

Psychophysiological and Modulatory Interactions in Neuroimaging

K. J. Friston,^{*,1} C. Buechel,^{*} G. R. Fink,^{*} J. Morris,^{*} E. Rolls,[†] and R. J. Dolan^{*,‡}

^{*}Wellcome Department of Cognitive Neurology, London, United Kingdom; [†]Department of Psychology, Oxford University, Oxford, United Kingdom; and [‡]Royal Free Hospital School of Medicine, London, United Kingdom

Received February 19, 1997

In this paper we introduce the idea of explaining responses, in one cortical area, in terms of an interaction between the influence of another area and some experimental (sensory or task-related) parameter. We refer to these effects as *psychophysiological interactions* and relate them to interactions based solely on experimental factors (i.e., psychological interactions), in factorial designs, and interactions among neurophysiological measurements (i.e., physiological interactions). We have framed psychophysiological interactions in terms of functional integration by noting that the degree to which the activity in one area can be predicted, on the basis of activity in another, corresponds to the contribution of the second to the first, where this contribution can be related to effective connectivity. A psychophysiological interaction means that the contribution of one area to another changes significantly with the experimental or psychological context. Alternatively these interactions can be thought of as a contribution-dependent change in regional responses to an experimental or psychological factor. In other words the contribution can be thought of as modulating the responses elicited by a particular stimulus or psychological process. The potential importance of this approach lies in (i) conferring a degree of functional specificity on this aspect of effective connectivity and (ii) providing a model of modulation, where the contribution from a distal area can be considered to modulate responses to the psychological or stimulus-specific factor defining the interaction. Although distinct in neurobiological terms, these are equivalent perspectives on the same underlying interaction. We illustrate these points using a functional magnetic resonance imaging study of attention to visual motion and a position emission tomography study of visual priming. We focus on interactions among extrastriate, inferotemporal, and posterior parietal regions during visual processing, under different attentional and perceptual conditions. © 1997 Academic Press

INTRODUCTION

This paper is about explaining regionally specific neuronal responses in terms of the interaction between influences from distal brain areas and sensory or task-related parameters. These interactions will be referred to as *psychophysiological interactions* as distinct from conventional interactions between experimentally manipulated factors and physiological interactions based purely on measurements of brain activity. Examples of the latter include modulatory interactions between two or more brain areas that could be inferred on the basis of measured neuronal or hemodynamic activity (Friston *et al.*, 1995a; Buechel and Friston, submitted for publication).

Although the techniques required to identify psychophysiological interactions in neuroimaging [functional magnetic resonance imaging (fMRI) and positron emission tomography (PET)] are well established and easy to implement, the conceptual issues of interpretation are not straightforward. Psychophysiological interactions bring together techniques normally associated with effective connectivity and those used to model designed effects that have been manipulated experimentally. By virtue of the integration of these physiological and experimental influences on regional responses, one is able to confer a degree of functional specificity when making inferences about functional integration or interactions between cortical areas. In this paper we will demonstrate this in terms of (i) effective connectivity that modulates *stimulus-specific* responses and (ii) *context-sensitive* changes in effective connectivity. The aim of this paper is to introduce and define psychophysiological interactions and address the issues of interpretation that arise. In this paper the term interaction is used in a specific way, to denote a measurable effect that can be modeled in terms of the interaction between two factors or variables and about which statistical inferences can be made.

The basic idea behind this paper is very simple: If one were to regress the activity of one region, on the activity of a second region, the slope of this regression would reflect the influence the second area could be exerting over the first. If one then repeated this regression,

¹ To whom correspondence should be addressed at Wellcome Department of Cognitive Neurology, Institute of Neurology, 12 Queen Square, London WC1N 3BG, UK. Fax: 44 (171) 813 1445. E-mail: k.friston@fil.ion.ucl.ac.uk.

using data acquired in a different context, then the slope might change. This change in slope is a psychophysiological interaction and what follows is an attempt to understand what such interactions mean.

The paper is divided into four sections. In the first section we present the background to psychophysiological interactions by reviewing effective connectivity, as it is employed in neuroimaging. Effective connectivity is usually understood in terms of connection strengths, by analogy with synaptic efficacy in electrophysiology (e.g., Gerstein and Perkel, 1969; Aertsen and Preissl, 1991). The second section reviews factorial designs and looks at how one generally assesses interactions between factors. The third section considers how interactions between physiological measurements can be understood in terms of effective connectivity, using concepts from the first section. This section includes an example of physiological interactions that speak to neuromodulatory effects in the visual pathway during attention to motion. The final section brings together the themes of the previous sections in terms of psychophysiological interactions and provides an illustrative example of the approach, as applied to the visual attention motion study and a study of priming-dependent perceptual processing.

