

# Learning to Relate Images

Roland Memisevic

**Abstract**—A fundamental operation in many vision tasks, including motion understanding, stereopsis, visual odometry, or invariant recognition, is establishing correspondences between images or between images and data from other modalities. Recently, there has been increasing interest in learning to infer correspondences from data using relational, spatiotemporal, and bilinear variants of deep learning methods. These methods use multiplicative interactions between pixels or between features to represent correlation patterns across multiple images. In this paper, we review the recent work on relational feature learning, and we provide an analysis of the role that multiplicative interactions play in learning to encode relations. We also discuss how square-pooling and complex cell models can be viewed as a way to represent multiplicative interactions and thereby as a way to encode relations.

**Index Terms**—Learning image relations, spatiotemporal features, mapping units, energy models, complex cells

## 1 INTRODUCTION

CORRESPONDENCE is arguably the most ubiquitous computational primitive in vision: *Motion estimation*, *action recognition*, and *tracking* amount to establishing correspondences between frames, *geometric inference* and *stereo vision* between multiple views of a scene, *invariant recognition* between images and invariant templates in memory, *visual odometry* between images and estimates of motion, and so on. In these and many other tasks, the relationship *between* images, not the content of a single image, carries the relevant information. Representing structures within a single image, such as contours, can be also considered as an instance of a correspondence problem, namely, between areas, or pixels, within an image. The fact that correspondence is such a common operation across vision suggests that the task of *representing relations* may have to be kept in mind when trying to build autonomous vision systems and when trying to understand biological vision.

A lot of progress has been made recently in building models that learn to solve tasks like object recognition from independent, static images. One of the reasons for the recent progress is the use of *local features*, which help deal with occlusions and small invariances. A central finding is that the right choice of features, not the choice of high-level classifier or computational pipeline, is what typically makes a system work well. Interestingly, some of the best performing recognition models are highly biologically consistent in that they are based on features that are learned unsupervised from data. Besides being biologically plausible, feature learning comes with various benefits, such as helping overcome tedious engineering, helping adapt to new domains, and allowing for some degree of *end-to-end*

*learning* in place of constructing, and then combining, a large number of modules to solve a task. The fact that tasks like object recognition can be solved using biologically consistent, learning-based methods raises the question whether understanding *relations* can be amenable to learning in the same way. If so, this may open up the road to learning-based and/or biologically consistent approaches to a much larger variety of problems than static object recognition, and perhaps also beyond vision.

In this paper, we review a variety of recent methods that address correspondence tasks by learning local features. We discuss how these methods are fundamentally based on *multiplicative interactions* between pixels or between filter responses. The idea of using multiplicative interactions in vision was introduced about 30 years ago under the terms “mapping units” [1] and “dynamic mappings” [2]. An illustration of mapping units is shown in Fig. 1: The three variables shown in the figure interact multiplicatively, and as a result, each variable (say,  $z_k$ ) can be thought of as dynamically *modulating* the connections between other variables in the model ( $x_i$  and  $y_j$ ). Likewise, the value of any variable (e.g.,  $y_j$ ) can be thought of as depending on the product of the other variables ( $x_i, z_k$ ) [1]. This is in contrast to common feature learning models like ICA, Restricted Boltzmann Machines (RBMs), autoencoder networks, and many others, all of which are based on bipartite networks that do not involve any three-way multiplicative interactions. In these models, independent hidden variables interact with independent observable variables such that the value of any variable depends on a weighted *sum* not product of the other variables. Closely related to models of mapping units are energy models (e.g., [3]), which may be thought of as a way to “emulate” multiplicative interactions by computing squares.

We shall show how both mapping units and energy models can be viewed as ways to learn and detect rotations in a set of *shared invariant subspaces* of a set of *commuting matrices*. Our analysis may help understand why action recognition methods seem to profit from using squaring nonlinearities and it predicts that the use of squaring

• The author is with the Department of Computer Science and Operations Research, University of Montreal. E-mail: memisevr@iro.umontreal.ca.

Manuscript received 8 Apr. 2012; revised 7 Oct. 2012; accepted 14 Jan. 2013; published online 6 Mar. 2013.

Recommended for acceptance by S. Bengio, L. Deng, H. Larochelle, H. Lee, and R. Salakhutdinov.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number TPAMSI-2012-04-0256.

Digital Object Identifier no. 10.1109/TPAMI.2013.53.

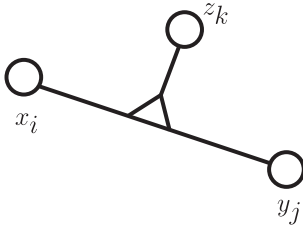


Fig. 1. A *mapping unit* [1]. The triangle symbolizes multiplicative interactions between the three variables  $z_k$ ,  $x_i$ , and  $y_j$ . The value of any one of the three variables is a function of the product of the others.

nonlinearities or multiplicative interactions will be essential in any tasks that involve representing relations.

## 2 MULTIVIEW FEATURE LEARNING

### 2.1 Feature Learning

We briefly review standard feature learning models in this section and we discuss relational feature learning in Section 2.2. We discuss extensions of relational models and how they relate to complex cells and to energy models in Section 3.

Practically all standard feature learning models can be represented by a graphical model like the one shown in Fig. 2 (top). The model is a bipartite network that connects a set of unobserved, latent variables  $z_k$  with a set of observable variables (e.g., pixels)  $y_j$ . The weights  $w_{jk}$ , which connect pixel  $y_j$  with hidden unit  $z_k$ , are learned from a set of training images  $\{\mathbf{y}^\alpha\}_{\alpha=1,\dots,N}$ . The vector of latent variables  $\mathbf{z} = (z_k)_{k=1,\dots,K}$  in Fig. 2 (top) is considered to be unobserved, so one has to infer it separately for each training case, along with the model parameters for training. The graphical model shown in the figure represents how the dependencies between components  $y_i$  and  $z_k$  are parameterized, but it does not define a model or learning algorithm. A wide variety of models and learning algorithms can be parameterized as in the figure, including principal components analysis (PCA), mixture models, k-means clustering, or RBMs [4]. Each of these can in principle be used as a feature learning method (see, e.g., [5] for a recent quantitative comparison of several models in a recognition task).

For the hidden variables to extract useful structure from the images, their capacity needs to be constrained. The simplest form of constraining it is to let the dimensionality  $K$  of the hidden variables be smaller than the dimensionality  $J$  of the images. Learning in this case amounts to performing dimensionality reduction. It has become increasingly obvious recently that it is more useful in most applications to use an *overcomplete* representation, that is,  $K > J$ , and to constrain the capacity of the latent variables instead by forcing the hidden unit activities to be *sparse*. In Fig. 2, and in what follows, we use  $K < J$  to symbolize the fact that  $\mathbf{z}$  is capacity constrained, but it should be kept in mind that capacity can be (and often is) constrained in other ways. The most common operations performed by a trained model are: Inference (or “Analysis”): Given an image  $\mathbf{y}$ , compute  $\mathbf{z}$ ; and Generation (or “Synthesis”): Invent a latent vector  $\mathbf{z}$ , then compute  $\mathbf{y}$ .

A simple way to train this type of model, given training images, is by minimizing the squared reconstruction error

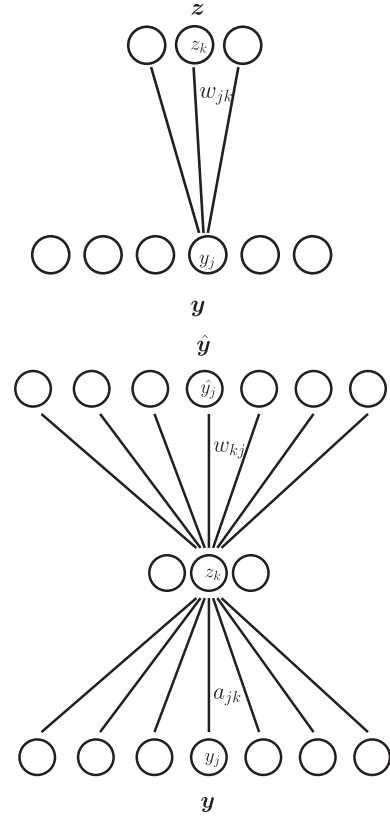


Fig. 2. *Top*: Feature learning graphical model. *Bottom*: Autoencoder network.

combined with a sparsity term for the hidden variables (e.g., [6]):

$$\sum_{\alpha} \sum_j \left( y_j^{\alpha} - \sum_k w_{jk} z_k^{\alpha} \right)^2 + \lambda \sum_k |z_k^{\alpha}|. \quad (1)$$

Optimization is with respect to both  $\mathbf{W} = (w_{jk})_{j=1,\dots,J,k=1,\dots,K}$  and all  $\mathbf{z}^{\alpha}$ . It is common to alternate between optimizing  $\mathbf{W}$  and optimizing all  $\mathbf{z}^{\alpha}$ . After training, inference then amounts to minimizing the same expression w.r.t.  $\mathbf{z}$  for test images (with  $\mathbf{W}$  fixed).

To avoid iterative optimization during inference, one can eliminate  $\mathbf{z}$  from (1) by defining it implicitly as a function of  $\mathbf{y}$ . A common choice of function is  $\mathbf{z} = \text{sigmoid}(\mathbf{A}\mathbf{y})$ , where  $\mathbf{A}$  is a matrix and  $\text{sigmoid}(a) = (1 + \exp(-a))^{-1}$  is a squashing nonlinearity which confines the values of  $\mathbf{z}$  to reside in a fixed interval. This model is the well-known autoencoder (e.g., [7]) and it is depicted in Fig. 2 (bottom). Learning amounts to minimizing reconstruction error with respect to both  $\mathbf{A}$  and  $\mathbf{W}$ , with gradients that are usually computed using back-prop. In practice, it is common to define the autoencoder in a symmetric fashion by setting  $\mathbf{A} = \mathbf{W}^T$  to reduce the number of parameters and for consistency with other feature learning models.

In practice, it is common to encourage sparse hidden activities by adding an appropriate sparsity penalty during training. Alternatively, it has been shown that a similar effect can be achieved by training the autoencoder to denoise corrupted versions of its inputs [7]. To this end, one feeds in noisy inputs during training (e.g., by adding Gaussian noise to the input, or by randomly setting

individual dimensions of the input to zero) and minimizes reconstruction error with respect to the original (not noisy) data. This turns the autoencoder into a “denoising autoencoder,” which shows properties similar to common sparse coding methods, but inference, like in a standard auto-encoder, is a simple feed-forward mapping [7]. In the rest of the paper, we shall use the term autoencoder to refer to denoising autoencoders; in other words, we shall always assume that inputs are corrupted for training.

A technique similar to the autoencoder is the RBM [4], [8]: RBMs define the joint probability distribution:

$$p(\mathbf{y}, \mathbf{z}) = \frac{1}{Z} \exp(-E(\mathbf{y}, \mathbf{z})), \quad (2)$$

$$\text{with } E(\mathbf{y}, \mathbf{z}) = - \sum_{jk} w_{jk} y_j z_k \quad (3)$$

$$\text{and } Z = \sum_{\mathbf{y}, \mathbf{z}} \exp(-E(\mathbf{y}, \mathbf{z})). \quad (4)$$

From the joint, one can derive

$$p(z_k | \mathbf{y}) = \text{sigmoid}\left(\sum_j w_{jk} y_j\right), \quad (5)$$

$$p(y_j | \mathbf{z}) = \text{sigmoid}\left(\sum_k w_{jk} z_k\right). \quad (6)$$

This shows that inference, again, amounts to a linear mapping plus nonlinearity. Learning amounts to maximizing the average log-probability  $\frac{1}{N} \sum_{\alpha} \log p(\mathbf{y}^{\alpha})$  of the training data. Since the derivatives with respect to the parameters are not tractable (due to the normalizing constant  $Z$  in (2)), it is common to use approximate Gibbs sampling to approximate them. This leads to a Hebbian-like learning rule known as contrastive divergence training [8]. As with autoencoders, it is common to enforce sparsity of the hidden units during training (e.g., [9]).

Another common feature learning method is independent components analysis (ICA) (e.g., [10]). One way to train an ICA-model that is complete (that is, where the dimensionality of  $\mathbf{z}$  is the same as that of  $\mathbf{y}$ ) is by encouraging latent responses to be sparse while preventing weights from becoming degenerate [10]:

$$\min_W \|W^T \mathbf{y}\|_1 \quad (7)$$

$$\text{s.t. } W^T W = I. \quad (8)$$

The constraint can be inconvenient in practice, where it is commonly enforced by repeated orthogonalization using an eigendecomposition.

For most feature learning models, inference and generation are variations of the two linear functions:

$$z_k = \sum_j w_{jk} y_j, \quad (9)$$

$$y_j = \sum_k w_{jk} z_k. \quad (10)$$

The set of model parameters  $W_k$  for any  $k$  is typically referred to as “features” or “filters” (although a more appropriate term would be “basis functions”; we shall use these interchangeably). Practically all methods yield Gabor-like features when trained on natural images. An advantage of nonlinear models, such as RBMs and autoencoders, is that stacking them makes it possible to learn feature hierarchies (deep learning) [11].

