

---

# Learning a Convolutional Bilinear Sparse Code for Natural Videos

---

**Dimitrios C. Gklezakos**  
Paul G. Allen School of CSE  
University of Washington  
gklezd@cs.washington.edu

**Rajesh P. N. Rao**  
Paul G. Allen School of CSE &  
Center for Neurotechnology  
University of Washington  
rao@cs.washington.edu

## Abstract

1 In contrast to the monolithic deep architectures used in deep learning today for  
2 computer vision, the visual cortex processes retinal images via two functionally  
3 distinct but interconnected networks: the ventral pathway for processing object-  
4 related information and the dorsal pathway for processing motion and transforma-  
5 tions [8]. Inspired by this cortical division of labor and properties of the magno-  
6 and parvocellular systems [5], we explore an unsupervised approach to feature  
7 learning that jointly learns object features and their transformations from natu-  
8 ral videos. We propose a new convolutional bilinear sparse coding model that  
9 (1) allows independent feature transformations and (2) is capable of processing  
10 large images. Our learning procedure leverages smooth motion in natural videos.  
11 Our results show that our model can learn groups of features and their transfor-  
12 mations directly from natural videos in a completely unsupervised manner. The  
13 learned "dynamic filters" exhibit certain equivariance properties, resemble corti-  
14 cal spatiotemporal filters, and capture the statistics of transitions between video  
15 frames. Our model can be viewed as one of the first approaches to demonstrate  
16 unsupervised learning of primary "capsules" (proposed by Hinton and colleagues  
17 for supervised learning) and has strong connections to the Lie group approach to  
18 visual perception.

## 19 1 Motivation

20 During early development, the brain learns a general-purpose internal representation of objects from  
21 unlabeled image sequences. This representation is compositional and leverages the decomposition  
22 of objects into parts, sub-parts, and features, along with their relative transformations. In contrast,  
23 modern object recognition systems based on deep learning require thousands of labeled examples  
24 and typically discard information about transformations (via pooling) in order to achieve invari-  
25 ance. Information about transformations is critical for tasks such as movement planning and spatial  
26 reasoning.

27 Current unsupervised models produce representations that either lack interpretability or hierarchical  
28 depth. Variational autoencoders and generative adversarial networks (GANs) typically produce non-  
29 interpretable features that do not match the object/parts hierarchy inherent in natural visual scenes.  
30 Because they do not explicitly model transformations, these models have difficulty generalizing to  
31 the vast range of viewing conditions that objects can appear in. Sparse coding and its variants can  
32 learn interpretable features from unlabeled images: these features resemble the localized oriented  
33 (Gabor) receptive fields found in the primary visual cortex. However, these models again do not  
34 model transformations and have been difficult to generalize to deeper hierarchies due to the combi-  
35 natorial explosion of possible features.

36 We propose a new model for unsupervised learning motivated by the idea that the combinatorial  
 37 explosion problem can be mitigated by a neural architecture that processes the identity (“what”) and  
 38 the pose (“where”) of objects and their parts separately. Such an architecture acknowledges the ven-  
 39 tral/dorsal processing dichotomy in the visual cortex: the first is mostly responsible for processing  
 40 content and identity of objects while the latter is responsible for processing motion and transforma-  
 41 tions.

42 We introduce a new bilinear sparse coding model that builds on previous bilinear generative models  
 43 by (1) allowing each feature to have its own transformation and (2) accommodating large images  
 44 via transposed convolutions. Furthermore, emulating the slower response times of the parvo path-  
 45 way compared to the magno pathway, we assume that at short time scales, object identities at each  
 46 location will remain the same, modeling any fast changes as changes in object transformation val-  
 47 ues. We demonstrate our model by using short natural video sequences to learn features and their  
 48 transformations. The resulting collection of “steerable” filters can be viewed as dynamic features  
 49 resembling the spatiotemporal receptive fields reported in the primary visual cortex. Our model is  
 50 also one of the first to apply ideas from sparse coding to solve the problem of unsupervised learning  
 51 of “primary capsules”<sup>1</sup> previously proposed by Hinton and colleagues for supervised learning [4].

