

# View-based Models of 3D Object Recognition: Invariance to Imaging Transformations

Thomas Vetter,<sup>1</sup> Anya Hurlbert,<sup>2</sup> and Tomaso Poggio<sup>3</sup>

<sup>1</sup> Max-Planck Institut für Biologische Kybernetik, 72076 Tübingen, Germany, <sup>2</sup> Physiological Sciences, University of Newcastle-upon-Tyne, Newcastle-upon-Tyne NE2 4HH, United Kingdom, and <sup>3</sup> Center for Computational and Biological Learning, and Department of Brain Sciences, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139

This report describes the main features of a view-based model of object recognition. The model does not attempt to account for specific cortical structures; it tries to capture general properties to be expected in a biological architecture for object recognition. The basic module is a regularization network (RBF-like; see Poggio and Girosi, 1989; Poggio, 1990) in which each of the hidden units is broadly tuned to a specific view of the object to be recognized. The network output, which may be largely view independent, is first described in terms of some simple simulations. The following refinements and details of the basic module are then discussed: (1) some of the units may represent only components of views of the object—the optimal stimulus for the unit, its “center,” is effectively a complex feature; (2) the units’ properties are consistent with the usual description of cortical neurons as tuned to multidimensional optimal stimuli and may be realized in terms of plausible biophysical mechanisms; (3) in learning to recognize new objects, preexisting centers may be used and modified, but also new centers may be created incrementally so as to provide maximal view invariance; (4) modules are part of a hierarchical structure—the output of a network may be used as one of the inputs to another, in this way synthesizing increasingly complex features and templates; (5) in several recognition tasks, in particular at the basic level, a single center using view-invariant features may be sufficient.

Modules of this type can deal with recognition of specific objects, for instance, a specific face under various transformations such as those due to viewpoint and illumination, provided that a sufficient number of example views of the specific object are available. An architecture for 3D object recognition, however, must cope—to some extent—even when only a single model view is given. The main contribution of this report is an outline of a recognition architecture that deals with objects of a *nice* class undergoing a broad spectrum of transformations—due to illumination, pose, expression, and so on—by exploiting prototypical examples. A *nice* class of objects is a set of objects with sufficiently similar transformation properties under specific transformations, such as viewpoint transformations. For *nice* object classes, we discuss two possibilities: (1) class-specific transformations are to be applied to a single model image to generate additional *virtual* example views, thus allowing some degree of generalization beyond what a single model view could otherwise provide; (2) class-specific, view-invariant features are learned from examples of the class and used with the novel model image, without an explicit generation of virtual examples.

In the past three years we have been developing systems for 3D object recognition that we labeled *view based* (or memory based; see Poggio and Hurlbert, 1993) since they require units tuned to views of specific objects or object classes.<sup>1</sup> Our work has led to artificial systems for solving toy problems such as the recognition of paperclips as in Figure 3 (Poggio and Edelman, 1990; Brunelli and Poggio, 1991), as well as more real problems such as the recognition of frontal faces (Brunelli and Poggio, 1993; Gilbert and Yang, 1993) and the recognition of faces in arbitrary pose (Beymer, 1993). We have discussed how this approach may capture key aspects of the cortical architecture for 3D object recognition (Poggio, 1990; Poggio and Hurlbert, 1993), we have tested successfully with

psychophysical experiments some of the predictions of the model (Bülthoff and Edelman, 1992; Edelman and Bülthoff, 1992; Schyns and Bülthoff, 1993), and recently we have gathered preliminary evidence that this class of models is consistent with both psychophysics and physiology [specifically, of inferotemporal (IT) cortex] in alert monkeys trained to recognize specific 3D paperclips (Logothetis and Pauls, 1995).

This report is a short summary of some of our theoretical work; it describes work in progress and it refers to other reports that treat in more detail several aspects of this class of models. Some of these ideas are similar to those of Perrett et al. (1989), though they were developed independently; they originate instead from applying regularization networks to the problem of visual recognition and noticing an intriguing similarity between the hidden units of the model and the tuning properties of cortical cells. The main problem this report addresses is that of how a visual system can learn to recognize an object after exposure to only a *single* view, when the object may newly appear in different views corresponding to a broad spectrum of image transformations. Our main novel contribution is the outline of an architecture capable of achieving invariant recognition for a single model view, by exploiting transformations learned from a set of prototype objects of the same class.

We will first describe the basic view-based module and illustrate it with a simple simulation. We will then discuss a few of the refinements that are necessary to make it biologically plausible. The next section will sketch a recognition architecture for achieving invariant recognition. In particular, we will describe how it may cope with the problem of recognizing a specific object of a certain class from a single model view. Finally, we will describe an hypothetical, secondary route to recognition—a *visualization* route—in which (1) class-specific RBF-like modules estimate parameters of the input image, such as illumination, pose, and expression; (2) other modules provide the appropriate transformation from prototypes and synthesize a “normalized” view from the input view; and (3) the normalized input view is compared with the model view in memory. Thus, analysis and synthesis networks may be used to close the loop in the recognition process by generating the “neural” imagery corresponding to a certain interpretation and eventually comparing it to the input image. In the last section we will outline some of the critical predictions of this class of biological models and discuss some of the existing data.

## The Basic Recognition Module

Figure 1 shows our basic module for object recognition. As Poggio and Hurlbert (1993) have argued, it is representative of a broad class of memory-based modules (MBMs). Classification or identification of a visual stimulus is accomplished by a network of units. Each unit is broadly tuned to a particular view of the object. We refer to this optimal view as the center of the unit. One can think of it as a template to which the input is compared. The unit is maximally excited when

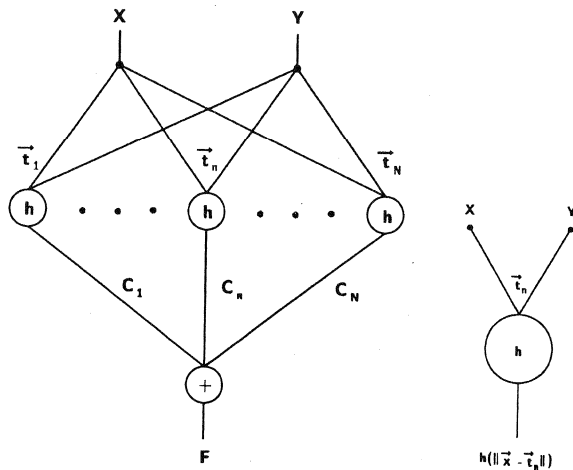


Figure 1. An RBF network for the approximation of 2D functions (left) and its basic "hidden" unit (right).  $x$  and  $y$  are components of the input vector that is compared via the RBF  $h$  at each center  $t$ . Outputs of the RBFs are weighted by the  $c_i$  and summed to yield the function  $F$  evaluated at the input vector.  $N$  is the total number of centers.

the stimulus exactly matches its template but also responds proportionately less to similar stimuli. The weighted sum of activities of all the units represents the output of the network.