In practice, it is common to add bias terms such that inference and generation ((9) and (10)) are affine, not linear, functions; for example,  $y_j = \sum_k w_{jk} z_k + b_j$  for some parameter  $b_j$ . We shall refrain from adding bias terms to avoid clutter, noting that, alternatively, one may think of  $\mathbf{y}$  and  $\mathbf{z}$  as being in “homogeneous” coordinates, containing an extra, constant 1-dimension.

Feature learning is typically performed on small image patches of size between around  $5 \times 5$  and  $50 \times 50$  pixels. One reason for this is that training and inference can be computationally demanding. More importantly, local features make it possible to deal with images of different size, and to deal with occlusions and local object variations. Given a trained model, two common ways to perform invariant recognition on test images are:

“*Bag-of-features*”: Crop patches around interest points (such as SIFT or Harris corners) compute latent representation  $\mathbf{z}$  for each patch, collapse (add up) all representations to obtain a single vector  $\mathbf{z}^{\text{Image}}$ , and classify  $\mathbf{z}^{\text{Image}}$  using a standard classifier. There are several variations of this scheme, including using an extra clustering step before collapsing features, or using a histogram similarity in place of euclidean distance for the collapsed representation.

“*Convolutional*”: Crop patches from the image along a regular grid; compute  $\mathbf{z}$  for each patch; concatenate all descriptors into a very large vector  $\mathbf{z}^{\text{Image}}$ ; classify  $\mathbf{z}^{\text{Image}}$  using a standard classifier. One can also use combinations of the two schemes (see, e.g., [5]).

Local features yield highly competitive performance in object recognition tasks [5]. In the next section, we discuss recent approaches to extending feature learning to encode relations between, as opposed to content within, images.

## 2.2 Learning Relations

We now consider the task of learning relations between two images  $\mathbf{x}$  and  $\mathbf{y}$  as illustrated<sup>1</sup> in Fig. 4, and we discuss the role of multiplicative interactions when learning relations.

### 2.2.1 The Need for Multiplicative Interactions

A naive approach to modeling relations between two images would be to perform standard feature learning on the concatenation of the images. Obviously, a hidden unit in such a model would receive as input the sum of two projections, one from each image. To detect a particular transformation, the two receptive fields would need to be defined such that one receptive field is the other modified by the transformation that the hidden unit is supposed to detect. The net input that the hidden unit receives would then tend to be high for image pairs showing the transformation. Unfortunately, however, the net input would be equally dependent on the images themselves,

1. Face images are taken from the database described in [12].

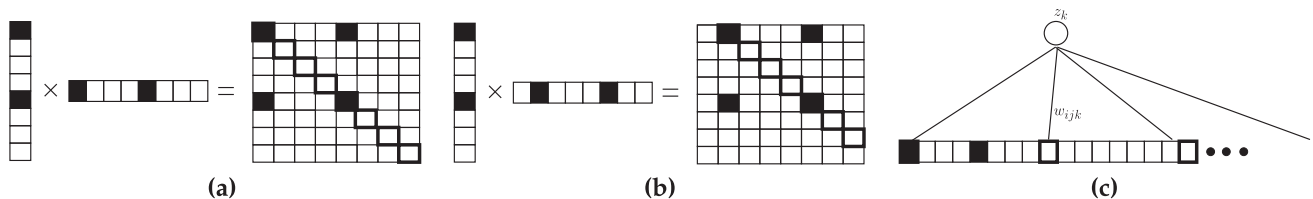


Fig. 3. (a) The diagonal of  $L := \mathbf{xy}^T$  contains evidence for the identity transformation. (b) The secondary diagonals contain evidence for shifts. (c) A hidden unit that pools over one of the diagonals can detect transformations. Such a hidden unit would need to compute a *sum over products* of pixels.

not just the transformation: If both images change but not the transformation between them, the hidden activity would also change. Another way to see this is by noting that hidden variables act like logical “OR”-gates, which can only accumulate the information from their receptive fields [13].

It is straightforward to build a content-independent detector, however, if we allow for *multiplicative interactions* between the variables. In particular, consider the outer product  $L := \mathbf{xy}^T$  between two one-dimensional, binary images, as shown in Fig. 3. Every component  $L_{ij}$  of this matrix constitutes evidence for exactly one type of transformation (translation, in the example). The components  $L_{ij}$  act like AND-gates that can detect coincidences. Since a component  $L_{ij}$  is equal to 1 only when both corresponding pixels are equal to 1, a hidden unit that pools over multiple components (Fig. 3c) is much less likely to receive spurious activity that depends on the image content rather than on the transformation. Note that pooling over the components of  $L$  amounts to computing the *correlation* of the output image with a transformed version of the input image. The same would be true for real-valued images.

When a variable,  $z$ , depends linearly on the product of two other variables, that is,  $z = a(xy)$ , then both  $x$  and  $y$  can be thought of as *gating* the connection between the other variable and  $z$  because  $z = (ax)y = (ay)x$  (see Fig. 1). In other words, commutativity allows us to think of  $x$  as modulating the parameter  $a$  that connects  $z$  and  $y$  (and vice versa).

Based on this observation, a variety of feature learning models that encode transformations have been suggested (see, e.g., [14], [15], [16]). *The idea is to let the pixels in one*

*image gate the parameters of a feature learning model applied to another image. This is equivalent to letting hidden variables encode the product of pixel  $x_i$  in one image and pixel  $y_j$  in the other image.* If every pixel in the first image is allowed to independently gate every parameter in the model of the other image, then the total number of parameters is (number of hidden variables)  $\times$  (number of input pixels)  $\times$  (number of output pixels). It is common to think of the parameters as populating a 3-way-tensor  $W$  with components  $w_{ijk}$ .

Fig. 5 shows two illustrations of this type of model (adapted from [16]). The left subfigure shows a feature learning model whose parameters are modulated by another image. Each input pixel of that other image can be thought of as blending in a slice  $w_{i..}$  of the parameter tensor. This turns the model into a kind of *predictive* or *conditional* feature learning model [14], [16].

Fig. 5 (right) shows an alternative visualization of the same type of model, where each hidden variable can blend in a slice  $w_{..k}$  of the parameter tensor. Each slice in turn is a matrix connecting an “input” pixel from image  $x$  to an “output” pixel from image  $y$ . We can think of this matrix as performing linear regression in the space of stacked gray-value intensities. A linear transformation in “pixel-space” is commonly known as a “warp.” Although linear in pixel-space, a warp can be a highly nonlinear transformation in image coordinates. Note also that a warp has a very large number of parameters in comparison to an affine image transformation. The model in Fig. 5 (right) may be thought of as defining a *mixture of warps*.

In both cases, hidden variables take on the roles of dynamic mapping units [1], [2], which encode the relationship not the content of the images. Each unit in the model

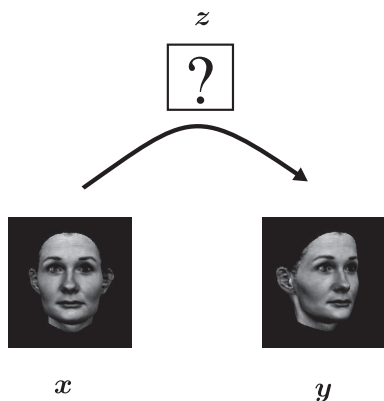


Fig. 4. Learning to encode relations: We consider the task of learning latent variables  $z$  that encode the relationship between images  $x$  and  $y$ , independently of their content.

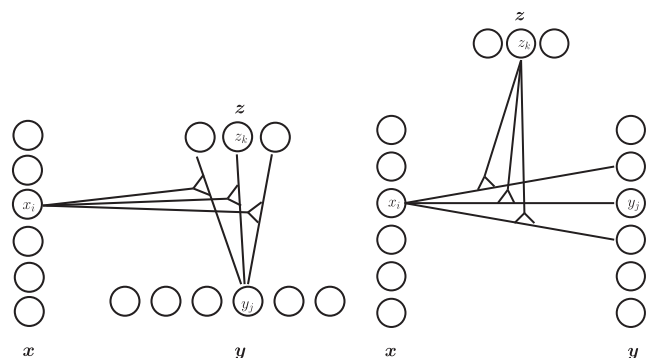


Fig. 5. Relating images using multiplicative interactions. Two views of the same model. Left: Input-modulated feature learning. Right: Mixture of warps.

can gate connections between other variables in the model. We shall refer to this type of model as a “gated feature learning,” “cross correlation,” or “multiview feature learning” model in the following.

Like in a standard feature learning model, one needs to include biases in practice. The set of model parameters thus consists of the three-way parameters  $w_{ijk}$ , as well as of single-node parameters  $w_i$ ,  $w_j$ , and  $w_k$ . One could also include “higher-order-biases” [16] like  $w_{ik}$ , which connect two groups of variables, but it is not common to do so. Like before, we shall drop all bias terms in what follows to avoid clutter. Both simple biases and higher-order biases can be implemented by introducing appropriate constant-1 dimensions to the images or the hidden variables.

### 2.3 Inference

The graphical model for gated feature learning is tripartite. That of a standard feature learning model is bipartite. As a result, inference can be performed in almost the same way as in a standard feature learning model whenever two out of three groups of variables have been observed, as we show now.

Consider the task of inferring  $\mathbf{z}$ , given  $\mathbf{x}$  and  $\mathbf{y}$ . Recall that for a standard feature learning model, we have  $z_k = \sum_j w_{jk} y_j$  (up to component-wise nonlinearities like the sigmoid). Formally, we may think of the gated feature learning model as turning the weights into a *linear function* of  $\mathbf{x}$ :

$$w_{jk}(\mathbf{x}) = \sum_i w_{ijk} x_i \quad (11)$$

so that the inference equation becomes

$$z_k = \sum_j w_{jk}(\mathbf{x}) y_j = \sum_j \left( \sum_i w_{ijk} x_i \right) y_j = \sum_{ij} w_{ijk} x_i y_j, \quad (12)$$

which amounts to computing, for each hidden  $z_k$ , a quadratic form in  $\mathbf{x}$  and  $\mathbf{y}$ . The quadratic form is defined by the weight tensor  $w_{\cdot k}$ . Thus, we can think of inference in a gated feature learning model either as computing a (quadratic) function of two images or as standard inference for a single image,  $\mathbf{y}$ , where parameters are linearly dependent on another image,  $\mathbf{x}$ . We shall refer to the latter as “*predictive coding*,” because we can think of the model as predicting  $\mathbf{y}$  from  $\mathbf{x}$  via  $\mathbf{z}$ . The fact that inference may be interpreted in multiple ways is common in models with bilinear dependencies [17].

Note that (12) is symmetric in  $\mathbf{x}$  and  $\mathbf{y}$ . We could therefore, in principle, switch their roles in (11) and (12) and define the linear parameters (11) as a function of  $\mathbf{y}$  rather than  $\mathbf{x}$ . During training, however, it has been common to drop this symmetry and to declare one of the two images as the gating “input” and the other as the “output,” as we shall show.

The *meaning* of the hidden variables differs from that in standard feature learning models despite the similarity of inference: In standard feature learning,  $\mathbf{z}$  constitutes a representation of the input image; in gated feature learning  $\mathbf{z}$  represents the *transformation* that takes  $\mathbf{x}$  to  $\mathbf{y}$ .

Inferring  $\mathbf{y}$ , given  $\mathbf{x}$  and  $\mathbf{z}$ , yields the analogous expression:

$$y_j = \sum_k w_{jk}(\mathbf{x}) z_k = \sum_k \left( \sum_i w_{ijk} x_i \right) z_k = \sum_{ik} w_{ijk} x_i z_k, \quad (13)$$

which again amounts to computing a quadratic form. The meaning of  $\mathbf{y}$  is now “ $\mathbf{x}$  transformed according to the known transformation  $\mathbf{z}$ .” Likewise, we could compute  $\mathbf{x}$  given  $\mathbf{z}$  and  $\mathbf{y}$  using an analogous equation, but, again, this would not be common if the model was trained asymmetrically (see Section 2.4.1).

It is important to note that for any given transformation  $\mathbf{z}$ ,  $\mathbf{y}$  is a linear function of  $\mathbf{x}$ , so it can be written

$$\mathbf{y} = L(\mathbf{z})\mathbf{x}. \quad (14)$$

From (13) it follows that the entries of matrix  $L(\mathbf{z})$  are given by

$$L_{ij}(\mathbf{z}) = \sum_k w_{ijk} z_k. \quad (15)$$

When  $\mathbf{x}$  and  $\mathbf{y}$  are images, the linear function is a *warp*. So the hidden variables can be thought of as composing a warp from “basis warps”  $w_{\cdot k}$ , in exactly the same way that feature learning models can be thought of as composing an image from basis images (features)  $w_{\cdot k}$  (see (10)). Whereas inferring a component  $z_k$  in a standard feature learning model amounts to computing the inner product between image  $\mathbf{y}$  and feature  $w_{\cdot k}$ , for gated feature learning, it amounts to computing the inner product between  $\mathbf{x}\mathbf{y}^T$  and the basis warp  $w_{\cdot k}$  (see (9) and (12)).

