

# Neural and Psychophysical Analysis of Object and Face Recognition

Irving Biederman and Peter Kalocsai

Department of Psychology and Computer Science  
University of Southern California  
Los Angeles, California 90089, U.S.A.  
{ib, kalocsai}@selforg.usc.edu

**Abstract.** A number of behavioral phenomena distinguish the recognition of faces and objects, even when members of the set of objects are highly similar. Because faces have the same parts in approximately the same relations, individuation of faces typically requires specification of the metric variation in a holistic and integral representation of the facial surface. The direct mapping of a hypercolumn-like pattern of activation onto a representation layer that preserves relative spatial filter values in a 2D coordinate space, as proposed by C. von der Malsburg and his associates (Lades et al., 1993; Wiskott, et al., 1997), may account for many of the phenomena associated with face recognition. An additional refinement, in which each column of filters (termed "a jet") is centered on a particular facial feature (or fiducial point), allows selectivity of the input into the holistic representation to avoid incorporation of occluding or nearby surfaces. The initial hypercolumn representation also characterizes the first stage of object perception, but the image variation for objects at a given location in a 2D coordinate space may be too great to yield sufficient predictability directly from the output of spatial kernels. Consequently, objects can be represented by a structural description specifying qualitative (typically, nonaccidental) characterizations of an object's parts, the attributes of the parts, and the relations among the parts, largely based on orientation and depth discontinuities (e.g., Hummel & Biederman, 1992). A series of experiments on the name priming or physical matching of complementary images (in the Fourier domain) of objects and faces (See Kalocsai & Biederman, this volume) documents that whereas face recognition is strongly dependent on the original spatial filter values, object recognition evidences strong invariance to these values, even when distinguishing among objects that are as similar as faces.

**Keywords.** Face recognition, object recognition

**Acknowledgements.** This research was supported by ARO NVESD grant DAAH04-94-G-0065. This paper is excerpted, with minor modifications, from Biederman and Kalocsai (1997).

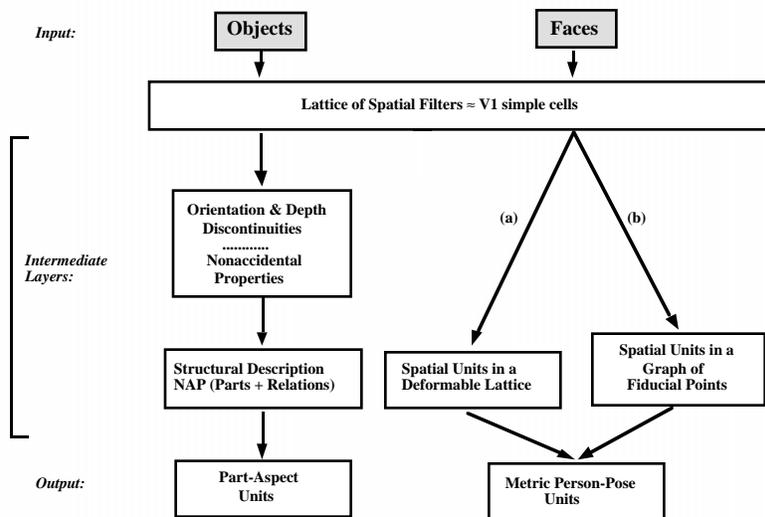
## 1 Introduction

We propose a theoretical account of the neural, perceptual, and cognitive differences that are apparent in the individuation of faces and the entry- and

subordinate-level classification of objects. After a general theoretical overview, we review some of the behavioral and neural phenomena by which face and object recognition can be contrasted and then present a neurocomputational account of these differences, with particular attention to the perceptual representation of faces. Last, original experiments testing a key assumption of this account are described.

## 2 A Theoretical Overview of Face and Object Recognition

The basic theoretical differences that we will propose are diagrammed in Figure 1. The object model follows that of Hummel and Biederman (1992) and only a brief overview will be presented here. Specification of the edges at an object's orientation and depth discontinuities in terms of nonaccidental properties (NAPs) is employed to activate units that represent simple, viewpoint invariant parts (or *geons*), such as bricks, cones, and wedges. Other units specify a geon's attributes, such as the geon's approximate orientation (e.g., HORIZONTAL) and aspect ratio, and still other units specify the relative relations of pairs of geons to each other, such as TOP-OF, LARGER-THAN, END-TO-MIDDLE-CONNECTED.



**Figure 1.** Relations between presumed models of object and face recognition. Both start with a lattice of columns of spatial filters characteristics of V1 hypercolumns. The object pathway is modeled after Biederman(1987) and Hummel and Biederman (1992) and computes a geon structural description (GSA) which represents the parts and their relations in a view of an object. Both face pathways retain aspects of the original spatial filter activation patterns. In the (a) pathway, modeled after Lades et al. (1993), the default position of the columns (termed "jets") of filters is a lattice similar to that of the input layer but which can be deformed to provide a best match to a probe image. In the (b) pathway, modeled after Wiskott, Fellous, Krüger, & von der Malsburg (1997), the jets are centered on a particular facial feature, termed a fiducial point.

The separate units associated with a given geon, its attributes, and its relations, are bound (though correlated firing) to a unit termed a geon feature assembly, GFA. A unit representing a geon structural description, *GSD*, specifying the geons and their relations in a given view of the object can then self-organize to the activity from a small set of GFAs.

Differences in GFAs are usually sufficient to distinguish entry level classes and most subordinate level distinctions that people can make quickly and accurately in their everyday lives. Sometimes the GSDs required for subordinate level distinctions are available at a large scale, as in distinguishing a square from a round table. Sometimes they are at a small scale, as when we use a logo to determine the manufacturer of a car.

Although there are some person individuation tasks that can be accomplished by the information specified by a GSD ("Steve is the guy wearing glasses"), generally we will focus on cases where such easy information as a distinctive GSD or texture field ("Steve is the guy with freckles") is insufficient. We will argue that the information required for general purpose face recognition is holistic, surface-based, and metric, rather than parts-based, discontinuous, and nonaccidental (or qualitative), as it is with objects. A representation that preserves the relative scale of the original spatial filter values in a coordinate space that normalizes scale and position may allow specification of the metric variation in that region for determining the surface properties of a face. The coordinate system is preserved because the location of facial characteristics are highly predictable from a given pose of a face. For objects they are not. (What is in the upper, right hand part of an object?) Relative (cycles/face) rather than absolute (cycles/degree) allows invariance over size changes of the face.

We consider two recent proposals by C. von der Malsburg and his associates for face representation. The first, labeled (a) in figure 1, is described by Lades et al. (1993). This system maps columns (or "jets") of V1-like spatial filter activation values to images of faces or objects. The jets are arranged in a hypercolumn-like lattice where they are stored. This stored lattice serves as a representation layer and is then matched against probe faces or objects by correlating the filter values of the original lattice against a new lattice that has been allowed to deform to achieve its own best match. The second model, labeled (b), proposed by Wiskott et al. (1997), positions each of the jets not on the vertices of a rectangular lattice but to assigned "fiducial points" on a face, such as the left corner of the mouth. These face models will be considered in more detail in a later section.

### **3 Distinguishing Face and Object Recognition: Empirical Results**

One problem with an effort to distinguish face and object recognition is that there are a large number of tasks that can be loosely described as "recognition." This problem will be examined in more detail below but for the present purposes we will consider the identification of an image of a face to the criterion of individuation and that of an object with its assignment to its basic level or common subordinate level class.

### 3.1 Behavioral Differences.

Table 1 lists seven behavioral differences between face and object recognition. These will be considered in turn with respect to the different properties that should be captured by a particular representation.