## 52 2 Model

### 53 2.1 Independent Bilinear Sparse Coding

54 In bilinear sparse coding [3, 1], an image patch is modeled as a combination of features  $B_{ij}$  with  
 55 two sets of coefficients  $r_i$  (object coefficients) and  $x_j$  (transformation coefficients) that interact  
 56 multiplicatively:

$$I \simeq \sum_i \sum_j r_i x_j B_{ij} \quad (1)$$

57 Let  $\sum_j x_j B_{ij} = B_i(\mathbf{x})$  where  $\mathbf{x}$  represents the transformation vector consisting of  $x_j$ ’s. Then  
 58  $I \simeq \sum_i r_i B_i(\mathbf{x})$ , which is the standard linear generative model used in sparse coding, PCA, ICA  
 59 etc. The  $r_i$  coefficients correspond to the degree to which each feature exists in the input. The  
 60  $x_j$  coefficients linearly combine a set of similar features to produce a dynamic “steerable” feature  
 61  $B_i(\mathbf{x})$ . The goal is for these dynamic features to capture an equivariance class centered around  
 62 an underlying feature  $B_i$ . As a result, the  $r$  coefficients remain invariant regardless of the specific  
 63 instantiation of the features, the variation being accounted for by  $\mathbf{x}$ . To learn sparse part-like features  
 64 of objects, sparsity is enforced on either  $r$  or both  $r$  and  $x$  via some appropriate sparsity penalty.

65 Typically bilinear sparse coding models are trained using pairs of video frames  $I_{t+1}$  and  $I_t$ , with  $r$   
 66 fixed and  $x$  inferred separately to account for the difference between frames:

$$\Delta I = I_{t+1} - I_t \simeq \sum_i r_i \sum_j (x_{t+1,j} - x_{t,j}) B_{ij} = \sum_i r_i \sum_j \Delta x_{t,j} B_{ij} \quad (2)$$

67 There is a strong connection to the Lie group approach to vision [2] where two consecutive frames  
 68 are modelled as  $I_{t+1} = T(\Delta x)I_t$  where  $T$  is a transformation operator. The first-order Taylor  
 69 series approximation of the Lie model [7, 6] is given by:  $I_{t+1} = I_t + \sum_j \Delta x_{t,j} \nabla x_j I_t$  which  
 70 means that  $\Delta I = \sum_j \Delta x_{t,j} \nabla x_j I_t$ . Suppose  $I_t \simeq \sum_i r_i U_i$  where  $U_i \in \mathbb{R}^{d \times 1}$  form an un-  
 71 derlying feature set. Replacing  $\nabla x_j$  with the transformation matrix  $G_j \in \mathbb{R}^{d \times d}$ , we obtain:  
 72  $\Delta I \simeq \sum_j \Delta x_{t,j} G_j \sum_i r_i U_i = \sum_i r_i \sum_j \Delta x_{t,j} G_j U_i$ . Comparing with Equation 2 above, we see  
 73 that  $B_{ij} = G_j U_i$ .

74 We build on this model by allowing features to have independent pose parameters  $x_{ij}$  so that features  
 75 can transform independently from frame to frame. We also go beyond image patches to modeling  
 76 large images by using transposed convolutions ( $*^T$ ), resulting in a new bilinear model for images:

$$I \simeq \sum_i r_i \sum_j x_{ij} *^T (G_j U_i) = \sum_i r_i *^T B_i(\mathbf{x}_i) \quad (3)$$

77 To distinguish our model from past models, we refer to traditional bilinear sparse coding as BSC  
 78 and our independent bilinear sparse coding model as IBSC.

<sup>1</sup>Primary capsules are capsules in the first layer of processing that convert the image into a collection of activations and poses.