Besides computing hidden unit activities or generating fantasy data, a third common operation in many feature learning models is to compute a confidence value for new input data which quantifies how well that data can be represented by the model. For this number to be useful, it has to be “calibrated,” which is typically achieved by using a probabilistic model. In contrast to standard feature learning, training a probabilistic gated feature learning model can be slightly more complicated because of the dependencies between  $\mathbf{x}$  and  $\mathbf{y}$  conditioned on  $\mathbf{z}$ . We shall discuss this issue in more detail in Section 2.4.2.

## 2.4 Learning

### 2.4.1 Predictive Training

The training data for a multiview feature learning model consists of pairs of data points  $(\mathbf{x}^\alpha, \mathbf{y}^\alpha)$ . Training is similar to standard feature learning, but there are some important differences. In particular, recall that the gated model may be viewed as a feature learning model whose input is the vectorized outer product  $\mathbf{x}\mathbf{y}^T$ . Standard learning criteria such as squared reconstruction error are obviously not appropriate for the outer product.

It is, however, possible to use squared error for training, albeit not on the products. To this end, consider the perspective from predictive coding, where the parameters of a feature learning model on  $\mathbf{y}$  are modulated by an input image  $\mathbf{x}$ . This suggests deriving the learning criterion from the task of predicting  $\mathbf{y}$  from  $\mathbf{x}$  [15], [14], [16]. We shall first discuss this learning approach from the “inference-free” perspective (1), and we shall subsequently discuss how using the linear inference equations (see Section 2.1) can simplify learning in direct analogy to standard feature learning.

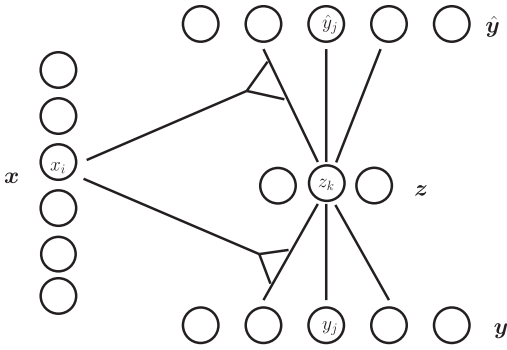


Fig. 6. A gated autoencoder is an autoencoder that learns to represent an image,  $\mathbf{y}$ , using parameters that are modulated by another image,  $\mathbf{x}$ . This makes it possible to learn relationships between  $\mathbf{x}$  and  $\mathbf{y}$  with gradient-based learning and back-prop.

The modulation of parameters in (11) is *case dependent*, that is, each input example leads to a different model for  $\mathbf{y}$ . Learning can therefore be viewed as feature learning with case-dependent weights. In analogy to (1), we can write the reconstruction error that data-case  $(\mathbf{x}^\alpha, \mathbf{y}^\alpha)$  contributes as

$$\sum_j \left( y_j^\alpha - \sum_{ik} w_{ijk} x_i^\alpha z_k^\alpha \right)^2. \quad (16)$$

By using (11) and by adding a sparsity penalty for  $\mathbf{z}$ , we can write the cost over the whole dataset also as

$$\sum_\alpha \sum_j \left( y_j^\alpha - \sum_k w_{jk}(\mathbf{x}^\alpha) z_k^\alpha \right)^2 + \lambda \sum_k |z_k^\alpha|, \quad (17)$$

which highlights the similarity with (1). Differentiating (16) with respect to  $w_{ijk}$  is the same as in a standard feature learning model. In particular, the model is still linear w.r.t. the parameters. Predictive learning is therefore possible with gradient-based optimization similar to standard feature learning (see Section 2.1).

However, in analogy to standard feature learning, it can be useful to use a feed-forward inference function to simplify inference and learning. A simple way to eliminate the hidden variables is by using an encoder network as defined in (12) (possibly followed by a sigmoid nonlinearity) and by using a set of decoder parameters  $a_{ijk}$  to compute reconstructions as defined in (13). Like in standard feature learning, one may tie decoder and encoder parameters by setting  $a_{ijk} = w_{ijk}$ .

This type of model is known as *gated autoencoder* [18], [19], and it is depicted in Fig. 6. Learning is similar to learning a standard autoencoder. This becomes obvious by plugging (12) into (16) and by noting that  $\mathbf{x}$  is *fixed* in each training example. It is also possible to add multiple layers and to apply nonlinearities to the hidden layers (in which case, of course, the derivatives will no longer be linear w.r.t. the parameters). Conditioning on  $\mathbf{x}$  always ensures that the model is a directed acyclic graph, so one can use standard back-prop to compute derivatives.

As a second example of a gated feature learning model, we obtain the *gated Boltzmann machine* (GBM) by changing the energy function into the three-way energy [16]:

$$E(\mathbf{x}, \mathbf{y}, \mathbf{z}) = - \sum_{ijk} w_{ijk} x_i y_j z_k. \quad (18)$$

Exponentiating and normalizing yields the conditional probability:

$$p(\mathbf{y}, \mathbf{z} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(-E(\mathbf{x}, \mathbf{y}, \mathbf{z})), \quad (19)$$

$$Z(\mathbf{x}) = \sum_{\mathbf{y}, \mathbf{z}} \exp(-E(\mathbf{x}, \mathbf{y}, \mathbf{z})). \quad (20)$$

Note that the normalization is over  $\mathbf{y}$  and  $\mathbf{z}$  only, which ensures that we obtain a conditional model of the output image  $\mathbf{y}$ . While it is possible to define a joint model over both images, this makes training more difficult (see Section 2.4.2).

Like in a standard RBM, maximum likelihood and contrastive divergence training involve sampling  $\mathbf{z}$  and  $\mathbf{y}$ . In the GBM, samples are drawn from the *conditional* distributions  $p(\mathbf{y} | \mathbf{z}, \mathbf{x})$  and  $p(\mathbf{z} | \mathbf{y}, \mathbf{x})$ . Training the GBM is like training a standard RBM if we again utilize the fact that every input training pair defines a standard RBM whose parameters are defined as a (case-dependent) function of the input. For both the GBM and the gated autoencoder, one can include weight-decay terms and penalty terms to encourage hidden variable responses to be sparse.

#### 2.4.2 Relational Training

Modeling the *joint* distribution over two images, rather than the conditional distribution of one image given the other, can make it possible to perform image matching by allowing us to quantify how compatible any two images are under to the trained model [20].

Formally, modeling the joint amounts simply to changing the normalization constant of the GBM to  $Z = \sum_{\mathbf{x}, \mathbf{y}, \mathbf{z}} \exp(-E(\mathbf{x}, \mathbf{y}, \mathbf{z}))$  (see (20)). Learning is more complicated as a result, however, because the view of input-dependent parameters no longer holds. Susskind et al. [20] show that it is possible to use a “three-way” version of contrastive divergence learning, where each iteration involves sampling  $\mathbf{x}$  from  $p(\mathbf{x} | \mathbf{z}, \mathbf{y})$  and sampling  $\mathbf{y}$  from  $p(\mathbf{y} | \mathbf{z}, \mathbf{x})$ .

Another potential advantage of learning a symmetric model is that it allows us to model higher-order features for a single image; in other words, features that encode products of pixel intensities within the image. See, for example, [21], who train second-order features by using a joint version of a GBM with  $\mathbf{x} = \mathbf{y}$ . In contrast to [20], they use hybrid Monte Carlo for learning.

Alternatively, a gated autoencoder can be turned into a symmetric model by defining the cost as the sum of two symmetric reconstruction costs:

$$\sum_j \left( y_j^\alpha - \sum_{ik} w_{ijk} x_i^\alpha z_k^\alpha \right)^2 + \sum_i \left( x_i^\alpha - \sum_{jk} w_{ijk} y_j^\alpha z_k^\alpha \right)^2. \quad (21)$$

This makes it possible to learn higher-order features in a nonprobabilistic way using gradient descent [19]. For further approaches to learning higher-order within-image features see [22] and [23].



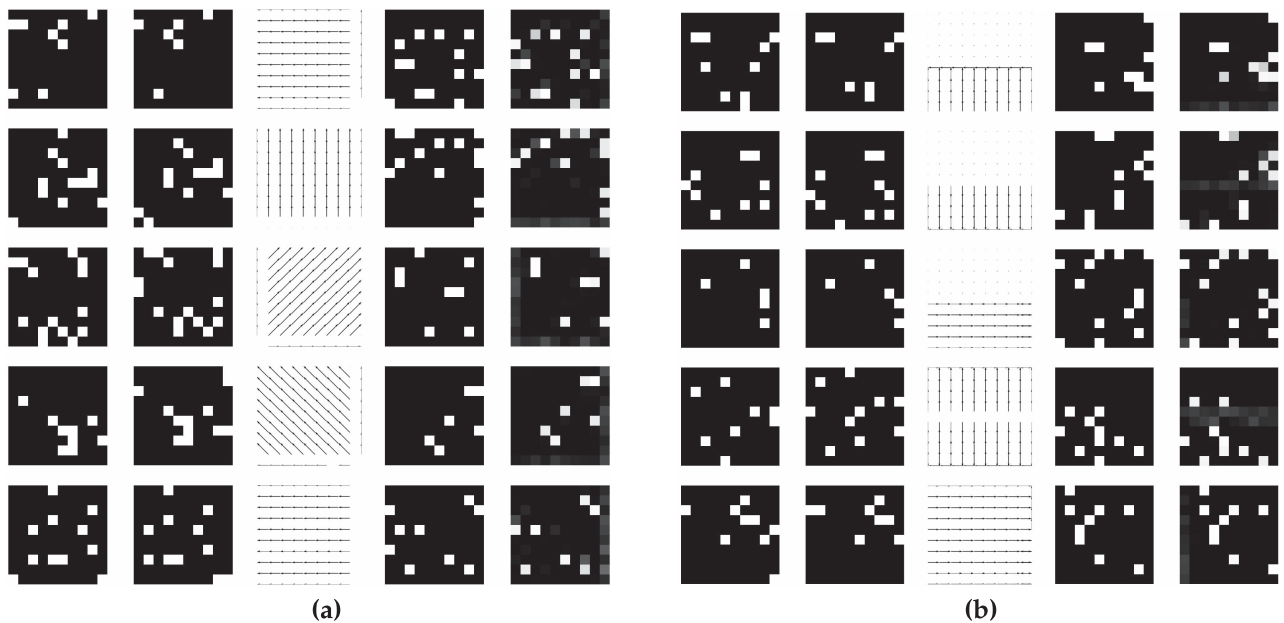


Fig. 7. Inferring transformations from test data. (a) Coherent motion across the whole image. (b) “Factorial motion” that is independent in different image regions. In both plots, the meaning of the five columns is as follows (left-to-right): Random test images  $x$ , random test images  $y$ , inferred flow-field, new test-image  $\hat{x}$ , inferred output  $\hat{y}$ .

## 2.5 Toy Example: Motion Extraction and Analogy Making

An example of a GBM applied to a motion inference task is shown in Fig. 7. We trained a GBM on binary image pairs containing random dots such that the output image  $y$  is a translated copy of the input image  $x$ . Some example image pairs are shown in the two left-most columns of Fig. 7a. The center column of Fig. 7a visualizes the corresponding inferred transformations as vector fields. To generate the vector field, we first infer the linear warp from an image pair using (12) and (15). Subsequently, we find for each input pixel the output pixel to which it is most strongly connected according to the inferred linear transformation, and we draw an arrow pointing from the input pixel to the output pixel. The plot shows that up to unpredictable edge effects, the model can correctly infer the translations inherent in the image pairs after being trained on shifts.

The two right-most columns in Fig. 7 show how the inferred transformation can be applied to new images not seen during training *by analogy*. To this end, we apply the inferred linear transformation to the input test image using (13).

Fig. 7b shows a variation of this task, where the transformations are “split-screen” translations, that is, translations that are independent in the top half versus the bottom half of the image. This example demonstrates how the model is able to decompose transformations into independent constituting transformations. This ability of the model is crucial for encoding natural videos, which contain a multitude of transformations as a result of a combination of many local transformations. We shall discuss the learning of natural video data in more detail below.

## 2.6 A Brief History of Gating

Shortly after mapping units were introduced in 1981, energy models [3] received a lot of attention. Energy

models apply *squaring nonlinearities* to features and have therefore also been referred to as “square-pooling” models. Energy models have also been common as models of *complex cells* [10].

Since the square of a multiview (e.g., binocular) feature can be shown to implicitly encode pairwise products between simple (monocular) features, energy models are closely related to multiplicative feature learning models. In fact, energy models can be used to implement mapping units and vice versa. We discuss this relationship in detail in Sections 3 and 4. Early work on energy models suggested these as a way to encode motion by relating time frames in a video [3], and to perform stereo vision by relating images from different viewpoints [24], [25], [26]. An approach to learning-based disparity estimation was introduced later by Becker and Hinton [27].