Here we consider as an example of such a structure an RBF network that we originally used as a learning network (Poggio and Girosi, 1989) for object recognition while discovering that it was biologically appealing (Poggio and Girosi, 1989; Poggio, 1990; Poggio and Edelman, 1990; Poggio and Hurlbert, 1993) and representative of a much broader class of network architectures (Girosi et al., 1993).

### RBF Networks

Let us review briefly RBF networks. RBF networks are approximation schemes that can be written as (see Fig. 1; Poggio, 1990; Poggio and Girosi, 1990b)

$$f(\mathbf{x}) = \sum_{i=1}^N c_i h(\|\mathbf{x} - \mathbf{t}_i\|) + p(\mathbf{x}). \quad (1)$$

The Gaussian case,  $h(\|\mathbf{x} - \mathbf{t}\|) = \exp(-\|\mathbf{x} - \mathbf{t}\|^2/2\sigma^2)$ , is especially interesting: (1) each "unit" computes the distance  $\|\mathbf{x} - \mathbf{t}\|$  of the input vector  $\mathbf{x}$  from its center  $\mathbf{t}$  and (2) applies the function  $h$  to the distance value; that is, it computes the function  $h(\|\mathbf{x} - \mathbf{t}\|)$ ; (3) in the limiting case of  $h$  being a very narrow Gaussian, the network becomes a *look-up* table; (4) centers are like *templates*.

The simplest recognition scheme we consider is the network suggested by Poggio and Edelman (1990) to solve the specific problem of recognizing a particular 3D object from novel views. This is a problem at the *subordinate* level of recognition; it assumes that the object has already been classified on the *basic* level but must be discriminated from other members of its class. In the RBF version of the network, each center stores a sample view of object, and acts as a unit with a Gaussian-like recognition field around that view. The unit performs an operation that could be described as "blurred" template matching. At the output of the network the activities of the various units are combined with appropriate weights, found during the learning stage.

Consider how the network "learns" to recognize views of the object shown in Figure 3. In this example the inputs of the network are the  $x, y$ -positions of the vertices of the object images and four training views are used. After training, the network consists of four units, each one tuned to one of the four views as in Figure 2. The weights of the output connec-

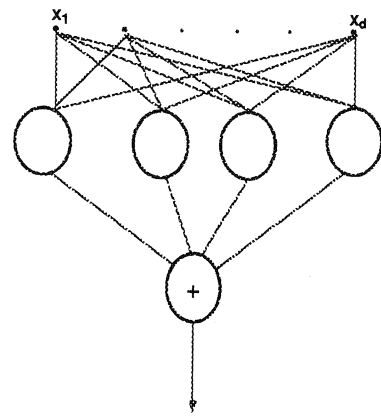


Figure 2. An RBF network with four units each tuned to one of the four training views shown in Figure 3. The tuning curve of each unit is also shown in Figure 3. The hidden units are view dependent but selective relative to distractors of the same type. The output unit is in this case view invariant for rotations of the object around the vertical axis, examples of which are represented by the centers of the four hidden units.

tions are determined by minimizing misclassification errors on the four views and using as negative examples views of other similar objects ("distractors").

Figure 3 shows the tuning of the four units for images of the "correct" object. The tuning is broad and centered on the training view. Somewhat surprisingly, the tuning is also very selective: the dotted line shows the average response of each unit to 300 similar distractors (paperclips generated by the

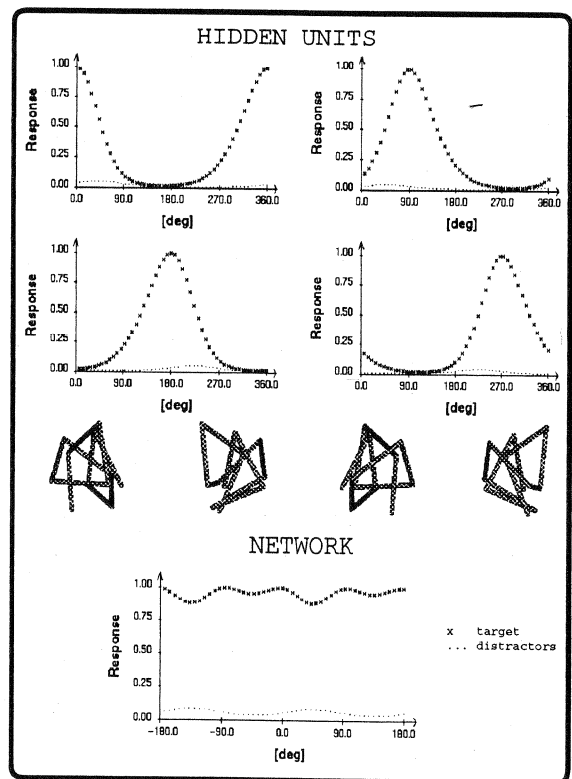


Figure 3. Tuning of each of the four hidden units of the network of the previous figure for images of the "correct" 3D objects. The tuning is broad and selective: the dotted lines indicate the average response to 300 distractor objects of the same type. The bottom graph shows the tuning of the output of the network after learning (i.e., after the computation of the weights  $c_i$ ): it is view invariant and object specific. Again, the dotted curve indicates the average response of the network to the same 300 distractors.

same mechanisms as the target; for further details about the generation of paperclips, see Edelman and Bühlhoff, 1992). Even the maximum response to the best distractor is in this case always less than the response to the optimal view. The output of the network, being a linear combination of the activities of the four units, is essentially view invariant and still very selective. Notice that each center is the *conjunction* of all the features represented: the Gaussian can in fact be decomposed into the product of 1D Gaussians, one for each input component. The activity of the unit measures the global similarity of the input vector to the center: for optimal tuning all features must be close to the optimum value. Even the mismatch of a single component of the template may set to zero the activity of the unit. Thus, the rough rule implemented by a view-tuned unit is the conjunction of a set of predicates, one for each input feature, measuring the match with the template. On the other hand, the output of the network is performing an operation more similar to the "OR" of the output of the units. Even if the output unit may have a sigmoidal nonlinearity (see Poggio and Girosi, 1990), its output does not need to be zero when one or more of the hidden units are inactive, provided there is sufficient activity in the remaining ones. In general, one expects a nonlinear decision-like mechanism (see Logothetis and Pauls, 1995) operating on the output signal of the network.