## 79 2.2 Inference

80 The reconstruction-based loss function for consecutive frames of a video is given by:

$$L(r, x_t) = \sum_t \|I_t - \sum_i \sum_j (r_i x_{ijt}) *^T P_{\ell_2, 1.0}(G_j U_i)\|_2^2 + \gamma |r|_1 + \lambda_G \sum_j \|G_j\|_2^2 + \lambda_U \sum_i \|U_i\|_2^2 \quad (4)$$

81 with  $r, x \geq 0$ . The first term is the mean-squared reconstruction error. The other terms include a  
82 sparsity penalty on  $r$  and weight decay for  $G$  and  $U$ . To stabilize learning we project each  $B_{ij} =$   
83  $G_j U_i$  to unit  $\ell_2$  norm ( $P_{\ell_2, 1.0}$ ).

84 Inference for BSC is typically performed by initializing  $x$  to some canonical vector and then alter-  
85 natively optimizing  $r$  and  $x$  [3]. One of the issues with this approach is that the canonical vector  
86 might be a poor approximation to the true underlying pose parameters, especially in the case of  
87 independent features as in our model. We convolve each feature  $B_{ij}$  with the image to produce a  
88 feature map  $\alpha_{ijt} = B_{ij} * I_t$ . We then project onto some appropriately chosen norm ball to compute  
89  $x_{ijt} = P_{\ell, \rho}(\alpha_{ijt})$ .<sup>2</sup> Inference proceeds by alternatively optimizing  $r$  and  $x$  until convergence. To  
90 optimize  $r$ , we use iterative thresholding, while  $x$  is optimized by projected gradient descent. Both  
91 sets of coefficients are forced to be non-negative, using a rectifier for  $r$  and projecting on the positive  
92 part of the norm ball for  $x$ .

## 93 3 Experiments

94 For our experiments, we used  $1920 \times 1080$  resolution YouTube videos converted to gray scale  
95 and scaled down to  $236 \times 176$  pixels per frame. The frames were normalized using subtractive  
96 normalization<sup>3</sup>. We extracted sequences of 5 consecutive frames, with  $r$  assumed to be constant  
97 for each sequence during training. We excluded sequences in the largest 5% of Euclidean norm  
98 difference between frames to exclude sudden camera changes or changes between scenes. We used  
99 a stride of half the size of the kernel for transposed convolutions.

100 Our model learns localized oriented Gabor-like features similar to those seen in sparse coding. Fig-  
101 ure 1 shows a subset of the learned  $12 \times 12$  pixel features: each column shows  $B_i$ : corresponding  
102 to different transformed versions of the same underlying feature. Note that not only translations but  
103 other transformations are learned as well, e.g., rotations and warping. The learned bilinear features  
104 allow accurate reconstruction, as seen for an example input in Figures 2(a) and 2(b). All feature sets  
105 were  $2 \times$  overcomplete.

106 To test whether each  $B_i(\mathbf{x}_i)$  corresponds to a “steerable” filter, we visualize in Figures 3(a-e) a  
107 subset of the different instantiations (with different  $\mathbf{x}_i$ ’s) of each feature across different inputs  
108 and image locations from our natural videos. Note that the model captures a wide range of such  
109 instantiations. To determine whether each  $B_i(\mathbf{x}_i)$  captures the progression of a single underlying  
110 feature across frames, we visualized the evolution of features across sequences of frames. As seen  
111 in Figure 4, the learned features evolve across frames in a manner similar to spatiotemporal filters  
112 in the visual cortex, e.g., direction-selective Gabor filters moving in a particular direction.

## 113 4 Conclusion & Future Work

114 We extend the bilinear sparse coding model to handle large images and independent feature transfor-  
115 mations. Our model learns to group similar features together, leveraging the smoothness of natural  
116 videos. Perhaps the most interesting direction for future work is that of extending this approach  
117 hierarchically.