In the early work on energy models, hard-wired Gabor features were used as the linear receptive fields instead of features that are learned from data [26], [25], [28]. The focus on Gabor features has somewhat biased the analysis of energy models to focus on the *Fourier-spectrum* as the main object of interest (see, e.g., [28], [26]). As we shall discuss in Section 3, Fourier-components arise just as the special case of one transformation class, namely, *translation*, and many important properties of these models apply also to other transformation classes.

Energy models based on Gabor features have also been applied to a single image. In this case, they encode features independently of the Fourier-phase of the input and, as a result, their responses are invariant to small translations as well as to contrast variations of the input (see, e.g., [10]).

Shortly after energy and cross-correlation models emerged, there was some interest in learning invariances with higher-order neural networks, which are neural networks trained on polynomial basis expansions of their inputs [29]. Higher-order neural networks can be composed

of computational units that compute sums of products. These units are sometimes referred to as “Sigma-Pi-units” [30] (where “Pi” stands for product and “Sigma” for sum). At the same time, multiplicative interactions have also been explored as an approach to building distributed representations of symbolic data (e.g., [31], [32]).

In 1995, Kohonen introduced the “Adaptive Subspace Self-Organizing Map” (ASSOM) [33], which computes sums over squared filter responses to represent data. Like the energy model, the ASSOM is based on the idea that the sum of squared responses is invariant to various properties of its inputs. In contrast to the early energy models, the ASSOM is trained from data. Inspired by the ASSOM, “Independent Subspace Analysis” (ISA) was introduced by Hyvärinen and Hoyer [23], who place the idea in the context of more conventional feature learning models. Shortly thereafter, extensions of this work showed how the grouping of squared filter responses can be used to learn topographic feature maps [34], [35].

In a parallel line of work, bilinear models were proposed at approximately the same time as an approach to learning in the presence of multiplicative interactions [17]. The early work on bilinear models used these as global models trained on whole images rather than using local receptive fields. In contrast to more recent approaches to learning with multiplicative interactions, training involved filling a two-dimensional grid with data that shows two types of variability (referred to as “style” and “content”). The purpose of bilinear models is then to untangle the two degrees of freedom in the data. More recent work does not make this distinction, and the purpose of multiplicative hidden variables is merely to capture the multiple ways in which two images can be related. The work by the authors of [15], [14], or [16], for example, shows how multiplicative interactions make it possible to model the multitude of relationships between frames in natural videos, or between artificially transformed images [16]. An earlier multiplicative interaction model that is also related to bilinear models is the “routing-circuit” [36].

Multiplicative interactions have also been used to model structure within static images, which can be thought of as modeling higher-order relations and, in particular, pairwise products between pixel intensities (e.g., [22], [23], [37], [38], [39], [40]).

### 3 FACTORIZATION AND ENERGY MODELS

In the following, we discuss the close relationship between gated feature learning models and energy models. To this end, we first describe how parameter factorization makes it possible to preprocess input images and thereby reduce the number of parameters.

#### 3.1 Factorizing the Gating Parameters

The number of gating parameters is roughly cubic in the number of pixels if we assume that the number of constituting transformations is about the same as the number of pixels. It can easily be more for highly overcomplete hiddens. One way to reduce that number is by factorizing the parameter tensor  $W$  into three matrices

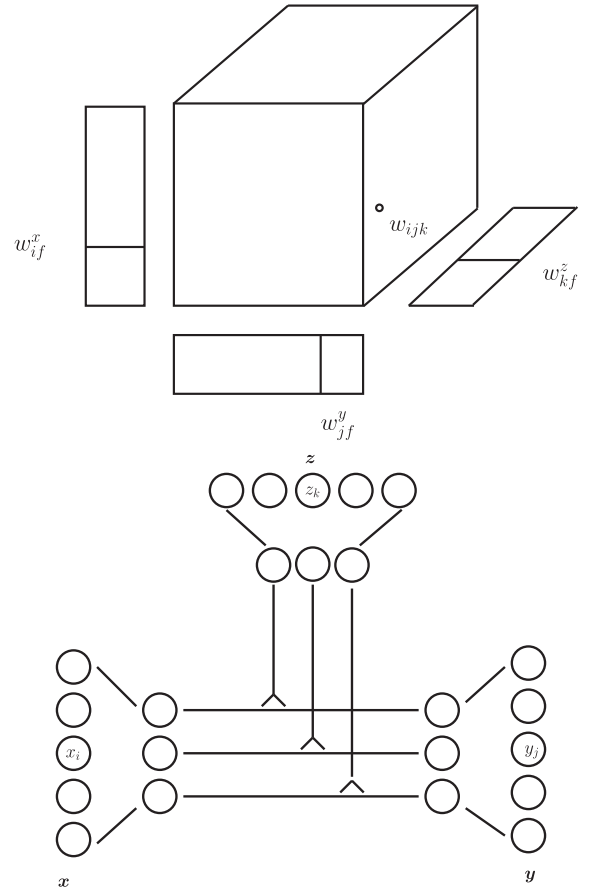


Fig. 8. Top: Factorizing the parameter tensor. Bottom: Interpreting factorization as filter matching.

such that each component  $w_{ijk}$  is given by a “three-way inner product” [41]:

$$w_{ijk} = \sum_{f=1}^F w_{if}^x w_{jf}^y w_{kf}^z. \quad (22)$$

Here,  $F$  is a number of hidden “factors,” which, like the number  $K$  of hidden units, has to be chosen by hand or by cross validation. The matrices  $w^x$ ,  $w^y$ , and  $w^z$  are  $I \times F$ ,  $J \times F$ , and  $K \times F$ , respectively.

An illustration of this factorization is given in Fig. 8 (top). It is interesting to note that, under this factorization, the activity of output-variable  $y_j$ , by using the distributive law, can be written

$$\begin{aligned} y_j &= \sum_{ik} w_{ijk} x_i z_k = \sum_{ik} \left( \sum_f w_{if}^x w_{jf}^y w_{kf}^z \right) x_i z_k \\ &= \sum_f w_{jf}^y \left( \sum_i w_{if}^x x_i \right) \left( \sum_k w_{kf}^z z_k \right). \end{aligned} \quad (23)$$

Similarly, for  $z_k$  we have

$$\begin{aligned} z_k &= \sum_{ij} w_{ijk} x_i y_j = \sum_{ij} \left( \sum_f w_{if}^x w_{jf}^y w_{kf}^z \right) x_i y_j \\ &= \sum_f w_{kf}^z \left( \sum_i w_{if}^x x_i \right) \left( \sum_j w_{jf}^y y_j \right). \end{aligned} \quad (24)$$



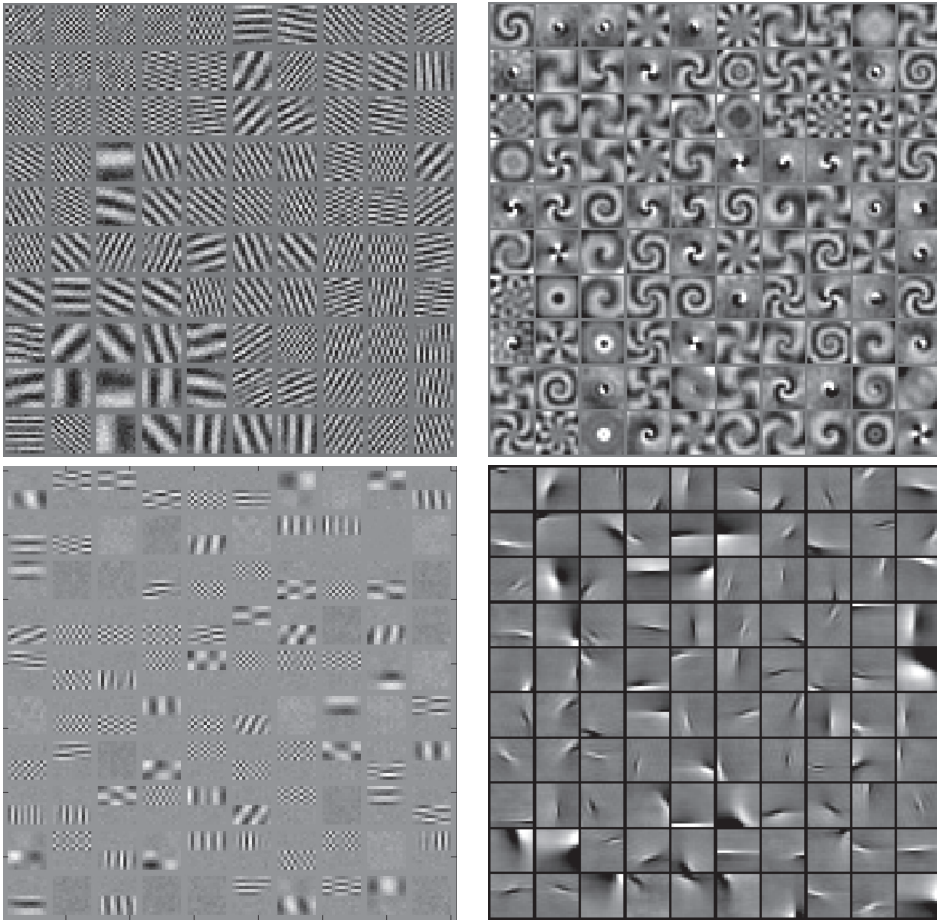


Fig. 9. Input filters learned from various types of transformation. Top-left: translation, Top-right: rotation, Bottom-left: split-screen translation, Bottom-right: natural videos. See Fig. 10 for corresponding output filters.

One can obtain a similar expression for the energy in a GBM. Equation (24) shows that factorization can be viewed as *filter matching*: For inference, each group of variables  $\mathbf{x}$ ,  $\mathbf{y}$ , and  $\mathbf{z}$  are projected onto linear basis functions, which are subsequently multiplied, as illustrated in Fig. 8 (bottom).

It is important to note that the way factorization reduces parameters is *not* necessarily by projecting data into a lower-dimensional space before computing the multiplicative interactions—a claim that can be found frequently in the literature. In fact, frequently  $F$  is chosen to be *larger* than  $I$  and/or  $J$ . The way that factorization reduces the number of parameters is by restricting the *three-way connectivity*, that is, with factorization, the number of pairwise products is equal to the number of factors rather than equal to the number of pixels squared. Learning then amounts to finding basis functions that can deal with this restriction optimally.

All gated feature learning models can be subjected to this factorization. Training is similar to training an unfactored model, which can be seen by using the chain rule and differentiating (22). An example of a factored gated autoencoder is described in [19]. Virtually all factored models that were introduced use the restriction of single multiplicative interactions (22). An open research question is to what degree a less restrictive connectivity—equivalently, using a nondiagonal core-tensor in the factorization—would be advantageous.

Factored models have empirically been shown to learn filter pairs that optimally represent transformation classes such as Fourier-components for translations and a polar variant of Fourier-components for rotations [41]. In contrast to the dictionaries learned with standard feature learning methods, the filters always come in pairs. These may be referred to as “predictionary” as they are often learned using predictive training (see Section 2.4).

Figs. 9 and 10 show examples of predictionaries learned from translations, affine transformations, split-screen translations, which are independent translations in the top and bottom half of the image, and natural video. For training the filters in the top rows and on the bottom right, we used datasets described in [41] and [19] and the model described in [19]. The filters resemble receptive fields found in various cells in visual cortex [42]. To obtain split-screen filters (bottom left) we generated a dataset of split-screen translations and trained the model described in [41]. In Section 4, we provide an analysis that sheds some light onto why the filters take on this form.

