

2

The Perfect Gestalt: Infinite Dimensional Riemannian Face Spaces and Other Aspects of Face Perception

James T. Townsend, Bruce Solomon,
and Jesse Spencer Smith

Indiana University

A number of papers in recent years have demonstrated that global aspects of faces can be extremely important in face perception and memory (e.g., Baenninger, 1994; Biederman & Kalocsai, 1998; Cottrell, Dailey, Padgett, & Adolphs, chap. 9, this volume; Farah, Wilson, Drain, & Tanaka, 1998; Tanaka & Sengco, 1997). Recently, Farah et al. (1998) adduced evidence that there are holistic properties of face cognition that go beyond configurational (i.e., relational) properties of features and other landmarks of faces. A longtime student of perception and philosophy of science in psychology, William Uttal, has repeatedly called for *mathematics and related psychological theories* that are suitable for capturing holistic aspects of perception (Uttal, 1988, chap. 12, this volume; see also Cottrell et al., chap. 9, this volume; Wenger & Townsend, chap. 7, this volume). It is becoming increasingly clear that no one approach could ever suffice for all aspects of face perception (e.g., Uttal, chap. 12, this volume; Wenger & Townsend, 2000). Nevertheless, we contend in this chapter that a quite natural theory immediately yields the quintessence of holism. This theory is constituted by our Riemannian face space. It is eminently holistic because each face in the theory is the entire function that is a perfect description of the perceptual

object. Each is more than the sum of the parts, in that in the space, each face is a unique point, in an analogous sense to a finite feature description that leads to a unique finite vector space representation.

The space is infinite dimensional and yet we show that the space bears potential for a number of standard and useful geometric properties. For instance, we devote considerable effort to showing that such notions as angle and distance may attend these seemingly esoteric spaces. There are various metrics that appear to be appropriate for different perceptual and cognitive tasks we discuss. Other global and local aspects of such spaces (e.g., morphing and low-dimensional subspaces) are considered. The final discussion relates our work, in a qualitative way, to the important notions of templates, prototypes, and similar concepts in categorization and identification models and experiments.

Naturally, the theory is quite general, but we believe this may be an advantage given the relatively impecunious knowledge available about these complex psychological processes at this point in time. For instance, certain simpler spaces, such as finite dimensional Euclidean spaces (e.g., Valentine, 1991; see also Busey, chap. 5, this volume; Steyvers & Busey, chap. 4, this volume; Valentine, chap. 3, this volume) or even present versus absent feature spaces (e.g., Townsend, Hu, & Kadlec, 1988), arise as special cases of our theory. As a prelude, experiments that define relevant versus irrelevant feature sets or construct alphabets based on the presence or absence of certain features (e.g., Rumelhart & Siple, 1974; Townsend, Hu, & Ashby, 1981) might eventuate in such dimensional reductions. Obviously, such a theory as we are posing here is meant to be an approximation to reality, not reality itself, in the same sense that any theoretical device is.

The representation aspects of the theory, on which we focus here, overlap some contemporary accounts in the sense that we believe three-dimensional shape to be important in face processing. A number of other theories and approaches have dealt with two-dimensional or three-dimensional actual faces. Many have contributed to the artificial intelligence aspects of face cognition, and some have been able to make predictions about human behavior. Nevertheless, geometric characterizations of face spaces almost always assume a finite dimensional, often Euclidean coordinate system. Needless to say, such approaches have led to considerable increase in our knowledge about face cognition, and Busey (chap. 5, this volume), Steyvers and Busey (chap. 4, this volume), and Valentine (chap. 3, this volume) give outstanding examples of this claim. As noted, such spaces may be realistic reductions of the more general settings.

The position that a spatial representation with any metric is appropriate for faces is worthy of quick review. Aside from the evidence that metric spatial models of face representations have been successful in modeling human performance (Johnston, Kanazawa, Kato, & Oda, 1997; Johnston, Milne, Williams, & Hosie, 1997; Valentine, 1991; Valentine & Endo, 1992), it is important to note that spaces even more general than metric spaces (e.g., uniform spaces) permit the imposition of a "pseudo-metric" in which all conditions of a metric hold, except that $d(x_1, x_2)$ may be 0 when $x_1 \neq x_2$, an argument made by Townsend and Thomas (1993). They also pointed out the interesting mathematical result that an arbitrary topology gives rise to something very close to a metric.

Even though our almost total emphasis in this study is on face processing, we do wish to perform some modest proselytization on behalf of investigation of more complex kinds of geometric and topological spaces in psychological milieus than is usually seen. There are a number of reasons to suspect that the fields of psychology and cognitive science will ultimately need a more general setting for face and other complex psychological object processing. One reason is that it seems very likely to us that general psychological spaces may require non-Euclidean (and even nonpower; see Tversky & Krantz, 1969) metrics, and possibly nonmetric descriptions (Baird, 1997; Suppes, Krantz, Luce, & Tversky, 1989). Theoretical examples are nongeometric featural spaces (e.g., Townsend & Ashby, 1982; Tversky, 1977), Riemannian spaces (e.g., Boothby, 1975; in cognitive science, Dzhafarov & Colonius, 1999; Lindman & Caelli, 1978; Townsend & Thomas, 1993), or spaces containing regularities without being so constrained that they necessarily admit a metric, such as spaces possessing affine connections (e.g., Synge & Schild, 1949; in cognitive science, D. N. Levin, 2000).

Moreover, there are quite general spaces that do admit the imposition of metrics and that allow all points to be pathconnected (e.g., in psychology, Beals & Krantz, 1967). Even where no useful metric exists, the powerful notion of ordinal similarity regulated by nested neighborhoods in some natural topology (e.g., mathematically, Kelley, 1955; Munkres, 1975) might be useful. It would seem almost inconceivable if it turned out that all of cognition required only finite, orthogonal dimensional spaces (e.g., the Euclidean or power metric), or even including such metrics as the ultrametric, associated with certain featural dominions, especially when physics itself, the most elegant of empirical sciences, is now founded on Riemannian and quasi-Riemannian (e.g., special relativity) spaces. In particular, face space, emotion space, semantic space, personality space, and so on, will

almost certainly demand more latitude of description than the usual geometries considered, especially when allied with dynamic process notions. We believe, moreover, that even though considerable progress has been made using the more circumscribed spatial tools, to continue to focus almost entirely on these would lead to unfortunate confinement of both theory and experimental designs to accommodate those relatively limited tools. Analogies abound in the voluntary incarceration within the confines of traditional statistics and hypothesis testing, to which most of us psychologists submit (e.g., Loftus, 1995, 1996; Loftus & Masson, 1994; Townsend, 1994).

It is useful to compare our situation with that of signal communication theory and control theory, which require complex, often infinite dimensional spaces. The dimensionality and geometry of such spaces are usually not obvious in elementary treatments in engineering and undergraduate mathematical courses, but they always lie in the background. The fact is that the stimulus space from which one begins the study of sensation and perception is often of infinite dimension; for instance, acoustic space, visual form and color space, and so on. Note that even color space starts out with infinite dimensionality and is then reduced to a lower number of dimensions (e.g., Suppes et al., 1989).

Furthermore, outside of the popular regions of multidimensional scaling (which again typically are limited to orthogonal or at least straight-line coordinate spaces), investigators have pretty much ignored the psychophysics of multidimensional objects. The latter usually scale a very small set of sub-dimensions of the objects. We know almost nothing about how objects like two- or three-dimensional forms as a whole are mapped into psychological space, or about their mutual discriminability. For example, does something like Weber's law hold on several dimensions simultaneously, but probably with a different constant than for dimensions taken one at a time? How do multidimensional sensitivity (e.g., Weber) functions relate to the more macroscopic sensation functions of the original stimuli (e.g., see Dzhafarov & Colonius, 1999; D. N. Levin, 2000)?

The physical characteristics of the shapes of faces can be, in fact arguably have to be, represented as surfaces in a three-dimensional space. We take as a reasonable starting point the proposition that some type of fairly smooth function or map takes the retinal impression into an object that is itself a surface in three-dimensional space. It is becoming increasingly clear that quite sophisticated and relatively global properties of three-dimensional objects must be coded very early, even retinally (e.g., Lappin, Ahlstrom, Craft, & Tschantz, 1995; Lappin & Craft, 1997). We are

aware that there is an ongoing controversy concerning the viability of three-dimensional representations as models of object perception.¹ However, on the side of the three-dimensional, we feel, is the fact that the human visual system indubitably does employ a variety of so-called monocular (e.g., linear perspective) and binocular (e.g., retinal disparity) cues to perceive a three-dimensional world. Although there is a splitting apart of the image anatomically when it reaches the visual striate cortex (Brodman's Area 17), somehow the neural connections are such that no such split is consciously available. Furthermore, perhaps a kind of teleological argument on the basis of the purpose of stereoscopic vision might not be entirely out of order. In any case, we base our discussion on the notion of a kind of isomorphic (again, in the intuitive rather than mathematical, sense) representation as a function space. We do not believe that this is the only kind of representation or processing that occurs; far from it. We believe that the brain can accommodate multiple representations simultaneously, including potential featural characterizations (e.g., the gross anatomical aspects such as eyes, nose, mouth, etc.), relational measurements (e.g., is the width separating the eyes broader than the mouth, etc.), and even the entire Gestalt. What receives further processing depends on predispositional characteristics and environmental demands.

CONSTRUCTING AN INFINITE DIMENSIONAL FACE SPACE

We now explore the idea of representing face space by a space of functions. This kind of approach provides the implicit backdrop, for instance, of the work of O'Toole and colleagues (e.g., O'Toole, Vetter, Volz, & Salter, 1997) and that of Edelman and colleagues (e.g., Edelman & Duvdevani-Bar, 1997). We might formulate O'Toole's framework for example, as follows:

A face can be represented in cylindrical coordinates by a function $r(\theta, z)$, which gives the distance of the facial point at height z and angle θ from the central vertical axis of the head. In other words, the

¹Shepard and Cermak (1973) and others (e.g., recently, Edelman, 1998) have argued in favor of a second order isomorphism, rather than a first order isomorphism. This topic goes beyond the present domain.

face is seen as the cylindrical “graph” of the function $r(\theta, z)$, defined over the rectangular domain

$$D := (0, 2\pi) \times (-1, 1),$$

which parameterizes the θ and z coordinates respectively; we have normalized the z coordinates so that no face extends more than 1 unit up or down vertically from height $z = 0$.

This representation lets us regard face space as a certain set Ω of functions on D , and we adopt this viewpoint later. Before doing so, however, we wish to stress a couple of points.

First, not all functions on D —in fact, not all continuous, or even infinitely differentiable functions on D —correspond to faces. For example, there are functions on D whose graphs (as set forth earlier) would be perceived by most of us as bas relief realizations of the Warholian message, “Campbell’s Tomato Soup.” Although perfectly respectable as functions on D , these would not belong to our face space Ω . More prosaically, neither would the zero function. In particular, face space is a proper subset of the much larger space of all functions on D , and although the latter is a vector space under the usual operations of addition and scalar multiplication of functions, our face space Ω sits inside it in some potentially complicated way, and not as a vector subspace.

Second, we emphasize that substantially different representations of face space as a function space are clearly possible—and perhaps even preferable—for particular problems. For instance, one could represent faces as “spherical” graphs, in contrast to our earlier cylindrical formulation. Quite generally, one would regard faces as differential-geometric surfaces in R^3 ; that is subsets of R^3 expressible near any point as the image of a vector function

$$F(u, v) = (x(u, v), y(u, v), z(u, v))$$

satisfying certain differentiability and nondegeneracy conditions. Indeed, both the cylindrical and spherical strategies can be seen as special cases of this more general situation, and all substantive aspects of our following explorations carry through in that generality. To simplify our exposition, however, we generally assume that some such model (e.g., the cylindrical model already outlined) has been selected, and that face space is represented by a set of functions on the rectangular domain D .

Central to our thinking about this model is the following reflection: We suspect that there is a natural topology of face space, and we would like to understand how to express it in this model.

Topologies and Metrics

Mathematically, a *topology* on any set of objects (here, the functions we use to represent faces) is simply a designation of certain subsets of these objects as so-called open subsets (see Munkres, 1975, for an excellent introduction). Openness is not a property of any particular subset in isolation; it only denotes membership in the distinguished class of “open” subsets, and the entire class must meet certain requirements to qualify as a topology. Chiefly, all unions and finite intersections of open subsets must also be open; that is, members of the distinguished class. Secondly, the entire set and the empty set must be open. Although very general, these simple axioms already let one make primitive “proximity” judgments. For instance, two objects may be considered “close” if there are comparatively few open sets that contain one, but not the other. Indeed, this much structure actually suffices for the definition of numerous critical notions: continuity, connectedness, compactness, and convergence of sequences among them.