EFFECTIVE CONNECTIVITY—EFFICACY AND CONTRIBUTION

In this section we review the notions of functional connectivity, effective connectivity, efficacy, and contribution, in the context of neuroimaging, and show how they can be assessed with the general linear model as employed by statistical parametric mapping. In short we will suggest that the significance of the regression of activity at any voxel on the activity at a reference voxel can be understood in terms of the contribution of the reference region to the voxel in question.

Functional Integration and Effective Connectivity

The principle of functional specialization in the brain is now well established, particularly in visual neuroscience (e.g., Zeki *et al.*, 1992). Functional integration refers to the interactions among specialized areas or neuronal populations and how these interactions depend upon the sensorimotor or cognitive context. Functional specialization and integration are not exclusive but complementary, with one making sense only in the context of the other. From the perspective of neuroimaging, functional specialization calls for the identification of *regionally specific effects* that can be attributed to changing stimuli or task conditions. Functional integration on the other hand is usually assessed by examining the correlations among activity in different brain areas or by trying to explain the activity in one area in

relation to others (e.g., Friston *et al.*, 1993a,b; McIntosh *et al.*, 1994a,b; Paus *et al.*, 1996). These analyses are usually framed in terms of *effective connectivity* (the influence that one neuronal system exerts over another). There is a fundamental distinction between demonstrating effective connectivity (in relation to some model of neuronal interactions) and simply observing correlated activity: The term *functional connectivity* is usually reserved for simple correlations between areas. Correlations can arise from many sources that do not reflect teleologically meaningful interactions (e.g., stimulus-evoked transients in two neuronal populations that are not connected, or the modulation of two cortical areas by a common subcortical input). One key aspect of effective connectivity is that it is always predicated on some model of the influence that one neuronal system exerts over another and therefore attempts to disambiguate correlations of a spurious sort from those mediated by direct or indirect neuronal interactions. The parameters (usually the connection strengths) of the model are then identified as those which allow the model to emulate, as closely as possible, the observed regional activities (or interregional correlations).

Contribution and Regression Analyses

Effective connectivity is as diverse as the models employed to model the influences among different regions. Clearly the measurements of effective connectivity that obtain are directly related to the validity of the model used and the implicit assumptions. In this paper we are less concerned with models of effective connectivity per se than how psychophysiological interactions can be understood in relation to them. Although there is current interest in nonlinear models, the model considered here is a simple linear one. Let the activity at the i th unit or voxel be modeled as

$$x_i(t) = \sum_j C_{ij} \cdot x_j(t), \quad (1)$$

with summation over j . This model says that the activity $x_i(t)$ is the sum of influences $C_{ij} \cdot x_j(t)$ from all other units or regions (note that there is no self-connection in this model, i.e., $C_{ii} = 0$). The parameters C_{ij} correspond to the connection strengths from j to i . Unlike correlations effective connectivity is not symmetric, i.e., C_{ij} is not the same as C_{ji} . In electrophysiology C_{ij} would be known as the *efficacy* and $x_i(t)$ would be the postsynaptic response to presynaptic inputs $x_j(t)$. In this very simple model the inputs sum linearly and there is no opportunity for the inputs to interact; furthermore the connection strengths are fixed, precluding time-dependent changes. These are strong and unrealistic assumptions but they have proved useful in making the analysis of effective connectivity more

robust and tractable, particularly in relation to multiple regression approaches and structural equation modeling (e.g., Friston *et al.*, 1993b; McIntosh *et al.*, 1994a).

Consider the statistical model associated with Eq. (1), which is a multiple linear regression and, in matrix format, looks like

$$\mathbf{x}_i = [\mathbf{x}_1 \cdots \mathbf{x}_k \cdots \mathbf{x}_K][\beta_{i1} \cdots \beta_{ik} \cdots \beta_{iK}]^T + \mathbf{G} \cdot \beta_G + \mathbf{e}_i \quad (2a)$$

where T denotes transpose. \mathbf{x}_i is a column vector of response variables at each voxel i (similarly for \mathbf{x}_k). \mathbf{G} is a matrix whose columns contain a collection of uninteresting effects or confounds, such as global activity, time effects, etc. β_{ik} and β_G are the parameter estimates and \mathbf{e}_i is a well-behaved error term. Here the parameter estimates β_{ik} can be identified with the effective connection strengths C_{ik} because we have included all the regions $1 \dots K$ in the model (Friston *et al.*, 1993b). However, in some situations only a few regions might be included and in this instance the parameter estimates and estimates of effective connectivity would not be the same. In this paper we consider the simplest case where only one region is included in the statistical model (although this is not a necessary constraint on what follows):

$$\mathbf{x}_i = \mathbf{x}_k \beta_{ik} + \mathbf{G} \cdot \beta_G + \mathbf{e}_i \quad (2b)$$