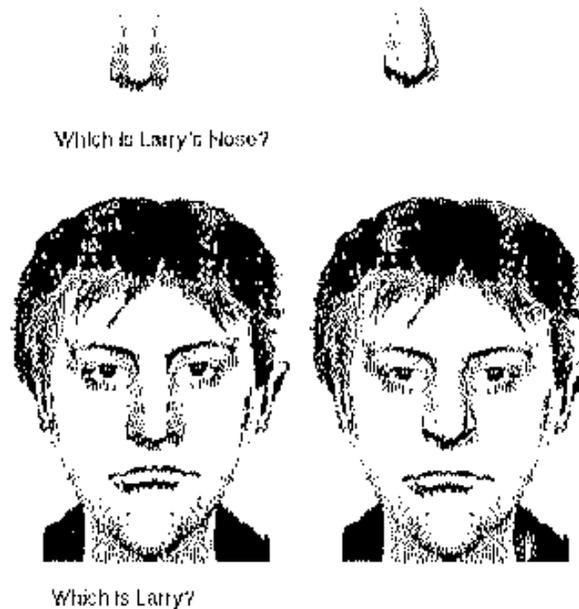
	FACES	OBJECTS
Configural Effects?	YES	NO
Basis of Expertise?	Holistic Representation	Feature Discovery
Differences Verbalizable?	NO	YES
SENSITIVE TO:		
Contrast Polarity?	YES	NO
Illumination Dir?	YES	NO
Metric Variation?	YES	SLIGHTLY
Rotation in Depth?	YES	NO, within part aspects ( $\approx 60^\circ$ )
Rotation in the Plane?	YES	SLIGHTLY

**Table 1.** Some Differences in the recognition of Faces and Objects

1. *Configural Effects.* Tanaka and Farah (1993) trained their subjects to recognize a set of Identikit faces, each with a different eyes, nose, and mouth. In testing, they presented pairs of images that differed in the shape of a single face part, the eyes, nose, or mouth (Figure 2). In one condition, only a pair of face parts was shown, for example, two slightly different noses. In the other, the stimuli were part of a context of a whole face, one with one of the noses the other with the other nose. The subjects did not know which face part might differ when they viewed a complete face. Remarkably, the context of the face facilitated detection of the difference. The facilitation from the presence of the context was not found for non-face objects, such as a house, or when the faces were inverted.

2. *Expertise.* Good face recognizers use the whole face, although with unfamiliar faces, the overall external shape and hairline receive extremely high weight (Young, Hay, McWeeny, Flude, & Ellis, 1985). When asked to describe a picture of a person's face, these individuals will often refer to a famous person, perhaps with some modification in the descriptions (Cesa, 1994). Poor recognizers tend to pick a single feature or small set of distinctive features. As people age, face recognition performance declines. This decline is marked by a qualitative shift in the representation such that older people, like poor face

recognizers in general, search for distinctive features. Prosopagnosics often report a distinctive feature strategy as well (Davidoff, 1988).



**Figure 2.** Sample stimuli from Tanaka & Farah's (1993) single feature and whole face conditions. In the single feature condition, subjects were presented with, for example, the upper pair of noses and were to judge, "Which is Larry's nose?" In the whole face condition, the subjects were presented with a pair of faces whose members were identical except they differed in a single feature, the one shown in the feature condition, and they had to judge, "Which is Larry?" Used with permission.

In contrast to the holistic processing of faces, expertise in the identification of an object from a highly similar set of objects is most often a process of discovery or instruction as to the location and nature of small differences that reliably distinguish the classes (Gibson, 1947; Biederman & Shiffrar, 1988). If such features are not present then performance is often slow and error prone (Biederman & Subramaniam, 1997). Gibson (1947) described the consequences of attempting to teach aircraft identification during World War II by "total form" versus distinctive features of the parts:

"Two principle observations made by the instructors who took part in the experiment are of some bearing on the question of the two methods under consideration. The impression was obtained by all three of the instructors, at about the time the course was two-thirds completed, that the group taught by emphasis on total form was definitely 'slipping' in comparison with the other group. The second observation was that a single question was insistently and repeatedly asked by the cadets in the group taught by emphasis on total form.

This question was 'How can I distinguish between this plane and the one which resembles it closely (e.g., the C-46 and the C-47)?' (Gibson, 1947, p. 120.)

Whether still more extensive training on non-face stimuli can lead to face-like processing is an open issue. Gauthier and Tarr (In press) provided extensive training to some of their subjects, termed "experts," in distinguishing among a family of "greebles," a set of stimuli composed of three rounded parts--a base, body, and head--one on top of the other with protrusions that are readily labeled penis, nose, and ears. These rounded, bilaterally symmetrical creatures, closely resemble humanoid characters, such as the Yoda (in Return of the Jedi). Despite Gauthier and Tarr's conclusions that they were able to mimic certain aspects of face processing with their training, none of the expected face results were obtained. Some were clearly inconsistent with face-like processing. For example, the identification of the parts (Is this Pimo's quiff?) was unaffected by inversion or scrambling of the greebles. Closer analysis of the stimuli suggest that the invariance to 2D orientation and the lack of a configural effects might have been a consequence of geon differences among the parts, rather than the metric variation in smooth surfaces required for face processing. Another shortcoming of the greebles as stimuli for the study of face perception is their resemblance to people. The parts that were tested are readily identified as ears, nose, and penis, so even if only metric variation in the surfaces of the parts had been varied, it would be unclear whether the stimuli engaged face or body processing because of their physical resemblance to people to because of the training.

3. *Differences Verbalizable?* People find it exceedingly difficult to express verbally the differences between two similar faces. This fact is well known to the chagrin of police investigators interviewing witnesses. When asked to describe an object, however, people readily name its parts and provide a characterization of the shape of these parts in terms of NAPs (Tversky & Hemenway, 1984; Biederman, 1987). Within highly similar shape classes, such as Western U. S. male Quail, people will spontaneously employ local shape features that closely correspond to those specified--verbally--by the bird guides (Biederman, 1997). Gibson (1947) concluded that the problem of training aircraft spotters was best solved by informing them as to the nonaccidental differences in the shapes of parts. It was a simple matter for Gibson to construct an outline--in words--providing this information.

4. *Sensitivity to contrast polarity and illumination direction?* Whereas people have great difficulty in identifying a face from a photographic negative or when illuminated from below (Johnston, Hill, & Carmen, 1992), there is little, if any, effect of reversing the polarity of contrast of a picture of an object (Subramaniam & Biederman, 1997). Viewing an object at one polarity provides essentially the same information as to the structure of the object as does the other polarity. A major reason for this difference between faces and objects is that, as noted previously, object recognition is largely based on distinguishable parts based on differences in NAPs of edges marking orientation and depth discontinuities. The position of these edges and their nonaccidental values (e.g., straight or curved) are unaffected by contrast reversal. Individuating faces typically requires metric differences that may be specified in terms of the convexities and