This example is clearly a caricature of a view-based recognition module but it helps to illustrate the main points of the argument. Despite its gross oversimplification, it manages to capture some of the basic psychophysical and physiological findings, in particular the existence of view-tuned and view-invariant units and the shape of psychophysically measured recognition fields. In the next section we will list a number of ways in which the network can be made more plausible.

### Toward More Biological Recognition Modules

The simple model proposed in the previous section contains view-centered hidden units.<sup>2</sup> More plausible versions allow for the centers and corresponding hidden units to be view invariant if the task requires. In a biological implementation of the network, we in fact expect to find a full spectrum of hidden unit properties, from view centered to view invariant. View-centered units are more likely in the case of subordinate level recognition of unfamiliar *not nice* objects (for the definition of a *nice* class, see below); view-invariant units would appear for the basic level recognition of familiar objects. We will now make a number of related observations, some of which can be found in Poggio and Hurlbert (1993), which point to necessary refinements of the model if it is to be biologically plausible.

(1) In the previous example each unit has a center that is effectively a full training view. It is much more reasonable to assume that most units in a recognition network should be tuned to *components* of the image, that is, to conjunctions of some of the elementary features but not *all* of them. This should allow for sufficient selectivity (the above network performs better than humans) and provide for significant robustness to occlusions and noise (see Poggio and Hurlbert, 1993). This means that the "AND" of a high-dimensional conjunction can be replaced by the "OR" of its components—a face may be recognized by its eyebrows alone, or a mug by its color. Notice that the disjunction (corresponding to the weighted combination of the hidden units) of conjunctions of a small number of features may be sufficient (each conjunction is implemented by a Gaussian center that can be written as the product of 1D Gaussians). To recognize an object, we may use not only templates (i.e., centers in RBF terminology) comprising all its features, but also, and in some

cases solely, subtemplates, comprising subsets of features (which themselves constitute "complex" features). This is similar in spirit to the technique of supplementing whole-face templates with several smaller templates in the Brunelli-Poggio work on frontal face recognition (see also Beymer, 1993).

(2) The units tuned to complex features mentioned above are similar to IT cells described by Fujita and Tanaka (1992) and could be constructed in a hierarchical way from the output of simpler RBF-like networks. They may avoid the correspondence problem, provided that the system has built-in invariance to image-plane transformations, such as translation, rotation, and scaling. Thus, cells tuned to complex features are constructed from a hierarchy of simpler cells tuned to incrementally larger conjunctions of elementary features. This idea—common among physiologists (see Perrett and Oram, 1993; Tanaka, 1993)—can immediately be formalized in terms of Gaussian radial basis functions, since a multidimensional Gaussian function can be decomposed into the product of lower dimensional Gaussians (Marr and Poggio, 1976; Ballard, 1986; Poggio and Girosi, 1990; Mel, 1992).

(3) The features used in the example of Figure 3 ( $x, y$ -coordinates of paperclip vertices) are biologically implausible. We have also used other more natural features such as orientation of lines. An attractive feature of this module is its recursive nature: detection and localization of a line of a certain orientation, say, can be thought of as being performed by a similar network with centers being units tuned to different examples of the desired line type. An eye detector can localize an eye by storing in its units templates of several eyes and using as inputs more elementary features such as lines and blobs. A face recognition network may use units tuned to specific templates of eyes and nose and so on. A homogeneous, recursive approach of this type in which not only object recognition is view based but also feature localization is view based has been successfully used in the Beymer-Poggio face recognizer (see Beymer, 1993). Both feature detection and face recognition depend on the use of several templates, the "examples."

(4) In this perspective there are probably elementary features such as blobs and oriented lines and center-surround patterns, but there is then a continuum of increasingly complex features corresponding to centers that are conjunctions of more elementary ones. In this sense a center is simply a more complex feature than its inputs and may in turn be the input to another network with even more complex center-features.

(5) The RBF network described in the previous sections is the simplest version of a more general scheme (hyperbasis functions, HBF) given by

$$f^*(\mathbf{x}) = \sum_{\alpha=1}^n c_{\alpha} G(\|\mathbf{x} - \mathbf{t}_{\alpha}\|_{\mathbf{W}}^2) + p(\mathbf{x}), \quad (2)$$

where the centers  $\mathbf{t}_{\alpha}$  and coefficients  $c_{\alpha}$  are unknown, and are in general fewer in number than the data points ( $n \leq N$ ). The norm is a *weighted norm*:

$$\|\mathbf{x} - \mathbf{t}_{\alpha}\|_{\mathbf{W}}^2 = (\mathbf{x} - \mathbf{t}_{\alpha})^T \mathbf{W}^T \mathbf{W} (\mathbf{x} - \mathbf{t}_{\alpha}), \quad (3)$$

where  $\mathbf{W}$  is an unknown square matrix and the superscript  $T$  indicates the transpose. In the simple case of diagonal  $\mathbf{W}$  the diagonal elements  $w_i$  assign a specific weight to each input coordinate, determining in fact the units of measure and the importance of each feature (the matrix  $\mathbf{W}$  is especially important in cases in which the input features are of a different type and their relative importance is unknown). During learning, not only the coefficients  $c$  but also the centers  $\mathbf{t}_{\alpha}$ , and the elements of  $\mathbf{W}$  are updated. Whereas the RBF technique is similar to and similarly limited as template matching, HBF networks perform a generalization of template matching

in an appropriately linearly transformed space, with the appropriate metric. As a consequence, HBF networks may "find" view-invariant features when they exist (Bricolo, unpublished observation). There are close connections between HBF networks, multilayer perceptrons, and regularization (see Girosi et al., 1993).