In this case there is only a tenuous relationship between β_{ik} and effective connectivity because the influences of all the remaining regions have been ignored when estimating β_{ik} . We will refer to β_{ik} as the “contribution” to make it clear that this regression slope is not necessarily an estimate of effective connectivity. Interestingly a test of the significance of the regression (i.e., with the null hypothesis that $\beta_{ik} = 0$) corresponds to a test for the significance of the correlation between \mathbf{x}_i and \mathbf{x}_k , in other words a test for the functional connectivity between regions i and k . However, this is not generally the case: if we add more areas the contribution will change, whereas the correlation will not. As more areas are added the contribution progressively approximates the estimate of effective connectivity. A map of contributions from region k to the rest of the brain can be created by designating activity at k as an explanatory variable in the general linear [model Eq. (2)] and creating a statistical parametric map (SPM) testing for the significance of this regression.

In conclusion the contribution from voxel or unit k to i can be thought of in terms of the regression that obtains by designating the activity in k as the explanatory variable and the activity in i as a response variable. We will use this interpretation below.

FACTORIAL DESIGNS AND PSYCHOLOGICAL INTERACTIONS

Factorial designs involve combining two or more factors within a task or tasks. The idea is to look at the interaction between these factors, or the effect that one factor has on the responses due to the other. An early and perhaps the simplest factorial design in neuroimaging involved an interaction between motor activation and time, which was interpreted in terms of physiological plasticity or adaptation (Friston *et al.*, 1992). The interaction here was simply an effect of time on the activation due to motor performance. Generally, interactions can be thought of as a difference in responses to one factor brought about by another factor or processing demand. In other words, in changing the context of a particular task, one can modulate the activation and examine the interaction between the response and the context employed. Dual-task interference paradigms are a clear example of this approach (Fletcher *et al.*, 1995). Of course, the factorial approach can be used in a parametric context. A simple example of this might involve examining the brain responses to increasing frequency of stimulus presentation in different contexts and looking for a differential sensitivity to increasing presentation rate. In Frith and Friston (1997), we presented auditory and visual stimuli while varying the presentation rate of the auditory stimuli. The subjects attended to either the auditory or the visual stimuli, over the complete range of presentation rates. We tested for a differential sensitivity to increasing rate, by specifying the statistical model

$$\mathbf{x}_i = \mathbf{g}_r \times \mathbf{g}_a \cdot \beta_i + [\mathbf{g}_r \mathbf{g}_a \mathbf{G}] \cdot \beta_G + \mathbf{e}_i \quad (3a)$$

where \mathbf{g}_r and \mathbf{g}_a are explanatory variables pertaining to the experimental conditions under which the scans were obtained. \mathbf{g}_r represents a mean corrected vector of presentation rates and \mathbf{g}_a is a similarly corrected vector of dummy variables taking the values 1 or -1 , denoting whether attention was directed to the auditory or visual modality, respectively. These are the *main effects* or factors in this study. $\mathbf{g}_r \times \mathbf{g}_a$ is the interaction term, modeling a difference in regression slopes (of response on presentation rate) under the different attentional sets and is simply the element by element product of \mathbf{g}_r and \mathbf{g}_a . Significant interactions were in fact found in the thalamus [see Frith and Friston (1997) for more details]. The nature of these interactions meant that thalamic activity increased with presentation rate when and only when the auditory stimuli were being attended to.