In practice, it is common to utilize a set of tricks to improve stability and efficiency of learning. It is common, for example, to normalize output filter matrices  $w^x$  and  $w^y$  during learning such that all filters  $w_{\cdot f}^x$  and  $w_{\cdot f}^y$  grow slowly and maintain roughly the same length as learning progresses. This is typically achieved by renormalizing filters after each parameter update (see, e.g., [20], [39]). It is also

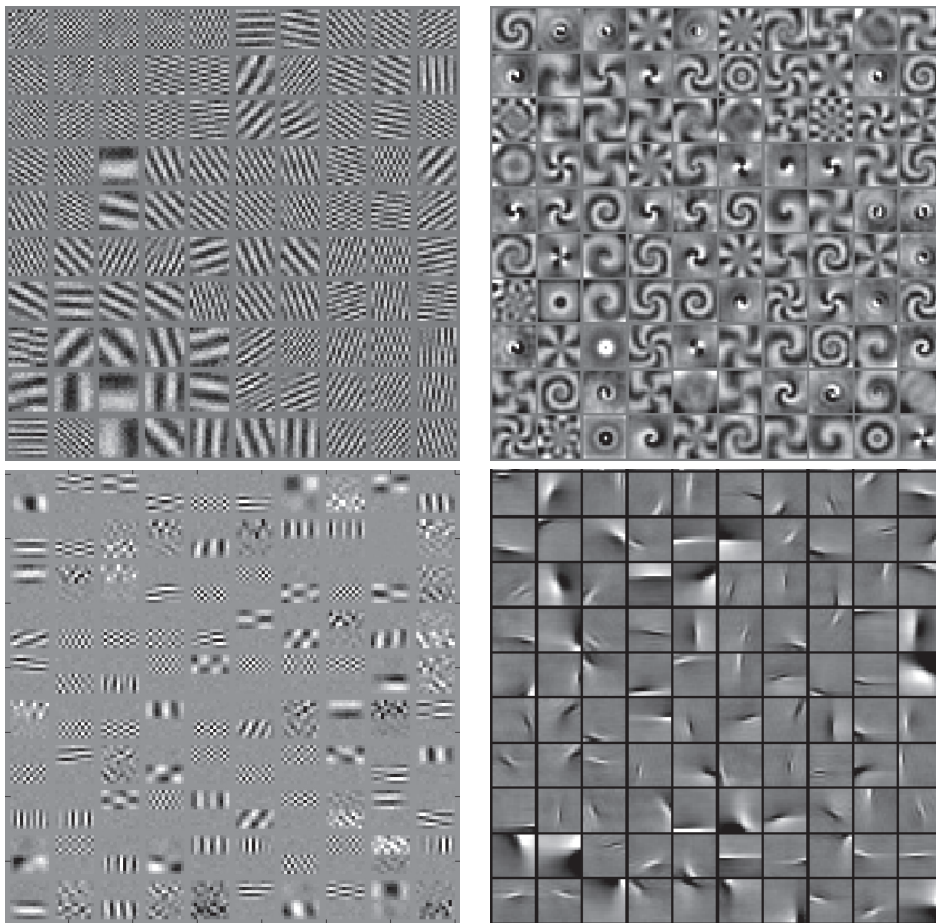


Fig. 10. Output filters learned from various types of transformation. Top-left: translation, Top-right: rotation, Bottom-left: split-screen translation, Bottom-right: natural videos. See Fig. 9 for corresponding input filters.

common to connect hidden units locally to the factors, rather than using full connectivity. A slightly more complicated approach is to let all hidden units populate a virtual “grid” in a low-dimensional space (e.g., 2D) and to connect hidden units to factors such that neighboring hidden units are connected to the same or to overlapping sets of factors. This typically leads to topographic organization of filters, an example of which is the set of shift filters shown in Figs. 9 and 10. This approach is also common for learning energy models on still images (e.g., [34], [35]). Finally, it is common to train the models using image patches that are DC centered and contrast normalized, and usually also whitened. For a quantitative comparisons of several variations of gated feature learning models see [43].

### 3.2 Energy Models

Energy models [3], [24] are an alternative approach to modeling image motion and disparities, and they have been deployed monocularly, too. A main application of energy models has been the detection of small translational motion in image pairs. This makes them suitable as biologically plausible mechanisms of both local motion estimation and binocular disparity estimation. Energy models detect motion by projecting two images onto two phase-shifted Gabor functions each (for a total of four basis function responses). The two responses *across* the images are added and

squared. The sum of these two squared, spatiotemporal responses then yields the response of the energy model.

The rationale behind the energy model is that because each within-image Gabor filter pair can be thought of as a localized spatiotemporal Fourier-component, the sum of the squared components yields an estimate of spectral energy, which is not dependent of the *phase*—and thus to some degree not dependent on the *content*—of the input images. The two filters within each image need to be sine/cosine pairs, which is commonly referred to as being “in quadrature.”

A detector of local shift can be built by using a set of energy models tuned to different frequencies. To turn a set of energy models into an estimate of local translation, one can, for example, read off the shift from the model with the strongest response [25], [26] or use pooling to get a more stable estimate [28].

In their early approach to *learning* energy models from training data, Hyvärinen and Hoyer [23] suggest extending a standard feature learning model by introducing an elementwise squaring operation and adding a linear pooling layer. For learning, they suggest adapting ICA by forcing latent variable responses (which are now sums of squared basis function responses) to be sparse, while keeping the filters orthogonal to avoid degenerate solutions. This approach is known as “ISA” [23]. ISA was introduced initially to model single images not pairs, but it is possible

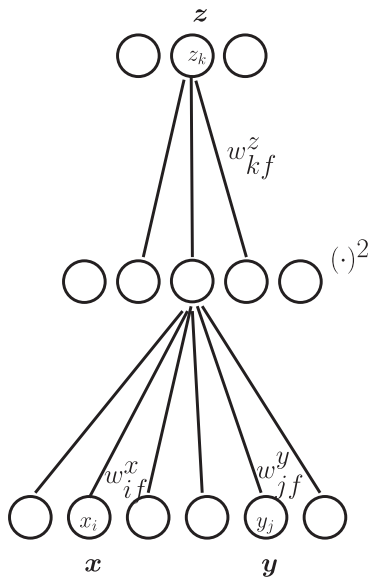


Fig. 11. Illustration of ISA applied to an image pair  $(\mathbf{x}, \mathbf{y})$ .

to apply it to the concatenation of multiple images, like the early versions of the energy model [3], [24]. In contrast to the early models, it is common to use ISA models that pool over more than two filters, and pooling weights can be learned along with the filters, instead of being fixed to one. Fig. 11 shows an illustration of ISA applied to an image pair. As the figure shows, the model can be viewed as a two-layer network, with a hidden layer that uses an elementwise squaring nonlinearity. The first hidden layer of an energy model is closely to the latent “factors” of a factored GBM as we shall show. Both ISA and factored GBMs were recently shown independently to yield state-of-the-art performance in various motion recognition tasks [44], [45].

### 3.3 Relationship between Gated Feature Learning and Energy Models

Learning energy models such as ISA on the concatenation of two inputs  $\mathbf{x}$  and  $\mathbf{y}$  is closely related to learning gated

feature learning models. Let  $w_{\cdot f}^x$  ( $w_{\cdot f}^y$ ) denote the set of weights connecting part  $\mathbf{x}$  ( $\mathbf{y}$ ) of the concatenated input with factor  $f$  (see Fig. 11). The activity of hidden unit  $z_k$  in the energy model is given by

$$\begin{aligned} z_k &= \sum_f w_{kf}^z (w_{\cdot f}^x \mathbf{x} + w_{\cdot f}^y \mathbf{y})^2 \\ &= \sum_f w_{kf}^z (2(w_{\cdot f}^x \mathbf{x})(w_{\cdot f}^y \mathbf{y}) + (w_{\cdot f}^x \mathbf{x})^2 + (w_{\cdot f}^y \mathbf{y})^2). \end{aligned} \quad (25)$$

Up to the quadratic terms in (25), hidden unit activities are the same as in a gated feature learning model (see (24)). As we shall discuss in detail below, the quadratic terms do not have a significant effect on the meaning of the hidden units. The hidden units in an energy model may therefore be interpreted as a way to implement mapping units which encode relations. See also [28] for a discussion of this relationship in the context of the traditional energy models.

## 4 RELATIONAL CODES AND SIMULTANEOUS EIGENSPPACES

We now show that hidden variables learn to detect subspace-rotations when they are trained on transformed image pairs. In Section 2.3 (14), we showed that transformation codes  $\mathbf{z}$  represent linear transformations,  $L$ , that is  $\mathbf{y} = L\mathbf{x}$ . We shall restrict our attention in the following to transformations that are orthogonal, that is,  $L^T L = L L^T = I$ , where  $I$  is the identity matrix. In other words,  $L^{-1} = L^T$ . Note that practically all relevant spatial transformations, like translation, rotation, or local shifts, can be expressed approximately as an orthogonal warp because orthogonal transformations subsume, in particular, all *permutations* (“shuffling pixels”).

An important fact about orthogonal matrices is that the eigendecomposition  $L = UDU^T$  is complex where eigenvalues (diagonal of  $D$ ) have absolute value 1 [46]. Multiplying by a complex number with absolute value 1 amounts to performing a rotation in the complex plane, as illustrated in Fig. 12 (left). Each eigenspace associated with  $L$  is also

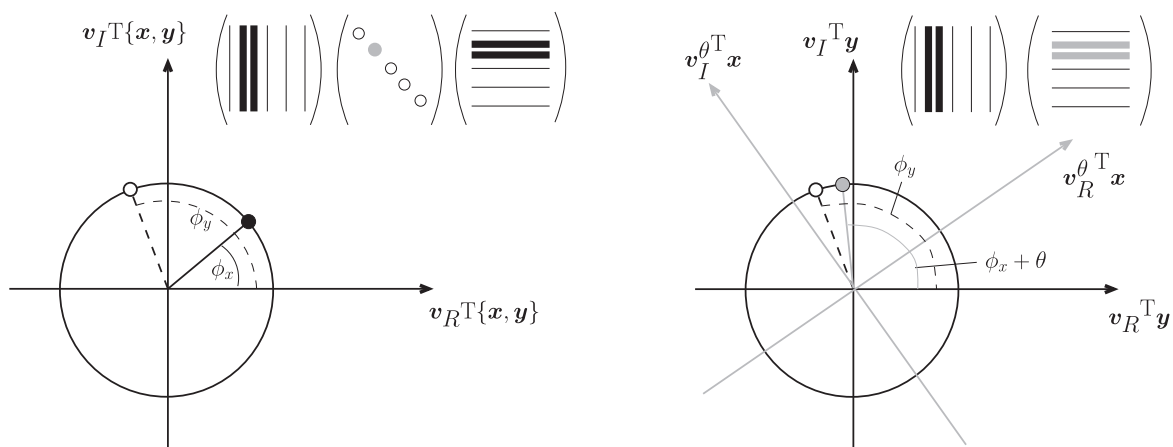


Fig. 12. Left: Inference in a gated feature learning model is equivalent to extracting rotation angles from two-dimensional invariant subspaces. Right: By absorbing eigenvectors into the eigenvectors, a mapping unit can learn to detect rotations by a particular angle (its preferred angle): The inner product between the projections of the two images  $\mathbf{x}$  and  $\mathbf{y}$  in the figure will be maximal when  $\phi_y = \phi_x + \theta$ , that is, when the rotation that the detector applies to  $\mathbf{x}$  has the effect of aligning  $\mathbf{x}$  with  $\mathbf{y}$ .

referred to as the *invariant subspace* of  $L$  (as the application of  $L$  will keep the eigenvectors within the subspace).

Applying an orthogonal warp is, thus, equivalent to 1) projecting the image onto *filter pairs* (the real and imaginary parts of each eigenvector), 2) performing a rotation within each invariant subspace, and 3) projecting back into the image-space. In other words, we can decompose an orthogonal transformation into a set of independent, 2D rotations. The most well-known examples are translations: A 1D-translation matrix contains ones along one of its secondary diagonals, and it is zero elsewhere.<sup>2</sup> The eigenvectors of this matrix are Fourier-components [47], and the rotation in each invariant subspace amounts to a phase-shift of the corresponding Fourier-feature. This leaves the norm of the projections onto the Fourier-components (the power spectrum of the signal) constant, which is a well-known property of translation.

It is interesting to note that the imaginary and real parts of the eigenvectors of a translation matrix correspond to sine and cosine features, respectively, reflecting the fact that Fourier-components naturally come in *pairs*. These are commonly referred to as *quadrature pairs* in the literature. In the special case of Gabor features, the importance of quadrature pairs is that they allow us to detect translations independently of the local content of the images [26], [28]. However, the property that eigenvectors come in *pairs* is not specific to translations. It is shared by all transformations that can be represented by an orthogonal matrix so that they can be composed from 2D rotations. Bethge et al. [48] use the term “*generalized quadrature pair*” to refer to the eigenfeatures of these transformations.

#### 4.1 Commuting Warps Share Eigenspaces

An observation that is central to our analysis is that eigenspaces can be *shared* among transformations. When eigenspaces are shared, then the only way in which two transformations differ is in the angles of rotation within the eigenspaces. That way, shared eigenspaces allow us to represent *multiple transformations with a single set of features*. An example of a shared eigenspace is the Fourier-basis, which is shared among translations. This well-known observation follows from the fact that the set of all circulant matrices (which are 1D translation matrices) of the same size have the Fourier-basis as eigenbasis [47]. However, eigenspaces can be shared between other transformations. An obvious generalization is local translation, which may be considered the constituting transformations of natural videos. Another, less obvious generalization is spatial rotation. Formally, two matrices  $A, B$  share eigenvectors if they *commute*, that is, if  $AB = BA$  holds [46].<sup>3</sup> As an example, consider translations: Translating an image by  $a$  pixels to the left and then by  $b$  pixels upward yields the same result as first translating it upward and then to the left, showing that translations commute.

2. To be exactly orthogonal it has to contain an additional one in another place so that it performs a rotation with wrap-around.

3. This can be seen by considering any two matrices  $A$  and  $B$  with  $AB = BA$  and with  $\lambda, v$  an eigenvalue/eigenvector pair of  $B$  with multiplicity one. It holds that  $BAv = ABv = \lambda Av$ . Therefore,  $Av$  is also an eigenvector of  $B$  with the same eigenvalue.