---

<sup>2</sup>This allows us to use the features themselves to derive a suitable pose vector. For the projection of  $x$  we use the simplex  $S_\rho : \sum_j |x_j| \leq \rho, |x_j| \geq 0$ ; the radius  $\rho$  determines how sparse the coefficients will be.

<sup>3</sup>A Gaussian kernel is used to estimate the mean intensity around each pixel, which is then subtracted from the pixel value.

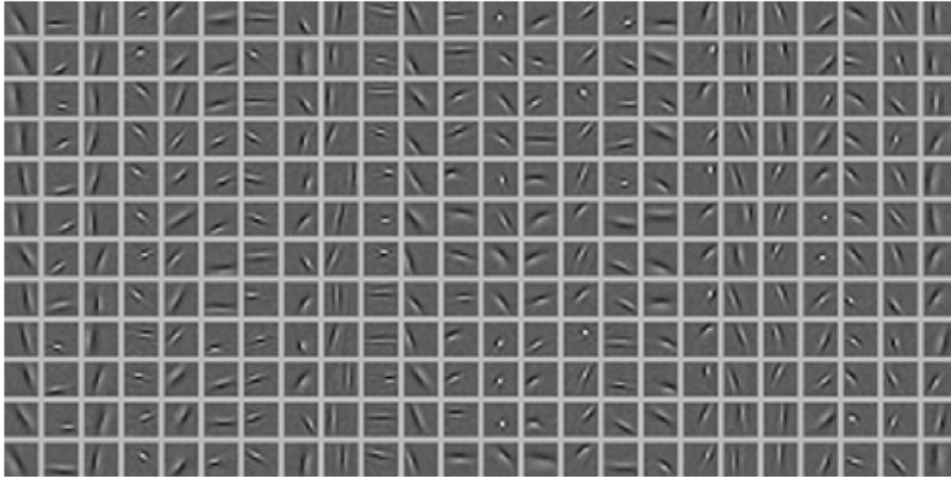


Figure 1: **Independent Bilinear Sparse Coding for Natural Videos.**  $12 \times 12$  pixel features  $B_{ij}$ : each column shows a feature  $i$  for different  $j$ 's.

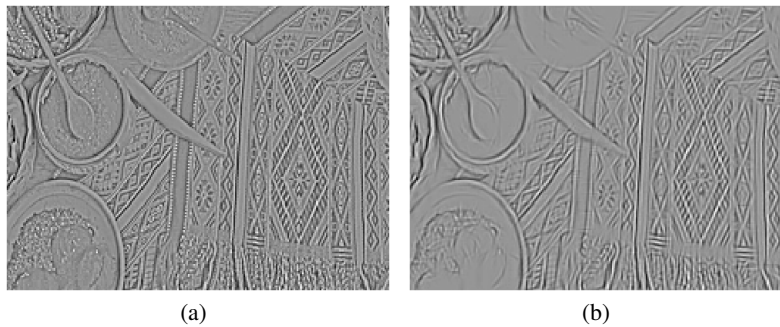


Figure 2: **Example Frame Reconstruction.** (a) Original image and (b) its reconstruction using the learned bilinear features.

## 118 5 Acknowledgments

119 This work was supported by NSF grant no. EEC-1028725, CRCNS/NIMH grant  
 120 no. 1R01MH112166-01, and a grant from the Templeton World Charity Foundation  
 121 (TWCF).

122

## 123 References

124 1

- 125 [1] Jack Culpepper David K Warland Bruno A. Olshausen, Charles Cadieu.  
 126 Bilinear models of natural images, 2007.
- 127 [2] Peter C. Dodwell. The lie transformation group model of visual per-  
 128 ception. *Perception & Psychophysics*, 34(1):1–16, 1983.
- 129 [3] David B. Grimes and Rajesh P. N. Rao. Bilinear sparse coding for  
 130 invariant vision. *Neural Computation*, 17(1):47–73, 2005.