(6) It is also plausible that some of the center-features are "innate," synthesized by evolution or by early experience of the individual or more likely by both. We assume that the adult system has at its disposal a *vocabulary* of simple as well as increasingly more complex center-features. Other centers are synthesized on demand in a task-dependent way. This may happen in the following way. Assume that a network such as the one in Figure 2 has to learn to recognize a new object. It may attempt to do so by using some of the outputs in the pool of existing networks as its inputs. At first no new centers are allocated and only the linear part of the network is used, corresponding to the term  $p(\mathbf{x})$  in Equation 1 and to direct connections between inputs and output (not shown in Fig. 2). This of course is similar to a simple OR of the input features. Learning may be successful in which case only some of the inputs will have a nonzero weight. If learning is not successful—or sufficiently weak—a new center of minimal dimension may be allocated to mimic a component of one of the training views. New centers of increasing dimensionality—comprising subsets of components, up to the full view—are added while old centers are continually pruned until the performance is satisfactory. Centers of dimension 2 effectively detect conjunctions of pairs of input features (see also Mel, 1992). It is not difficult to imagine learning strategies of this type that would select automatically centers, that is, complex features, that are as view invariant as possible (this can be achieved by modifying the associated parameters  $c$  and/or  $w$  in the  $\mathbf{W}$  matrix). Such features may be global—such as color—but we expect that they will be mostly local and perhaps underlie recognition of geon-like components (see Biederman, 1987; Edelman, 1991). View-invariant features may be used in basic-level more than in subordinate-level recognition tasks.

(7) One essential aspect of the simplest (RBF) version of the model is that it contains units that are viewer centered, not object centered. This aspect is independent of whether the model is 2D or 3D, a dichotomy that is not relevant here. Each center may consist of a set of features that may mix 2D with 3D information, by including shading, occlusion, or binocular disparity information, for example. The features that depend on the image geometry will necessarily be viewpoint dependent, but features such as color may be viewpoint independent. As we mentioned earlier, in situations in which view-invariant features exist (for basic- as well as for subordinate-level recognition) centers may actually be view independent.

(8) The network described here is used as a classifier that performs *identification*, or subordinate-level recognition: matching the face to a stored memory, and thereby labeling it. A similar network with a different set of centers could perform also basic-level recognition: distinguishing objects that are faces from those that are not.

#### Virtual Views and Invariance to Image Transformations: Toward a Recognition Architecture

In the example given above, the network learns to recognize a particular 3D object from novel views and thereby achieves one crucial aim in object recognition: viewpoint invariance. But recognition does not involve solely or simply the problem of recognizing objects in hitherto unseen poses. Hence, as Poggio and Hurlbert (1993) emphasize, the cortical architecture for recognition cannot consist simply of a collection of

the modules of Figures 1 and 3, one for each recognizable object. The architecture must be more complex than that cartoon, because recognition must be achieved over a variety of image transformations, not just those due to changes in viewpoint, but also those due to translation, rotation, and scaling of the object in the image plane, as well as non-image-plane transformations, such as those due to varying illumination. In addition, the cortex must also recognize objects at the basic as well as subordinate level.

In the network described above, viewpoint invariance is achieved by exploiting several sample views of the specific object. This strategy might work to obtain invariance under other types of transformations also, provided sufficient examples of the object under sample transformations are available. But suppose that example views are not available. Suppose that the visual system must learn to recognize a given object under varying illumination or viewpoint, starting with only a single example view. This is the problem that we will focus on in the next few sections, that of subordinate level recognition under non-image-plane transformations, given only a single model view.

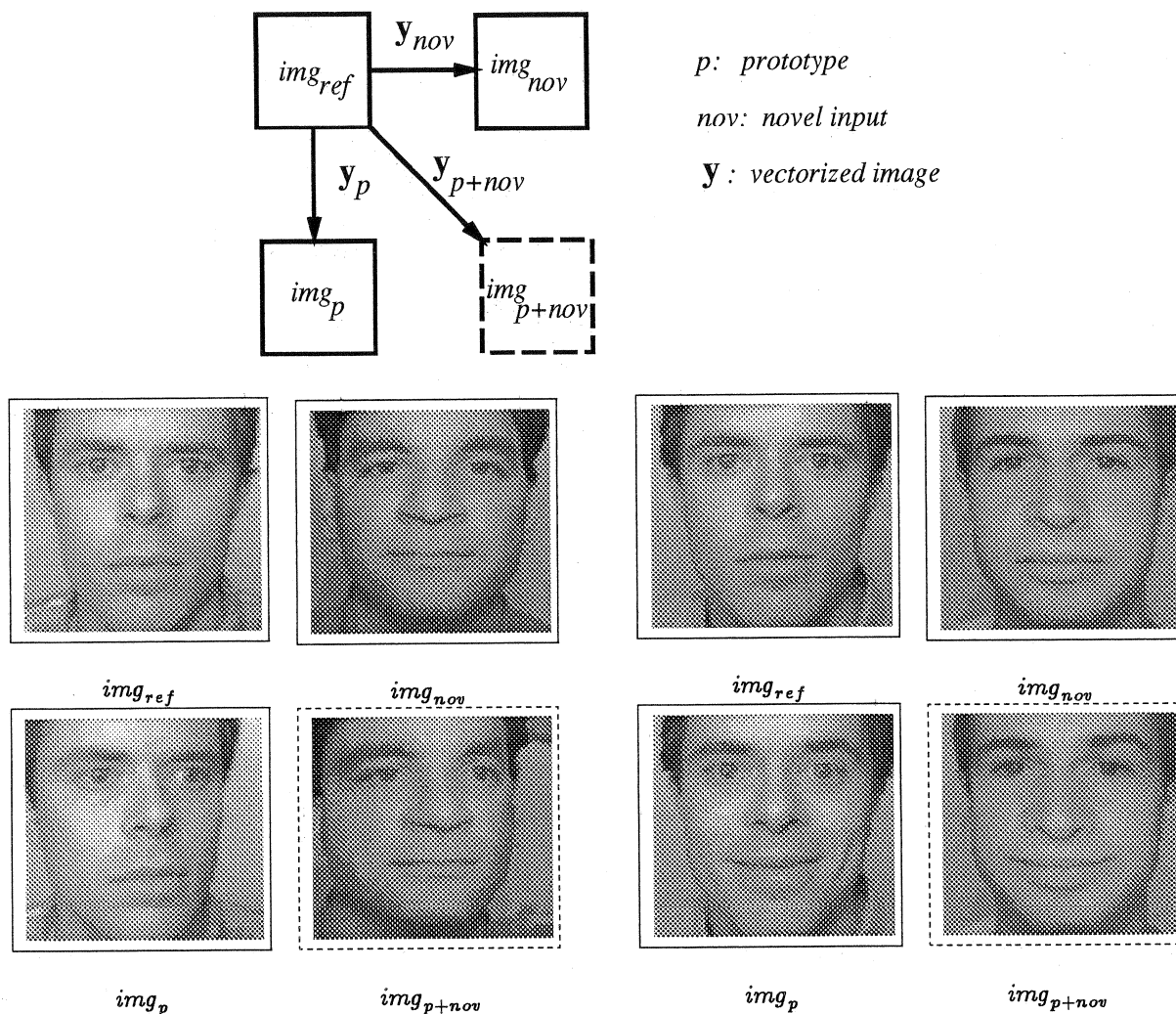
Probably the most natural solution is for the system to exploit certain invariant features, learned from examples of objects of the same class. These features could supplement the information contained in the single model view. Here we will put forward an alternative scheme that, although possibly equivalent at a computational level, may have a very different implementation. Our proposal is that when sample images of the specific object under the relevant transformations are not available, the system may generate virtual views of that object, using *image-based* transformations that are characteristic of the corresponding class of objects (Poggio and Vetter, 1992). We propose that the system learns these transformations from prototypical example views of other objects of the same class, with no need for 3D models. The idea is simple but it is not obviously clear that it will work. Figure 4 provides a plausibility argument.