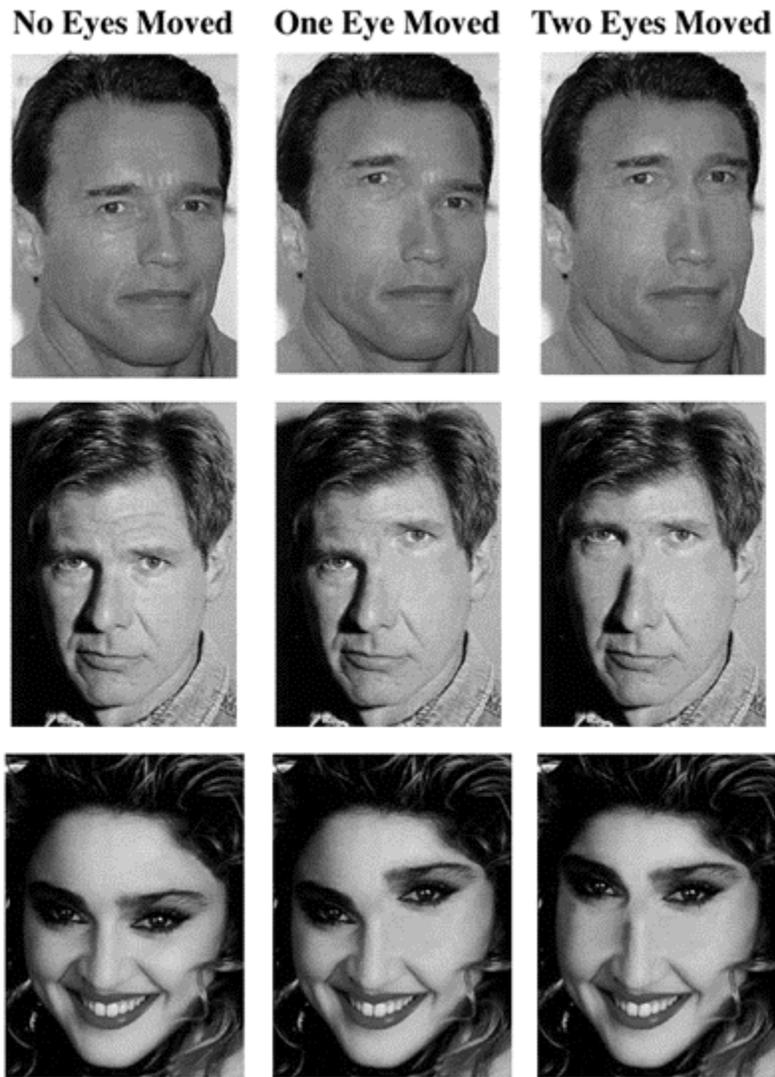
concavities that characterize a facial structure. A change in contrast polarity would reverse the interpretation of the luminance and shadow gradients that are employed to determine the convexity or concavity of a smooth surface. A similar explanation may account for some of the increased difficulty in identifying faces when they are illuminated from below as this would violate the strong assumption that illumination is from above.

5. *Metric variation?* Metric properties are those such as aspect ratio or degree of curvature that vary with the orientation of the object in depth. Such properties are to be contrasted with NAPs, such as whether an edge is straight or curved, which are only rarely affected by slight changes in viewpoint of an object. Other NAPs are the vertices that are formed by coterminating lines and whether pairs of edges are approximately parallel or not, given edges that are not greatly extended in depth.

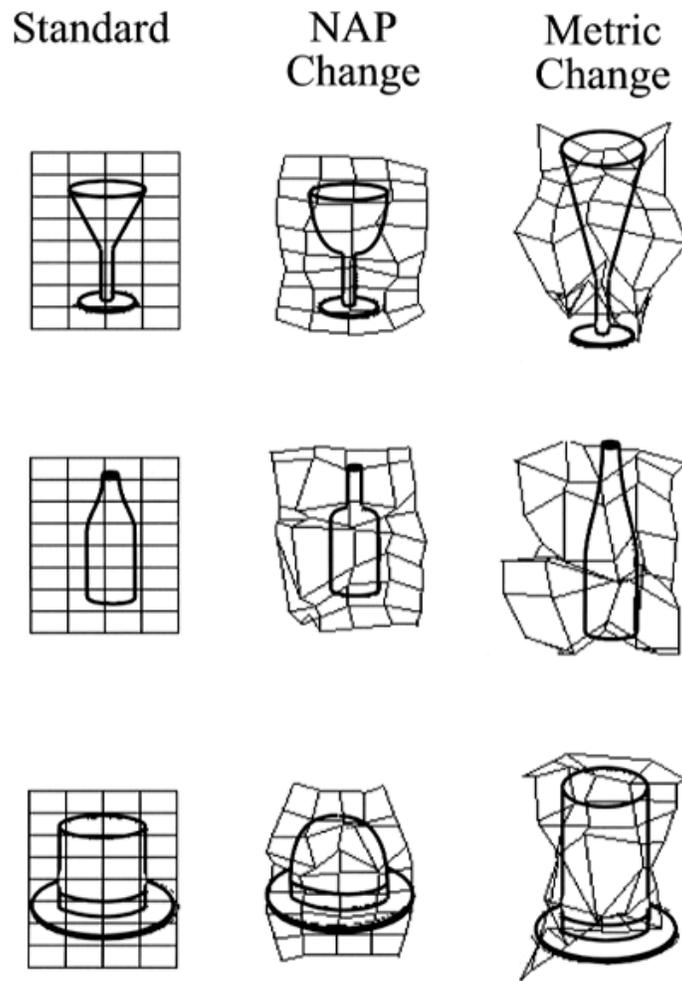
Before looking at figure 3 (from Cooper & Wojan, 1996), please cover the left and center columns. In looking at the right column, the reader can assess for himself or herself how modest variation in the metrics of a face can result in marked interference in the recognition of that face (see also Hosie, Ellis, & Haig, 1988). In these images of celebrities, the eyes have been raised. A similar variation in the length (and, hence, aspect ratio) of an object part, as illustrated in figure 4, has little or no effect in the assignment of objects to classes. As long as the relative relations, such as LARGER-THEN or ABOVE, between parts are not changed by altering a part's length, the effects of the variation appear to be confined to that part, rather than affecting the object as a whole. Unlike what occurs with the holistic effects with faces, there is little effect of the variation on a metric attribute of a part in the recognition of objects. Biederman and Cooper (1993) presented two images of simple, two-part objects (illustrated in Figure 4) sequentially. Subjects had to judge whether the two objects had the same name. When the objects differed in the aspect ratio of a part, RTs and error rates were only slightly elevated compared to when the images were identical. A change in a NAP produced a much larger interfering effect on the matching.

6. *Rotation in depth.* If objects differ in NAPs, then little or no cost is apparent when they are rotated in depth, as long as the same surfaces are in view (Biederman & Gerhardstein, 1993). In contrast, when the differences are in metric properties, such as aspect ratio or degree of curvature, then marked rotation costs are observed (e.g., Edelman, 1995). The robustness of the detection of nonaccidental differences under depth rotation is not simply a function of greater discriminability of NAPs compared to metric properties. Biederman and Bar (1995) equated the detectability of metric and nonaccidental part differences in a sequential same-different matching task with novel objects. Presenting the objects at different orientations in depth had no effect on the detectability of nonaccidental differences. When easy nonaccidental cues are eliminated, such as glasses, facial hair, and the hairline, even modest rotations of faces, from 20° left to 40° right, as illustrated in figure 7 (middle row), can result in marked increases in RTs and error rates in their matching (Kalocsai, Biederman, & Cooper, 1994).

7. *Rotation in the plane.* Recognizing an upside-down face is extremely difficult relative to identifying an upside-down object, such as a chair (e.g., Yin, 1969; Johnston, et al., 1992; Jolicoeur, 1985). According to the Hummel and



**Figure 3.** Sample stimuli from Cooper and Wojan (1996). Subjects were much worse at identifying the celebrities in the third column, where both eyes were raised, compared to those in the second column where only one eye was raised, despite the greater difficulty in judging the latter as a face. Copyright Eric E. Cooper. Used with permission.



**Figure 4.** Sample object stimuli from Cooper and Biederman (1993). Given the standard object on the left, a NAP of only a single part was changed in the objects in the middle column (NAP condition) and that same part was lengthened in the Metric condition illustrated by the objects in the third column. The magnitude of the Metric changes were slightly larger than the NAP changes, according to the Lades et al. (1993) model. Whereas the difference between Metric and Standard images were more readily detected when performing a simultaneous physical identity matching task (Are the objects identical?) , in a sequential object matching task (Do the objects have the same name?), a change in a NAP resulted in far more disruption than a change in a metric property.