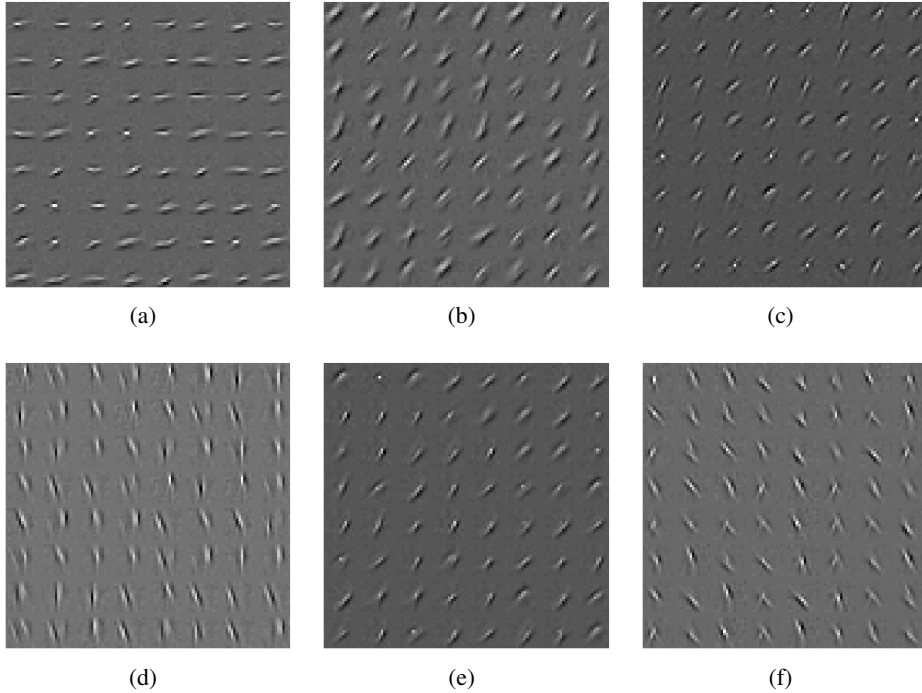


Figure 3: **Feature Equivariance Classes.** (a,b,c,d,e,f) Feature equivariance classes (each plot shows different transformations of the same underlying feature).

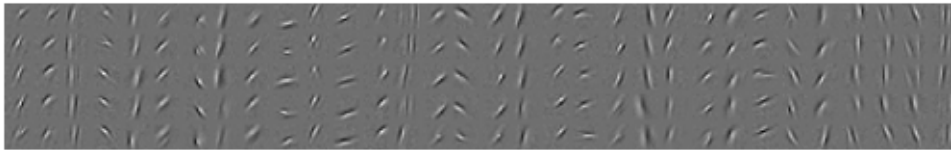


Figure 4: **Learned Feature Dynamics.** Feature dynamics between frames. Each column corresponds to a distinct instance of a dynamic spatiotemporal filter, resembling cortical spatiotemporal filters.

- 131 [4] Geoffrey E Hinton, Sara Sabour, and Nicholas Frosst. Matrix capsules  
 132 with EM routing. In *International Conference on Learning Representations*, 2018.  
 133
- 134 [5] Victor Lamme and Pieter Roelfsema. The distinct modes of vision of-  
 135 ferred by feedforward and recurrent processing. *Trends in neurosciences*,  
 136 23:571–9, 12 2000.
- 137 [6] Xu Miao and Rajesh P. N. Rao. Learning the Lie Groups of Visual  
 138 Invariance. *Neural Comput.*, 19(10):2665–2693, October 2007.
- 139 [7] Rajesh Rao and D Ballard. Development of localized oriented recep-  
 140 tive fields by learning a translation-invariant code for natural images\*.  
 141 *Network (Bristol, England)*, 9:219–34, 06 1998.
- 142 [8] L. G. Ungerleider and L. Pessoa. What and where pathways. *Scholar-*  
 143 *pedia*, 3(11):5342, 2008. revision #91940.