

176 The resolution to this conflict involves permitting the representation to mirror input similarity at the  
 177 most granular level, while enforcing the even code principle at a larger scale in the transformed space.  
 178 We will detail this method in subsequent sections.

## 179 5.1 Loss Function

180 To promote even distribution, we incorporate a loss function that compels input samples to repel each  
 181 other in the transformed space. This repulsive force diminishes with increasing distance, as described  
 182 by the following equation:

$$E = \langle -\ln |\mathbf{y}_s - \mathbf{y}_{s'}| \rangle_{\langle s, s' \rangle}. \quad (14)$$

183 Here,  $-\ln |\mathbf{y}_s - \mathbf{y}_{s'}|$  represents the potential energy due to the repulsive force, which is proportional  
 184 to the inverse of the Minkowski distance between the representations of samples  $s$  and  $s'$  in the  
 185 transformed space. Alternative forms of potential energy and distance measures could also be  
 186 applicable.  $\langle \rangle_{\langle s, s' \rangle}$  denotes the average over all sample pairs.

187 Should numerous samples converge at one point in the transformed space, they will exert a strong  
 188 repulsive force in the surrounding area, thereby discouraging other samples from occupying nearby  
 189 positions. To prevent samples from pushing each other infinitely far apart, we restrict the represen-  
 190 tation values to be within the range  $[0, 1]$ . With this constraint, the repulsive force pushes samples  
 191 towards the vertices of the unit hypercube, effectively reducing the representations from real vectors  
 192 to binary vectors. As a result, an even distribution is achieved on a larger scale in the transformed  
 193 space, which consists solely of the vertices.

194 In the context of binary vectors, the collection of output nodes can be viewed as a vocabulary,  
 195 and activated nodes by an input image patch act as its representative tokens. Unlike fixed-length  
 196 representations with real-valued vectors, binary representations can employ fewer tokens for more  
 197 common image patches (e.g., homogeneous patches), and more tokens for less common, structurally-  
 198 rich patches. This can be accomplished by introducing a second term,  $\langle |\mathbf{y}_s| \rangle_s$ , to the loss function,  
 199 which echoes the sparsity regularization term found in various studies [9, 16, 26, 27, 29, 4]. The  
 200 updated loss function becomes:

$$E = \langle -\ln |\mathbf{y}_s - \mathbf{y}_{s'}| \rangle_{\langle s, s' \rangle} + \alpha \langle |\mathbf{y}_s| \rangle_s, \quad (15)$$

201 where  $\alpha$  is a free parameter to adjust sparsity.

202 In practice, we add a small value  $\epsilon = 10^{-38}$  to the distance, allowing slightly different samples  
 203 to share the same representation and enhancing numerical stability. Another approach to improve  
 204 numerical stability involves using a theoretically equivalent form of the loss function, which instead  
 205 of allowing samples to repel each other in the output space, we enable nodes to repel one another,  
 206 encouraging output nodes to be as independent as possible [24].

## 207 5.2 Experiments

208 In the following experiments, we use either a single MLP with  $N$  outputs, or  $N$  MLPs each with  
 209 one output, as the IPU to approximate the transformation function  $\mathbf{y} = f(\mathbf{x})$  and model  $p(\mathbf{x})$ . The  
 210 last layer of the MLP is a sigmoid layer, ensuring the output value ranges between 0 and 1. Our  
 211 training data comprises random image patches extracted from the COCO 2017 image dataset [18] or  
 212 the ImageNet dataset [6]. No image preprocessing is used. Additional training details are provided in  
 213 Appendix E.

### 214 5.2.1 Output Statistics

215 First, we examine the statistics of the learned representation. Across all experiments, we observe  
 216 qualitatively similar output statistics, irrespective of the IPU architectures and training specifics,  
 217 provided the training has properly converged. For illustration, we present an example using a model  
 218 trained on  $5 \times 5$  color image patches. It uses 96 MLPs, each with one output node and a middle  
 219 layer of 48 nodes, as the IPU. Following training, the model is used to generate representations for 1  
 220 million random image patches for this analysis.