The importance of commuting transformations for our analysis is that since these transformations share an eigenbasis, they differ only w.r.t. the angle of rotation in the joint eigenspace. As a result, we may extract a particular transformation from a given image pair  $(\mathbf{x}, \mathbf{y})$  by recovering the angles of rotation between the projections of  $\mathbf{x}$  and  $\mathbf{y}$  onto the eigenspaces. To this end, consider the real and complex parts  $\mathbf{v}_R$  and  $\mathbf{v}_I$  of some eigenfeature  $\mathbf{v}$ . That is,  $\mathbf{v} = \mathbf{v}_R + i\mathbf{v}_I$ , where  $i = \sqrt{-1}$ . The real and imaginary coordinates of the projection of  $\mathbf{x}$  onto the invariant subspace associated with  $\mathbf{v}$  are given by  $\mathbf{v}_R^T \mathbf{x}$  and  $\mathbf{v}_I^T \mathbf{x}$ , respectively. For the output image, they are  $\mathbf{v}_R^T \mathbf{y}$  and  $\mathbf{v}_I^T \mathbf{y}$ .

Let  $\phi_x$  and  $\phi_y$  denote the angles of the projections of  $\mathbf{x}$  and  $\mathbf{y}$  with the real axis in the complex plane. If we normalize the projections to have unit norm, then the cosine of the angle between the projections,  $\phi_y - \phi_x$ , may be written

$$\cos(\phi_y - \phi_x) = \cos \phi_y \cos \phi_x + \sin \phi_y \sin \phi_x,$$

by trigonometric identity. This is equivalent to computing the inner product between two normalized projections (see Fig. 12 (left)). In other words, to estimate the (cosine of) the angle of rotation between the projections of  $\mathbf{x}$  and  $\mathbf{y}$ , we need to *sum over the product of two filter responses*.

Note, however, that normalizing each projection to 1 amounts to dividing by the sum of squared filter responses, an operation that is highly unstable if a projection is close to zero. This will be the case whenever one of the images is almost orthogonal to the invariant subspace. This, in turn, means that the rotation angle *cannot be recovered from the given image* because the image is too close to the axis of rotation. One may view this as a subspace generalization of the well-known *aperture problem* beyond translation to the set of orthogonal transformations. Normalization would ignore this problem and provide the illusion of a recovered angle even when the aperture problem makes the detection of the transformation component impossible. In the next section, we discuss how gated feature learning overcomes this problem by allowing us to treat the problem as a *rotation detection* task instead.

#### 4.2 Representing Transformations by Detecting Subspace Rotations

For each eigenvector,  $\mathbf{v}$ , and rotation angle,  $\theta$ , define the complex filter

$$\mathbf{v}^\theta = \exp(i\theta)\mathbf{v},$$

which represents a projection and simultaneous rotation by  $\theta$ . This amounts to *absorbing* the rotation, as given by some eigenvalue, into the eigenvector itself, allowing us to define a subspace rotation-detector with preferred angle  $\theta$  as follows:

$$r^\theta = (\mathbf{v}_R^T \mathbf{y})(\mathbf{v}_R^{\theta T} \mathbf{x}) + (\mathbf{v}_I^T \mathbf{y})(\mathbf{v}_I^{\theta T} \mathbf{x}). \quad (26)$$

Like before, if projections would be normalized to length 1, we would have

$$\begin{aligned} r^\theta &= \cos \phi_y \cos(\phi_x + \theta) + \sin \phi_y \sin(\phi_x + \theta) \\ &= \cos(\phi_y - \phi_x - \theta), \end{aligned}$$

which would be maximal whenever  $\phi_y - \phi_x = \theta$ , thus when the observed angle of rotation,  $\phi_y - \phi_x$ , is equal to the preferred angle of rotation,  $\theta$ . An illustration of such a rotation detector is given in Fig. 12 (right).

Although normalizing projections is not a good idea due to the subspace aperture problem, normalization turns out not to be necessary if the task is defined as a detection task. In particular, note that if features and data are contrast normalized, then the subspace inner product (26) will yield a strong response if the following two conditions are met:

- the angle between the projections of  $\mathbf{x}$  and  $\mathbf{y}$  matches the detector's preferred angle,  $\theta$ ,
- the projections of the images onto the invariant subspace are large enough (in other words, the images are sufficiently well aligned with the subspace).

The second condition implies that the output of the detector defined in (26) factors in not only the presence of a transformation but also its ability to discern it. In other words, when the detector fires, we know its preferred transformation is present. When it does not fire, then either the transformation is not present or it is present but invisible to the detector because the images are not well aligned with its invariant subspace.

However, when a transformation that is present is invisible to a detector, then a different detector defined over a different invariant subspace may still be able to observe the transformation. Thus, a population of detectors can yield a robust representation of transformations. As an example, consider two images  $\mathbf{x}$  and  $\mathbf{y}$  related through vertical translation. The invariant subspaces of translations are spanned by Fourier-components. Vertical translation, in particular, is represented by phase-shifts of horizontal Fourier-components. Now consider a rotation detector defined in a horizontal Fourier-component of a particular frequency. If the images happen to lack this horizontal frequency, then the detector will be silent, even if its transformation is present. Since vertical translation, however, is detectable not only in one but in many subspaces (such as in the other horizontal Fourier-components of different frequencies), detectors defined over those subspaces may still be able to detect the transformation. In fact, if we assume that detectors for all subspaces (horizontal frequencies) are present, then the only way that no detector fires would be if the transformation is not present.

One way to combine the information from multiple detectors into a single representation of a transformation is by *pooling* because a sum (or weighted sum) over multiple detectors will be able to represent the event that any of the detectors fires. That way, a weighted sum over rotation detectors can yield a representation of transformations that is not dependent on the content of the images (such as the frequency content in the case of translation).

Formally, if we stack imaginary and real eigenvector pairs for the input and output images,  $\mathbf{v}$  and  $\mathbf{v}^\theta$ , column-wise in matrices  $V$  and  $U$ , respectively, we may define the representation  $t$  of a transformation, given two images  $\mathbf{x}$  and  $\mathbf{y}$ , as

$$t = W^T P (U^T \mathbf{x}) \cdot (V^T \mathbf{y}), \quad (27)$$

where  $W$  is an appropriate “across-subspace” pooling matrix, and  $P$  is a band-diagonal “within-subspace”

pooling matrix that defines the two-dimensional inner product in (26).<sup>4</sup>

Note that (27) takes exactly the same form as inference in a factored gated feature learning model (see (24)) if we absorb the within-subspace pooling matrix  $P$  into  $W$ .

This shows that we may interpret mapping units in a gated feature learning model as a way to encode transformations by *representing rotation angles in invariant subspaces*.

### 4.3 Learning as Simultaneous Diagonalization

The interpretation of mapping units as encodings of rotations relies on several assumptions:

- Images  $\mathbf{x}$  and  $\mathbf{y}$  are contrast-normalized.
- The columns of both  $U$  and  $V$  come in pairs, each of which spans a two-dimensional invariant subspace of a transformation class. (Formally, the pairs represent the real and imaginary components of some complex eigenvector of the transformation class.)
- Corresponding filter pairs in  $U$  and  $V$  are related through rotations only. In other words, for each filter pair  $\mathbf{v}_f$  in  $V$  there exists  $\theta$  such that the corresponding filter pair in  $U$  can be written  $\mathbf{v}_f^\theta = \exp(i\theta)\mathbf{v}_f$ .

While it would be possible in principle to define filter pairs with these properties by hand to represent a given transformation class, model parameters can be learned from image pairs, as we discussed in Sections 2 and 3. In particular, if we interpret the inference equations in a factored model (24), (27) as constraints under which we learn to represent the training image pairs, then it becomes clear that learning may be viewed as a way to find appropriate sets of filter pairs that are able to detect the rotations.

Learning a factored gated feature learning model thus has the effect of performing an *approximate simultaneous diagonalization of a set of transformations*. As we showed in Section 3, training on translations, for example, yields Fourier-features, which indeed represent the invariant subspaces of translation. In addition to finding representations of the invariant subspaces, learning involves finding an appropriate across-subspace pooling matrix  $W$ .

In [49], it is shown empirically that when a dataset contains more than one transformation class, learning can involve partitioning the set of observed transformations into commutative subsets, simultaneously diagonalizing each subset.

It is important to note that although complex arithmetic provides a convenient notational framework in the analysis, complex numbers are not typically used in any actual implementations. Learning yields filters which implicitly span the two-dimensional invariant subspaces of a transformation class. However, these subspaces do not need to be predefined as the real/imaginary components of a complex vector nor do filter pairs need to be exactly in quadrature as long as the learning can adapt them to represent two-dimensional rotations. As an analogy, consider training a linear autoencoder to perform PCA: While the autoencoder features will span the same subspace as the eigenvectors of the data covariance matrix, they will

4. To this end,  $P$  has to contain exactly 2 ones along each row and has to be zero elsewhere.



typically not be exactly the same (in particular, the features will typically not be orthogonal). By making use of (27) rather than the more common (24) to define mapping unit activations (and thereby keeping within-subspace pooling and across-subspace pooling separate), it is nevertheless possible to visualize the learned “real” and “imaginary” components of the learned filters [49].

It is interesting to note that diagonalizing a single transformation would amount to performing a kind of canonical correlations analysis (CCA), so training a multi-view feature learning model may be thought of as performing multiple (possibly nonlinear) canonical correlation analyses with tied features. Note that learning such a “mixture” of related CCA models, rather than a single CCA model, is crucial in practical applications because typically any image can potentially be transformed by many different transformations from the given class and not just by a single instance from that class. Similarly, modeling within-image structure by setting  $\mathbf{x} = \mathbf{y}$  (see, Section 2.4.2) is like learning a PCA mixture with tied weights.

#### 4.4 Relation to Energy Models

By concatenating images  $\mathbf{x}$  and  $\mathbf{y}$ , as well as filters  $\mathbf{v}$  and  $\mathbf{v}^\theta$ , we may approximate the subspace rotation detector (26) with the response of an energy detector:

$$\begin{aligned} r^\theta &= ((\mathbf{v}_R^T \mathbf{y}) + (\mathbf{v}_R^{\theta T} \mathbf{x}))^2 + ((\mathbf{v}_I^T \mathbf{y}) + (\mathbf{v}_I^{\theta T} \mathbf{x}))^2 \\ &= 2((\mathbf{v}_R^T \mathbf{y})(\mathbf{v}_R^{\theta T} \mathbf{x}) + (\mathbf{v}_I^T \mathbf{y})(\mathbf{v}_I^{\theta T} \mathbf{x})) \\ &\quad + (\mathbf{v}_R^T \mathbf{y})^2 + (\mathbf{v}_R^{\theta T} \mathbf{x})^2 + (\mathbf{v}_I^T \mathbf{y})^2 + (\mathbf{v}_I^{\theta T} \mathbf{x})^2. \end{aligned} \quad (28)$$

Equation (28) is equivalent to (26) up to the four quadratic terms, which represent the squared norms of the projections of  $\mathbf{x}$  and  $\mathbf{y}$  onto the invariant subspace. Thus, like the norm of the projections, they contribute information about the discernibility of transformations. By adding the square terms, the detector will have a different tuning behavior than the inner product detector defined in (26). However, the energy detector still attains its peak response exactly when both images reside within its invariant subspace and when their projections are rotated by the detector’s preferred angle  $\theta$ , so it is tuned to the same transformation as the inner product detector. This shows that energy models applied to the concatenation of two images are well suited to modeling transformations, too. The inner product detector (26) is commonly referred to as cross-correlation model in the literature (e.g., [28]).

#### 4.5 Representing Videos

Both energy models and cross-correlation models can be applied to sequences of more than two images. In a gated feature learning model, (26) may be modified to contain pairwise products across all time steps or all those that are deemed relevant (e.g., products between adjacent frames).

In the energy model, (28) may be modified to compute the square of the concatenation of more than two images. In particular, consider the concatenation (“vectorization”)  $\mathbf{X}$  of a sequence of  $T$  frames  $\mathbf{x}_t$ ,  $t = 1, \dots, T$ , projected onto the concatenation of  $T$  corresponding feature vectors  $\mathbf{v}_R^t$  and  $\mathbf{v}_I^t$ ,  $t = 1, \dots, T$ . The sum of squares of the projections onto these filters will yield an energy response of the form

$$\begin{aligned} r &= \left( \sum_t \mathbf{v}_R^{t T} \mathbf{x}_t \right)^2 + \left( \sum_t \mathbf{v}_I^{t T} \mathbf{x}_t \right)^2 \\ &= \Omega + \sum_{st} \left[ (\mathbf{v}_R^{s T} \mathbf{x}_s)(\mathbf{v}_R^{t T} \mathbf{x}_t) + (\mathbf{v}_I^{s T} \mathbf{x}_s)(\mathbf{v}_I^{t T} \mathbf{x}_t) \right], \end{aligned} \quad (29)$$

where we use  $\Omega$  to denote all square terms which, like before, represent subspace norms. Equation (29) shows that the energy response will implicitly contain the sum over *all* pairwise products between projected frames, or equivalently, the sum over inner products between all pairs of two-dimensional projections of the data. Thus, for  $r$  to take on a meaningful value, all angles implicit in these inner products have to be consistent; in other words, they have to be multiples of the first angle.