Note that we still model the main effects of rate and attention but, because we are not interested in these, they are relegated to the confound partition of the design matrix. The design matrix contains the explana-

tory variables in its columns. In Eq. (3a) the significance of the interaction is tested with the F ratio to create an SPM[F]. An equivalent inference could be obtained by using the t statistic and a contrast of parameter estimates that was 1 for β_i and zero elsewhere in the usual way (see Friston *et al.*, 1995b). The advantage of the t statistic is that it distinguishes between positive and negative contributions of the interaction term. In fact, in this instance, the F value is the t statistic squared. When using the t statistic, the main effects modeled by \mathbf{g}_r and \mathbf{g}_a can be treated as effects of interest, giving

$$\mathbf{x}_i = [\mathbf{g}_r \times \mathbf{g}_a \ \mathbf{g}_r \ \mathbf{g}_a] \cdot \boldsymbol{\beta}_i + \mathbf{G} \cdot \boldsymbol{\beta}_G + \mathbf{e}_i \quad (3b)$$

The SPM[t] is obtained using a contrast of [1 0 0]. Equations (3a) and (3b) are the same model, statistically speaking. The only difference is that to obtain the same inference one has to ensure that the interaction term is not confounded with the main effects, i.e., $[\mathbf{g}_r \times \mathbf{g}_a]^T [\mathbf{g}_r \ \mathbf{g}_a] = 0$. This was the case in the experiment described above and, when it is not, the interaction effect can be orthogonalized with respect to the main effects. In what follows we will variously use Eq. (3a) or Eq. (3b) depending on whether we want to leave the main effects in the adjusted data or remove them to look at the effects of the interaction in isolation. Adjusting the data simply means removing the effects of confounds by subtracting the estimate of $\mathbf{G} \cdot \boldsymbol{\beta}_G$. For simplicity, however, we will specify our models using the format of Eq. (3a).

Consider what would happen if we replaced the explanatory variables above with measured physiological activity in two brain regions (A and B) and identified a significant interaction in a third brain area (C). In this case activity in any significant region (C) could be explained in terms of activity in one area (A) in a way that depended on activity in the other (B). This would constitute a physiological interaction.

PHYSIOLOGICAL INTERACTIONS

An interesting extension of contribution, as described above, is the contribution of the interactions between two areas in explaining the variation in activity in a third. A powerful example of this is reported in Buechel and Friston (submitted for publication). This paper reports a fMRI study of attention to visual motion (wherein the subject was asked to detect changes in the speed of radially moving dots) that was designed to examine interactions among V5 (the human homologue of MT), posterior parietal cortex (PP), and the dorsolateral prefrontal cortex (PFC). On the basis of nonlinear structural equation modeling we were able to infer modulation of $V5 \rightarrow PP$ connections by PFC. The

question then was “how regionally specific was this modulation?” Physiological interactions were used to answer this question: In this context a modulatory effect of PFC on the efferent projections from V5 would be expressed as a contribution from V5 that depended on activity in PFC (consider how this relates to an influence of *rate* that is modulated by *attention* in the example of the previous section). By using the notation of Eq. (1) this modulatory effect can be characterized as $[C_{PP,PFC \times V5} \cdot x_{PFC}(t)] \cdot x_{V5}(t)$, where the *activity-dependent* efficacy from V5 to PP $[C_{PP,PFC \times V5} \cdot x_{PFC}(t)]$ embodies the modulatory effect of $x_{PFC}(t)$. There is another way of looking at this effect in which the coefficient $C_{PP,PFC \times V5}$ can be regarded as the *activity-independent* efficacy of the interaction term $x_{PFC}(t) \cdot x_{V5}(t)$, i.e., $[C_{PP,PFC \times V5} \cdot x_{PFC}(t)] \cdot x_{V5}(t) = C_{PP,PFC \times V5} \cdot [x_{PFC}(t) \cdot x_{V5}(t)]$. The corresponding contribution is, by analogy with Eq. (2) and Eq. (3), assessed with the statistical model:

$$\mathbf{x}_i = \mathbf{x}_{PFC} \times \mathbf{x}_{V5} \cdot \boldsymbol{\beta}_i + [\mathbf{x}_{PFC} \ \mathbf{x}_{V5} \ \mathbf{G}] \cdot \boldsymbol{\beta}_G + \mathbf{e}_i \quad (4)$$

The interaction term $\mathbf{x}_{PFC} \times \mathbf{x}_{V5}$ is simply the element by element product of the mean corrected vectors containing the activities in PFC and V5 (\mathbf{x}_{PFC} and \mathbf{x}_{V5}). If the modulatory effect is regionally specific, then we should see this second-order contribution effect in, and only in, the posterior parietal complex. This is exactly what we observed and we went on to replicate this regionally specific effect in three different subjects.