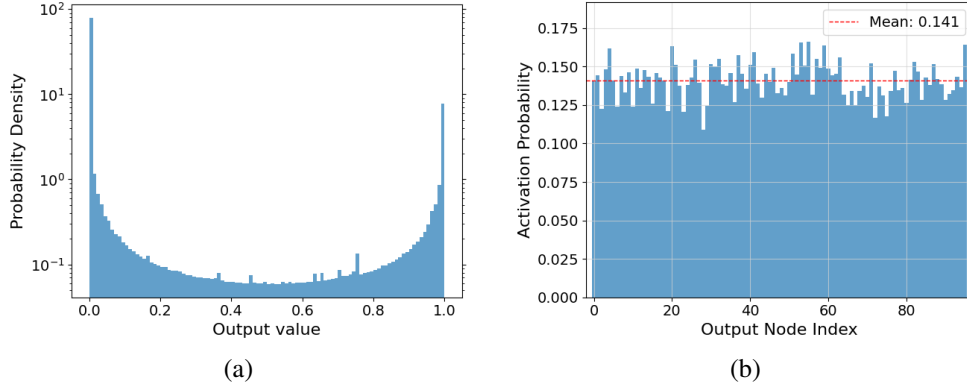


Figure 2: Statistical analysis of the learned representation using the loss function Eq. (15). (a) Histogram of the model’s output values on a log scale. (b) Probability of an output node being activated by a random image patch.

221 Fig. 2 (a) presents the histogram of output values on a logarithmic scale. The vast majority of output  
 222 values are either at 0 or 1. As such during inference, we can round the outputs to yield a binary  
 223 representation. Fig. 2 (b) illustrates the probability of an output node being activated by a random  
 224 image patch. All nodes demonstrate similar activation probabilities, indicating an even distribution at  
 225 this coarse scale. Further statistical analysis of the output representation is available in Appendix F.

### 226 5.2.2 Image Patch Similarity

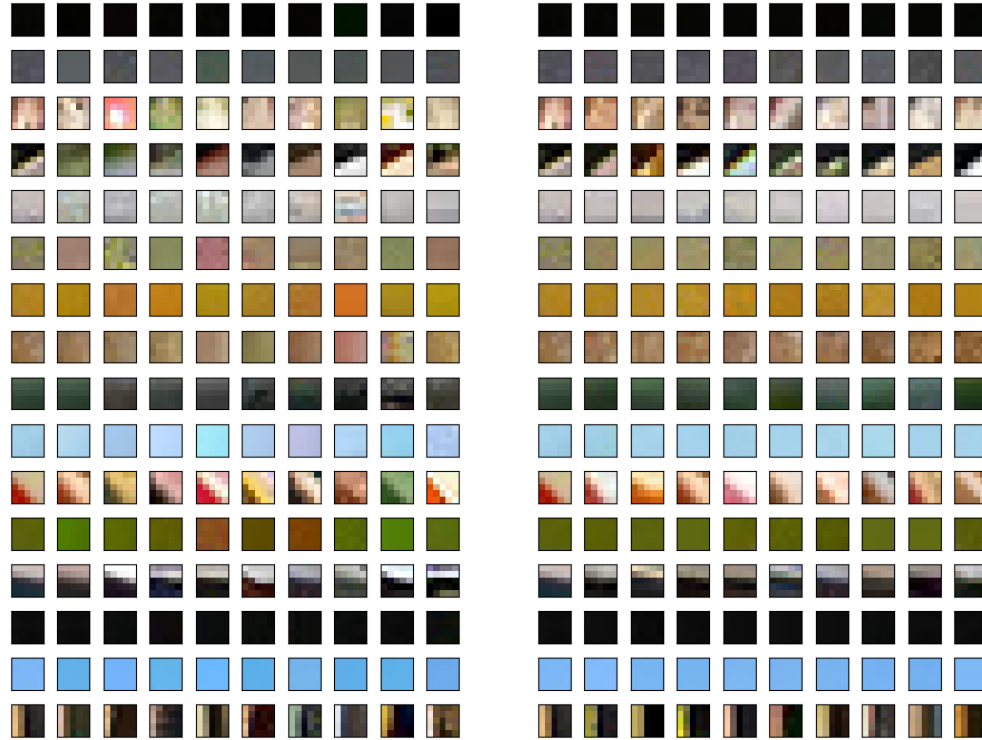
227 Next, to examine how the learned representation reflects the similarity between image patches, we  
 228 display 16 random image patches, each followed by 9 image patches similar to them in the binary  
 229 representation space, as shown in Fig. 3 (a). The same  $5 \times 5$  color image patch model is used. The  
 230 learned representation clearly captures perceptual similarity. The results shown in Fig. 2, Fig. 3 (a),  
 231 and additional results in Appendix E confirm that with the loss function Eq. (15), we can indeed learn  
 232 a sparse binary representation which reflects the image similarity while adhering to the even code  
 233 principle.