By far the most common (and intuitive) way to impose a topology on a set of objects is to give a *metric*, a numerical distance function $d(x, y)$, defined between pairs of objects (x and y) in the set. This distance function must behave like the common notion of physical distance by satisfying three properties:

- It must be order independent: $d(x, y) = d(y, x)$, the symmetry condition.
- It must assign a positive distance to any pair of distinct objects: $d(x, y) \geq 0$, and return $d(x, y) = 0$ if and only if $x = y$.
- The distance between any two objects must never exceed the sum of their individual distances to a third object: $d(x, y) \leq d(x, z) + d(z, y)$ for any x, y , and z in the set.

This final metric axiom is known as *triangle inequality*, because it corresponds to the familiar fact that each side of a triangle is shorter than the sum of the other two.

Now, given any metric $d(p, q)$ on a space X , one can define, for any point p in that space and any $r > 0$, the ball $B(p, r)$ of radius r and

center p :

$$B(p, r) := \{q \in X \mid d(p, q) < r\}$$

In the Euclidean plane, for instance, with the usual metric

$$d\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}\right) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2},$$

a ball by this definition is just a round open disc of radius r and center p .

Indeed, in any metric space, we may designate all such balls open, for they generate (via their unions and finite intersections) a topology, known as the *metric topology* on X . In fact, the open sets of a metric topology have a special property with regard to metric balls that we need later; namely, An arbitrary set O in X will be open relative to the topology determined by a metric if and only if every point p in O forms the center of a ball $B(p, r)$ lying entirely in O for some radius $r > 0$. Although we omit the argument, this fact does follow in short order from our definitions.

Simple mathematical examples show that many topologies are not generated by any metric.² But when the topology of face space is given by an easily expressible metric, a great deal of mathematical analysis becomes available, including the notion of a geodesic (i.e., “shortest” path) joining a pair of faces, which we explore in detail later. Admittedly, the existence of a “natural” metric on face space is, for the time being, suspect. An interesting test of this assumption could be made by attempting to measure whether the property expressed by the triangle inequality seems to hold for human perception of faces. To our knowledge, no such test has yet been carried out.

Nevertheless, we consider here some ramifications of the assumption that significant aspects of human face perception can be understood as arising from a metric on face space.

The following question then arises: What sort of metric might we expect? Perhaps the simplest plausible metric on Ω would be the so-called *sup* (pronounced *soup*) metric, which we can impose if we assume that all elements of face space are represented by continuous functions on D .

²For instance, on the set $\{a, b, c\}$, the subsets $\{\}, \{b\}, \{c\}, \{a, b\}, \{b, c\}$, and $\{a, b, c\}$ comprise a topology (check their unions and intersections). However, no metric generates this topology, because no ball centered at a excludes b and c —if it did, $\{a\}$ would constitute a ball, and hence be open.

Formulaically, this metric is given by

$$d(f, g) := \sup \{|f(x) - g(x)| \mid x \in D\}.$$

Intuitively, this just assigns the distance between two faces to be the largest pointwise difference between their defining functions. This is a simple notion of distance, but from an analytical standpoint, the sup metric is not amenable to the sorts of mathematical tools that have been used by O'Toole, Edelman, and others, because it is not compatible with any inner product. We elaborate on inner products later; they play a key role in our subsequent discussion.

Leaving inner products temporarily aside, however, we find a serious defect in the simple sup metric: Important distinctions between faces—gender and aging, for example—seem to result not from large localized (i.e., pointwise) differences, but rather from small differences that are more evenly distributed over the whole face.

Fortunately, there are simple metrics that unlike the sup metric, treat small but distributed differences as significant. In particular, we have in mind metrics of L^p type. The simplest exemplar of this class is given by the formula

$$d(f, g) := \left(\int_D |f(x) - g(x)|^p w(x) dx \right)^{1/p} \quad (1)$$

where $w(x)$ is a function that may weight the relative contribution of specific areas of the face. Because we do not expect the weighting to change precipitously over an infinitesimal distance, we may assume $w(x)$ is continuous. Normally we require $w(x) > 0$, but we note in passing that were $w(x)$ to vanish away from the main facial features (e.g., eyes, nose, mouth), the metric would then degenerate into a finite feature-weighted metric. This could model human perception performance in discrimination experiments where the task is to say “same” if and only if the two faces are identical in every one of the main features.

In such a setting, one can additionally measure the “size” of any individual function by determining its metric distance from the zero function, known technically as its *w-weighted L^p norm*:

$$\|f\| := \left(\int_D |f(x) - 0|^p w(x) dx \right)^{1/p} = \left(\int_D |f(x)|^p w(x) dx \right)^{1/p}$$

Notice that with this notation, the distance from f to g reduces to the norm

of $f - g$; that is,

$$d(f, g) = \|f - g\|$$

At this juncture, we find ourselves in the position to refine our earlier ideas. Until now, we have regarded face space Ω as a subset of the set of all functions on D . But the L^p metric just proposed helps us narrow our focus a bit; it only makes sense for functions on D that have finite p -power integrals; that is, functions in the space known as

$$L^p(D) := \left\{ \int_D |f(x)|^p w(x) dx < \infty \right\}$$

where f is a function of the type defined earlier.

The restriction we admit by redefining our universe of functions as $L^p(D)$ instead of all functions on D is really no restriction at all, however, because functions that do yield recognizable faces are likely to be nearly continuous. Indeed, except possibly for jumps across certain contours (e.g., the nostrils), one expects actual continuity. Wildly chaotic functions, or functions with huge spikes, and so forth, would seem to play no role in the theory: Surely we can safely exclude functions outside $L^p(D)$ from consideration.

Two more facts about $L^p(D)$, established in any graduate analysis text, will prove useful as we proceed. First, as is true for the space of all functions on D , $L^p(D)$ is a vector space, albeit of infinite dimension. That is, linear combinations of functions in $L^p(D)$ remain in $L^p(D)$. A bit more subtly, $L^p(D)$ is complete. That is, when a certain type of sequence of functions, all belonging to $L^p(D)$, converge to a limit using a weighted L^p metric, the limit will also belong to $L^p(D)$. This property—completeness—will play an especially important role next.

L^2 and Inner Products

In the special case $p = 2$, our functional “universe” $L^p(D)$ acquires an additional key feature that enables one to measure not only distances but “angles” as well, and brings a wealth of useful mathematical techniques into play. We refer to the *inner product* mentioned earlier, which becomes possible here because the product of two functions in $L^2(D)$ is guaranteed to be integrable. This fact—which fails for all $p \neq 2$ —guarantees that for any continuous, bounded, positive weighting function w , the integral

defining the following pairing converges:

$$\langle f, g \rangle := \int_D f(x)g(x)w(x) dx \quad (2)$$

When $w(x) = 1$, this pairing \langle, \rangle is called the *standard inner product on $L^2(D)$* . It constitutes an inner product for any strictly positive continuous weighting function $w(x)$ because it satisfies the following three axioms:

- It is *symmetric*; that is, $\langle f, g \rangle = \langle g, f \rangle$ for all f and g in $L^2(D)$.
- It is *bilinear*, meaning linear in both slots, for example,

$$\langle a_1 f_1 + a_2 f_2, g \rangle = a_1 \langle f_1, g \rangle + a_2 \langle f_2, g \rangle$$

for any scalars a_1 and a_2 and any f_1, f_2 , and g in $L^2(D)$.

- It is *positive definite*, meaning that $\langle f, f \rangle > 0$ for any nonzero f in $L^2(D)$.

These properties all follow very easily from the pairing's definition, and indeed, with regard to positive definiteness, we see that the ever-positive norm of a function f in $L^2(D)$ can be defined by the inner product:

$$\|f\| := \sqrt{\langle f, f \rangle}$$

Thus equipped, $L^2(D)$ becomes a *Hilbert space*—a complete inner product space. Indeed, L^2 spaces provide the most important examples of Hilbert spaces, for which a great deal of theory exists. We state here one fairly elementary result from this theory that will serve us well later on:

Theorem (Riesz Representation): *Any continuous linear function λ on a Hilbert space \mathcal{H} may be represented as “pairing with some vector” relative to the inner product on \mathcal{H} . That is, there is a vector v_λ in \mathcal{H} such that for all \bar{v} in \mathcal{H} , we have*

$$\lambda(\bar{v}) = \langle \bar{v}, v_\lambda \rangle.$$

The key application of this theorem for our purposes is the following:

Corollary: *Any inner product \langle, \rangle^* on $L^2(D)$ can be written in terms of the standard one \langle, \rangle , and some positive symmetric linear operator*

W as follows:

$$\langle f, g \rangle^* = \langle f, W(g) \rangle \quad \text{for all } f, g \text{ in } \mathbf{L}^2(D)$$

By positive symmetric operator, we mean a continuous linear transformation $W : \mathbf{L}^2(D) \rightarrow \mathbf{L}^2(D)$ that is symmetric ($\langle f, W(g) \rangle = \langle g, W(f) \rangle$) and positive definite (i.e., $\langle f, W(f) \rangle > 0$ for all nonzero f in $\mathbf{L}^2(D)$).

The proof of this corollary is short and easy, but we omit it in favor of a key example. Namely, we can represent any w -weighted inner product as in Equation 2 by using the operator $W(f) := wf$ (multiplication of $f(x)$ by $w(x)$). This multiplication-by- w operator W is a linear transformation, mainly because multiplication distributes over addition. It is clearly symmetric too, as

$$\begin{aligned} \langle W(f), g \rangle &:= \langle wf, g \rangle \\ &= \int_D w(x)f(x)g(x) dx = \int_D f(x)w(x)g(x) dx \\ &= \langle f, wg \rangle =: \langle f, W(g) \rangle \end{aligned}$$

Using the fact that $w(x) > 0$ for all x in D , the reader will also verify that the operator W is positive definite.

We next explain how an inner product yields a notion of the angle between two functions in $\mathbf{L}^2(D)$. Indeed, one simply defines the angle θ between two functions f and g in $\mathbf{L}^2(D)$ by the formula

$$\cos(\theta) := \frac{\langle f, g \rangle}{\|f\| \|g\|}$$

This definition exploits the *Cauchy–Schwarz inequality* $|\langle f, g \rangle| \leq \|f\| \|g\|$ to ensure that the right-hand side never exceeds 1 in absolute value, and hence forms the cosine of a unique angle $0 \leq \theta \leq \pi$. Angles defined in this manner are compatible with the metric notion of length in every expected way; for example, the usual law of cosines relating the angles and side lengths of triangles will hold, just as in the plane. If these facts and formulas seem familiar, it is because they all arise for the standard “dot product” for vectors in Euclidean space. Indeed, the dot product was the original phenomenon motivating inner product axioms.

From a face space perspective, the notion of angle would be useful, for instance, if we wanted to measure the orthogonality of gender and aging axes in face space.

POSSIBLE GLOBAL PROPERTIES
OF FACE SPACE AND RIEMANNIAN
METRICS IN FUNCTION SPACES

With these more precise notions about mathematical context in place, we return to our discussion of face space, and begin to consider the topological and geometric situation it has within $L^2(D)$. For a start, we assert:

Proposition: *Given any metric of the form (1), the face space Ω forms an open subset of $L^2(D)$.*

This proposition is really just a careful formulation of the heuristic notion that when many humans perceive the cylindrical graph of some f in $L^p(D)$ as a face, then surely many will similarly perceive the same for any function g in $L^p(D)$ that only differs from f by a tiny amount. To make this intuition precise, however, one has to measure the meaning of “tiny” here, using the metric under consideration. That is, we must argue that given any face f in Ω , one can find a sufficiently small radius $r > 0$ such that not only f , but the entire metric ball $B(f, r)$ belongs not just to $L^p(D)$, but in fact to Ω . However, this seems inherent in the nature of human perception; we are simply asserting the following:

If the graph of f is widely recognizable as a face, then there is a “noise tolerance” $r > 0$ such that any function of the form $f + \phi$ will be recognizable as a “noisy” image of the face f , provided the noise ϕ is sufficiently small, in the sense that $\|\phi\| < r$.

In more familiar terms, the openness claimed by this proposition amounts to the notion that face space sits inside the infinite dimensional vector space $L^p(D)$ as a “blob” of full dimension, as opposed to a “surface-like” object of lower dimension. By way of analogy, the air in a balloon in three-dimensional space forms an open set, as opposed to the lower dimensional surface of the balloon, which does not.

Such information carries very little detail as to the shape of face space, however, and one could ask for much more. For instance, one might hope, for the sake of simplicity, that face space is linearly convex. To explain this, we require the notion of a *line* in L^p . In Euclidean space, the line joining two points p and q can be defined as the set

$$\{tp + (1 - t)q \mid t \in \mathbf{R}\}$$

For any value of the real parameter t , one gets in this way a point on the line. In fact, p and q themselves occur for the parameter values 1 and 0, respectively, and the interval of t values between 0 and 1 construct precisely the line segment between p and q .