Biederman (1992) network, turning an object upside down would leave most of the units coding the structural description intact, affecting only the relations TOP-OF and BELOW. Consequently, only a small effect for objects would be

expected. Some of the large effect of inversion with face photos lies in the misinterpretation of luminance gradients where the light source is typically assumed to be coming from above. But when the light source is controlled, there still remains a large cost to viewing a face upside-down (Johnston, et al., 1992; Enns & Shore, 1997).

### 3.2 Neural Differences Between Faces and Objects

There are several neural differences distinguishing the representation of faces and objects. Only a brief summary will be presented here. (See Grüsser & Landis, 1991, for a comprehensive treatment of this general area.)

1. *Selective impairment: Prosopagnosia and object agnosias.* Prosopagnosia, the inability to recognize familiar faces but with a normal or near normal capacity for object recognition, is a well documented phenomenon, generally associated with lesions to the right, inferior mesial hemispheric (Grüsser & Landis, 1991), although some (e.g., Damasio, Damasio, Van Hoesen, 1985) have argued that the lesions must be bilateral. Farah (1990) theorized that the underlying continuum in visual recognition extended from holistic processing, which would be required for faces, to the capacity to represent multiple shapes (or parts), which would be typified by the integration of letters into words in reading. She surmised that right inferior occipital-temporal lesions affected holistic processing whereas bilateral lesions to the inferior temporal-occipital region (including the fusiform) resulted in a condition (ventral simultanagnosia) in which the patient could not simultaneously process multiple parts of an object or letters of a word (alexia). (Other authors, e.g., Behrmann & Shallice, 1995, have also argued that alexia is associated with lesions to the left hemisphere.) Object recognition, according to Farah, employs both types of processing so object agnosia should be accompanied by either prosopagnosia or alexia. Several recent cases, however, have described individuals manifesting strong object agnosias who are neither prosopagnosic nor alexic (Rumiati, Humphreys, Riddoch, & Bateman, 1994; Moscovitch, Winocur, & Behrmann, 1997). We interpret these findings as evidence that object recognition does not generally entail holistic processing and that the integration of letters into a word in reading may not necessarily be engaging the same mechanisms or representations that mediate face recognition.

2. *Imaging studies.* Recent fMRI studies in humans give clear evidence for object and shape specific regions in the occipital cortex. Tootell, Dale, Sereno, & Malach (1996) have documented an area just anterior to V4v and partly overlapping with regions of the fusiform, termed Lateral Occipital (LO), that gives vigorous responses to interpretable faces and objects even when they are unfamiliar, such as an abstract sculpture, but not to these stimuli when they have been rendered into textures as, for example, digitized blocks characteristic of the "Lincoln" illusion or in gratings, texture patterns, or highly jumbled object images. In contrast to LO, V4 does not show this specificity to objects as compared to textures. LO is thus sensitive to shapes--faces or objects--that have an interpretable structure rather than being characterizable as a texture pattern.

More anterior regions in the ventral pathway such as IT are sensitive to the familiarity of the objects, as described in the next section. That LO's responsivity is unaffected by familiarity suggests that it may be a region where shape descriptions--even novel ones--are created. A number of fMRI and PET studies have demonstrated that the processing of faces and objects activate different loci in or near LO. These areas are generally consistent with the results of the lesion work, showing greater posterior right hemisphere activity, particularly in the fusiform gyrus, for face processing and greater left hemisphere activity for object processing (Kanwisher, Chun, & McDermott, 1996; Sergent, Ohta, Macdonald, & Zuck. 1994; Sergent, Ohta, & Macdonald; 1994). The two Sergent et al. PET studies are noteworthy in showing virtually identical loci for the differential activity of judging whether a face was that of an actor. The control task was one of judging whether the orientation of a gratings was horizontal or vertical.

3. *Single unit recording.* It is well established that individual IT cells can be found that are differentially tuned either to faces or to complex object features, but not both (e.g., Bayliss, Rolls, & Leonard, 1987; Kobatake & Tanaka, 1994; Young & Yamane, 1992). However, as recently argued by Biederman, Gerhardstein, Cooper, & Nelson (1997), it is likely that these IT cells are not involved in the initial perceptual description of an image--which they suggest is accomplished by LO or in the area immediate anterior to it--but, instead, in coding episodic memories *following* perception. Because these experiences include contribution of the dorsal system in which position, size, and orientation of the stimulus is specified, it is not surprising to find cells that are tuned to the specific orientations and characteristics of the trained stimuli (e.g., Logothetis, Pauls, Bülthoff, & Poggio, 1994). That IT may not be involved in the perceptual recognition of a face or object is suggested by the requirement of an interval between stimulus presentation and testing in order to show any deficits in object processing of macaques who have undergone bilateral ablation of IT (Desimone & Ungerleider, 1989). However, the differential tuning of IT cells to faces and complex object features indicates that these two classes of stimuli are distinguished neurally. A given IT face cell does not fire in all-or-none fashion to a given face but participates in a population code to that face by which the firing of the cell is modulated by the specific characteristics of the face (Young & Yamane, 1992; Rolls, 1992). Young and Yamane showed that the code for macaques looking at pictures of men could be summarized by two dimensions, one coding the width of the face and one the distance of the pupil of the eye to the hairline. Somewhat remarkably, as noted earlier, these same two dimensions characterize human performance with unfamiliar faces. Recently, Scialoja, Wilson and Goldman-Rakic (1997) showed that the isolation of face and object processing extended to the prefrontal cortex where they found cells in the macaque that were tuned exclusively to faces and were quite unresponsive to objects, scrambled faces, or objects of interest such as food.

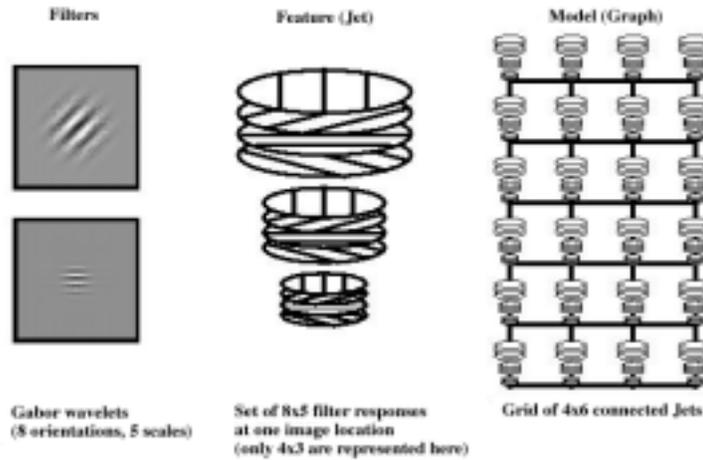
4. *Universal Classes of Facial Attributes.* All cultures appear to process faces in highly similar ways. Faces are not only processed for identity, but for the information they provide about emotion, age, sex, direction of gaze, and attractiveness. Different areas mediate at least some of these attributes. Cells

tuned to differences in emotional expression and direction of gaze are found in the superior temporal sulcus in the macaque, an area different from the IT locus of the units that contribute to a population code that can distinguish identity. Prosopagnosics can often readily judge these other attributes, e.g., sex, age, etc., as we have recently witnessed in our laboratory. To the extent that these areas are segregated from those for object recognition, we have additional evidence supporting the face-object distinction. However, it is not clear to what extent, if any, these classes contribute to face individuation.