234 For comparison with a traditional convolutional neural network, we present the results generated with  
 235 the first 10 layers of a VGG16 model [31] pre-trained on ImageNet in Fig. 3 (b). The image patch  
 236 representation from the first 10 layers of the VGG16 model is a float vector of size 128. The even  
 237 code model, with only 96 binary outputs, achieves results similar to the VGG16 model. These 96  
 238 binary outputs occupy the same storage space as a float vector of length 3 — just 1/42 of the VGG16  
 239 representation’s size, which underscores the exceptional efficiency of the even code method in image  
 240 patch representation.

### 241 5.2.3 Local Edge Detectors and Orientation-Selective Units

242 Biological visual systems’ initial stages are known to possess local edge detectors and orientation-  
 243 selective units [17, 2]. While CNNs have been successfully trained to detect boundaries via supervised  
 244 learning [20], their initial layers have not shown proficiency in edge detection [15]. Notably, prevalent  
 245 local edge detection algorithms, such as the Canny edge detector [5], still primarily rely on non-deep  
 246 learning methods.

247 Does the even code model, proposed as the initial stage of an optimal image processing system,  
 248 resemble biological systems more closely? To answer this, we trained an even code model on  $4 \times 4$   
 249 grayscale image patches and applied it to images with a stride of 1 pixel, generating feature maps for  
 250 each output node. The model comprises a MLP with 64 outputs and an intermediary layer with 100  
 251 nodes. Fig. 4 illustrates the feature maps of 4 output nodes of the even code model for 4 different  
 252 input images. It also shows edges generated by the Canny edge detector as comparison. Interestingly,



(a) Even code model with 96 binary outputs      (b) Early layers of VGG16 with 128 real outputs

Figure 3: Image patches with the shortest distance in the representation space to 16 randomly selected image patches. The first column presents the 16 random image patches, while the succeeding nine columns display patches that are closest to the first-column patches in the same row. (a) Distances are computed using using an even code model with 96 binary outputs. (b) Distances are computed using the first 10 layers of a pretrained VGG16 model with 128 real outputs.

253 with this simple network architecture, the even code model demonstrated a remarkable capability in  
 254 edge detection, rivaling the multi-stage Canny edge detector.

255 Furthermore, Fig. 5 shows the feature maps of 5 output nodes for a sample bike image. Spokes of  
 256 different orientations activate different nodes, indicating that these output nodes of the even code  
 257 model have varying orientation preferences, similar to orientation-selective units found in biological  
 258 systems.

## 259 6 Conclusion

260 In summary, this paper demonstrates that maximizing the information-carrying capacity of output  
 261 channels and modeling the input probability distribution are not mutually exclusive objectives and  
 262 can be pursued independently. Given a specific output resolution, the sole goal of an information  
 263 processing unit is to preserve as much information from the input as possible by ensuring an even  
 264 distribution of samples in the output space. We applied the even code principle to study the probability  
 265 distributions of two-pixel systems and image patches. For the two-pixel system, we learned orthogonal  
 266 bases with independent states to model its probability distribution. For image patches, the even code  
 267 approach naturally led to a nonlinear sparse binary representation. The even code model also shares  
 268 additional similarities with early visual systems, such as the presence of local edge-detecting and  
 269 orientation-selective units. These similarities suggest that the even code model could potentially  
 270 serve as a new representation for neurons in early visual systems.