Exactly the same construction can be made in our setting. Any two functions f and g in $L^p(D)$ determine a line segment, namely, the one-parameter family of functions

$$\text{seg}[f, g](t, x) := tg(x) + (1 - t)f(x), \quad 0 < t < 1.$$

A subset of $L^p(D)$ is now called *convex* if, whenever it contains two functions f and g , it also contains all of $\text{seg}[f, g]$. Notice how much this implies about the shape of a set. In three-dimensional Euclidean space, for instance, a solid object fails to be convex as soon as it has any indentations in its boundary, any holes drilled through it, or any internal “air pockets.” Solid spheres and cubes are convex; bowling balls, doughnuts, and Swiss cheese—indeed, most everyday objects—are not.

Nevertheless, it seems reasonable to conjecture that face space is convex as a subset Ω of $L^p(D)$. From a perceptual standpoint, the question boils down to the following: Suppose we have two “faces” $f, g \in \Omega$. Then, as we vary t from 0 to 1, $\text{seg}[f, g](t)$ would seem to provide a simple direct “morphing” from one face into the other, but does each of the intermediate functions $\text{seg}[f, g](t)$ on this segment look like a face? If the answer is yes (or at least yes subject to some simple caveats), it would say quite a bit about the geometry and topology of face space, as the earlier examples in Euclidean three-dimensional space suggest.

Whether or not Ω is convex in $L^p(D)$, the question arises as to whether line segments, economical as they are, form the most psychologically direct morphing from one face into another. Of course, this begs the question of what one means by “direct.” We propose an answer: Direct means shortest, as measured by a psychologically meaningful metric. Even this involves some subtlety, however. A metric allows one to directly compare distances between points, but how can one measure lengths of competing paths to declare one a shorter route from f to g than the other? To answer this, we turn to the mathematics of differential geometry, which has developed a substantial toolkit for precisely this type of problem.

The basic idea is simple. A path ϕ joining one face to another (say f to g) within Ω is a continuous one-parameter family of functions ϕ_t , with $\phi_t(x)$ defining a face for each $0 \leq t \leq 1$, and in particular, with $\phi_0 = f$ and $\phi_1 = g$. We can compute an approximate length for ϕ by choosing many

(say N) equally spaced intermediate times t_i in the parameter interval $[0, 1]$, that is, $0 = t_0 < t_1 < t_2 < \cdots < t_N = 1$, and then adding up the stepwise distances, *i.e.*

$$\text{length}(\phi) \approx \sum_{i=1}^N \|\phi_{t_i} - \phi_{t_{i-1}}\|$$

Under reasonable conditions, these approximations converge to a limit as $N \rightarrow \infty$. In particular, we feel it reasonable to assume that ϕ varies “smoothly” in the sense that it is differentiable, the velocity vector $\dot{\phi}_t$, defined for each x in D via

$$\dot{\phi}_t(x) := \frac{d\phi_t(x)}{dt} = \lim_{h \rightarrow 0} \frac{\phi_{t+h}(x) - \phi_t(x)}{h}$$

exists. This velocity vector is, like ϕ_t itself, a well-defined function of x , and belongs to $\mathbf{L}^p(D)$ for each t in the interval $[0, 1]$. In this case, the length approximations just described converge to an integral, providing the following length formula:

$$\text{length}(\phi) = \int_0^1 \|\dot{\phi}_t\| dt \quad (3)$$

To make this more concrete, we illustrate by computing the length of the simplest type of path we know, namely the line segment $\text{seg}[f, g]$ joining a face f to a face g in $\mathbf{L}^p(D)$. As in the Euclidean analogue, this curve has an unchanging velocity vector. For, using the formula given earlier for $\text{seg}[f, g]$, we note that $f(x)$ and $g(x)$ are, for any fixed x in D , constants relative to t , so that elementary calculus gives

$$\begin{aligned} \dot{\phi}_t(x) &= \frac{d}{dt}(\text{seg}[f, g](t)(x)) \\ &= \frac{d}{dt}(tf(x) + (1-t)g(x)) = f(x) - g(x) \end{aligned}$$

As this is true for every x in D , we conclude that $\dot{\phi}_t = f - g$, independently of t . That is, the velocity vector of our line segment from f to g is, at each time t , the difference from function $f - g$. And because $\mathbf{L}^p(D)$ is closed under addition, this velocity, like f and g themselves, lies

in $L^p(D)$. Hence line segments in $L^p(D)$ are always differentiable, and we can compute

$$\begin{aligned} \text{length}(\text{seg}[f, g]) &= \int_0^1 \|f - g\| dt \\ &= \|f - g\|, && \text{(because } \|f - g\| \text{ is constant} \\ & && \text{with respect to } t\text{).} \\ &= d(f, g) \end{aligned}$$

Trivial as it may seem, this computation does establish an important fact: Just as in Euclidean space, the length of a line segment exactly equals the metric distance between its endpoints. Moreover, again as in Euclidean space, no path ψ_t from f to g is shorter than a line segment. Indeed, the following computation shows that every such path ψ_t has length at least $d(f, g)$. For, suppose that ψ_t maps some interval $a \leq t \leq b$ into $L^p(D)$, with $\psi_a = f$ and $\psi_b = g$. Then

$$\begin{aligned} \text{length}(\psi_t) &= \int_a^b \|\dot{\psi}_t\| dt \\ &\geq \left\| \int_a^b \dot{\psi}_t dt \right\| && \text{(limit of the triangle inequality for} \\ & && \text{integrals)} \\ &= \|\psi_b - \psi_a\| && \text{Fundamental Theorem of Calculus} \\ &= \|g - f\| \\ &= d(f, g). \end{aligned}$$

In sum, we see that for the weighted L^p type metric given in Equation 1, line segments do provide the the shortest route from one face to another.

We suspect, however, that a good model for face space will require a more general type of metric, for which line segments may not provide the shortest routes. We have in mind here the *Riemannian metrics* that are fundamental to differential geometry and find wide application in physics at both the macro (general relativity) and micro (quantum dynamic) level. Indeed, it would be very surprising if a metric as simple as the weighted L^p metric (Equation 1) did accurately measure psychological distances in face space. We say this because individuals seem to measure small

distinctions between faces differently in some regions of face space than in others.

For example, consider two sets of identical twin girls. Represent the faces of one pair by the function ϕ_1 and ϕ_2 in Ω , and the faces of the other by g_1 and g_2 . We clearly expect that the distances $d(\phi_1, \phi_2)$ and $d(g_1, g_2)$ will both be very small as measured by a metric of the type (Equation 1) we have been discussing. Assuming the families are not acquainted with each other, however, we expect that the mother of the first pair will perceive a far more substantial difference between her own daughters' faces than she will between the other twins' faces. The perceived difference is clearly distorted relative to the physical and mathematical difference. Similarly, many Whites seem less sensitive to differences between Asian faces than they are to differences of a similar physical magnitude in White faces, and vice versa. Again, this indicates a perceptual metric that magnifies or shrinks relative to the simple metric (Equation 1) in different regions of face space.

The perceptual metric probably even varies differently in different directions even in the vicinity of any one particular face. For example, suppose f_{25}^m and f_{25}^w in Ω represent the faces of a 25-year-old female impersonator before and after preparing for a show. The success of this impersonator would seem to depend heavily on the possibility of inducing a large difference in an observer's perceptual response with the aid of small physical changes in the actual topography of his face. But consider f_{65}^m , which, say, represents the off-stage face of this same individual at age 65. As measured by a metric of the form in Equation 1, we can easily imagine that the physical distance $d(f_{25}^m, f_{65}^m)$ might substantially exceed $d(f_{25}^m, f_{25}^w)$ while an observer might perceive just the opposite.

A Riemannian metric, however, can easily model both these types of perceptual distortion. It makes the assignment of distances between faces in Ω secondary; rather, a Riemannian metric first norms all velocity vectors to paths through Ω . This device enables measurement of distances—turning Ω into a metric space—because it allows us to invoke the length formula in Equation 3, integrating velocity norms to assign lengths to paths. One then measures the distance between two points (two faces in our context) as the length of the shortest path connecting them.³ As mentioned earlier, these shortest paths generally do not follow straight line segments. They

³Actually, there may not be any "shortest" path joining two points, just as there is no "smallest" number in the interval $0 < x < 1$. In this case, the distance between two points is defined as the *infimum* (greatest lower bound) of all lengths of connecting paths. (For example, the infimum of the set $\{x \mid 0 < x < 1\}$ is 0.) We ignore this technical point in our exposition; doing so does not compromise our work here in any substantive way.

are called *geodesics*, and can be found by solving a differential equation that we derive later.

To do so, we need to describe more precisely the form a Riemannian metric takes on our face space Ω . As indicated previously, such a metric operates by norming velocity vectors. It accomplishes this by introducing, for each and every f in Ω , an inner product on the set of all velocity vectors based at f . The key new feature here is that this inner product can vary from one base point to another. This allows the metric to vary in the ways we have been discussing.

To carry this out, we begin by requiring $p = 2$, for the reasons mentioned in our earlier discussion of inner products. Recall that in this case, the velocity of a differentiable curve ϕ_t of faces that, say, passes through f at the time $t = 0$, is again a function in $L^2(D)$.⁴ It is easy to show that, no matter what f is, every function in $L^2(D)$ forms the velocity of some path through f . For example, if g is in $L^2(D)$, then g is the velocity of the path $f + tg$ at time $t = 0$ and where the velocity is evaluated at the point (face) f . So the set of all velocities of paths through f comprise an entire copy of $L^2(D)$ based at f , called the *tangent space to Ω at f* , denoted $T_f\Omega$.

We can now articulate our setup more clearly: A Riemannian metric on Ω is a differentially varying assignment of an inner product $\langle \cdot, \cdot \rangle_f$ to each tangent space $T_f\Omega$.

For instance, suppose we associate to each face f in Ω a weighting function, w_f , that we allow to vary in a continuous way, depending on the base point f . We can then measure the inner product of any pair of velocities (say ϕ and ψ) based at f —or the norm of either one—using formulas of the type seen previously (e.g., Equation 2). That is,

$$\begin{aligned} \langle \dot{\phi}, \dot{\psi} \rangle_f &:= \int_D \dot{\phi}(x)\dot{\psi}(x)w_f(x) dx & (4) \\ \|\dot{\phi}\|_f &:= \langle \dot{\phi}, \dot{\phi} \rangle_f \end{aligned}$$

Using this f -varying norm, we can still assign lengths to paths by introducing the position subscript into our previous length formula:

$$\text{length}(\phi_t) := \int_a^b \|\dot{\phi}_t\|_{\phi_t} dt$$

⁴To make this more explicit, we could use the notion $\phi(x, t) := \phi_t(x)$; i.e., $\phi(x, t)$ is, for each fixed t , the “face” given by ϕ_t at time t . In this notation, the velocity $\dot{\phi}(t)$ becomes $\dot{\phi}(x, t)$, which, at any fixed time t is just another function of x . Our definition of differentiability of ϕ_t requires that it belong to $L^2(D)$ at each time t , as opposed to being “just any” function.

In this way, our Riemannian metric $\langle \cdot, \cdot \rangle_f$ yields a metric space metric: We define the distance between any pair of faces f and g to be the infimum or greatest lower bound of lengths of all paths connecting them.

By way of illustration, consider this construction in light of the mother of twins mentioned earlier. In the vicinity of her daughters' and other family members' faces, we expect her to use relatively large weighting functions w_f . Perhaps she applies a somewhat smaller weighting function at the faces of her friends. Of course, we are not suggesting that such weighting is necessarily conscious, although this might be the case in certain situations. If so, she would be more likely to agree with statements like "This face looks like your friend Doris" than she would be to "This face looks like your daughter Michelle," given a physical (i.e., L^2) difference of a similar magnitude because mathematically similar paths joining faces to Doris would all seem longer than equivalent paths emanating from Michelle. Shown two similarly distinct faces belonging to members of a racial group with which she has little contact, she might even have trouble telling them apart, her weighting functions in that region of face space being much smaller still.

Before proceeding further, we note that besides making distance measurements available, a Riemannian metric allows measurements of angle, speed, and a variety of curvature quantities, including the fairly recondite tensor-theoretic curvatures associated with Einstein's general theory of relativity. We do not define these invariants here, and the extent to which they might prove useful in characterizing perceptual phenomena remains to be seen. However, recent work by D. N. Levin (2000) suggests that even tensor-theoretic curvatures may be measurable in a perceptual context, and even useful. In particular, Levin theorized that curvature can encode and quantify the ways in which the subjective experience of a given objective stimulus domain varies from individual to individual.