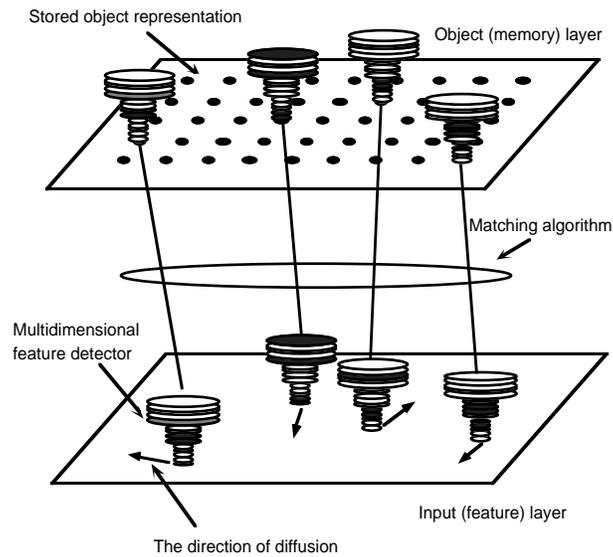
#### **4 A Theory of Perceptual Recognition of Faces**

A biologically inspired face recognition system developed by Christoph von der Malsburg and his associates (Lades, et al., 1993; Wiskott, et al., 1997) suggests a theoretical perspective from which many of the phenomena associated with face perception described in the previous section might be understood. The fundamental representation element is a column of multiscale, multiorientation spatial (Gabor) kernels with local receptive fields centered on a particular point in the image. Each column of filters is termed a "Gabor jet" and each jet is presumed to model aspects of the wavelet-type of filtering performed by a V1 hypercolumn. We will first consider the initial version of the model (Lades et al, 1993), which will be referred to as the lattice version. This model can be applied to the recognition of faces and objects so it has the potential to serve as a device for the scaling of both kinds of stimuli. A more recent version (Wiskott et al., 1997) , the "fiducial point" model, incorporates general face knowledge. We will ignore preprocessing stages by which a probe image is translated and scaled to achieve a normalized position and size. Overall illumination levels and contrast are similarly normalized.

As illustrated in figure 5, Lades et al. (1993) posited a two-layer network. The input layer is a rectangular lattice of Gabor jets. The pattern of activation of the 80 kernels (5 scales X 8 orientations X 2 phases, sine and cosine) in each of the jets is mapped onto a representation layer, identical to the input layer, that simply stores the pattern of activation over the kernels from a given image. An arbitrary large number of facial images can be stored in this way to form a gallery. Matching of a new image against those in the gallery is performed by allowing the jets (in either the probe or a gallery image) to independently diffuse (gradually change their positions) to determine their own best fit, as illustrated by the arrows on the jets in the input layer. The diffusion typically results in distortion of the rectangular lattice, as illustrated in Figures 6 and 7. The similarity of two images is taken to be the sum correlation in corresponding jets of the magnitudes of activation values of the 80 corresponding kernels. The correlation (range 0 to 1) for each pair of jets is the cosine of the angular difference between the vectors of the kernels in a 80 dimensional space. (If the values are identical, the angular difference will be 0 deg and the cosine will be 1. A 90 deg [orthogonal] difference in angles will be 0.00.) The correlations over the jets are summed to get a total similarity score. Figure 7 illustrates distortion of the lattice as a person changes expression, orientation, and both expression and orientation. Typically, the greater the deformation of the lattice, the lower the similarity of the match.



**Figure 5.** Illustration of the input layer to the Lades et al. (1993) network. The basic kernels are Gabor filters at different scales and orientations, two of which are shown on the left. The center figure illustrates the composition of a jet, with the larger disks representing lower spatial frequencies. The number of jets, scales, and orientation can be varied.



**Figure 6.** Schematic representation of the Lades et al. (1993) two-layer spatial filter model. The model first convolves each input image with a set of Gabor kernels at five scales and eight orientations and sine and cosine kernels arranged in a 5 x 9 lattice. These values can be varied. The set of kernels at each node in the lattice is termed a "Gabor jet". The activation values of the kernels in each jet along with their positions are stored for each of the images to form a "gallery". The figure shows the diameters of the receptive fields to be much smaller than actual size in that the largest kernels had receptive fields that were almost as large as the whole face.

Given a test image against a number of stored images, the most similar image, if it exceeds some threshold value, is taken to be the recognition choice.<sup>1</sup>

As noted earlier, the model does a good job at recognizing faces. Given modest changes in pose and expression, recognition accuracy can exceed 90 percent. How well does the model reflect the phenomena associated with faces listed in Table 1?

1. *Rotation Effects.* We will first consider the model's handling of rotation effects, particularly rotation in depth, as that is an extremely common source of image variation and we have assessed its effects under well controlled conditions.

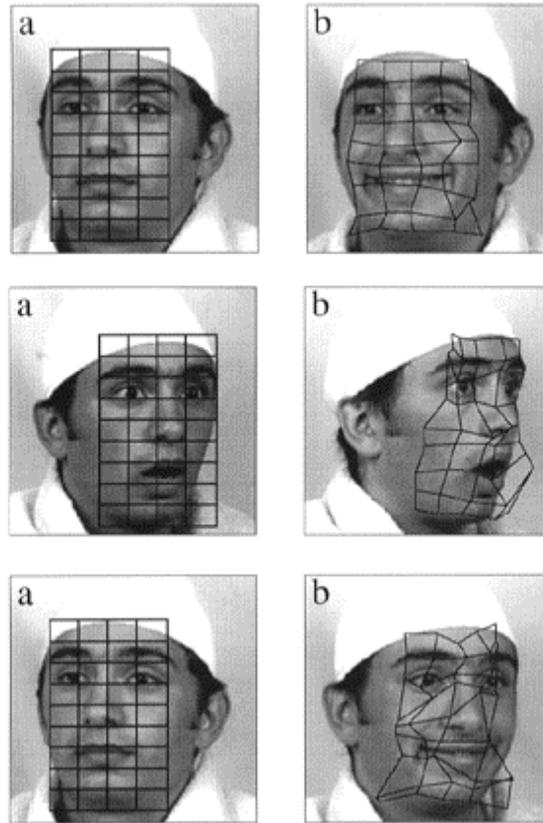
Kalocsai, Biederman, & Cooper (1994) had subjects judge whether two sequentially presented faces were of the same or different person. The faces could be at different orientations in depth and/or with a different expression, as shown in Figure 7. Easy cues, such as facial hair, clothing and the hairline (all stimulus models wore a white bathing cap) were eliminated. A change in the depth orientation of the two poses, such as that shown in the middle row of figure 7, increased RTs and error rates for 'same' trials. The magnitude of this cost was strongly and linearly correlated with the lattice model's similarity values for the pair of pictures, -.90 for RTs and -.82 for error rates. That is, the more dissimilar the two figures according to the model, the longer the RTs and error rates for judging them to be the same person. We can consider the effects of depth rotation as a yardstick for determining the model's adequacy for handling other effects.

Turning a face upside down would greatly reduce its similarity to that of the original image. Although it would be a simple matter, computationally, to rotate the coordinate space of the jets to eliminate the effects of planar rotation, the large cost to human recognition performance from inversion suggests that such a transformation is not available to human vision. Given a yardstick of depth rotation, it is an open question whether the same similarity function would also account for the cost of 2D inversion or other variables. That is, would a 60 deg rotation in depth (around the y-axis) result in as much cost as a 60 deg rotation in the plane? What would human subjects evidence?

Given that we have a scaling device (viz., the Lades et al. model), the analysis that could be undertaken to compare rotation in depth to rotation in the plane can be illustrated by Kalocsai et al.'s (1994) comparison of the effects of differences in depth orientation to the effects of differences in expression. Kalocsai et al. (1994) showed that when the degree of image dissimilarity of two images of the same person produced by differences in depth orientation (holding expression

---

<sup>1</sup>In terms of a current psychological theory of face recognition, the two-layer network would be an alternative to Bruce's "Face Recognition Units (or FRUs). Whereas FRUs are pose independent (Burton, 1994), the Lades et al. (1993) network has only modest capabilities to generalize over large rotations in depth, insofar as it starts with the facial image itself and the image is altered by even modest variations in pose, lighting direction, etc. It would be by associating different person-pose units (the output of the Lades et al. model) to the same Person Identification Node, or *PIN*, (Bruce, 1988) that the same semantic information about a person could be activated independent of the pose.