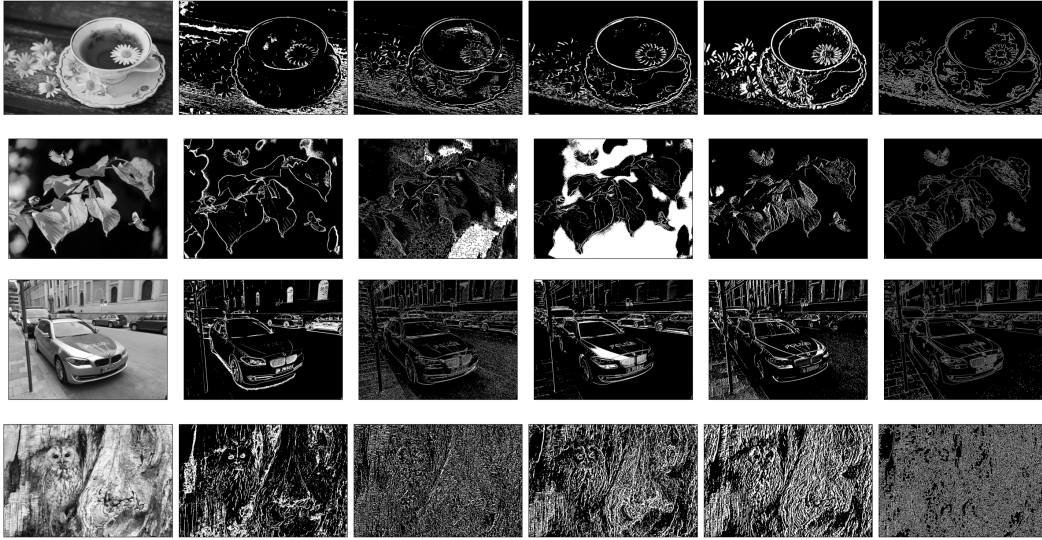


Figure 4: Feature maps of nodes resembling local edge detectors. The first column presents four grayscale test images. Each subsequent column, except the last one, displays the feature maps corresponding to the same output node for the four test images. The last column shows edges generated by the multi-stage Canny edge detector for comparison.

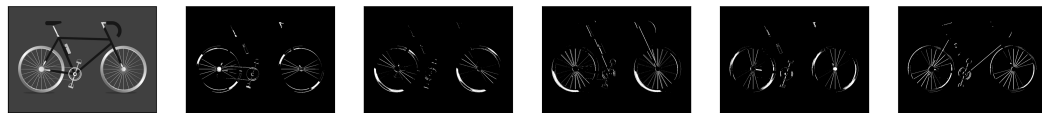


Figure 5: Feature maps of five orientation-selective nodes applied to a test image. Spokes in different orientations activate distinct nodes, illustrating the orientation-selectivity of these nodes.

271 There are several intriguing directions for future research. First, the even code model has been  
 272 applied to inputs ranging from as simple as one pixel to more complex color image patches. Can we  
 273 extend its application beyond the early stage of visual information processing? Second, the even code  
 274 model could be extended to videos by incorporating an additional time dimension alongside color,  
 275 width, and height dimensions. Investigating time-varying inputs, which produce spike train-like  
 276 outputs, and conducting an in-depth comparison with early visual systems would be very interesting.  
 277 Third, the even code model can also be extended to binocular vision data by adding another input  
 278 dimension of size two. Whether the model with binocular and/or video data can construct a 3D  
 279 model of the world based on data of two spatial dimensions is an intriguing question. Fourth, while  
 280 this paper focuses on visual information, the even code model is a general method that could be  
 281 applied to model other multivariate probability distributions as well. Lastly, on the application  
 282 side, the even code model has potential in various areas, including local edge detection, image and  
 283 video compression/denoising/retrieval, texture classification, and multispectral/hyperspectral image  
 284 processing.

## 285 References

- 286 [1] Joseph J. Atick. Could information theory provide an ecological theory of sensory processing?  
 287 *Network: Computation in neural systems*, 3(2):213–251, 1992.
- 288 [2] Tom Baden, Philipp Berens, Katrin Franke, Miroslav Román Rosón, Matthias Bethge, and  
 289 Thomas Euler. The functional diversity of retinal ganglion cells in the mouse. *Nature*,  
 290 529(7586):345–350, 2016.