THE SHORTEST PATHS ARE GEODESICS IN FACE SPACE

Let us return now to the exploration of shortest paths in face space, as manifest in the context of a Riemannian metric; that is, the geodesics mentioned earlier. As promised earlier, we now derive the differential equation defining geodesics in this setting.

In other words, we seek conditions under which a differentiable path of "faces" ϕ_t in Ω may be recognized as a shortest path connecting its

endpoints. To simplify the ensuing calculations, we assume that ϕ_t has *unit speed*; that is, that at each point along ϕ_t , the velocity vector $\dot{\phi}_t$ has unit norm $\|\dot{\phi}_t\| \equiv 1$. This assumption is a standard mathematical move; one can always realize it by using the one-dimensional Implicit Function Theorem to “reparametrize” the path, which, conceptually, simply amounts to smoothly slowing ϕ_t down or speeding it up to make the speed constantly 1 at all times. This may entail expanding or shrinking the time parameter interval, and once it is done, the length of the path and the length of the parameter interval—call it I —coincide, with the former being given by our length integral:

$$\text{length}(\phi_t) = \int_I \sqrt{\langle \dot{\phi}_t, \dot{\phi}_t \rangle_{\phi_t}} dt = \int_I \|\dot{\phi}_t\|_{\phi_t} dt$$

Note that the inner product $\langle \cdot, \cdot \rangle_{\phi_t}$ appearing in the integral here may itself be given by some complicated rule; for example, an integral against some $w_f = w_{\phi_t}$, as proposed earlier. Such details turn out to be irrelevant to the calculation we are about to exhibit, however. In fact, we can proceed virtually identically with the analogous calculation for characterizing shortest paths in a finite-dimensional physical space for this reason.

The key idea that gives a differential equation for geodesics is now this: We regard the length formula as defining a function (often called a *functional* in this context) on the enormous space of all paths connecting the endpoints of ϕ_t . If this length functional actually takes a minimum on our path ϕ_t , then—just as in single variable calculus—its “derivative” at ϕ_t must vanish. This is the standard approach via calculus of variations, which we make precise here as follows.

Let ϵ_t be any differentiable path in $L^2(D)$ that is defined on the same parameter interval I as ϕ_t , but vanishes at both its endpoints; that is, ϵ_t begins and ends at the zero function. We can use ϵ_t to perturb the path ϕ_t while leaving its endpoints fixed. Indeed, for any sufficiently small number s , the formula

$$\phi_{t,s} := \phi_t + s \epsilon_t$$

constructs a new path having the same endpoints as ϕ_t (because ϵ_t vanishes at the endpoint values of t), and elsewhere very near ϕ_t . With regard to this nearness, we can also be sure that for s sufficiently small, the new path $\phi_{t,s}$ stays in face space Ω . This follows from our earlier proposition, namely that Ω is open in $L^2(D)$, so that every point of ϕ_t forms the center of some open ball in Ω . For small enough s , addition of $s \epsilon_t$ will not perturb ϕ_t

outside that ball. Now, however, we have really boiled things down to a single variable calculus problem. The one-variable function

$$\lambda(s) := \text{length}(\phi_{t,s}) = \text{length}(\phi_t + s\epsilon_t)$$

clearly attains a minimum when $s = 0$, the unperturbed path ϕ_t being, by hypothesis, no longer than any other path having the same endpoints. Consequently, if ϕ_t is a geodesic, then

$$\left. \frac{\partial \lambda(s)}{\partial s} \right|_{s=0} = 0 \quad (5)$$

This relation encodes, although vaguely, the differential equation we seek. It remains to unravel it into a condition more explicitly concerned with the path ϕ_t and the Riemannian metric.

We make our first step in this direction by unpacking $\lambda(s)$. It represents, after all, the length integral, so Equation 5 becomes:

$$\left. \frac{\partial}{\partial s} \right|_{s=0} \int_I \sqrt{\langle \dot{\phi}_{t,s}, \dot{\phi}_{t,s} \rangle_{\phi_{t,s}}} dt = 0,$$

where

$$\dot{\phi}_{t,s} := \frac{\partial}{\partial t} \phi_{t,s} = \frac{\partial}{\partial t} (\phi_{t,s} + s\epsilon_t) = \dot{\phi} + s\dot{\epsilon}_t$$

At each point along the path $\phi_{t,s}$, our Riemannian metric provides the inner product $\langle \cdot, \cdot \rangle_{\phi_{t,s}}$ used there, as the subscript indicates. It benefits us to view this slightly differently, via the representation result mentioned in our earlier discussion of inner products. We therefore rewrite the inner product in terms of the representative symmetric operator $W_{\phi_{t,s}}(f)$ characterized by the equation

$$\langle f, g \rangle_{\phi_{t,s}} = \langle W_{t,s}(f), g \rangle,$$

where $\langle \cdot, \cdot \rangle$ denotes the standard L^2 inner product on $T_{\phi_{t,s}}\Omega$ for each time t . This rewrites the previous equation as

$$\left. \frac{\partial}{\partial s} \right|_{s=0} \int_I \sqrt{\langle W_{t,s}(\dot{\phi}_{t,s}), \dot{\phi}_{t,s} \rangle} dt = 0.$$

It is justifiable here to interchange the integration and differentiation, but doing so, taking the resulting derivative and using the product rule (which

works in an inner product just as well as with a product of ordinary functions), produces an unwieldy-looking result:

$$\int_I \frac{\left\langle \frac{\partial}{\partial s}(W_{t,s}(\dot{\phi}_{t,s})), \dot{\phi}_{t,s} \right\rangle + \left\langle W_{t,s}(\dot{\phi}_{t,s}), \frac{\partial}{\partial s} \dot{\phi}_{t,s} \right\rangle}{2\sqrt{\langle W_{t,s}(\dot{\phi}_{t,s}), \dot{\phi}_{t,s} \rangle}} \Big|_{s=0} dt = 0$$

Fortunately, our unit-speed assumption now simplifies matters, because when we set $s = 0$, as this formula instructs us to do, the inner product under the radical in the denominator is just the norm of $\dot{\phi}_t$, and hence identically 1. This fact simplifies the denominator to the scalar constant 2.

We can also simplify the numerator by noting that when $s = 0$, we have $\dot{\phi}_{t,s} = \dot{\phi}_t$, and $W_{t,s} = W_t$. Moreover,

$$\frac{\partial}{\partial s} \dot{\phi}_{t,s} \Big|_{s=0} = \frac{\partial}{\partial s} (\dot{\phi}_t + s\dot{\epsilon}_t) \Big|_{s=0} = \dot{\epsilon}_t$$

So our geodesic condition (Equation 5) reduces considerably; we now have

$$\frac{1}{2} \int_I \left[\left\langle \frac{\partial}{\partial s}(W_{t,s}(\dot{\phi}_t)), \dot{\phi}_t \right\rangle \Big|_{s=0} + \langle W_t(\dot{\phi}_t), \dot{\epsilon}_t \rangle \right] dt = 0$$

We omit the routine limiting argument, similar to the one that proves the product rule in elementary calculus, that evaluates the remaining s -derivative to yield

$$\frac{\partial}{\partial s}(W_{t,s}(\dot{\phi}_{t,s})) = \frac{\partial W}{\partial \epsilon_t}(\dot{\phi}_t) + W_t(\dot{\epsilon}_t),$$

where $\partial W/\partial \epsilon_t$ denotes the result of evaluating the differential ∂W on the vector ϵ_t at ϕ_t . At each time t , ∂W denotes the linear approximation to the operator field W at the corresponding point along ϕ_t . The domain of ∂W , at time t , is $T_{\phi_t}\Omega$, and its range is the space of all symmetric operators. We emphasize that it is a linear mapping; this is important later. In this particular instance, the vector of which it is a linear function is ϵ_t .

Alternatively, one may regard $\partial W/\partial \epsilon_t$ as the “directional derivative of W in the ϵ_t direction.” Note that this description makes sense: W is the field of symmetric operators that define the Riemannian metric at each point of Ω , the latter being a space of functions. We can thus differentiate this field in the direction of ϵ_t , which, at each time t is itself a function in $L^2(D)$,

and thereby a vector in the tangent space $T_\phi \Omega$, as we have identified the latter space as a copy of $L^2(D)$.

In any case, when we insert this result in our geodesic condition and exploit the symmetry of W_t , it becomes

$$\begin{aligned} & \frac{1}{2} \int_I \left[\left\langle \frac{\partial W}{\partial \epsilon_t}(\dot{\phi}_t), \dot{\phi}_t \right\rangle + \langle W_t(\dot{\epsilon}_t), \dot{\phi}_t \rangle + \langle W_t(\dot{\phi}_t), \dot{\epsilon}_t \rangle \right] dt \\ &= \int_I \left[\frac{1}{2} \left\langle \frac{\partial W}{\partial \epsilon_t}(\dot{\phi}_t), \dot{\phi}_t \right\rangle + \langle W_t(\dot{\phi}_t), \dot{\epsilon}_t \rangle \right] dt \\ &= 0 \end{aligned}$$

To finish implementing the variational strategy, we need to make the integrand here involve only the perturbation ϵ_t , as opposed to its velocity $\dot{\epsilon}_t$. To do this, we first compute

$$\frac{\partial}{\partial t} \langle W_t(\dot{\phi}_t), \epsilon_t \rangle = \left\langle \frac{\partial W}{\partial \dot{\phi}_t}(\dot{\phi}_t), \epsilon_t \right\rangle + \langle W_t(\ddot{\phi}_t), \epsilon_t \rangle + \langle W_t(\dot{\phi}_t), \dot{\epsilon}_t \rangle.$$

(We again omit the limiting argument that justifies applying product rule here.) The last term on the right here coincides with the last term under our previous integral, so we may replace the latter by the other three terms here; that is, integrate by parts:

$$\begin{aligned} & \int_I \frac{1}{2} \left\langle \frac{\partial W}{\partial \epsilon_t}(\dot{\phi}_t), \dot{\phi}_t \right\rangle + \frac{\partial}{\partial t} \langle W_t(\dot{\phi}_t), \epsilon_t \rangle - \left\langle \frac{\partial W}{\partial \dot{\phi}_t}(\dot{\phi}_t), \dot{\epsilon}_t \right\rangle \\ & \quad - \langle W_t(\ddot{\phi}_t), \epsilon_t \rangle dt = 0 \end{aligned}$$

But we can throw out the second term here, because by the Fundamental Theorem of Calculus, $\int_a^b \frac{du}{dt} dt = u(b) - u(a)$ for any differentiable function u , and in this situation, with the role of u played by $\langle W_t(\dot{\phi}), \epsilon_t \rangle$, we have $u(b) = u(a) = 0$, because our perturbation ϵ_t vanishes at both ends of I . So we are left with

$$\int_I \frac{1}{2} \left\langle \frac{\partial W}{\partial \epsilon_t}(\dot{\phi}), \dot{\phi} \right\rangle - \left\langle \frac{\partial W}{\partial \dot{\phi}_t}(\dot{\phi}), \epsilon_t \right\rangle - \langle W_t(\ddot{\phi}_t), \epsilon_t \rangle dt = 0,$$

and each of the three terms in the integrand now depend linearly on ϵ_t . This is perhaps less obvious in the first term, but as noted earlier, it does hold there.

We now invoke the Riesz Representation Theorem again, which says that any linear functional on $L^2(D)$ (or indeed on any Hilbert space) can be

represented as an inner product with some element of the space. We apply this to the first two terms, thereby defining a vector field $\vec{\Gamma}$ along the path ϕ_t , by requiring that for all ϵ in $T_\phi\Omega$,

$$\langle \vec{\Gamma}(\dot{\phi}_t), \epsilon \rangle = \frac{1}{2} \left\langle \frac{\partial W}{\partial \epsilon}(\dot{\phi}_t), \dot{\phi}_t \right\rangle - \left\langle \frac{\partial W}{\partial \dot{\phi}_t}(\dot{\phi}_t), \epsilon \right\rangle$$

That is, at each face along our path ϕ_t , the linear function of ϵ on the right is represented by taking the inner product of ϵ with the vector $\vec{\Gamma}(\dot{\phi}_t)$ on the left. We write $\vec{\Gamma} = \vec{\Gamma}(\dot{\phi}_t)$ to emphasize the dependence of $\vec{\Gamma}$ on $\dot{\phi}_t$, which, it is worth noting, is quadratic.

With this notation, our integral condition for a geodesic, written originally as Equation 5, takes the following far more specific form:

$$\int_I \langle \vec{\Gamma}(\dot{\phi}_t) - W_t(\ddot{\phi}_t), \epsilon_t \rangle dt = 0$$

We can now draw the conclusion we have been aiming toward: Because the integral vanishes for every differentiable perturbation ϵ_t on the interval I , the quantity pairing with ϵ_t under the integral must vanish identically.