The problem of achieving invariance to image plane transformations such as translation, rotation, and scaling, given only one model view, is also difficult, particularly in terms of biologically plausible implementations. But given a single model view, it is certainly possible to generate virtual examples for appropriate image-plane translations, scalings, and rotations *without specific knowledge about the object*. This is not the case for the non-image-plane transformations we will consider here, caused by, for example, changes in viewpoint, illumination, facial expression, or physical attitude of a flexible or articulated object such as a body.

Within the *virtual views* theory, there are two extreme ways in which virtual views may be used to ensure invariance under non-image-plane transformations. The first one is to precompute all possible "virtual" views of the object or the object class under the desired group of transformations and to use them to train a classifier network such as the one of Figure 1. The second approach—equivalent from the point of view of information processing—is instead to apply all the relevant transformations to the input image and to attempt to match the transformed image to the data base, which under our starting assumption, may contain only one view per object. These two general strategies may exist in several different variations and can also be mixed in various ways.

#### An Example

Consider as an example of the general recognition strategy we propose the following architecture for biological face recognition based on our own work on artificial face recognition systems (Beymer, 1993; Brunelli and Poggio, 1993; see also Gilbert and Yang, 1993).



**Figure 4.** A face transformation is "learned" from a prototypical example transformation. Here, face rotation and smiling transformations are represented by prototypes,  $y_p$ .  $y_p$  is mapped onto the new face image  $img_{nov}$ . The virtual image  $img_{p+nov}$  is synthesized by the system. In a biological implementation cell activities instead of gray levels would be the inputs and the outputs of the transformation. From Beymer et al. (1993).

First, the face has to be localized within the image and segregated from other objects. This stage might be template based, and may be equivalent to the use of a network like that in Figure 3, with units tuned to the various low-resolution images a face may produce. From the biological point of view, the network might be realized by the use of low-resolution face detection cells at each location in the visual field (with each location examined at a resolution dictated by the cortical map, in which the fovea of course dominates), or by connections from each location in, say, V1 to "centered" templates (or the equivalent networks) in IT, or by a routing mechanism to achieve the same result with fewer connections (see Olshausen et al., 1992). Of course, the detection may be based on disjunction of face components rather than on their conjunction in a full-face template.

The second step in our face recognizer is to normalize the image with respect to translation, scale, and image rotation. This is achieved by finding two anchor points, such as the eyes, again with a template-based strategy, equivalent to a network of the type of Figure 1 in which the centers are many templates of eyes of different types in different poses and expressions. A similar strategy may be followed by biological systems for both faces and other classes of objects. The ex-

istence of two stages would suggest that there are modules dedicated to detect certain classes of complex features—such as eyes—and other modules that use the result to normalize the image appropriately. Again, there could be eye detection networks at each location in the visual field or a routing of relevant parts of the image—selected through segmentation operations—to a central representation in IT.

The third step in our face recognizer is to match the localized, normalized face to a data base of individual faces while at the same time providing for invariance of view, expression, and illumination. If the data base contains several views of each particular face, the system may simply compare the normalized image to each item there (Beymer, 1993); this is equivalent to classifying the image using the network of Figure 1, one for each person. But if the data base contains only a single model view for each face, which is the problem we consider here, virtual examples of the face may be generated using transformations—to other poses and expressions—learned from examples of other faces (see Poggio and Brunelli, 1992; Poggio and Vetter, 1992; Beymer et al., 1993). Then the same approach as for a multiexample data base may be followed, but in this case most of the centers will correspond to "virtual examples."



### **Transformations and Virtual Examples**

In summary, our proposal is to achieve invariance to non-image-plane transformations by using a sufficient number of views of the specific objects for various transformation parameters. If real views are available they should be used directly; if not, virtual views can be generated from the real one(s) using image-based transformations learned from example views of objects of the same class.

#### *Transformation Networks*

How can we learn class-specific transformations from prototypical examples? There are several simple technical solutions to this problem, as discussed by Poggio (1991), Poggio and Brunelli (1992), and Poggio and Vetter (1992). The proposed schemes can "learn" approximate 3D geometry and underlying physics for a *sufficiently restricted* class of objects—a *nice class*.<sup>3</sup> We define informally here *nice classes* of objects as sets of objects with sufficiently similar transformation properties. A class of object is *nice* with respect to one or more transformations. Faces are a nice class under viewpoint transformations because they typically have a similar 3D structure. The paperclip objects used by Poggio and Edelman (1990) and Bülthoff and Edelman (1992; Bülthoff et al., 1994 this issue) and by Logothetis and Pauls (1995) are *not nice* under viewpoint transformation because their global 3D structures are different from each other. Poggio and Vetter (1992) describe a special set of nice classes of objects—"linear classes." For linear classes, linear networks can learn appropriate transformations from a set of prototypical examples. Figure 4 shows how Beymer et al. (1993) used the even simpler technique (linear additive) of Poggio and Brunelli (1992) for learning transformations due to face rotation and change of expression.