In Fig. 1 we present an analysis based on the data from one of the above subjects that asked whether PP, in mediating changes in attentional set, could be shown to modulate $V1 \rightarrow V5$ connections. In this instance we created an SPM of the t statistic testing the null hypothesis $\beta_i = 0$ based on the following statistical model:

$$\mathbf{x}_i = \mathbf{x}_{PP} \times \mathbf{x}_{V1} \cdot \boldsymbol{\beta}_i + [\mathbf{x}_{V1} \ \mathbf{x}_{PP} \ \mathbf{G}] \cdot \boldsymbol{\beta}_G + \mathbf{e}_i \quad (5)$$

The location of the voxel in V1 was 12, -81 , -4 mm (Talairach and Tournoux, 1988). These data included a contribution from V2 and “V1” is really a euphemism for the V1/V2 complex. The location of the voxel in PP was 21, -57 , 60 mm. As one might have predicted there was indeed evidence for a regionally specific interaction in the vicinity of V5. Figure 1 shows the regression implicit in Eq. (5) by plotting \mathbf{x}_i from a voxel at 54, -72 , -12 mm (adjusted for confounds, including the main effects of \mathbf{x}_{V1} and \mathbf{x}_{PP}) against the interaction term $\mathbf{x}_{PP} \times \mathbf{x}_{V1}$. Removing the effects of \mathbf{x}_{V1} and \mathbf{x}_{PP} on \mathbf{x}_i is important because, by virtue of the contribution from both V1 and PP to V5, a substantial amount of variance in \mathbf{x}_i can be explained by \mathbf{x}_{V1} and \mathbf{x}_{PP} and this would obscure the contribution of the interaction term.

The markedly positive regression slope suggests that

a significant ($Z = 5.77$, $P < 0.001$ corrected) amount of the variation in this inferior V5 satellite can be explained by the contribution of the interaction term. This is consistent with a modulation of V1 inputs by PP projections (but see below). It will be noted that the region expressing this positive modulation is somewhat posterior and inferior to V5 proper. This region is the lower right-hand focus in the inset in Fig. 1. More superior and anterior regions of cortex in this area expressed a negative interaction. Figure 2 shows an identical analysis for a voxel only 24 mm away. This somewhat unexpected juxtaposition of very significant but opposite interactions within V5 and proximate regions may reflect a functional segregation within this

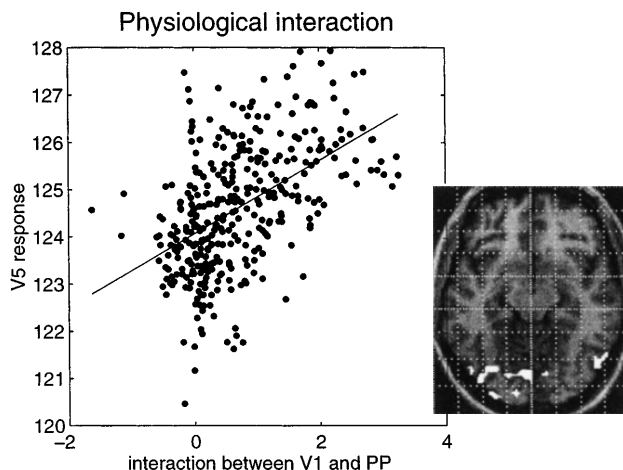


FIG. 1. Regression of (adjusted) activity near right V5 on the product of activity in V1 and posterior parietal cortex (PP). This regression shows the substantial contribution of the interaction effect to V5 responses and can be interpreted as evidence for a positive modulation of V1 \rightarrow V5 influences by PP projections. This is an example of a physiological interaction. The data come from an fMRI experiment in which a single female subject was asked to view radially moving dots under two attentional states. In one condition the subject was asked to attend to small changes in the speed of the dots (which did not actually occur) and in the second the subject was simply instructed to look at the stimuli. These data were acquired using T₂*-weighted fMRI at 2 T. One thousand volume images were acquired comprising isotropic 3-mm voxels. A volume image was acquired every 6 s. Radially moving dots on a black background were displayed in epochs of 10 scans followed by 10 scans where only a fixation dot was presented. On alternate presentations of the moving dots the subject was asked to attend to changes in their speed. The data were realigned, spatially normalized, and analyzed using the general linear model as implemented in SPM96 (Wellcome Department of Neurology; Friston *et al.* 1995c,d, 1996). In this instance the confounds (G) comprised the main effects of \mathbf{x}_{V1} and \mathbf{x}_{PP} , low-frequency time effects and global activity. The voxel from which the data were taken was the most significant voxel ($Z = 5.77$, $P < 0.001$ corrected) in the lower right focus in the inset. The location of this voxel was 54, -72 , -12 mm according to the atlas of Talairach and Tournoux (1988). The white regions in the inset correspond to regions in the SPM that survived a height threshold of $P = 0.001$ (uncorrected) and a volume threshold of $P = 0.05$ (corrected). The resulting clusters of voxel have been superimposed on a structural T₁-weighted MRI.