**Figure 7.** Sample images from the Kalocsai, Biederman, and Cooper (1994) experiment with the Lades et al. (1993) lattice deformations superimposed over different pairs of images of the same person. The positioning of the lattice over an original image is shown in the left hand column (a) and the deformed lattice is shown in the right column (b). Top, middle, and bottom rows show changes in expression, orientation (60°), and both expression and orientation, respectively. The similarities as determined by the Lades et al. (1993) model correlated highly with performance in matching a pair of images when there were at different orientations and expressions (Kalocsai et al., 1994).

constant) and expression differences (holding depth orientation constant) were equated, the increase in RTs and error rates in responding "same" were three times greater when the dissimilarity was produced by expression differences than when produced by depth rotation. They modeled this effect by assuming that a classifier for expression, which was also highly correlated with Gabor similarity, would signal a mismatch to a decision stage [same vs. different person?] between two face images that differed in expression, even though the images were of the same

person. That mismatch signal resulted in the increased cost for faces differing in expression.

2. *Configural and verbalization effects.* Contrast variation within any small region of the face would affect all those kernels whose receptive fields included that region. The pattern of activation of the kernels implicitly contains a holistic or configural representation in that the shape of all facial features and their positions with respect to each other are implicitly coded by the activation of the kernels. Indeed, the representation if run with sufficient jets would be equivalent to a picture of a face and so it does not distinguish contrast variation arising from the shape of facial features from contrast variation arising from translation of those features. It would be impossible to move a region or a feature or to change a feature without affecting the coding of a number of kernels from a number of jets. The representation thus becomes integral (Shepard, 1964) or nonanalytical (Garner, 1966) in that it is not decomposed into readily perceivable independent attributes. This spatially distributed population code of activation values of many kernels of varying scales and orientations in a number of different jets thus captures many of the characteristics of what is generally meant by "holistic representations." Consistent with human performance, this spatially distributed code would be extraordinarily difficult to verbalize.

3. *Lighting, and Contrast Reversal Effects.* Although the model's normalization routines allows its performance to be invariant to overall lighting and contrast levels, a change in the direction of lighting would result in a cost in similarity for the lattice model. It is not clear whether changing the light source vertically, from top to bottom, would result in a greater reduction in similarity, than a right to left change, nor would the cost of contrast reversal necessary be as severe as that evidenced in human performance when compared to, say, rotation in depth. There is nothing in the model, at present, that would identify regions on the surface as convex or concave.

4. *Metric sensitivity.* Metric variation such as that performed by Cooper and Wojan (1996) in raising the eyes in the forehead would alter the pattern of activation values in the lattice. Although the distortion of the lattice might be sufficient to account for the effects on recognition performance of such an operation, it is not obvious how lattice distortion would handle the much smaller effect of moving only one eye. In this case, the relation between the eyes would be disrupted, although one half of the lattice would, most likely, not be affected. We will return to this problem when we consider the incorporation of fiducial points.

Another result that is not obviously derived from the lattice model is the extraordinary difficulty in recognizing the components of a face where the upper half is of one famous person and the lower half another, with the upper and lower halves smoothly aligned to constitute a single face (Young, Hellawell, & Hay, 1987). When the upper and lower halves are offset it is much easier to identify the component individuals.

A third result is that we experience little distortion of other regions when a face is partially occluded as, for example, when a person holds his chin with his

hand. The hand is not seen as part of the face but instead is regarded as another object, with the occluded regions contributing little, if anything, to the perception of the face.

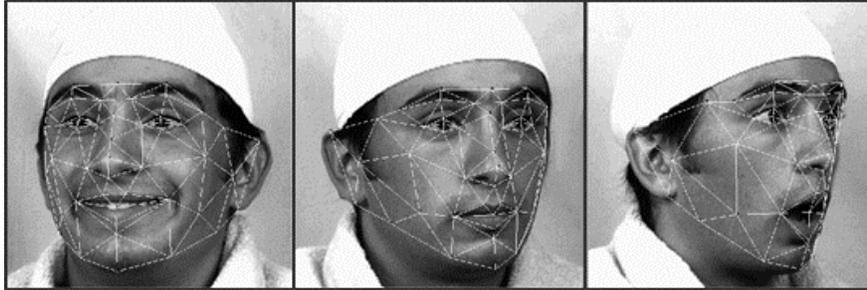
5. *Direct tests of filter-based matching in face but not object recognition.* A series of experiments on the name priming or physical matching of complementary images (in the Fourier domain) of objects and faces (See Kalocsai & Biederman, this volume) documents that whereas face recognition is strongly dependent on the original spatial filter values, object recognition evidences strong invariance to these values, even when distinguishing among objects that are as similar as faces.

## 5 Beyond a lattice of spatial features: Fiducial points

We now consider the fiducial point version of the face recognition system so that we can appreciate the potential gains in making facial features explicit by centering designated jets onto salient feature points. We will also consider two other possible extensions of the model: The explicit use of spatial distances and normative coding by which a face is represented in terms of its deviations from a population norm.

In the fiducial point model (Wiskott, Fellous, Krüger, & von der Malsburg, 1997), the jets are not initially arranged in a rectangular lattice but, instead, each jet is centered on a particular landmark feature of the face, termed a *fiducial point*, such as the corner of the right eye. This step has been implemented and was achieved by centering each of 45 jets (by hand) on a particular fiducial feature, e.g., the outside corner of the right eye, for a "learning set" of 70 faces, which differed in age, sex, expression, depth orientation, etc. Figure 8 shows some of the fiducial points on a face at different orientations and expressions. The 70 jets for each of the 45 points are stored as a "bunch graph." When a new face is presented to the system, not the mean *but the closest fitting* of the 70 jets for each feature is taken as a basis for refining the position by undergoing local diffusion. For example, if the right eye in the probe image is blinking, then a best match might be an eye that is blinking, rather than the mean. A jet on the center of the chin might come from another face. Once a sufficiently large set of faces are included in the bunch graphs ( $\approx 50$ ), it is possible to automatically add new fiducial points. After the matching jet from the bunch graph finds its optimal position, the actual pattern of activation for a jet at that fiducial point is taken to be one of the jets representing that particular face.

The fiducial points, in addition to potentially allowing better resolution in matching, can readily be employed to reject inappropriate image information, such as would occur if the face were partially occluded by a hand. When none of the jets for a given fiducial point in the bunch graph can match their feature to some confidence level in a circumscribed region (constrained in part, by the neighboring jets), that jet is simply not employed in the matching phase. In this way partial occlusion can be made to exact a much smaller cost on recognition than it would if the occluder were incorporated into the representation of a face. Although not implemented, it may be possible to suppress the activity from parts of the receptive fields of jets that lie outside of the bounding contours of the face so they



**Figure 8.** Illustration of the mapping of jets onto fiducial points (the vertices of the triangles) on three images of the same person at different orientations and expressions.

do not contribute to the representation as well. Young et al.'s (1987) finding that offsetting the upper and lower halves of a composite face resulted in much better performance in recognizing the component individuals might be handled by a similar application of a fiducial point model. In this case the fiducial points in the upper and lower halves of the face were not in their expected locations so their activation pattern would not be included in matching one half of the face to the other half. It is possible, of course, that beyond the offset of the fiducial points, the matched cusps provide strong evidence of separate parts and this evidence could also enter into the easier retrieval of the offset face.

It will be recalled that in the Cooper and Wojan (1996) experiment, better recognition was obtained for faces in which one eye was raised, rather than both of them, despite the former stimuli looking less like a face. If the expected locations of the fiducial points for the eye on the opposite side of the head differed for the left and right halves of the face, then each face half might not have been integrated the fiducial points of the eye in the opposite half. Consequently, the original half could vote for the correct face, without incorporation of the distorted region.