- 291 [3] Horace B Barlow et al. Possible principles underlying the transformation of sensory messages.  
292 *Sensory communication*, 1(01):217–233, 1961.
- 293 [4] Michael Beyeler, Emily L. Rounds, Kristofor D. Carlson, Nikil Dutt, and Jeffrey L. Krichmar.  
294 Neural correlates of sparse coding and dimensionality reduction. *PLOS Computational Biology*,  
295 15(6):e1006908, jun 2019.
- 296 [5] John Canny. A computational approach to edge detection. *IEEE Transactions on pattern*  
297 *analysis and machine intelligence*, (6):679–698, 1986.
- 298 [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-  
299 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*  
300 *recognition*, pages 248–255. Ieee, 2009.
- 301 [7] James J DiCarlo, Davide Zoccolan, and Nicole C Rust. How does the brain solve visual object  
302 recognition? *Neuron*, 73(3):415–434, 2012.
- 303 [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
304 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.  
305 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*  
306 *arXiv:2010.11929*, 2020.
- 307 [9] David J. Field. What Is the Goal of Sensory Coding? *Neural Computation*, 6(4):559–601, jul  
308 1994.
- 309 [10] Cheng Guo and Felix Berkhahn. Entity embeddings of categorical variables. *arXiv preprint*  
310 *arXiv:1604.06737*, 2016.
- 311 [11] Torkel Hafting, Marianne Fyhn, Sturla Molden, May Britt Moser, and Edvard I. Moser. Mi-  
312 crostructure of a spatial map in the entorhinal cortex. *Nature* 2005 436:7052, 436(7052):801–  
313 806, jun 2005.
- 314 [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
315 recognition. arxiv 2015. *arXiv preprint arXiv:1512.03385*, 14, 2015.
- 316 [13] G. E. Hinton and R. R. Salakhutdinov. Reducing the Dimensionality of Data with Neural  
317 Networks. *Science*, 313(5786):504–507, 2006.
- 318 [14] Simon Laughlin. A Simple Coding Procedure Enhances a Neuron’s Information Capacity.  
319 *Zeitschrift für Naturforschung C*, 36(9-10):910–912, oct 1981.
- 320 [15] Minh Le and Subhradeep Kayal. Revisiting edge detection in convolutional neural networks. In  
321 *2021 International Joint Conference on Neural Networks (IJCNN)*, pages 1–9. IEEE, 2021.
- 322 [16] Ann B Lee, Kim S Pedersen, and David Mumford. The Nonlinear Statistics of High-Contrast  
323 Patches in Natural Images. *International Journal of Computer Vision*, 54(5413):83–103, 2003.
- 324 [17] William R Levick. Receptive fields and trigger features of ganglion cells in the visual streak of  
325 the rabbit’s retina. *The Journal of physiology*, 188(3):285, 1967.
- 326 [18] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays,  
327 Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common  
328 objects in context, 2014.
- 329 [19] Ralph Linsker. Self-organization in a perceptual network. *Computer*, 21(3):105–117, 1988.
- 330 [20] David R Martin, Charless C Fowlkes, and Jitendra Malik. Learning to detect natural image  
331 boundaries using local brightness, color, and texture cues. *IEEE transactions on pattern analysis*  
332 *and machine intelligence*, 26(5):530–549, 2004.
- 333 [21] Richard H. Masland. The Neuronal Organization of the Retina, oct 2012.

- 334 [22] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word  
335 representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- 336 [23] Guido Montúfar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio. On the Number of  
337 Linear Regions of Deep Neural Networks. *Advances in Neural Information Processing Systems*,  
338 4(January):2924–2932, feb 2014.
- 339 [24] Jean Pierre Nadal and Nestor Parga. Nonlinear neurons in the low-noise limit: a factorial code  
340 maximizes information transfer. [http://dx.doi.org/10.1088/0954-898X\\_5\\_4\\_008](http://dx.doi.org/10.1088/0954-898X_5_4_008), 5(4):565–581,  
341 2009.
- 342 [25] Benjamin Odermatt, Anton Nikolaev, and Leon Lagnado. Encoding of Luminance and Contrast  
343 by Linear and Nonlinear Synapses in the Retina. *Neuron*, 73:758–773, 2012.
- 344 [26] Bruno A Olshausen. sparse codes and spikes. *Probabilistic models of the brain*, page 257, 2002.
- 345 [27] Bruno A Olshausen and David J Field. Sparse coding of sensory inputs. *Current opinion in*  
346 *neurobiology*, 14(4):481–487, 2004.
- 347 [28] Bruno A. Olshausen and David J. Field. Emergence of simple-cell receptive field properties by  
348 learning a sparse code for natural images. *Nature*, oct 2015.
- 349 [29] Marc’ Aurelio Ranzato, Christopher Poultney, Sumit Chopra, and Yann Cun. Efficient learning of  
350 sparse representations with an energy-based model. *Advances in neural information processing*  
351 *systems*, 19, 2006.
- 352 [30] Joshua R. Sanes and S. Lawrence Zipursky. Design Principles of Insect and Vertebrate Visual  
353 Systems, 2010.
- 354 [31] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale  
355 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 356 [32] JH van Hateren. A theory of maximizing sensory information. *Biol. Cybern*, 68:23–29, 1992.