For, if $\vec{\Gamma}(\dot{\phi}) - W_t(\ddot{\phi}_t)$ did not vanish identically, we could choose a perturbation ϵ_t that coincided with it over the entire interval I , excluding an arbitrarily small neighborhood of the endpoints I , where our perturbations are required to vanish. The integrand would then be (almost completely) the inner product of a nonzero function with itself, hence strictly positive, because inner products are positive definite. This would make the integral positive—in particular, nonzero—that by assumption happens for no perturbation whatsoever.

We can therefore state one version of the conclusion we seek: A unit-speed path ϕ_t in face space is geodesic only if

$$W_t(\ddot{\phi}_t) - \vec{\Gamma}(\dot{\phi}_t) = 0 \quad \text{for all } t$$

This is a second order differential equation for the path ϕ_t . The quantity on the left-hand side is called the *geodesic curvature* (or *geodesic acceleration*) vector of ϕ_t . It measures the deviation from perceptual straightness perceived by someone watching ϕ_a morph into ϕ_b along the path ϕ_t , given that they use the Riemannian metric defined by the operator field W to process their experience. So our conclusion amounts to the following: When geodesic curvature of a unit-speed path ϕ_t vanishes, that path is a geodesic, perceived as being straight, and having constant speed, by the observer.

EXPRESSING FACES AND GEODESICS IN TERMS OF BASIS FUNCTIONS

The preceding formulae are admittedly abstract, and we now seek to bring them down to earth. In practical work with function spaces—as with all vector spaces—one generally selects a complete set of *basic* functions $\{u_i\}$, where the indexes i run over a set I that may be finite, countably infinite, or even uncountably infinite, according to the dimension of the function space in question. Every other function f in the space can then be written as a linear combination of the basis “vectors”: $f = \sum_i c_i u_i$. The vector (c_1, c_2, \dots) of coefficients thus encodes f in a unique and useful way. If the index set is as numerous as the real numbers we would write t in terms of an integral, $f = \int_\alpha c_\alpha u_\alpha d\alpha$, rather than a sum.

For instance, if we were to represent faces as two-dimensional images on a rectangular $n \times m$ monochrome pixel display as is common in several approaches to object and face processing (e.g., Abdi, Valentin, & Edelman, & O’Toole, 1995; Kersten, 1987; Walton & Bower, 1993), each “face” would be encoded as a vector of $n \times m$ numbers corresponding to pixel intensities on the display rectangle D . The space of all such vectors form an $n \times m$ -dimensional approximation to $L^2(D)$, and the $n \times m$ functions that light single pixels to some fixed intensity and leave all other pixels blank, form a basis for the resulting finite-dimensional function space.

Alternatively, one can single out an *eigenface basis* (e.g., Abdi et al., 1995; Turk & Pentland, 1991). Obtained by doing a principal component analysis on a representative collection of faces in $L^2(D)$, these eigenfaces provide a basis that itself encodes information about “expected” facial configuration. Because of this, any psychologically meaningful metric is likely to take a much simpler form with respect to such an eigenface basis, as compared with the former raw pixel basis.

In the presence of an inner product on $L^2(D)$, one can ask for an *orthonormal* basis; a basis such that $\langle u_i, u_j \rangle = 1$ or 0 , according to whether $i = j$ or not, respectively. Standard linear algebra then shows that one can compute the coordinates of any other function contained in the subspace spanned by the u_i simply by taking inner products: If $f = \sum c_i u_i$, then we have $c_i = \langle f, u_i \rangle$. One can easily arrange for an eigenspace basis to be orthonormal with respect to the standard L^2 inner product on $L^2(D)$, for instance, and then find out “how much” of each eigenface u_i is in a given face f by simply computing the inner product $\langle f, u_i \rangle$. (Note that this is much more interesting from a psychological standpoint than computing “how much the i th pixel is lighted” in representing that same face.)

In any event, using coordinates relative to an orthonormal (e.g., eigenface) basis $\{u_i\}$ in $L^2(D)$, we can give our earlier geodesic equation a much more concrete form. First of all, we can express our curve Φ_t relative to this basis as

$$\phi_t = \sum_i c_i(t)u_i, \quad (6)$$

where each coefficient $c_i(t)$ is now a simple numeric function of our t parameter interval. We then assume the basis functions u_i are time independent, so we then have

$$\dot{\phi}_t = \sum_i \dot{c}_i(t)u_i, \quad \ddot{\phi}_t = \sum_i \ddot{c}_i(t)u_i \quad (7)$$

The field of positive symmetric operators W_t that define our Riemannian metric along the path—and their inverses, which we shall need—can likewise be expressed relative to our basis. Given any particular basis element u_i we can expand $W_t(u_i)$ relative to the entire basis using some coefficients, which we call g_{ij} (or, for the inverse operator, g^{ij}):

$$W_t(u_i) = \sum_j g_{ij}u_j, \quad W_t^{-1}(u_i) = \sum_j g^{ij}u_j \quad (8)$$

To avoid complicating our notation more than it already is, we do not indicate here the fact that because W_t varies from point to point in face space, so do the g_{ij} values. Indeed they must, because they now encode our Riemannian metric that (presumably) varied from point to point.

We spare the reader the mathematically standard debauch of index manipulation that now ensues to provide coordinates for the differential equation for geodesics we derived earlier; we present only the result—it is complex enough. To state it succinctly, one first needs to define the traditional Christoffel symbols

$$\Gamma_{ij}^k := -\frac{1}{2} \sum_l \left(\frac{\partial g_{ij}}{\partial u_l} - \frac{\partial g_{il}}{\partial u_j} - \frac{\partial g_{je}}{\partial u_i} \right) g^{lk}, \quad (9)$$

in which $\partial/\partial u_i$ means partial differentiation with respect to the u_i coordinate. Given these Christoffel symbols, one can put the geodesic equation

into fairly simple form; it becomes a *system* of second order differential equations for the coefficient functions ϕ_i^k :

$$\ddot{c}_k(t) = - \sum_{i,j} \Gamma_{ij}^k \dot{c}_i(t) \dot{c}_j(t), \quad i, j, k, \text{ in } I \quad (10)$$

We now discuss two examples to illustrate the application of this system.

Example 1: Consider the simple case in which the inner product does not vary from point to point. That is, suppose that the Riemannian metric is the same everywhere, a single positive operator W determining the norms of velocities at every point f in face space:

$$\langle \phi, \psi \rangle_f := \langle \phi, W(\psi) \rangle$$

The unsubscripted pairing \langle, \rangle here denotes the standard $L^2(D)$ inner product, as usual. In this situation, the constancy of W implies the constancy of all the g_{ij} s, whose partial derivatives therefore vanish, causing the Christoffel symbols all to vanish in turn. The differential equation for geodesics, as given by Equation 10 then becomes simply

$$\ddot{c}_k(t) \equiv 0 \quad \text{for all } k$$

Of course, any function whose second derivative vanishes identically must be linear. So for every index k , we have

$$c_k(t) = a_k + b_k t \quad \text{for all } k,$$

where a_k and b_k are constants. The reader can easily check that this set of equations defines ϕ_t as the line seg $[f, g]$, where f is the face whose coordinates relative to our basis are (a_1, a_2, \dots) , and g is the face whose coordinates are $(a_1 + b_1, a_2 + b_2, \dots)$.

In other words; if the Riemannian metric is constant relative to the standard L^2 inner product, all geodesics are straight lines, and vice versa.

This example is admittedly simple to a fault; it does verify the geodesic equation in an easy case, but our whole point in introducing the Riemannian metric was to allow a varying metric. Among other things, our next example illustrates this.

A TWO-DIMENSIONAL PLANE OF FACES

Example 2: Our infinite dimensional function space model provides a metric context for the finite-dimensional approximations to face space that typically arise in practice. To show how, we consider now a two-dimensional face space: a “plane” determined by three specific faces. We let the Riemannian metric vary in a simple way on this plane.

Assume, therefore, that we have data representing the face of a particular woman (Kay) in three states, giving us faces f_0 , f_1 , and f_2 , each a function on the rectangle D , as we have been discussing. Perhaps f_0 and f_1 represent Kay’s face at ages 20 and 65 with emotionally neutral expressions, so that the line $\text{seg}[f_0, f_1]$ forms an “aging” axis, and f_2 represents Kay’s face at age 20 expressing astonishment, so that $\text{seg}[f_0, f_2]$ provides a simple one-dimensional “emotion” axis.

Defining $u_i := f_i - f_0$ for $i = 1, 2$, we parameterize the “face plane” Ω_K containing f_0 , f_1 , and f_2 as follows:

$$\Omega_K := \{f_0 + c_1 u_1 + c_2 u_2 \mid c_1, c_2 \text{ arbitrary}\}$$

Note that the c_i s provide coordinates in this plane, in such a way that f_0 corresponds to the origin ($c_1 = c_2 = 0$), and f_1 and f_2 correspond to the coordinate vectors $(1, 0)$ and $(0, 1)$ respectively. We use square brackets in what follows (e.g., $[f]$), to denote this “coordinatization” of a function f . Thus, $[f_0] = (0, 0)$, and $[f_1] = (0, 1)$, for instance.

In this model, we may regard f_0 as the young, emotionally neutral origin face, and the basis vector u_1 in the tangent space $T_{f_0}\Omega_K$ then represents a movement from the origin toward the aged, emotionally neutral face f_1 . Similarly, u_2 in $T_{f_0}\Omega_K$ represents a movement toward the young but astonished face f_2 . It seems psychologically reasonable to provisionally regard these two movements as perceptually orthogonal, although it is, naturally, an empirical question. Because the coordinate vectors for u_2 and u_1 , respectively $(0, 1)$ and $(1, 0)$, are indeed orthogonal with respect to the standard dot product on \mathbf{R}^2 , one encodes precisely this idea by using the standard dot product to compare small displacements from the origin in Ω_K . There is a simply defined positive operator W that, when inserted in the standard L^2 inner product as discussed already, constructs an inner product in which u_1 and u_2 are orthonormal. To show that our infinite dimensional model contains the two-dimensional case in this way, we now expose this relation.

Compute the four L^2 inner products

$$\gamma_{ij} := \int_D u_i(x) u_j(x) dx, \quad i, j = 1, 2,$$

and form the matrix

$$[W] := \begin{pmatrix} \gamma_{11} & \gamma_{12} \\ \gamma_{21} & \gamma_{22} \end{pmatrix}^{-1}$$

Given any face f in Ω_K , we can then transform its coordinate vector $[f]$ to the coordinate vector of some other face in Ω_K via $[f] \rightarrow [W][f]$, and we can call the new face indicated by these transformed coordinates Wf . Routine calculations (simply expand out $[W]$ and the resulting integrals) then show that for any displacements ϕ and ψ in $T_{f_0}\Omega_K$, we have

$$\int_D \phi(x) (W\psi)(x) dx = [\phi] \cdot [\psi]$$

In other words, W modifies the standard L^2 inner product so that it coincides precisely with the result of just “dotting together coordinate vectors” of face space displacements.

Of course, when it comes to the actual business of doing analysis on the face plane Ω_K , or on any other finite-dimensional approximation to our full theoretical face space Ω , a coordinate representation of \mathbf{R}^n like the one given earlier for Ω_K is by far the more practical approach. Our infinite dimensional function space model, however, provides a universal theoretical context for all such approximations. As we show next, theoretical results (e.g., the geodesic differential) derived in the infinite dimensional model provide useful corresponding results in the finite-dimensional setting.

To accomplish this, we now define a simple family of nontrivial Riemannian metrics of Ω_K by taking the coordinate “dot-product” on $T_{f_0}\Omega_K$ given earlier, and duplicating it at $T_f\Omega_K$ for every f in Ω_K , simultaneously multiplying it by a positive weighting factor w that depends on the base point f ; that is, $w = w(f)$. That is, for velocities ϕ_f and ψ_f based at some face f in Ω_K , we set⁵

$$\langle \phi_f \psi_f \rangle_f := w(f)([\phi_f] \cdot [\psi_f]) = w(f) \int_D \phi_f(x) (W_f \psi_f)(x) dx$$

⁵Note that this function W , defined on Ω_K still does not make use of the full generality inherent in the w_f s we inserted into L^2 inner products in our earlier discussion; there, for each f in face space, $w_f(x)$ defined a function on D . Here, for each f in face space $w_f(x)$ is merely a constant function on D , varying with f , but not with x in D , so we shall write $w(f)$ instead of $w_f(x)$.

We henceforth calculate using only the first identity in this formula; we have displayed the integral version only to place these efforts more clearly in the general context of our earlier discussion.