In summary, in addition to greater accuracy in recognizing faces over a wider range of conditions, the great value in employment of a fiducial point representation is that it allows selective attention to be exercised over a holistic representation of the face.

### 5.1 The use of topological relations

A second modification of the filter model would be the incorporation of the *distances* between the jets. This could be done either with the original lattice or with the fiducial points. Figures 7 and 8 show both arrangements with the nodes of the lattice (upper) connected to its nearest nodes and the fiducial points (lower) connected to their nearest fiducial points to form a set of triangles. A change in the image of a face produced by changes in orientation and expression, as in figures 7 and 8, results in distortion of the lattice or the triangles. A potentially important representational problem is whether the distances among the jets (or the distortions of these distances) should be incorporated into the representation or whether the jet similarities are sufficient to account for the accuracy of the model's

performance in modeling human face recognition. Many issues remain about the possible inclusion of an explicit measure of distance (e.g., the sum of the squares of the differences in corresponding distances) as a component of similarity in the matching phase. The fiducial point model has a strong potential for serving as a research platform for addressing these and a number of the other issues in face recognition, such as norm based coding.

## 5.2 Norm Based Coding?

In the current versions of the model, the match of a probe face to a face stored in the gallery is only a function of the similarity between the two. An alternative basis for matching could be to include not only the similarity of the two faces but their distances from the norms of a population of faces. There are several effects that would suggest some role of such norm based coding in face recognition. Caricatures can be created by enhancing deviations (e.g., by 50%) of points on a particular face from the population values (see Rhodes & Tremewan, 1994, for a recent review). Moreover, for famous faces the recognition accuracy of such caricatures does not suffer in comparison to--and can sometimes be found to exceed-- the recognition accuracy of the original face (Rhodes & Tremewan, 1994). Carey (1992) and Rhodes (1994) tested whether the caricature gains its advantage in recognition (or resists a loss) because of the increased "distinctiveness" of the distortions in face space. They showed that "lateral" caricatures, in which the distortions were made in a direction orthogonal to the direction of the deviation of a point, were recognized less well than 50% characters, which were recognized as well as the original, and even less well than *anticaricatures*, faces where the distortion was reduced by 50% towards the norm. Thus, it is not merely *any* distortion that produces an advantage, but only those that enhance the deviations from the norm.

The fiducial point model of Wiskott et al. (1997) would seem to be particularly well designed to incorporate norm-based coding. Whether the perception of caricatures differs from that of non caricatured faces can be assessed with such a representation. A caricature matched against its original image will have a lower similarity value with the standard matching routines in the Wiskott et al. system. But it would be a simple matter to include deviations of both the jet locations and the kernel activation values from a normed face. One can also ask whether the advantage of the caricature is one of deviations from the norm or deviations from near neighbors? In general these two measures will covary. An explicit model also offers the possibility of more detailed tests of how caricatures function. When performed over a set of faces, would it be possible to predict which faces would enjoy a caricature advantage and which not? Should greater weight in matching be given to kernels in proportion to their departure from their normed activation value? This last question raises a possible issue with respect to caricatures. People typically realize that they are looking at a caricature and not the original face. Is it possible that caricature perception alters the way in which faces are coded or matched? Specifically, do models that predict the distinctiveness of uncaricatured faces also serve to predict the distinctiveness of caricatured faces?

## 6 Conclusion

A number of differences are apparent in the behavioral and neural phenomena associated with the recognition of faces and objects. Readily recognizable objects can typically be represented in terms of a geon structural description which specifies an arrangement of viewpoint invariant parts based on a nonaccidental characterization of edges at orientation and depth discontinuities. The parts and relations are determined in intermediate layers between the early array of spatially distributed filters and the object itself and they confer a degree of independence between the initial wavelet components and the representation. The units in a structural description of an object allow ready verbalization. The nonaccidental characterization of discontinuities endows the representation with considerable robustness over variations in viewpoint, lighting, and contrast variables. Last, object experts discover mapping of small nonaccidental features. Individuation of faces, by contrast, requires specification of the fine metric variation in a holistic representation of a facial surface. This can be achieved by storing the pattern of activation over a set of spatially distributed filters. Such a representation will evidence many of the phenomena associated with faces such as holistic effects, nonverbalizability, and great susceptibility to metric variations of the face surface, as well as to image variables such as rotation in depth or the plane, contrast reversal, and direction of lighting. Face experts represent the whole face. A series of experiments demonstrated that the recognition or matching of objects is largely independent of the particular spatial filter components in the image whereas the recognition or matching of a face is closely tied to these initial filter values.

## References

- Baylis, G. C., Rolls, E. T., and Leonard, C. M. (1987). Functional subdivisions of the temporal lobe neocortex. *Journal of Neuroscience*, 7, 330-342.
- Behrmann, M., Winocur, G., & Moscovitch, M. (1992). Dissociation between mental imagery and object recognition in a brain-damaged patient. *Nature*, 359, 636-637.
- Behrmann, M., & Shallice, T. (1995). Pure alexia: A nonspatial visual disorder affecting letter activation. *Cognitive Neuropsychology*, 12, 409-454.
- Biederman, I. (1987). Recognition-by-components: A theory of human image understanding. *Psychological Review*, 94, 115-147.
- Biederman, I. (1995). Visual object recognition. In S. F. Kosslyn and D. N. Osherson (Eds.). *An Invitation to Cognitive Science*, 2nd edition, Volume 2., *Visual Cognition*. MIT Press. Chapter 4, pp. 121-165.
- Biederman, I., & Bar, M. (1995). One-Shot Viewpoint Invariance with Nonsense Objects. Paper presented at the Annual Meeting of the Psychonomic Society, 1995, Los Angeles, November.
- Biederman, I. & Cooper, E. E. (1991). Priming contour-deleted images: Evidence for intermediate representations in visual object recognition. *Cognitive Psychology*, 23, 393-419.

- Biederman, I., & Cooper, E. E. (1992). Size invariance in visual object priming. *Journal of Experimental Psychology: Human Perception and Performance*, 18, 121-133.
- Biederman, I., & Gerhardstein, P. C. (1995). Viewpoint-dependent mechanisms in visual object recognition: Reply to Tarr and Bülthoff (1995). *Journal of Experimental Psychology: Human Perception and Performance*, 21, 1506-1514.
- Biederman, I. & Gerhardstein, P. C. (1993). Recognizing depth-rotated objects: Evidence and conditions for three-dimensional viewpoint invariance. *Journal of Experimental Psychology-Human Perception and Performance*, 19(6), 1162-1182.
- Biederman, I., Gerhardstein, P.C. , Cooper, E. E., & Nelson, C. A. (1997). High Level Object Recognition Without an Anterior Inferior Temporal Cortex. *Neuropsychologia*, 35, 271-287.
- Biederman, I., & Kalocsai, P. (1997). Neurocomputational bases of object and face recognition. *Philosophical Transactions of the Royal Society of London B*, 352, 1203-1219.
- Biederman, I., & Subramaniam, S. (1997). Predicting the shape similarity of objects without distinguishing viewpoint invariant properties (VIPs) or parts. *Investigative Ophthalmology & Visual Science*, 38, 998.
- Biederman, I., Subramaniam, S., Kalocsai, P, and Bar, M. (1998). Viewpoint-invariant information in subordinate-level object classification. In D. Gopher & A. Koriat (Eds.) *Attention and Performance XVII. Cognitive Regulation of Performance: Interaction of Theory and Application*. Cambridge, MA: MIT Press.
- Bruce, V. (1988). *Recognizing Faces*. Hove and London, UK: Erlbaum.
- Bruce, V., & Humphreys, G. W. (1994). Recognizing objects and faces. *Visual Cognition*, 1, 141-180.
- Cesa, I. L. (1994). Of attractive librarians and lop-sided faces: Development and testing of a training procedure for improving the accuracy of eyewitnesses. Unpublished Doctoral Dissertation, Department of Psychology, Univ. of Southern California.
- Cooper, E. E., & Biederman, I. (1993). Metric versus viewpoint-invariant shape differences in visual object recognition. *Investigative Ophthalmology & Visual Science*, 34, 1080.
- Cooper, E. E., & Wojan, T. J. (1996). Differences in the coding of spatial relations in faces and objects. *Investigative Ophthalmology & Visual Science*, 37, 177.
- Damasio, A. R., Damasio, H., & Van Hoesen, G. E. (1982). Prosopagnosia: Anatomic basis and behavioral mechanisms. *Neuropsychologia*, 2, 237-246.
- Davidoff, J. B. (1988). Prosopagnosia: A disorder of rapid spatial integration. In G. Denes, C. Semenza, & P. Bisiachi (Eds/), *Perspectives on Cognitive Neuropsychology*, (pp. 297-309). Hillsdale, NJ: Erlbaum.
- Desimone, R., & Ungerleider, L. G. (1989). Neural mechanisms of visual processing in monkeys. (Chapter 14, Pp. 267-299). In F. Boller & J. Grafman (Eds.) *Handbook of neuropsychology*, Vol. 2. Amsterdam: Elsevier.
- Edelman, S. (1995). Representation of similarity in 3D object discrimination. *Neural Computation*, 7, 407-422.