A limitation of the energy detector for modeling video data is therefore that it can only detect the repeated application of a single transformation.

##### 4.5.1 Example: Implementing an Energy model via a Cross-Correlation Model

The close relation between energy models and gated feature learning models makes it possible to implement one via the other. Fig. 13 shows example filters from an energy model trained on concatenated frames from videos showing moving random dots.<sup>5</sup> We trained a gated autoencoder with  $F = 256$  factors and  $K = 128$  mapping units, where  $\mathbf{x} = \mathbf{y}$  is given by the concatenation of 10 frames. Filters are constrained such that  $w_{if}^x = w_{if}^y$ . Each 10-frame input shows random dots moving at a constant speed. Speed and direction vary across movies.

Since the gated autoencoder, a cross-correlation model, multiplies two sets of filter responses which are the same, it effectively computes a square, and thus implements an energy model. In the absence of any within-image structure, all filters learn to represent only across-image correlations. Thus, as predicted by (25) the energy model in turn effectively implements a cross-correlation model.

Fig. 13 depicts, separately, the 10 sets of 256 filters corresponding to the 10 time frames. It shows that the model learns spatiotemporal Fourier-features that are selective for speed, frequency, and orientation.

## 5 DISCUSSION

While energy and cross-correlation models have been well known for almost 20 years, they have only recently begun to be applied successfully in many vision tasks, such as motion understanding and action recognition, where they now achieve state-of-the-art performance. We shall discuss some selected applications in the next section. Most of the applications make use of multilayer or deep learning extensions of multiview feature learning. We shall discuss the relationship between gating and deep learning in detail in Section 5.2, and we shall point to some open research problems in Section 5.3.

### 5.1 A Tour of Recent Applications

*Action recognition* is a prototypical example of a feature learning application that requires the ability to extract

5. Data and code are available at <http://learning.cs.toronto.edu/~rfm/relational>.

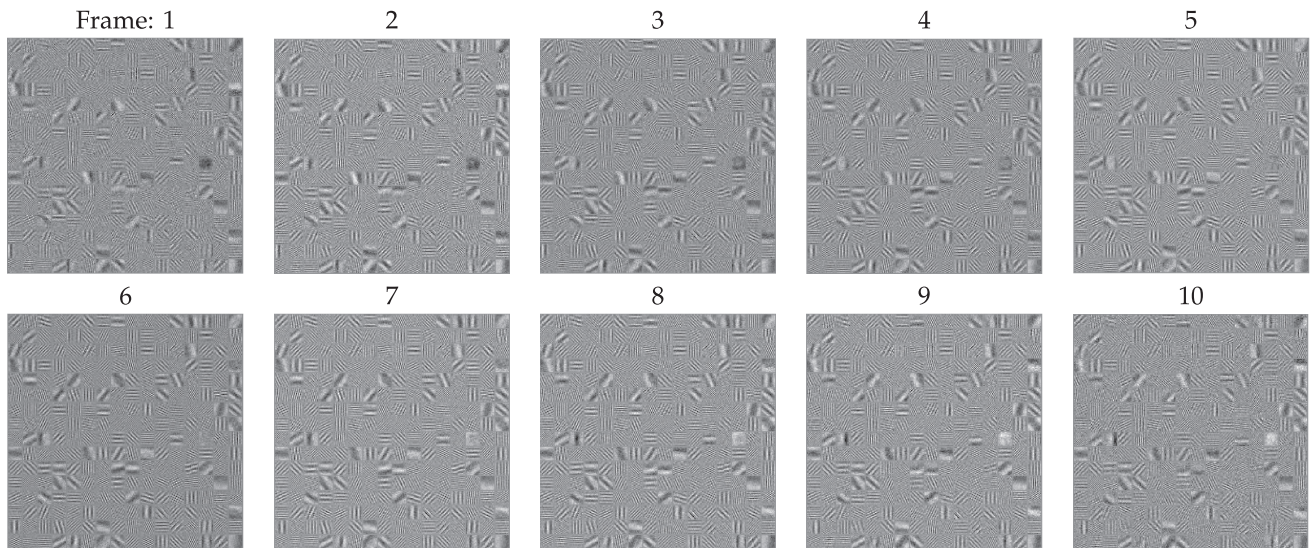


Fig. 13. Implementing a *cross-correlation model* via an *energy model* via a *cross-correlation model*. Sequence of filters learned from the concatenation of 10 frames of moving random dots.

*interframe* representations to extract motion from videos. The state-of-the-art approaches to action recognition are based on multiview feature learning, either by making use of gating connections [45] or of energy models [44]. To yield state-of-the-art performance, both types of model utilize multiple layers of features, so they combine the use of deep architectures with multiplicative gating. An empirical investigation of the usefulness of multiple layers is presented in [44].

*Learning invariant features from images.* It is interesting to note that invariant object recognition itself can be viewed as a correspondence problem where the goal is to match an input observation to an invariant template in memory. An approach to invariant recognition based on a variation of a multiview feature learning model is proposed in [50]. The model can be considered as an approach to invariant recognition through modeling mappings that take images to class labels. The input to this model is an image, the output is an orthogonal encoding of a class label, and prediction amounts to marginalizing over the set of possible mappings. The model transforms an input into a canonical pose so that it can be matched with a template, which itself represents the object in some canonical pose. As shown in [50], “swirly features” similar to the rotation features in Figs. 9 and 10 emerge when learning to perform rotationally invariant recognition. The model was shown to achieve state-of-the-art performance in a variety of invariant recognition tasks.

*Learning invariant features from videos.* Common object recognition systems differ from the way humans learn about objects in that they are trained on static views instead of movies, which show objects moving around or changing their pose. A model that computes squares or cross products can automatically learn to associate object identity with 3D structure or other object properties simply by being trained on multiple, concatenated frames. A recent application of this idea is the learning invariant features from videos [51]. In that work, a deep autoencoder is combined with an energy model, which makes it possible to learn

image features whose subspace energies stay constant across frames (see also [52] for a similar, probabilistic formulation of the same idea). A multilayer version of the model was shown to yield state-of-the-art results in object recognition by using subspace energies as features. A similar approach to learning invariances was proposed in [49]. In contrast to [51], the work in [49] uses phase-differences rather than subspace energies to represent objects.

*Tracking and pose estimation.* While in this paper we focus on the use of gating units to learn about image relations, there has been some recent work on applying multiplicative interactions in other tasks and in other domains. Examples include [53], [54], who use multiplicative interactions to gate hidden state transitions for tracking. The model by Denil et al. [54] is based on using three-way interactions to model the relationship between image regions [55]. Factored gating units have also been applied to the task of learning human pose dynamics from MOCAP data [56], and an extension of that work has been applied to tracking [57].

There exist several further recent applications of multiplicative interactions in a variety of domains. An example is the work by Tang et al. [58], who use gating connections to model occlusion with applications in denoising. Another example is the work by Sutskever et al. [59], who use factored gating connections to let input data modulate the hidden state transitions in a recurrent neural network. They report strong improvements in character-level language modeling, both over conventional recurrent networks and over other state-of-the-art models on these tasks.

A general explanation of these recent successes of multiplicative interactions, which goes beyond our analysis on image relations, may be that they provide a simple way to increase model capacity: When a very large amount of training data is available relative to the size of the models that may be trained within a reasonable amount of time, multiplication and unconventional nonlinearities can provide a computationally and memory-efficient way to increase the expressiveness of models.

## 5.2 Multiview Features and Deep Learning

A gated feature learning module is inherently a multilayer architecture because it uses pooling over multiplicative combinations (or squares) of features. Building deep models by stacking multiple layers of this type of module has been shown to further benefit applications, similar to the learning of image features, as we discussed in Section 5.1. There are several subtle but important differences in the role of deep learning for learning multiview representations, however, and we shall discuss these in the following.

Our analysis in Section 4 raises the question of to what degree our ability to learn about transformations is hampered by deviating from exactly commuting transformations and exactly orthogonal matrices. Existing experiments (e.g., in [16], [41]) suggest that there is some robustness w.r.t. deviations from exact orthogonality, but there has been no quantitative analysis. It is conceivable that one could preprocess input examples such that they can be related through approximately orthogonal matrices to make them amenable to gating or energy models. This would obviously require additional processing layers. In that way, deep learning can play an enabling role for building certain relational feature learning models. What the exact requirements of extra layers should be to improve the learnability of transformations is an important research question. It seems conceivable that high-dimensional and sparse representations as well as topographic filter organization may play an important role as they may make it easier to express relations in terms of *sets of local translations*.

A direct implication for learning of (nonlinear) visual feature hierarchies can be derived for spatial (e.g., geometric) transformations of images: The fact that the commutativity and orthogonality of the set of learned transformations must be preserved throughout the hierarchy imposes strong constraints on the learning of that hierarchy. One condition which seems to be well suited (if not necessary) to preserving this structure is *retinotopy*. In a feature hierarchy that is composed of retinotopic maps, a local shift in the image, for example, will map to a corresponding local shift in a higher layer feature map. While retinotopic organization of features is well known to hold in visual cortex, it has been much less common to explicitly enforce it in deep learning. Yet it may play a crucial role in keeping spatial relations intact.

A further implication of our analysis for the learning of deep architectures is that it suggests unconventional nonlinearities, such as squaring or rectification, to be useful for building models and for extending their applicability beyond still images. From a biological point of view, multiplicative interactions may also be viewed as a conceptually simple approximation to more complex *dendritic computations* [60] than the common neuron abstraction used in practically all deep learning models. The use of more elaborate models of dendritic computation in deep learning is a wide-open field, although uncommon nonlinearities, such as rectification, have begun to get an increasing amount of attention recently (e.g., [61]).

## 5.3 Directions for Further Research

The current standard approach to solving correspondence tasks in vision is by matching similar image regions and by using robust statistics to deal with the typically large

number of false positives [62]. Multiview feature learning differs from this type of approach in that it learns dense correspondences, albeit only in small image regions. The advantage of learning is that noise and false matches are dealt with automatically because hidden units jointly clean up and make sense of imperfect correspondences in the data. The disadvantage is that, for computational reasons, learning has to be restricted to very small image regions. An interesting research question is whether there is a bridge between multiview feature learning and solving sparse, global correspondence tasks in vision. Two ways in which multiview feature learning may be brought to use in such tasks are the following: 1) It may be possible to improve sparse matching by using similarity metrics derived from gated feature learning. This could help improve matching accuracy, and it could make it possible to deal with geometrical and other nuisance transforms through learning rather than hand-coding [63]. 2) It may be possible to increase the efficiency of gated feature learning on the scale of larger images by combining local receptive fields (factors) with mapping units that pool across factors defined at widely different image locations. To further reduce the connectivity, one could use prior knowledge about likely transformations.

Using multiplicative interactions can also be related to analogy making [41]. It can be argued that analogy making is at the heart of many cognitive capabilities [64]. An interesting question is in what way may an analogy-making module be a useful building block to solve tasks that are too difficult for the common deep architectures, which are composed of bipartite feature learning modules. Since gated feature learning and energy models can be trained with standard, even Hebbian-like, learning (see Section 2.2), analogy making does not require any uncommon or unusual machinery besides multiplicative interactions.

Squaring can be approximated using other nonlinearities (e.g., [13]). A possible research question is what type of approximations of computing squares or cross products may be advantageous computationally and/or more plausible biologically. Of course, squares could be approximated using a feed-forward network with sigmoid hidden unit activations [65]. However, the abundance of matching and correspondence tasks in vision does seem to provide an inductive bias in favor of genuine multiplicative interactions and/or squaring nonlinearities.

## ACKNOWLEDGMENTS

The author thanks Felix Bauer, Yoshua Bengio, Taco Cohen, Georgios Exarchakis, Geoffrey Hinton, Kishore Konda, Christoph von der Malsburg, Joshua Susskind and Christopher Zach for useful discussions and the reviewers for their useful suggestions. This work was supported in part by the German Federal Ministry of Education and Research (BMBF) in the project 01GQ0841 (BFNT Frankfurt).

## REFERENCES

- [1] G.F. Hinton, "A Parallel Computation That Assigns Canonical Object-Based Frames of Reference," *Proc. Seventh Int'l Joint Conf. Artificial Intelligence*, vol. 2, pp. 683-685, 1981.