In any case, a sufficient number of prototype transformations—which may involve shape, color, texture, shading, and other image attributes by using the appropriate features in the vectorized representation of images—should allow the generation of more than one virtual view from a single "real" view. The resulting set of *virtual* examples can then be used to train a classification network. The argument so far is purely on the computational level and is supported only by preliminary and partial experiments. It is totally unclear at this point how IT cortex may use similar strategies based on learning class-specific prototypical transformations. The alternative model in which virtual examples are not explicitly generated and instead view-invariant features are learned is also attractive. Since networks such as multilayer perceptrons and HBF networks may "find" some view-invariant features, the two approaches may actually be used simultaneously.

#### **An Alternative Visualization Route?**

As we hinted earlier, an alternative implementation of the same approach to invariant recognition from a single model view is to transform the (normalized) input image using the learned transformations and compare each one of the resulting virtual views to the available real views (in this case only one per specific object). As pointed out by Ullman (1991), the cortex may perform the required search by generating simultaneously transformations of both the input image and the model views until a match is found.

The number of transformations to be tested may be reduced by first estimating the approximate pose and expression parameters of the input image. The estimate may be provided by an RBF-like network of the "analysis" type in which the centers are generic face prototypes (or face parts) spanning different poses, expressions, and possibly illuminations.<sup>4</sup> They can be used if trained appropriately to do the *analysis* task of estimating state parameters associated with the image

of the object such as its pose in space, its expression (if a face), its illumination, and so on (see Poggio and Edelman, 1990; Beymer et al., 1993).

The corresponding transformation will then be performed by networks (linear or of a more general type).<sup>5</sup> Analysis-type networks may help reduce dramatically the number of transformations to be tried before successful recognition is achieved. A particular version of the idea is the following.

Assume that the data base consists of single views of different, say, faces in a "zero" pose. Then in the visualization route the analysis network provides an estimate of "pose" parameters; a synthesis network (Librande, 1992; Poggio and Brunelli, 1992; Beymer et al., 1993) generates the corresponding view of a prototype; the transformation from the latter prototype view to the reference view of the prototype is computed and applied to the input array to obtain its "zero" view; finally, this corrected input view is compared with the data base of single views. Of course, the inverse transformation could be applied to each of the views in the data base, instead of applying the direct transformation to the input image. We prefer the former strategy because of computational considerations but mixtures of both strategies may be suitable in certain situations.

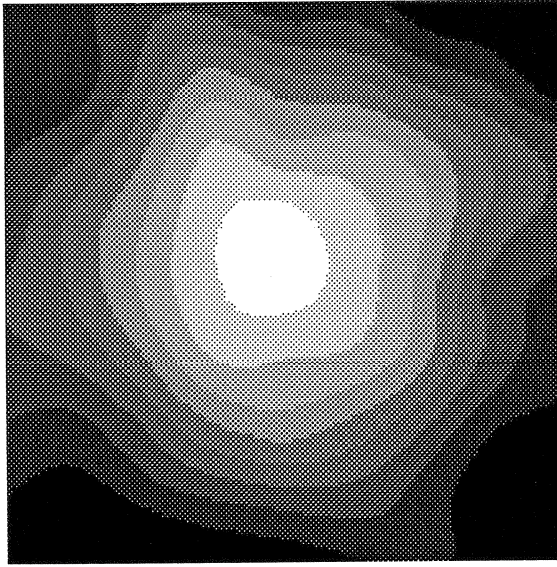
This estimation-transformation route (which may also be called analysis-synthesis) leads to an approach to recognition in which parameters are estimated from the input image, and then used to "undo" the deformation of the input image and "visualize" the result, which is then compared to the data base of reference views. A "visualization" approach of this type can be naturally embedded in an iterative or feedback scheme in which discrepancies between the visualized estimate and the input image drive further cycles of analysis-synthesis and comparison (see Mumford, 1992). It may also be relevant in explaining a role in mental "imagery" of the neurons in IT (see Sakai and Miyashita, 1991). A few remarks follow.

(1) Transformation parameters may be estimated from images of objects of a class; some degree of view invariance may therefore be achievable for new objects of a known class (such as faces or bilaterally symmetric objects; see Poggio and Vetter, 1992). This should be impossible for unique objects for which prior class knowledge may not be used (such as the paperclip objects; Bülthoff and Edelman, 1992).

(2) From the computational point of view it is possible that a "coarse" 3D model—rather like a marionette—could be used successfully to compute various transformations typical for a certain class of objects (such as faces) to control 2D representations of the type described earlier for each specific object. Biologically, this coarse 3D model may be implemented in terms of learned transformations characteristic for the class.

(3) We believe that the classification approach—the one summarized by Figures 1 and 3, as opposed to the visualization approach—is the main route to recognition, which should be used with real example views when a sufficient number of training views is available. Notice that this approach is memory based and in the extreme case of many training views should be very similar to a lookup table. When only one or very few views of the specific object are available, the classification approach may still suffice, if either (1) view-invariant features are discovered and then used or (2) virtual examples generated by the transformation approach are exploited. But this is possible only for objects belonging to a familiar class (such as faces). The analysis-synthesis route may be an additional, secondary strategy to deal with only one or very few real model views.<sup>6</sup>

(4) We have assumed here a supervised learning framework. Unsupervised learning may not be of real biological interest because various natural cues (object constancy, sen-



**Figure 5.** The generalization field associated with a single training view. Whereas it is easy to distinguish between, say, tubular and amoeba-like 3D objects, irrespective of their orientation, the recognition error rate for specific objects within each of those two categories increases sharply with misorientation relative to the familiar view. This figure shows that the error rate for amoeba-like objects, previously seen from a single attitude, is viewpoint dependent. Means of error rates of six subjects and six different objects are plotted versus rotation in depth around two orthogonal axes (Bülthoff et al., 1991; Edelman and Bülthoff, 1992). The extent of rotation was  $\pm 60^\circ$  in each direction; the center of the plot corresponds to the training attitude. Shades of gray encode recognition rates, at increments of 5% (white is better than 90%, black is 50%). From Bülthoff and Edelman (1992). As predicted by our model, viewpoint independence can be achieved by familiarizing the subject with a sufficient number of real training views of the 3D object. For objects of a "nice" class the generalization field may be broader because of the possible availability of virtual views of sufficient quality.

sorimotor cues etc.) usually provide the equivalent of supervised learning. Unsupervised learning may be achieved by using either a bootstrap approach (see Poggio et al., 1992) or an appropriate cost-functional for learning or special network architectures.