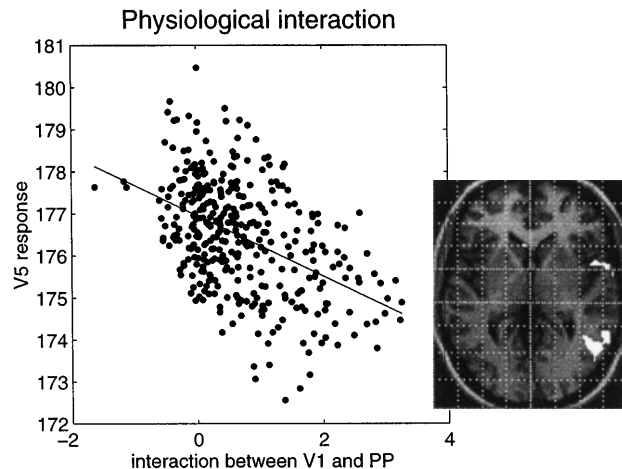


FIG. 2. As for Fig. 1 but for a voxel more anteriorly and superiorly placed in relation to V5 (57, -57 6 mm). In this instance the regression slope is negative, suggesting a negative modulation of V1 \rightarrow V5 influences by PP activity [$Z = 5.32$, $P = 0.007$ (corrected)].

complex; wherein only certain components of the motion-sensitive area are implicated in detecting change in speed or are dedicated to “optical flow.” If this were the case, then the PP modulation could be seen as turning off some aspects of visual motion processing and highlighting others in a highly regionally specific fashion. This of course is just speculation, but an interesting one.

An important aspect of this analysis, which we will return to later, is that there is an entirely equivalent and symmetric interpretation of the physiological interactions above; namely that they reflect a modulation of PP \rightarrow V5 connections by V1 activity. This is because an interaction can be construed as either a modulation of the effects of the first factor by the second or equivalently a modulation of the second's effects by the first. There is no formal distinction between what is an effect and what is a modulatory factor. This has interesting implications for psychophysiological interactions.

Physiological Interactions and Nonlinear Models of Effective Connectivity

Before turning the psychophysiological interactions it is worthwhile considering the model of effective connectivity implicit in Eq. (5) that includes the interaction between two areas when explaining the activity of a third. This model is a second-order model of the form

$$\mathbf{x}_i(t) = \sum C_{ij} \cdot \mathbf{x}_j(t) + \sum \sum C_{ijk} \cdot \mathbf{x}_j(t) \cdot \mathbf{x}_k(t), \quad (6)$$

where we have selected just two first-order effects (\mathbf{x}_{V1} and \mathbf{x}_{PP}) and one second-order effect ($\mathbf{x}_{PP} \times \mathbf{x}_{V1}$) to make a simplified statistical model corresponding to Eq. (5). Clearly Eq. (5) could be embellished with further second-order terms (e.g., $\mathbf{x}_{V1} \times \mathbf{x}_{V1}$ and

$\mathbf{x}_{pp} \times \mathbf{x}_{pp}$) and indeed effects from other areas (e.g., \mathbf{x}_{V3a} and $\mathbf{x}_{V3a} \times \mathbf{x}_{pp}$). The more terms we include in the design matrix, the more comprehensive the statistical model becomes and the more it resembles the model of effective connectivity [Eq. (6)]. Indeed the use of second-order models of this sort is the subject of current work (Buechel and Friston, submitted for publication). The key conceptual point here is that nonlinear or high-order models of effective connectivity allow for context-sensitive effects because they model interactions among inputs to a particular area. This is crucial for characterizing activity—and implicitly time-dependent effective connectivity. The inclusion of these interaction terms in simple statistical models, such as those employed in this paper, represents a step in this direction.