Our specific goal now is to answer the question, what geodesics on Ω_K are relative to the Riemannian metric we have just defined. Of course, if $w(f)$ does not actually vary with f , then neither does the metric, and (according to our previous example) the geodesics will all be straight lines. However, when $w(f)$ does vary, this will not generally be true, as follows.

First, note that just as in our earlier general computations, we may regard our Riemannian metric as a composition of the standard dot product with the “multiplication” operator W_f that simply multiplies tangent vectors based at f by the weight $w(f)$; that is, $W_f(\phi_f) := w(f)\phi_f$. So we can express our Riemannian metric in terms of the operator field W_f :

$$\langle \phi_f, \psi_f \rangle_f = [\phi_f] \cdot [W_f(\psi_f)]$$

In particular, we have $W_f(u^1) = w(f)u^1$ and $W_f(u^2) = w(f)u^2$. Similarly, W_f^{-1} corresponds to multiplication by $1/w(f)$, so that in the language of Equation 8,

$$\begin{aligned} g_{11} = g_{22} = w(f), & \quad \text{while} \quad g_{12} = g_{21} = 0 \\ g^{11} = g^{22} = \frac{1}{w(f)}, & \quad \text{while} \quad g^{12} = g^{21} = 0 \end{aligned}$$

Using these facts in Equation 9 to compute the Christoffel symbols, we immediately notice that because $g^{kl} = 0$ unless $k = l$, in which case it equals $1/w(f)$, the sum there only contains the $k = l$ term, and hence

$$\Gamma_{ij}^k = -\frac{1}{2} \left(\frac{\partial g_{ij}}{\partial u_k} - \frac{\partial g_{ik}}{\partial u_j} - \frac{\partial g_{jk}}{\partial u_i} \right) \cdot \frac{1}{w(f)}$$

Similarly, we can avail ourselves of the fact that $g_{ij} = w(f)\delta_{ij}$, where δ_{ij} is the Kronecker delta that equals either 1 or 0 according as $i = j$ or not. Moreover, because, for example, $\frac{1}{2w(f)} \frac{\delta w(f)}{\delta u_i} = \frac{1}{2} \frac{\delta}{\delta u_i} \ln(w(f))$, we define $\omega(f) := \frac{1}{2} \ln(w(f))$. These notational moves simplify the Christoffel formula to

$$\Gamma_{ij}^k = \left(\delta_{ij} \frac{\partial \omega}{\partial u_k} - \delta_{ik} \frac{\partial \omega}{\partial u_j} - \delta_{jk} \frac{\partial \omega}{\partial u_i} \right)$$

We therefore have

$$\begin{aligned}\Gamma_{11}^1 &= -\Gamma_{22}^1 = \frac{\partial\omega}{\partial u_1}, & \Gamma_{12}^1 &= \Gamma_{21}^1 = \frac{\partial\omega}{\partial u_2} \\ \Gamma_{22}^2 &= -\Gamma_{11}^2 = \frac{\partial\omega}{\partial u_2}, & \Gamma_{12}^2 &= -\Gamma_{21}^2 = \frac{\partial\omega}{\partial u_1}\end{aligned}$$

Using these formulae for the Christoffel symbols into the geodesic differential system as given by Equation 10, we see that, relative to the Riemannian metric under consideration, a path of faces

$$\phi_t := f_0 + c_1(t)u_1 + c_2(t)u_2$$

in Ω_K is geodesic if and only if its coordinates $(c_1(t), c_2(t))$ satisfy the following pair of second order differential equations:

$$\begin{aligned}\ddot{c}_1 &= \frac{\partial\omega}{\partial u_1}((\dot{c}_1)^2 - (\dot{c}_2)^2) + 2\frac{\partial\omega}{\partial u_2}\dot{c}_1\dot{c}_2 \\ \ddot{c}_2 &= \frac{\partial\omega}{\partial u_2}((\dot{c}_2)^2 - (\dot{c}_1)^2) + 2\frac{\partial\omega}{\partial u_1}\dot{c}_1\dot{c}_2\end{aligned}\tag{11}$$

The Riemannian metric enters into these equations through the coefficients $\frac{\partial\omega}{\partial u_i}$. For instance, suppose we consider not the entire face plane Ω_K (which, after all, contains unrecognizable faces such as “enormously aged faces” like $f_0 + 1,000u_1$), but only the ball

$$\Omega_K^R := \{f \mid [f] = (c_1, c_2) \text{ with } c_1^2 + c_2^2 < R^2\},$$

for some (presumably large) radius $R > 0$. We can then put a Poincaré disc metric (a standard example from differential geometry) on Ω_K^R by taking

$$w(f) = \frac{4}{(R^2 - c_1^2 - c_2^2)^2}, \quad ((c_1, c_2) = [f])$$

Note that $w(f) \rightarrow \infty$ as $c_1^2 + c_2^2 \rightarrow R^2$. This has the effect of magnifying distances between faces more and more the nearer they are to the bounding circle $c_1^2 + c_2^2 \rightarrow R^2$ —and hence are (presumably) less and less “familiar” looking. Although our aim in selecting this particular metric is more mathematical than psychological, the effect seems quite reasonable from a perceptual viewpoint.

In any event, though somewhat complicated to show using Equation 11, this choice of w (and thereby of $w = \frac{1}{2} \ln(w)$) makes every geodesic in Ω_K^R follow a circular arc that meets the circle $c_1^2 + c_2^2 = R^2$ at right angles. In particular, the geodesics here are not straight lines in the Euclidean sense.⁶

Although nonstraightness will hold quite typically, the coefficients $\frac{\partial \omega}{\partial u_i}$ in the (nonlinear) differential system (Equation 11) usually vary in such a way that (in contrast to the example here) the system cannot generally be solved in closed form. On the other hand, Equation 11 does lend itself to straightforward numerical solution, and hence is quite serviceable from a practical standpoint.

Using the method already outlined, one can construct a metric that provides an account of the “other-race” effect (Brigham & Barkowitz, 1978; Brigham & Malpass, 1985; Goldstein & Chance, 1985; Valentine & Bruce, 1986), in which an individual’s performance on a face recognition task is better for faces of his or her own race than for faces of other races. Valentine (1991) proposed a multidimensional space framework that accounted for this and other effects. In this framework, faces are represented as either points or as vectors in a multidimensional space. The point representation is consonant with our current approach. The main goal of Valentine’s study was to demonstrate that a spatial representation could account for the various effects—it was not necessary to posit a complex process operating on some simpler representation. Although this model employed a Euclidean metric, it is worth noting that Valentine stated that this assumption was made for simplicity and in the absence of evidence for another metric. In his closing discussion, he stated that the Euclidean assumption, “is almost certainly an oversimplification” (Valentine, 1991, p. 201). We demonstrate that employing Riemannian metrics can provide a concise, metric-based interpretation of the other-race effect, while also accounting for D. T. Levin’s (1996) finding that people are faster at classifying the race of other-race than same-race faces.

Example 3: We model the psychological face space of a White male observer whose perceptual experience has been primarily with White faces. His ability to discriminate individual within-race faces is superior for White faces as compared with African American faces,

⁶Actually, geodesics that pass through the origin face f_0 will follow circular arcs of infinite radius, which are, in fact, straight; however, a random geodesic has zero probability of passing through the origin, so this is a very exceptional case. Also, neither the circular nor the straight geodesics follow their paths with constant speed (as measured by the Euclidean metric we usually use in the plane).

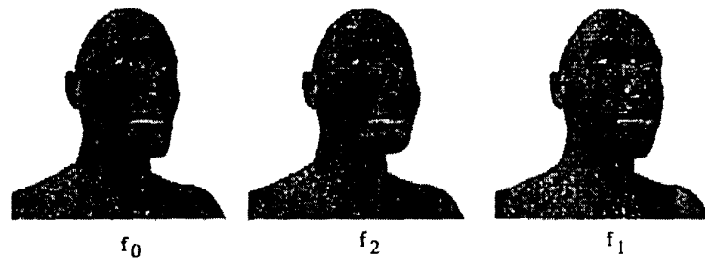


FIG. 2.1. The above faces were used in the "other-race" example. Faces f_1 and f_2 are based on f_0 ; f_2 differs by a single parameter change, whereas f_1 differs by five parameter changes. Note that the skin texture on all figures is identical. The three-dimensional face rendering was accomplished in Metacreations Poser 4.

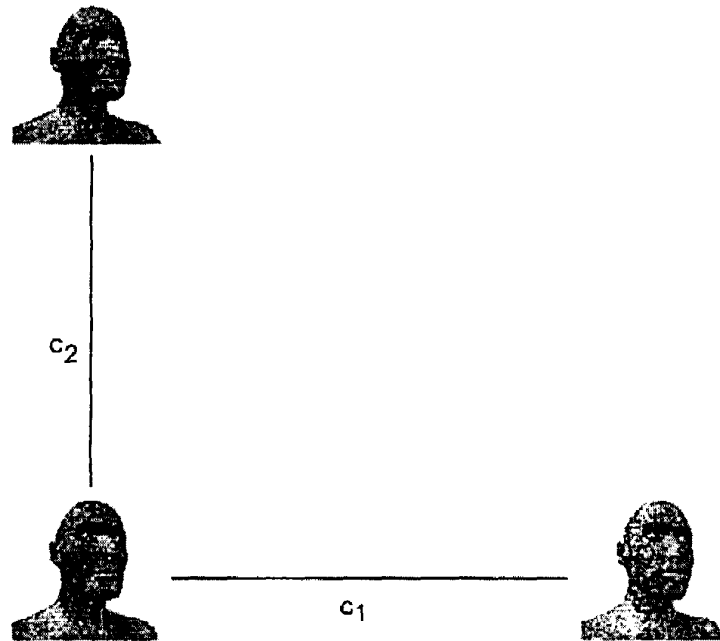


FIG. 2.2. Three faces, f_0 , f_1 , and f_2 , define a two-dimensional plane in an infinite-dimensional face space. Face f_0 is located at $c_1 = 0, c_2 = 0$, or $(0, 0)$; face f_1 at $c_1 = 0, c_2 = 1$ or $(0, 1)$; and face f_2 at $c_1 = 1, c_2 = 0$ or $(1, 0)$.

whereas he is faster to make a race classification for African American faces than he is for White faces. For our example, we use three-dimensional facial information derived from realistic African American and White characters in an animation and rendering program (Fig. 2.1). Three faces, two White and one African American, define a plane in our infinite-dimensional face space (Fig. 2.2).

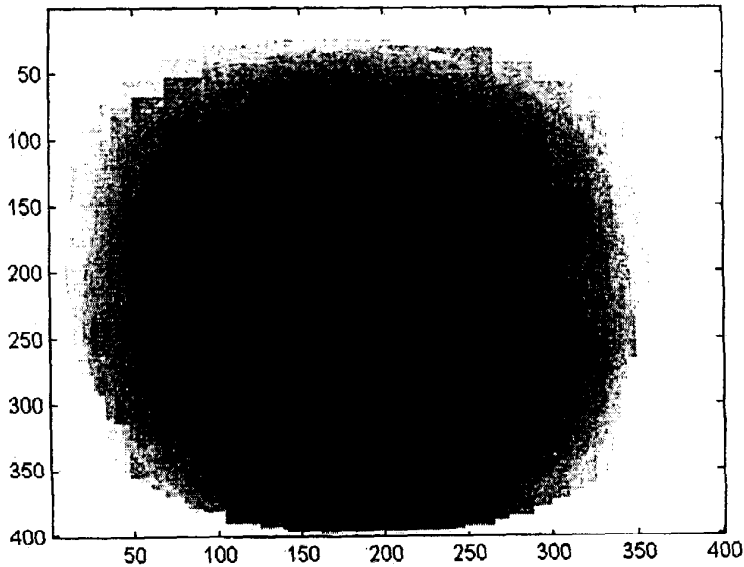


FIG. 2.3. The height matrix for face f_1 is shown as a 400×400 grayscale image, with height 0 represented as white.

Using the method outlined in Example 2, we define:

$$\Omega_K := \{f_0 + c_1 u_1 + c_2 u_2 \mid c_1, c_2 \text{ arbitrary}\}$$

where f_0 is a White face that serves as the origin in our plane, u_1 represents a vector from f_0 to f_1 (given by $f_0 - f_1$), and u_2 represents a vector from f_0 to f_2 . The parameters c_1 and c_2 serve as coordinates in the face plane. At $c_1 = 1$ and $c_2 = 0$, for example, we have the African American face, f_1 . Previously, faces have been represented as functions, f_i , defined over a rectangle, D . For this example, we equate each face, f_i , with a 400×400 matrix representing the height of each point on the surface of the face from a coronal plane positioned in front of the ears (Fig. 2.3). All calculations are based on these height matrices.