### Critical Predictions and Experimental Data

In this section we list a few points that may lead to interesting experiments both in psychophysics and physiology.

#### Predictions

##### Viewer-centered and Object-centered Cells

Our model (see the module of Fig. 2) predicts the existence of viewer-centered cells (in the "hidden" layer) and object-centered cells (the output of the network). Evidence pointing in this direction in the case of face cells in IT is already available. We predict a similar situation for other 3D objects. It should be noted that the module of Figure 2 is only a small part of an overall architecture. We expect therefore to find other types of cells, such as for instance pose-tuned, expression-tuned and illumination-tuned cells. Very recently Logothetis and Pauls (1995) have succeeded in training monkeys to the same objects used in human psychophysics and in reproducing the key results of Bülthoff and Edelman (1992). As we mentioned above, they also succeeded in measuring generalization fields of the type shown in Figure 5 after training on a single view. We believe that such a psychophysically measured generalization field corresponds to a group of cells tuned in a Gaussian-like manner to that view. We conjecture (though this is not a critical prediction of the theory) that

the step of creating the tuned cells, that is, the centers, is unsupervised; in other words, it would be sufficient to expose the monkeys to the objects without actually training them to respond in specific ways.

##### Cells Tuned to Full Views and Cells Tuned to Parts

As we mentioned, we expect to find high-dimensional as well as low-dimensional centers, corresponding to full templates and template parts. Physiologically this corresponds to cells that require the whole object to respond (say, a face) as well as cells that respond also when only a part of the object is present (say, the mouth).

Computationally, this means that instead of high-dimensional centers any of several lower-dimensional centers are often sufficient to perform a given task. This means that the "AND" of a high-dimensional conjunction can be replaced by the "OR" of its components—a face may be recognized by its eyebrows alone, or a mug by its color. To recognize an object, we may use not only templates comprising all its features, but also subtemplates, comprising subsets of features. Splitting the recognizable world into its additive parts may well be preferable to reconstructing it in its full multidimensionality, because a system composed of several independently accessible parts is inherently more robust than a whole simultaneously dependent on each of its parts. The small loss in uniqueness of recognition is easily offset by the gain against noise and occlusions and the much lower requirements on system connectivity and complexity.

##### View-Invariant Features

For many objects and recognition tasks there may exist features that are invariant at least to some extent (color is an extreme example). One would expect this situation to occur especially in basic level recognition tasks (but not only). In this case networks with one or very few centers and hidden units—each one being invariant—may suffice. One or very few model views may suffice.

##### Generalization from a Single View for "Nice" and "Not Nice" Object Classes

An example of a recognition field measured psychophysically for an asymmetric object of a "not nice" class after training with a single view is shown in Figure 5. As predicted from the model (see Poggio and Edelman, 1990), the surface of the recognition errors is bell shaped and is centered on the training view. If the object belongs to a familiar and "nice" class of objects—such as faces—the generalization from a single view is expected to be better and broader because information equivalent to additional virtual example views can be generated from familiar examples of other objects of the same class. Ullman et al. (1993) report evidence consistent with this view. They use two "nice" classes of objects, one familiar—upright faces—and one unfamiliar—inverted faces. They find that generalization from a single training view over a range of viewpoint and illumination transformations is perfect for the familiar class and significantly worse for the unfamiliar inverted faces. They also report that generalization in the latter case improved with practice, as expected in our model.

Notice again that instead of creating virtual views the system may discover features that are view invariant for the given class of objects and then use them.

##### Generalization for Bilaterally Symmetric Objects

Bilaterally symmetric objects—or objects that may seem bilaterally symmetric from a single view—are a special example of nice classes. They are expected from the theory (Poggio and Vetter, 1992) to have a generalization field with additional peaks. The prediction is consistent with old and new psycho-

physical (Vetter et al., 1994) and physiological data (Logothetis and Pauls, 1995).

## Notes

1. Of course the distinction between view-based and object-centered models makes little sense from an information processing perspective: a very small number of views contains full information about the visible 3D structure of an object (compare Poggio and Edelman, 1990). Our view-based label refers to an overall approach that does not rely on an explicit representation of 3D structure and in particular to a biologically plausible implementation in terms of view-centered units.

2. A computational reason for why a few views are sufficient can be found in the results (for a specific type of features) of Ullman and Basri (1991). Shashua (1992a,b) describes an elegant extension of these results to achieve illumination as well as viewpoint invariance.

3. The linear classes definition of Poggio and Vetter (1992) may be satisfactory, even if not exact, in a number of practically interesting situations such as viewpoint invariance and lighting invariance for faces.

4. Invariance to illumination can be *in part* achieved by appropriate preprocessing.

5. Of course, in all of the modules described above the centers may be parts of the face rather than the full face.

6. It turns out that the RBF-like classification scheme and its implementation in terms of view-centered units is quite different from the linear combination scheme of Ullman and Basri (1990). On the other hand a regularization network used for synthesis—in which the output is the image  $y$ —is similar to their linear combination scheme (though more general) because its output is always a linear combination of the example views (see Beymer et al., 1993).

We thank Heinrich Buelthoff, Amnon Shashua, Philip Schyns, Emanuela Bricolo, and especially Nikos Logothetis for discussions and useful suggestions. This report describes research done within the Center for Biological and Computational Learning in the Department of Brain and Cognitive Sciences, and at the Artificial Intelligence Laboratory. This research is sponsored by grants from the Office of Naval Research under Contracts N00014-91-J-1270 and N00014-92-J-1879, by a grant from the National Science Foundation under Contract ASC-9217041 (funds provided by this award include funds from DARPA provided under the HPCC program), and by a grant from the National Institutes of Health under Contract NIH 2-S07-RR07047. Additional support is provided by the North Atlantic Treaty Organization, ATR Audio and Visual Perception Research Laboratories, Mitsubishi Electric Corporation, Sumitomo Metal Industries, and Siemens AG. Support for the A.I. Laboratory's artificial intelligence research is provided by ONR Contract N00014-91-J-4038.