PSYCHOPHYSIOLOGICAL INTERACTIONS

Psychophysiological interactions have been so dubbed by analogy with psychopharmacological interactions. In psychopharmacological experiments we are interested in the interaction between the sensorimotor or cognitive-evoked responses and some pharmacological or neurotransmitter manipulation. In psychophysiological interactions we are trying to explain the physiological response in one part of the brain in terms of an interaction between prevalence of a sensorimotor or cognitive process and activity in another part of the brain. For example, by combining information about activity in the parietal region, mediating attention to a particular stimulus, and information about the stimulus, can we identify regions that respond to that stimulus when, and only when, activity in the parietal region is high? If such an interaction exists, then one might infer that the parietal area is modulating responses to the stimulus for which the area is selective. This has clear ramifications in terms of the top-down modulation of specialized cortical areas by higher brain regions.

The statistical model for psychophysiological interactions is

$$\mathbf{x}_i = \mathbf{x}_k \times \mathbf{g}_p \cdot \beta_i + [\mathbf{x}_k \mathbf{g}_p \mathbf{G}] \cdot \beta_G + \mathbf{e}_i \quad (7)$$

The term $\mathbf{x}_k \times \mathbf{g}_p \cdot \beta_i$ represents the psychophysiological interaction between the physiological activity in region k and some psychological or experimental parameter of the experimental design \mathbf{g}_p and is constructed by multiplying the two effects as in the previous sections. Consider the above example concerning the contribution of V1 to V5. We have shown that this is increased when PP activity is high and we infer that this is a modulatory effect of PP that mediates attention. However, we can explicitly test the hypothesis that the contribution from V1 is higher under attention to the visual motion in parts of V5 by looking for the regional

specificity of a psychophysiological interaction between attention and V1 activity and ensuring that this is maximally expressed in V5. In other words, we replace PP activity in the previous section by a psychological variable that denotes differences in attentional set. Figure 3 shows the variables employed: \mathbf{x}_k was the activity in V1 (corrected to a minimum of 0) and \mathbf{g}_p corresponds to attentional set, containing elements of 1 when the subject was expecting to detect changes in motion and -1 when she was not. Changes in attentional set were instantiated by a verbal cue several seconds before the presentation of the moving stimuli. The interaction effect is shown in the lower panel. The most significant effects were indeed seen near V5. Figure 4 shows significant ($P < 0.05$ corrected) interactions at this level of the brain and the activity and the adjusted time-series for a voxel at $42 -78 -9$ (upper

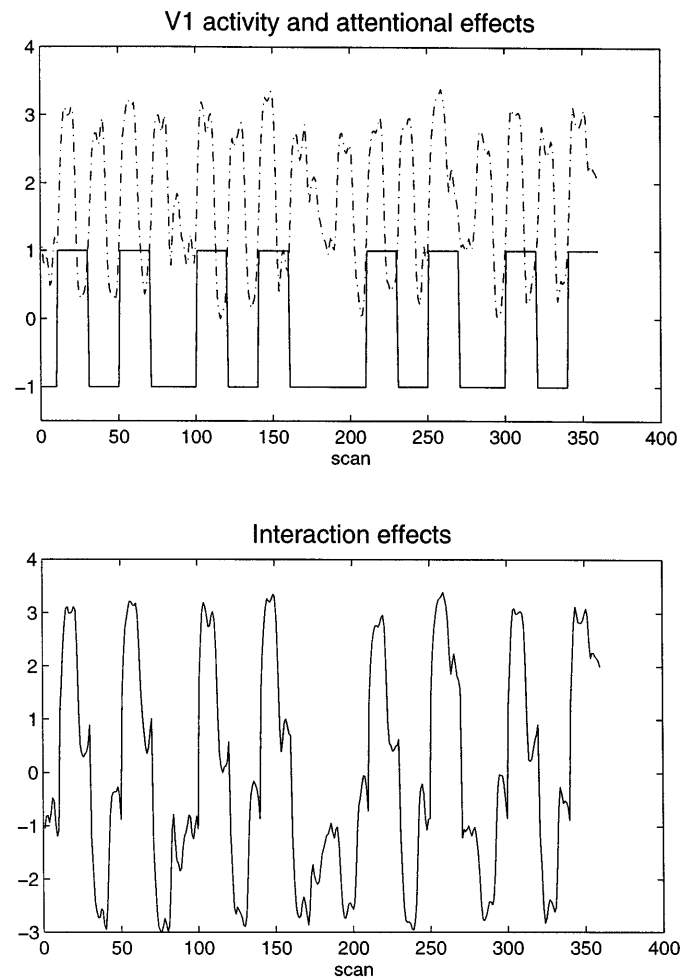


FIG. 3. The effects used to model psychophysiological interaction between activity in V1 (12, -81 , -4 mm) and attentional set. These variables (top) correspond to \mathbf{x}_k (dashed and dotted line) and \mathbf{g}_p (solid line) in the main text. The bottom panel shows the element-wise product of these effects and corresponds to the interaction effect ($\mathbf{x}_k \times \mathbf{g}_p$).

panel). In this context the interaction can be seen as a significant difference in the regression slopes of V5 activity on V1 activity when assessed under the two attentional conditions (lower panel of Fig. 4). Using the concept of contribution, this means that attention to visual motion significantly modulated the contribution of V1 to V5.