We can compute distances in our face plane using methods described earlier. To compute the distance between two faces in matrix form using the \mathbb{L}^2 metric, we calculate item-by-item differences between the two matrices:

$$d(f_0, f_1) \approx \sum_{i=1}^{400} \sum_{j=1}^{400} (|H_0(i, j) - H_1(i, j)|^2)^{1/2}$$

where H_0 and H_1 are the height matrices for f_0 and f_1 , respectively, and i and j are indexes for the matrices. Using the L^2 metric (strictly speaking, we are employing the Euclidean approximation to the L^2 metric based on the function space—however, no harm is done here, as the Euclidean metric can be considered as a special case of the L^2 metric), the distance between f_0 , the White face that serves as the origin in the face plane, and f_1 , the African American face, is 173.3. The distance between f_0 and f_2 is 259.2. In an informal sampling of 10 people, however, the two White faces were rated as the most similar. Each of these viewers rated f_0 and f_2 more similar than f_0 and f_1 .⁷ In the L^2 space, however, assuming distance in the psychological space is related to similarity, f_0 and f_1 are more similar than f_0 and f_2 . It might be possible to generate a process-oriented account for this reversal, but it would be more satisfying for our metric to address this disparity. With the L^2 metric, we also find that faces near the African American face f_1 will be equally discriminable as faces near the two White faces, f_0 and f_2 . Our metric should account for this as well.

Using our Riemannian metric, we can account for the higher similarity between the White faces, at the same time providing an account of the other-race effect and the race classification speed data. We devise a metric on the plane in which movements from the White face at the point $(0, 0)$ to the White face at $(0, 1)$ are unweighted—displacements in the u_2 direction using our new metric will be the same as in the L^2 metric (e.g., the distance between the point $(0, 0)$ and the point $(0, 1)$ is the same under either metric). Displacements from $(0, 0)$ to the African American face at $(1, 0)$ are weighted relative to the L^2 metric. Close to $(0, 0)$, displacements in the u_1 direction are magnified relative to the L^2 metric, whereas closer to $(1, 0)$, displacements in the u_1 direction are reduced relative to the L^2 metric. The metric is fully specified by the weighting function:

$$w(c_1, c_2) = e^{R-c_1}$$

where c_1 and c_2 are coordinates in the face plane, and R is a parameter between 0 and 1 that is related to where the labeling shift (White to African

⁷In an unanticipated finding, observers also reported that the shading differed between the faces. The African American faces were rated by each of our observers as being either darker or as having a smoother shading than the White faces. The faces had been rendered with the same skin tone—only the structural information was changed. Additionally, when first presented with the faces, observers immediately stated that face f_0 and f_1 , the White face at the origin and the African American face, were most similar. After a closer look, all observers changed their rating, stating that the two White faces were more similar to each other, and that the African American face was less similar.

American) occurs. R is stated as a threshold here to ease computation. R could be stated as a function of c_1 to more closely match the pattern one might likely find in data—that the shift does not occur at a clearly defined boundary. In either interpretation, R serves to change the scaling of measurements in the u_1 direction relative to the L^2 metric. For our computations, $R = 0.90$.

To calculate distance using this new metric, we perform a similar calculation on the face matrices as outlined earlier:

$$d_A(f_0, f_1) \approx \sum_{k=0}^N w(k/N + 1) \times \left(\sum_{i=1}^{400} \sum_{j=1}^{400} (|H_{k/N+1}(i, j) - H_{k/N}(i, j)|^2)^{1/2} \right)$$

where d_A is the distance using our alternate metric, w is the weighting function of the alternate metric, N is a large number, and $H_{k/(N+1)}$, $H_{k/N}$ are the height matrices of faces along the path of the geodesic from f_0 and f_1 . As N is increased, the accuracy of the approximation improves.

Using our new metric to measure distances, we find that the distance between the two White faces f_0 and f_2 is the same as when we measure using the L^2 metric—259.2. The distance between f_0 , the white face at the origin, and f_1 , the African American face, is now 265.5. Our distances now match our informal sample of similarity measures. Also, we find that changes in the u_1 direction to faces in the vicinity of $(0, 1)$, our African American face, have less effect than the same changes in faces near $(1, 0)$ or $(0, 0)$, our White faces. Relative to White faces, African American faces appear to be more similar to one another in our observer's psychological space.

Because structural changes in our African American face result in less movement away from $(0, 1)$, we find that African American faces, although distant from White faces, also form a cluster near $(0, 1)$. If race identification is facilitated by having clearly defined clusters, then our new metric also provides a natural account of D. T. Levin's (1996) findings.

Our change of metric is of a simple form. It is conformal, meaning that angle measures made using the L^2 metric agree with those made in the new metric. The weighting function that determines the metric is a function of only c_1 , meaning only movements in the u_1 direction are affected. The weighting function does not weight any area within a single face; all areas within a single face are equally weighted. This, along with the fact that the weighting function is of the form e^u , results in a greatly simplified equation

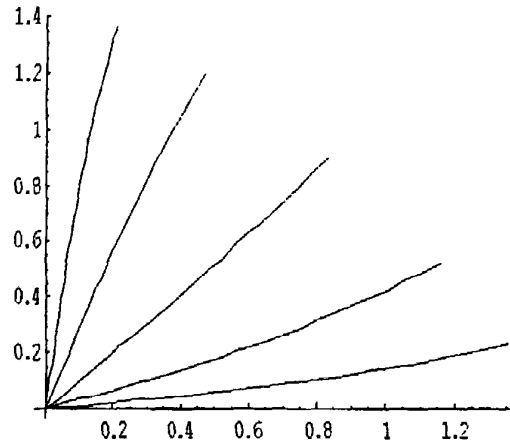


FIG. 2.4. The exponential map for the alternate metric for Example 3 reveals that not all the geodesics are straight. The geodesics from f_0 to f_1 and from f_0 to f_2 lie on the c_1 and c_2 axis, respectively.

for the geodesics. A curve, $(u(t), v(t))$ in our face plane is a geodesic if and only if:

$$\begin{aligned} u'' - 1/2 u e^{R-u} (v')^2 &= 0 \\ v'' - u e^{R-u} u' v' &= 0 \end{aligned}$$

One can numerically solve this equation to draw geodesics with the point $(0, 0)$ as a starting point in what is known as the exponential map (Fig. 2.4).

It is important to note that models such as the one described here are eminently falsifiable. Although it might well be possible to construct an arbitrary metric to account for any finite set of psychological phenomena, in doing so one fixes global properties of the proposed perceptual metric space. In particular, the geodesics for the proposed space are determined by the metric, as described previously. In addition, the curvature tensor is also determined by the metric; empirical methods exist that provide measures for this property (D. N. Levin, 2000; Lindman & Caelli, 1978).

CONCLUSION: TEMPLATES, PROTOTYPES, AND RELATED CONCEPTS

We have presented a general metatheory that provides for a geometry of perceived faces or other objects that can lie in an infinite dimensional space. We now discuss some important psychological aspects of such

spaces in a general, relatively nontechnical way. The set of all smooth three-dimensional objects comprises such a space and the set of faces, recognizable as such, form a subspace of the more general space. We demonstrated that such spaces could possess familiar notions of paths between two objects represented as points in the space and distances as being the shortest paths between two objects. Indeed, the present developments provide not only a natural rigorous "home" for concepts like morphing, but also introduce more novel concepts, some of which may possess psychological implications. One example is the question of whether all faces on a straight line (which we called geodesic in the general case) themselves are perceived or remembered as faces (the notion of convexity). It was argued that the future of many areas in cognitive science will require the general ideas and powerful machinery associated with modern geometry (e.g., infinite dimensional Riemannian manifolds) and topology, concepts and tools that have so far been visited by only a very small number of research domains. It should be emphasized that the claim is not being made that all properties of, say, a Riemannian manifold are enjoyed by face and object perception—all are empirical questions. We do think it probable that the present outlook could help lead to other enrichments of the theoretical terrain in addition to those discussed in this investigation.

The Gestalt psychologists intended that various measures of aspects of perceptual objects associated with "goodness of form" lie on a continuum. An object could possess varying degrees of these measured aspects and thus any figure might, in principle if not in practice, be assigned a number reflecting its goodness of form or measure of "Gestalthood." Our concept of a perfect Gestalt is a little different. We mean simply a unique mathematical description that would permit anyone to replicate the object.⁸ Although our tendering of the expression perfect Gestalt was done with part of our collective tongues-in-cheek, a unique mathematical description does seem to imply at least a configularity of all parts lying in a specified relation to one another. Patently such configularity could be aided by several of the traditional Gestalt aspects, including grouping mechanisms.

Some investigators are distinguishing between configularity wherein some degree of independence or separability could remain versus a more palpable and complete holism (e.g., Farah et al., 1998). Clearly, a strong

⁸We naturally choose to ignore philosophical quandaries associated with whether any potential description satisfies some ultimate definition of completeness. For instance beyond a unique mathematical specification of the implied facial surface, is a given "face" constructed from plastic, biological tissue, a holograph, or what?

form of holism could imply some type of configularity without configularity implying a total kind of holism. Another related issue regards where holism comes from. There is evidence that both perceptual learning and innate mechanisms could both play strategic roles in binding components (e.g., features) together to form good Gestalt forms (for a useful survey and update on the long neglected topic of perceptual learning, see Goldstone, 1998). Our total kind of perceptual form could be an approximation to an ideal that is not quite ever completed, especially with regard to components that start out relatively independent of one another. In addition, even most components such as features require a complex specification such as through function theory analogous to our approach, for a complete and relatively unique description. Hoffman and Bennett (1986) and Lappin (e.g., 1990) deftly employed differential geometry and topology in developing mathematical machinery for producing object segregation, metric structure of objects, and related visual properties; for instance, those involving motion or that might result in geon-like entities (Biederman, 1987).

In addition, our notion of a perceptual object is naturally related to other holistic ideas in perception and cognition. Two of these are template and prototype. As with many, if not most concepts in sciences relating to thought, these have never been rigorously defined, except perhaps in individual investigations. Nevertheless, they have both supported useful theoretical and experimental tracks in cognitive science and everyone seems to know one when one sees it." The major informal defining aspects of both concern holism, uniqueness, and, especially with prototype, connotations of being a kind of centroid or average of a population of figures. In practice, templates have pertained primarily to perception, especially identification, whereas prototypes have been employed primarily in categorization.

Very strict notions of templates as being rigid and perceived in an all-or-none fashion (due partly to the thesis that a percept would be either perfectly congruent with a template or not, with no in-between similarity permitted) are readily falsified (e.g., Neisser, 1967; see also the all-or-none model of Townsend, 1971). However, a continuous measure of template overlap of uppercase English letters provided a quite successful parameter of similarity in the latter identification experiment. Engineers and computer scientists have long employed more flexible notions of templates. One relatively invariant demand has been that what corresponds to early sensory input is required to undergo various normalization (e.g., centering, size adjustment, and the like) procedures before being matched against a set of templates denoting the memory set of stimulus patterns. However, our notion of a Gestalt percept or memory as well as more sophisticated ideas

of templates are compatible with potential deformation due to noise, low energy displays, illusionary effects, and so on. In addition, they include the possibility of physical similarity effects on identification (something not possible in true all-or-none congruency testing), effectuated for instance, through graded perceptual and memory matching results. For instance, dissimilarity could act through the geodesic length between two figures or, a measure of similarity is gained as noted earlier, in the angle associated with the inner product (in informal language, a kind of correlation), at least locally. Such generalized similarity effects are obviously a straightforward extension of featural overlap among stimulus patterns leading to interitem confusion being an increasing function of the degree of that overlap. Models or artificial intelligence routines based on sophisticated template processing are apparently making a comeback (Goldstone, 1998; Hinton, Williams, & Revow, 1992; Poggio & Edelman, 1990; Tarr, 1995; Ullman, 1989).

In categorization, a central theme has been whether individuals compound their experiences into a prototype or store individual examples of the stimuli, known as *exemplars*. Although there are aspects of performance in categorization experiments that appear to support prototype kinds of predictions, Nosofsky (1991) showed that his generalized context model, a model based on exemplars and generalized from the Medin and Schaffer (1978), Context model, can encompass many of those predictions. The Medin and Schaffer model was, in turn, founded on the Shepard-Luce similarity choice model for identification experiments (Luce, 1963; Shepard, 1957; see also Townsend & Landon, 1982). Although such models as Nosofsky's are certainly falsifiable, a potential problem for total experimental discrimination of prototype versus exemplar theories arises because in principle, an exemplar theory could preserve all information from an experimental sequence of events, including everything related to a person's perception and storage of each stimulus. Hence, again in principle, an inventive theoretician could simply degrade the "perfect" exemplar theory down to where it more or less closely approximates the data. The most general prototype models seem much more constrained than this. For instance, a natural constraint would seem to be that all stimulus information up to, say, the present experimental trial n is mapped into the single prototype.