Correspondence should be addressed to Tomaso Poggio, Center for Computational Learning, and Department of Brain Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139.

## References

- Ballard DH (1986) Cortical connections and parallel processing: structure and function. *Behav Brain Sci* 9:67-120.
- Beymer D, Shashua A, Poggio T (1993) AI memo 1431, Example based image analysis and synthesis. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Beymer DJ (1993) AI memo 1461, Face recognition under varying pose. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Biederman I (1987) Recognition by components: a theory of human image understanding. *Psychol Rev* 94:115-147.
- Brunelli R, Poggio T (1991) HyperBF networks for real object recognition. In: *Proceedings IJCAI*. Sydney.
- Brunelli R, Poggio T (1993) Face recognition: features versus templates. *IEEE Trans Pattern Anal Machine Intell* 15:1042-1052.
- Bülthoff HH, Edelman S, Sklar E (1991) Mapping the generalization space in object recognition. *Invest Ophthalmol Vis Sci [Suppl]* 32:996.
- Bülthoff HH, Edelman S (1992) Psychophysical support for a 2-D view interpolation theory of object recognition. *Proc Natl Acad Sci USA* 89:60-64.
- Edelman S (1991) Report CS-TR 91-10, Features of recognition. Rehovot: Weizmann Institute of Science.
- Edelman S, Bülthoff HH (1992) Orientation dependence in the recognition of familiar and novel views of 3D objects. *Vision Res* 32:2385-2400.
- Fujita I, Tanaka K (1992) Columns for visual features of objects in monkey inferotemporal cortex. *Nature* 360:343-346.
- Gilbert JM, Yang W (1993) A real-time face recognition system using custom vlsi hardware. In: *IEEE work on computer architectures for machine vision*. IEEE.
- Girosi F, Jones M, Poggio T (1993) AI memo 1430, Priors, stabilizers and basis functions: from regularization to radial, tensor and additive splines. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Librande S (1992) Example-based character drawing. MS thesis, School of Architecture and Planning, MIT.
- Logothetis NK, Pauls J (1995) Psychophysical and physiological evidence for viewer-centered object representations in the primate. *Cereb Cortex*, in press.
- Marr D, Poggio T (1976) Cooperative computation of stereo disparity. *Science* 194:283-287.
- Mel BW (1992) The clusteron: toward a simple abstraction for a complex neuron. In: *Neural information processing systems* (Hanson S, Moody J, Lippmann R, eds). San Mateo, CA: Kaufmann.
- Mumford D (1992) The computational architecture of the neocortex. *Biol Cybern* 66:241-251.
- Olshausen B, Anderson C, Van Essen D (1992) CNS memo 18, A neural model of visual attention and invariant pattern recognition. *Computation and Neural Systems Program*, Pasadena: California Institute of Technology.
- Perrett DI, Mistlin AJ, Chitty AJ (1989) Visual neurones responsive to faces. *Trends Neurosci* 10:358-364.
- Perrett DI, Oram MW (1993) The neurophysiology of shape processing. *Image Vision Comput* 11:317-333.
- Poggio T (1990) A theory of how the brain might work. In: *Cold Spring Harbor symposia on quantitative biology* pp 899-910. Cold Spring Harbor, NY: Cold Spring Harbor Laboratory.
- Poggio T (1991) Technical report 9107-02, 3D object recognition and prototypes: one 2D view may be sufficient. Povo, Italy: IRST.
- Poggio T, Brunelli R (1992) AI memo 1354, A novel approach to graphics. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Poggio T, Edelman S (1990) A network that learns to recognize three-dimensional objects. *Nature* 343:263-266.
- Poggio T, Edelman S, Fahle M (1992) Fast perceptual learning in visual hyperacuity. *Science* 256:1018-1021.
- Poggio T, Girosi F (1989) AI memo 1140, A theory of networks for approximation and learning. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Poggio T, Girosi F (1990a) AI memo 1167, Extension of a theory of networks for approximation and learning: dimensionality reduction and clustering. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Poggio T, Girosi F (1990b) Regularization algorithms for learning that are equivalent to multilayer networks. *Science* 247:978-982.
- Poggio T, Girosi F (1990c) Networks for approximation and learning. *Proceedings IEEE* 78:1481-1497.
- Poggio T, Hurlbert A (1993) AI memo 1404, Sparse observations on cortical mechanisms for object recognition and learning. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Poggio T, Vetter T (1992) AI memo 1347, Recognition and structure from one 2D model view: observations on prototypes, object classes and symmetries. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Sakai K, Miyashita Y (1991) Neural organization for the long-term memory of paired associates. *Nature* 354:152-155.
- Schyns PG, Bülthoff HH (1993) AI memo 1432, Conditions for viewpoint dependent face recognition. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Shashua A (1992a) Geometry and photometry in 3D visual recognition. PhD thesis AI-TR-1401, Artificial Intelligence Laboratory, MIT.
- Shashua A (1992b) Illumination and view position in 3D visual recognition. In: *Advances in neural information processing systems 4* (Hanson SJ, Moody JE, Lippmann RP, eds), pp 404-411. San Mateo, CA: Kaufmann.
- Tanaka K (1993) Neuronal mechanisms of object recognition. *Science* 262:685-688.
- Ullman S (1991) AI memo 1311, Sequence-seeking and counter-



- streams: a model for information flow in the cortex. Cambridge, MA: Artificial Intelligence Laboratory, MIT.
- Ullman S, Basri R (1991) Recognition by linear combinations of models. *IEEE Trans Pattern Anal Machine Intell* 13:992-1006.
- Ullman S, Moses Y, Edelman S (1993) Generalization across changes in illumination and viewing position in upright and inverted faces. *Perception [Suppl ECVF]* 22:25.
- Vetter T, Poggio T, Bülthoff H (1994) The importance of symmetry and virtual views in three dimensional object recognition. *Curr Biol* 4:18-23.