- [2] C. von der Malsburg, "The Correlation Theory of Brain Function," *Models of Neural Networks II*, E. Domany, J.L. van Hemmen, and K. Schulten, eds., chapter 2, pp. 95-119, Springer-Verlag, 1994.
- [3] E. Adelson and J. Bergen, "Spatiotemporal Energy Models for the Perception of Motion," *J. Optical Soc. Am. A*, vol. 2, no. 2, pp. 284-299, 1985.
- [4] P. Smolensky, "Information Processing in Dynamical Systems: Foundations of Harmony Theory," *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, D.E. Rumelhart, J.L. McClelland, and C. PDP Research Group, eds., vol. 1, pp. 194-281, MIT Press, 1986.
- [5] A. Coates, H. Lee, and A.Y. Ng, "An Analysis of Single-Layer Networks in Unsupervised Feature Learning," *Proc. 14th Int'l Conf. Artificial Intelligence and Statistics*, 2011.
- [6] B. Olshausen and D. Field, "Emergence of Simple-Cell Receptive Field Properties by Learning a Sparse Code for Natural Images," *Nature*, vol. 381, no. 6583, pp. 607-609, June 1996.
- [7] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, "Extracting and Composing Robust Features with Denoising Autoencoders," *Proc. 25th Int'l Conf. Machine Learning*, 2008.
- [8] G.E. Hinton, "Training Products of Experts by Minimizing Contrastive Divergence," *Neural Computation*, vol. 14, no. 8, pp. 1771-1800, 2002.
- [9] G. Hinton, "A Practical Guide to Training Restricted Boltzmann Machines," technical report, 2010.
- [10] A. Hyvärinen, J. Hurri, and P.O. Hoyer, *Natural Image Statistics: A Probabilistic Approach to Early Computational Vision*. Springer Verlag, 2009.
- [11] G.E. Hinton, S. Osindero, and Y.-W. Teh, "A Fast Learning Algorithm for Deep Belief Nets," *Neural Computation*, vol. 18, pp. 1527-1554, July 2006.
- [12] N. Troje and H. Bülthoff, "Face Recognition under Varying Poses: The Role of Texture and Shape," *Vision Research*, vol. 36, no. 12, pp. 1761-1771, 1996.
- [13] C. Zetsche and U. Nuding, "Nonlinear and Higher-Order Approaches to the Encoding of Natural Scenes," *Network*, vol. 16, no. 2/3, pp. 191-221, 2005.
- [14] B. Olshausen, C. Cadieu, J. Culpepper, and D. Warland, "Bilinear Models of Natural Images," *Proc. SPIE Human Vision Electronic Imaging XII*, vol. 6492, 2007.
- [15] D. Grimes and R. Rao, "Bilinear Sparse Coding for Invariant Vision," *Neural Computation*, vol. 17, no. 1, pp. 47-73, 2005.
- [16] R. Memisevic and G. Hinton, "Unsupervised Learning of Image Transformations," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2007.
- [17] J. Tenenbaum and W. Freeman, "Separating Style and Content with Bilinear Models," *Neural Computation*, vol. 12, no. 6, pp. 1247-1283, 2000.
- [18] R. Memisevic, "Non-Linear Latent Factor Models for Revealing Structure in High-Dimensional Data," PhD dissertation, Univ. of Toronto, 2008.
- [19] R. Memisevic, "Gradient-Based Learning of Higher-Order Image Features," *Proc. IEEE Int'l Conf. Computer Vision*, 2011.
- [20] J. Susskind, R. Memisevic, G. Hinton, and M. Pollefeys, "Modeling the Joint Density of Two Images under a Variety of Transformations," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [21] M. Ranzato, A. Krizhevsky, and G.E. Hinton, "Factored 3-Way Restricted Boltzmann Machines for Modeling Natural Images," *Proc. 13th Int'l Conf. Artificial Intelligence and Statistics*, 2010.
- [22] Y. Karklin and M.S. Lewicki, "Is Early Vision Optimized for Extracting Higher-Order Dependencies?" *Proc. Advances in Neural Information Processing Systems 18*, 2006.
- [23] A. Hyvärinen and P. Hoyer, "Emergence of Phase- and Shift-Invariant Features by Decomposition of Natural Images into Independent Feature Subspaces," *Neural Computation*, vol. 12, pp. 1705-1720, July 2000.
- [24] I. Ohzawa, G.C. Deangelis, and R.D. Freeman, "Stereoscopic Depth Discrimination in the Visual Cortex: Neurons Ideally Suited as Disparity Detectors," *Science*, vol. 249, no. 4972, pp. 1037-1041, Aug. 1990.
- [25] T. Sanger, "Stereo Disparity Computation Using Gabor Filters," *Biological Cybernetics*, vol. 59, pp. 405-418, 1988, doi: 10.1007/BF00336114.
- [26] N. Qian, "Computing Stereo Disparity and Motion with Known Binocular Cell Properties," *Neural Computation*, vol. 6, pp. 390-404, May 1994.
- [27] S. Becker and G.E. Hinton, "A Self-Organizing Neural Network That Discovers Surfaces in Random-Dot Stereograms," *Nature*, vol. 355, pp. 161-163, 1992.
- [28] D. Fleet, H. Wagner, and D. Heeger, "Neural Encoding of Binocular Disparity: Energy Models, Position Shifts and Phase Shifts," *Vision Research*, vol. 36, no. 12, pp. 1839-1857, June 1996.
- [29] C.L. Giles and T. Maxwell, "Learning, Invariance, and Generalization in High-Order Neural Networks," *Applied Optics*, vol. 26, no. 23, pp. 4972-4978, Dec. 1987.
- [30] D.E. Rumelhart, G.E. Hinton, and J.L. McClelland, "A General Framework for Parallel Distributed Processing," *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, vol. 1, chapter 2, pp. 45-76, MIT Press, 1986.
- [31] P. Smolensky, "Tensor Product Variable Binding and the Representation of Symbolic Structures in Connectionist Systems," *Artificial Intelligence*, vol. 46, pp. 159-216, 1990.
- [32] T. Plate, "Holographic Reduced Representations: Convolution Algebra for Compositional Distributed Representations," *Proc. Int'l Joint Conf. Artificial Intelligence*, pp. 30-35, 1991.
- [33] T. Kohonen, "The Adaptive-Subspace SOM (ASSOM) and Its Use for the Implementation of Invariant Feature Detection," *Proc. Int'l Conf. Artificial Neural Networks*, vol. I, EC2, pp. 3-10, 1995.
- [34] A. Hyvärinen, P.O. Hoyer, and M. Inki, "Topographic ICA as a Model of Natural Image Statistics," *Proc. First IEEE Int'l Workshop Biologically Motivated Computer Vision*, S.-W. Lee, H.H. Blthoff, and T. Poggio, eds., pp. 535-544, 2000.
- [35] M. Welling, G.E. Hinton, and S. Osindero, "Learning Sparse Topographic Representations with Products of Student-t Distributions," *Proc. Advances in Neural Information Processing Systems*, 2002.
- [36] B. Olshausen, "Neural Routing Circuits for Forming Invariant Representations of Visual Objects," PhD dissertation, Computation and Neural Systems, California Inst. of Technology, 1994.
- [37] M.J. Wainwright and E.P. Simoncelli, "Scale Mixtures of Gaussians and the Statistics of Natural Images," *Proc. Advances in Neural Information Processing Systems*, vol. 12, pp. 855-861, 2000.
- [38] P. Hoyer and A. Hyvärinen, "A Multi-Layer Sparse Coding Network Learns Contour Coding from Natural Images," *Vision Research*, vol. 42, pp. 1593-1605, 2002.
- [39] M. Ranzato and G.E. Hinton, "Modeling Pixel Means and Covariances Using Factorized Third-Order Boltzmann Machines," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 2551-2558, 2010.
- [40] A. Courville, J. Bergstra, and Y. Bengio, "A Spike and Slab Restricted Boltzmann Machine," *Proc. Conf. Artificial Intelligence and Statistics*, 2011.
- [41] R. Memisevic and G.E. Hinton, "Learning to Represent Spatial Transformations with Factored Higher-Order Boltzmann Machines," *Neural Computation*, vol. 22, no. 6, pp. 1473-92, 2010.
- [42] J.L. Gallant, J. Braun, and D.C.V. Essen, "Selectivity for Polar, Hyperbolic, and Cartesian Gratings in Macaque Visual Cortex," *Science*, vol. 259, pp. 1001-1004, 1993.
- [43] F. Bauer, "Motion Analysis Using Local Multiplicative Interactions," master's thesis, Institut für Informatik, 2012.
- [44] Q. Le, W. Zou, S. Yeung, and A. Ng, "Learning Hierarchical Spatio-Temporal Features for Action Recognition with Independent Subspace Analysis," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [45] G. Taylor, R. Fergus, Y. LeCun, and C. Bregler, "Convolutional Learning of Spatio-Temporal Features," *Proc. European Conf. Computer Vision*, 2010.
- [46] R.A. Horn and C.R. Johnson, *Matrix Analysis*. Cambridge Univ. Press, 1990.
- [47] R.M. Gray, "Toeplitz and Circulant Matrices: A Review," *Comm. Information Theory*, vol. 2, pp. 155-239, Aug. 2005.
- [48] M. Bethge, S. Gerwinn, and J. Macke, "Unsupervised Learning of a Steerable Basis for Invariant Image Representations," *Proc. SPIE Human Vision and Electronic Imaging XII*, pp. 1-12, Feb. 2007.
- [49] R. Memisevic, "On Multi-View Feature Learning," *Proc. Int'l Conf. Machine Learning*, June 2012.
- [50] R. Memisevic, C. Zach, G. Hinton, and M. Pollefeys, "Gated Softmax Classification," *Proc. Advances in Neural Information Processing Systems 22*, 2010.
- [51] W.Y. Zou, S. Zhu, A.Y. Ng, and K. Yu, "Deep Learning of Invariant Features via Tracked Video Sequences," *Proc. Advances in Neural Information Processing Systems 25*, 2012.

- [52] C.F. Cadieu and B.A. Olshausen, "Learning Intermediate-Level Representations of Form and Motion from Natural Movies," *Neural Computation*, vol. 24, no. 4, pp. 827-866, Dec. 2011.
- [53] D.A. Ross, S. Osindero, and R.S. Zemel, "Combining Discriminative Features to Infer Complex Trajectories," *Proc. 23rd Int'l Conf. Machine Learning*, pp. 761-768, 2006.
- [54] M. Denil, L. Bazzani, H. Larochelle, and N. de Freitas, "Learning Where to Attend with Deep Architectures for Image Tracking," *Neural Computation*, vol. 24, no. 8, pp. 2151-2184, [http://dx.doi.org/10.1162/NECO\\_a\\_00312](http://dx.doi.org/10.1162/NECO_a_00312), Aug. 2012.
- [55] H. Larochelle and G. Hinton, "Learning to Combine Foveal Glimpses with a Third-Order Boltzmann Machine," *Proc. Advances in Neural Information Processing Systems 23*, pp. 1243-1251, 2010.
- [56] G. Taylor and G. Hinton, "Factored Conditional Restricted Boltzmann Machines for Modeling Motion Style," *Proc. 26th Int'l Conf. Machine Learning*, L. Bottou and M. Littman, eds., pp. 1025-1032, June 2009.
- [57] G. Taylor, L. Sigal, D. Fleet, and G. Hinton, "Dynamic Binary Latent Variable Models for 3D Pose Tracking," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2010.
- [58] Y. Tang, R. Salakhutdinov, and G. Hinton, "Robust Boltzmann Machines for Recognition and Denoising," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2012.
- [59] I. Sutskever, J. Martens, and G. Hinton, "Generating Text with Recurrent Neural Networks," *Proc. 28th Int'l Conf. Machine Learning*, L. Getoor and T. Scheffer, eds., pp. 1017-1024, June 2011.
- [60] K.A. Archie and B.W. Mel, "A Model for Intradendritic Computation of Binocular Disparity," *Nature Neuroscience*, vol. 3, no. 1, pp. 54-63, Jan. 2000.
- [61] A. Krizhevsky, I. Sutskever, and G.E. Hinton, "Imagenet Classification with Deep Convolutional Neural Networks," *Proc. Advances in Neural Information Processing Systems 25*, 2012.
- [62] R.I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, second ed. Cambridge Univ. Press, 2004.
- [63] K. Mikolajczyk and C. Schmid, "A Performance Evaluation of Local Descriptors," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 10, pp. 1615-1630, Oct. 2005.
- [64] D.R. Hofstadter, "The Copycat Project: An Experiment in Nondeterminism and Creative Analogies," Technical Report AI Memo No. 755, MIT, 1984.
- [65] K. Funahashi, "On the Approximate Realization of Continuous Mappings by Neural Networks," *Neural Networks*, vol. 2, pp. 183-192, May 1989.



**Roland Memisevic** received the PhD degree in computer science from the University of Toronto, Canada, in 2008. Subsequently, he held positions as a research scientist at PNYLab LLC in Princeton and as a postdoctoral fellow at the University of Toronto and at ETH Zurich, Switzerland. From 2011 to 2012, he was on faculty at the Department of Computer Science at the University of Frankfurt, Germany. In 2012, he joined the Department of Computer Science

and Operations Research at the University of Montreal, Canada, as an assistant professor of computer science. His research interests are in machine learning and computer vision, in particular unsupervised learning and feature learning. His scientific contributions include relational and higher-order feature learning models, approaches to learning motion and transformation patterns from data, and approaches to learning invariant recognition. He presented his work at conferences such as NIPS, ICCV, ICML, CVPR, AAAI, and in journals such as the *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *Neural Networks*, and *Neural Computation*. He has served as a program committee member or reviewer for most of these and other conferences and journals in machine learning and computer vision.

► **For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).**