- Enns, J. T., & Shore, D. I. (1997). Separate influences of orientation and lighting in the inverted-face effect. *Perception & Psychophysics*, 59, 23-31.
- Farah, M. J. (1990). *Visual Agnosia: Disorders of Object Recognition and What They Tell Us About Normal Vision*. Cambridge, MA: MIT Press.
- Farah, M. J. (1995). Dissociable systems for visual recognition: A cognitive neuropsychology approach. In S. F. Kosslyn and D. N. Osherson (Eds.). *An Invitation to Cognitive Science*, 2nd edition, Volume 2., *Visual Cognition*. MIT Press. Chapter 3, pp. 101-119.
- Fiser, J., Biederman, I., & Cooper, E. E. (1997). To what extent can matching algorithms based on direct outputs of spatial filters account for human shape recognition? *Spatial Vision*, 10, 237-271.
- Garner, W. R. (1966). To perceive is to know. *American Psychologists*, 1966, 31, 11-19.
- Gauthier, I., & Tarr, M. J. (1977). Becoming a "Greeble" expert: Exploring mechanisms for face recognition. *Spatial Vision*, in press.
- Grüsser, O.-J., & Landis, T. (1991). *Visual Agnosias and Other Disturbances of Visual Perception and Cognition*. Boca Raton: CRC.
- Hosie, J. A., Ellis, H. D., & Haig, N. D. (1988). The effect of feature displacement on the perception of well-know faces. *Perception*, 17, 461-474.
- Hummel, J. E., & Biederman, I. (1992). Dynamic binding in a neural network for shape recognition. *Psychological Review*, 99, 480-517.
- Johnston, A., Hill, H., & Carman, N. (1992). Recognizing faces: effects of lighting direction, inversion, and brightness reversal. *Perception*, 21, 365-375.
- Jolicoeur, P. (1985). The time to name disoriented natural objects. *Memory & Cognition*, 13, 289-303.
- Kanwisher, N., Chun, M., M., & McDermott, J. (1996). fMRI in individual subjects reveals loci in extrastriate cortex differentially sensitive to faces and objects. *Investigative Ophthalmology & Visual Science*, 37, 193.
- Kalocsai, P., Biederman, I., & Cooper, E. E. (1994). To what extent can the recognition of unfamiliar faces be accounted for by a representation of the direct output of simple cells. *Investigative Ophthalmology & Visual Science*, 35, 1626.
- Kobatake, E., & Tanaka, K. (1994). Neuronal selectivities to complex object features in the ventral visual pathway of the macaque cerebral cortex. *Journal of Neurophysiology*, 71, 856-867.
- Lades, M., Vortbrüggen, J. C., Buhmann, J., Lange, J., von der Malsburg, C., Würtz, R. P., & Konen, W. (1993). Distortion Invariant Object Recognition in the Dynamic Link Architecture. *IEEE Transactions on Computers*, 42, 300-311.
- Logothetis, N. K., Pauls, J., Bülthoff, H. H., & Poggio, T. (1994). View-dependent object recognition by monkeys. *Current Biology*, 4, 401-414.
- Moscovitch, M., Winocur, G., Behrmann, M. What is special about face recognition? Nineteen experiments on a person with visual object agnosia and dyslexia but normal face recognition. *Journal of Cognitive Neuroscience*, 9, 555-604.

- Rolls, E. T. (1992). Neurophysiological mechanisms underlying face processing within and beyond the temporal cortical visual areas. *Philosophical Transactions of the Royal Society of London B*, 335, 11-21.
- Rhodes, G., & Tremewan, T. (1994). Understanding face recognition: Caricature effects, inversion, and the homogeneity problem. *Visual Cognition*, 1, 275-311.
- Rumiati, R., I., Humphreys, G. W., Riddoch, M., J., & Bateman, A. Visual object agnosia without prosopagnosia or alexia: Evidence for hierarchical theories of visual object recognition. *Visual Cognition*, 1, 181-225.
- Scalaidhe, P. Ó., Wilson, A. W., Goldman-Rakic, P. S. (1997). Areal segregation of face-processing neurons in prefrontal cortex. *Science*, 278, 1135-1138.
- Schiller, P. H. (1995). Effect of lesions in visual cortical area V4 on the recognition of transformed objects. *Nature*, 376, 342-344.
- Sergent, J., Ohta, S., & MacDonald, B. (1992). Functional neuroanatomy of face and object processing: A PET study. *Brain*, 115, 15-29.
- Sergent, J., Ohta, S., MacDonald, B., & Zuck, E. (1994). Segregated processing of facial identity and emotion in the human brain: A PET study. *Visual Cognition*, 1, 349-369.
- Shepard, R. N. (1964). Attention and the metric structure of the stimulus space. *Journal of Mathematical Psychology*, 1964, 1, 54-87.
- Shepard, R. N., & Cermak, G. W. (1973). Perceptual-cognitive explorations of a toroidal set of free-form stimuli. *Cognitive Psychology*, 4, 351-377.
- Subramaniam, S. & Biederman, I. (1997). Does contrast reversal affect object identification. *Investigative Ophthalmology & Visual Science*, 38, 998.
- Tanaka, J. W., & Farah, M. J. (1993). Parts and wholes in face recognition. *The Quarterly Journal of Experimental Psychology*, 46A, 225-245.
- Tootell, R. B. H., Dale, A. M., Sereno, M. I., Malach, R. (1996). New images from human visual cortex. *Trends in Neural Science*, 19, 481-489.
- Tversky, B., & Hemenway, K. (1984). Objects, parts, and categories. *Journal of Experimental Psychology: General*, 113, 169-193.
- Wiskott, L., Fellous, J-M., Krüger, N., & von der Malsburg, C. (1997). Face recognition by elastic bunch graph matching. *PAMI*, in press.
- Yin, R. K. (1969). Looking at upside down faces. *Journal of Experimental Psychology*, 81, 141-145.
- Young, A. W., Hay, D. C., McWeeny, K. H., Flude, B. M., & Ellis, A. W. (1985). Matching familiar and unfamiliar faces on internal and external features. *Perception*, 14, 737-746.
- Young, A. W., Hellawell, D., & Hay, D. C. (1987). Configural information in face perception. *Perception*, 16, 747-759.
- Young, M. P., & Yamane, S. (1992). Sparse population coding of faces in the inferotemporal cortex. *Science*, 256, 1327-1331.