The major quantitative theoretical alternative to Nosofsky's Generalized Context Model has undoubtedly been the Bounded Performance Model of Ashby and colleagues (e.g., Ashby & Gott, 1988).⁹ That model views

⁹Recently, other quantitative models, some by the same authors and some extensions of the earlier ones, have been proposed. Obviously, these lie beyond the present scope.

presentation of a stimulus as leading to an observation specified as a point in a multidimensional observation space. The space is carved up into a set of mutually exclusive and exhaustive regions, each of which is associated with a particular category and therefore, response. The bounds separate these regions and there is a multivariate probability distribution (e.g., normal) on the observation space for each stimulus pattern. If a percept falls into a region associated with, say, Category C_i , it is predicated that the individual will give the response tied to that category. There exist certain cases within the Bounded Performance Model where the multidimensional means of the various percept distributions can be viewed as a set of prototypes, with the category response determined by the minimum distance of the observation from the prototypes (e.g., Ashby, 1992). Interestingly, if the probability distributions are multivariate normal, this kind of model is equivalent to the matched filter model, well known in engineering circles (Townsend & Landon, 1983). Here too, each filter can be interpreted as a kind of template or prototype.

Finally, it seems clear that theories of well-founded Gestalts such as faces, require dynamics that permit their holistic aspects to either hurt or help performance (e.g., Kuehn & Jolicouer, 1994; Suzuki & Cavanagh, 1995). Such dynamics call again for research targeting interrelations and potential syntheses of geometric and topological representation theory with quantitative human information processing theory as stressed in the introduction to this volume.

REFERENCES

- Abdi, H., Valentin, D., Edelman, B., & O'Toole, A. J. (1995). More about the difference between men and women: Evidence from linear neural networks and the principal-component approach. *Perception, 24*(5), 539–562.
- Ashby, F. G. (Ed.). (1992). *Multidimensional models of perception and cognition*. Hillsdale, NJ: Lawrence Erlbaum Associates.
- Ashby, F. G., & Gott, R. E. (1988). Decision rules in the perception and categorization of multidimensional stimuli. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 14*(1), 33–53.
- Baenninger, M. A. (1994). The development of face recognition: Featural or configurational processing? *Journal of Experimental Child Psychology, 57*, 377–396.
- Baird, J. C. (1997). *Sensation and judgment: Complementarity theory of psychophysics*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Beals, R., & Krantz, D. H. (1967). Metrics and geodesics induced by order relations. *Mathematische Zeitschrift, 101*, 285–298.
- Biederman, I. (1987). Recognition-by-components: A theory of human image understanding. *Psychological Review, 94*, 115–117.

- Biederman, I., & Kalocsi, P. (1998). Neural and psychophysical analysis of object and face recognition. In H. Wechsler, P. J. Phillips, V. Bruce, F. F. Soulie, & T. Huang (Eds.), *Face recognition: From theory to applications* (pp. 3–25). New York: Springer-Verlag.
- Boothby, W. M. (1975). *An introduction to differentiable manifolds and Riemannian geometry*. New York: Academic.
- Brigham, J. C., & Barkowitz, P. (1978). Do "They all look alike"? The effect of race, sex, experience, and attitudes on the ability to recognize faces. *Journal of Applied Social Psychology*, 8(4), 306–318.
- Brigham, J. C., & Malpass, R. S. (1985). The role of experience and contact in the recognition of faces of own- and other-race persons. *Journal of Social Issues*, 41(3), 139–155.
- Dzhafarov, E. N., & Colonius, H. (1999). Fechnerian metrics in unidimensional and multidimensional stimulus spaces. *Psychonomic Bulletin and Review*, 6(2), 239–268.
- Edelman, S. (1998). Representation is representation of similarities. *Behavioral and Brain Sciences*, 21(4), 449–498.
- Edelman, S., & Duvdevani-Bar, S. (1997). Similarity, connectionism, and the problem of representation in vision. *Neural Computation*, 9(4), 701–720.
- Farah, M. J., Wilson, K. D., Drain, M., & Tanaka, J. N. (1998). What is "special" about face perception? *Psychological Review*, 105, 482–498.
- Goldstein, A. G., & Chance, J. E. (1985). Effects of training on Japanese face recognition: Reduction of the other-race effect. *Bulletin of the Psychonomic Society*, 23(3), 211–214.
- Goldstone, R. L. (1998). Perceptual learning. *Annual Review of Psychology*, 49, 585–612.
- Hinton, G., Williams, K., & Revow, M. (1992). Adaptive elastic models for handprinted character recognition. In J. Moody, S. Hanson, & R. Lippmann (Eds.), *Advances in neural information processing systems, IV* (pp. 341–376). San Mateo, CA: Morgan Kaufmann.
- Hoffman, D. D., & Bennett, B. M. (1986). The computation of structure from fixed-axis motion: Rigid structures. *Biological Cybernetics*, 54, 71–83.
- Johnston, R. A., Kanazawa, M., Kato, T., & Oda, M. (1997). Exploring the structure of multidimensional face-space: The effects of age and gender. *Visual Cognition*, 4(1), 39–57.
- Johnston, R. A., Milne, A. B., Williams, C., & Hosie, J. (1997). Do distinctive faces come from outer space? An investigation of the status of a multidimensional face-space. *Visual Cognition*, 4(1), 59–67.
- Kelley, J. L. (1955). *General topology*. New York: Springer-Verlag.
- Kersten, D. (1987). Predictability and redundancy of natural images. *Journal of the Optical Society of America A*, 4(12), 2395–2400.
- Kuehn, S. M., & Jolicouer, P. (1994). Impact of the quality of the image, orientation, and similarity of the stimuli on visual search for faces. *Perception*, 23, 95–122.
- Lappin, J. S. (1990). Perceiving the metric structure of environmental objects from motion, self-motion and stereopsis. In R. Warren & A. H. Wertheim (Eds.), *Perception and control of self-motion: Resources for ecological psychology*. (pp. 541–578). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Lappin, J. S., Ahlstrom, U. B., Craft, W. D., & Tschantz, S. T. (1995). Spatial primitives for seeing 3D shape from motion. In T. Pappathomas, C. Chubb, E. Kowler, & A. Gorea (Eds.), *Early vision and beyond*. Cambridge, MA: MIT Press.
- Lappin, J. S., & Craft, W. D. (1997). Definition and detection of binocular disparity. *Vision Research*, 37(21), 2953–2974.
- Levin, D. N. (2000). A differential geometric description of the relationship among perceptions. *Journal of Mathematical Psychology*, 44(2), 241–284.
- Levin, D. T. (1996). Classifying faces by race: The structure of face categories. *Journal of Experimental Psychology: Learning, Memory and Cognition*, 22(6), 1364–1382.
- Lindman, H., & Caelli, T. (1978). Constant curvature Riemannian scaling. *Journal of Mathematical Psychology*, 17, 89–109.
- Loftus, G. R. (1995). Data analysis as insight: Reply to Morrison and Weaver. *Behavior Research Methods, Instruments & Computers*, 27, 57–59.

- Loftus, G. R. (1996). Psychology will be a much better science when we change the way we analyze data. *Current Directions in Psychological Science*, 5, 161–171.
- Loftus, G. R., & Masson, M. E. J. (1994). Using confidence intervals in within-subject designs. *Psychonomic Bulletin & Review*, 1, 476–490.
- Luce, D. (1963). Detection and recognition. In R. D. Luce, R. R. Bush, & E. Galanter (Eds.), *Handbook of mathematical psychology* (vol. 1, pp. 103–190). New York: Wiley
- Medin, D. L., & Schaffer, M. M. (1978). Context theory of classification learning. *Psychological Review*, 85(3), 207–238.
- Munkres, J. R. (1975). *Topology: A first course*. Englewood Cliffs, NJ: Prentice-Hall.
- Neisser, U. (1967). *Cognitive psychology*. New York: Appleton-Century-Crofts.
- Nosofsky, R. M. (1991). Tests of an exemplar model for relating perceptual classification and recognition memory. *Journal of Experimental Psychology: Human Perception and Performance*, 17(1), 3–27.
- O'Toole, A. J., Vetter, T., Volz, H., & Salter, E. M. (1997). Three-dimensional caricatures of human heads: Distinctiveness and the perception of facial age. *Perception*, 26(6), 719–732.
- Poggio, T., & Edelman, S. (1990). A network that learns to recognize three-dimensional objects. *Nature*, 343, 263–266.
- Rumelhart, D. E., & Siple, P. (1974). Process of recognizing tachistoscopically presented words. *Psychological Review*, 81, 99–118.
- Shepard, R. N. (1957). Stimulus and response generalization: A stochastic model relating generalization to distance in psychological space. *Psychometrika*, 22, 325–345.
- Shepard, R. N., & Cermak, G. W. (1973). Perceptual–cognitive explorations of a toroidal set of free-form stimuli. *Cognitive Psychology* 4(3), 351–377.
- Suppes, P., Krantz, D. M., Luce, R. D., & Tversky, A. (1989). *Foundations of measurement: Vol. 2. Geometrical, threshold, and probabilistic representations*. San Diego, CA: Academic Press.
- Suzuki, S., & Cavanagh, P. (1995). Facial organization blocks access to low-level features: An object inferiority effect. *Journal of Experimental Psychology: Human Perception and Performance*, 21, 901–913.
- Syngé, J. L., & Schild, A. (1949). *Tensor calculus*. New York: Dover.
- Tanaka, J. N., & Sengco, J. A. (1997). Features and their configuration in face recognition. *Memory & Cognition*, 25, 583–592.
- Tarr, M. J. (1995). Rotating objects to recognize them: A case study on the role of viewpoint dependency in the recognition of three-dimensional objects. *Psychonomic Bulletin and Review*, 2, 55–82.
- Townsend, J. T. (1971). Alphabetic confusion: A test of models for individuals. *Perception & Psychophysics*, 9(6), 449–454.
- Townsend, J. T. (1994). Methodology and statistics in the behavioral sciences: The old and the new. *Psychological Science*, 5, 321–325.
- Townsend, J. T., & Ashby, F. G. (1982). Experimental test of contemporary mathematical models of visual letter recognition. *Journal of Experimental Psychology: Human Perception and Performance*, 8, 834–864.
- Townsend, J. T., Hu, G. G., & Ashby, F. G. (1981). Perceptual sampling of orthogonal straight line features. *Psychological Research*, 43, 259–275.
- Townsend, J. T., Hu, G. G., & Kadlec, H. (1988). Feature sensitivity, bias, and interdependencies as a function of intensity and payoffs. *Perception & Psychophysics*, 43, 575–591.
- Townsend, J. T., & Landon, D. E. (1983). Mathematical models of recognition and confusion in psychology. *Mathematical Social Sciences*, 4, 25–71.
- Townsend, J. T., & Landon, D. E. (1982). An experimental and theoretical investigation of the constant-ratio rule and other models of visual letter confusion. *Journal of Mathematical Psychology*, 25(2), 119–162.
- Townsend, J. T., & Thomas, R. (1993). On the need for a general quantitative theory of pattern similarity. In S. C. Masin (Ed.), *Foundations of perceptual theory* (pp. 297–368). Amsterdam: Elsevier.

- Turk, K. M., & Pentland, A. (1991). Eigenfaces for recognition. *Journal of Cognitive Neuroscience*, 3, 71–86.
- Tversky, A. (1977). Features of similarity. *Psychological Review*, 84, 327–352.
- Tversky, A., & Krantz, D. H. (1969). Similarity of schematic faces: A test of interdimensional additivity. *Perception & Psychophysics*, 5, 124–128.
- Ullman, S. (1989). Aligning pictorial descriptions: An approach to object recognition. *Cognition*, 32(3), 193–254.
- Uttal, W. R. (1988). *On seeing forms*. Hillsdale, NJ: Lawrence Erlbaum Associates.
- Valentine, T. (1991). A unified account of the effects of distinctiveness, inversion, and race in face recognition. *Quarterly Journal of Experimental Psychology: Human Experimental Psychology*, 43A(3) 161–204.
- Valentine, T., & Bruce, V. (1986). The effect of race, inversion and encoding activity upon face recognition. *Acta Psychologica*, 61(3), 259–273.
- Valentine, T., & Endo, M. (1992). Towards an exemplar model of face processing: The effects of race and distinctiveness. *Quarterly Journal of Human Experimental Psychology*, 44A(4), 671–703.
- Walton, G. E., & Bower, T. G. (1993). Newborns form “prototypes” in less than 1 minute. *Psychological Science*, 4(3), 203–205.
- Wenger, M. J., & Townsend, J. T. (2000). Spatial frequencies in short-term memory for faces: A test of three frequency-dependent hypotheses. *Memory & Cognition*, 28(1), 125–142.