We now consider an alternative perspective on this interaction. Recall that an interaction can be equivalently formulated in terms of the modulation of an effect of the first factor on the responses to the second or vice versa. This means that the V5 responses can reflect either (i) a modulation of the contribution from V1 by

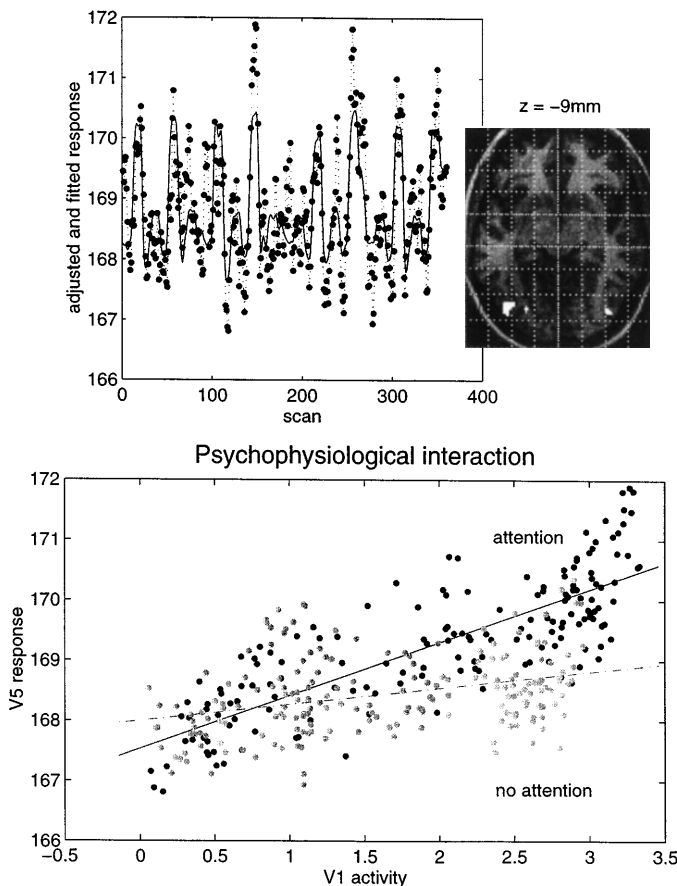
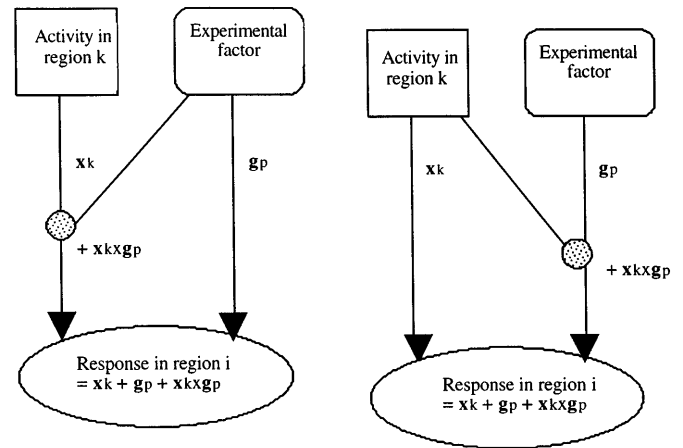


FIG. 4. (Top) SPM thresholded at $P = 0.001$ (height - uncorrected) and $P = 0.05$ (volume - corrected) superimposed on a structural T_1 -weighted image. This SPM tested for a significant psychophysiological interaction between activity in V1 and attention to visual motion. The most significant effects are seen in the vicinity of V5 (white region on the bottom right) and the associated time series of the most significant ($Z = 4.46$) voxel in this region is shown in the top panel (fitted data, line; adjusted data, dots). (Bottom) Regression of V5 activity (at 42, -78, -9 mm) on V1 activity when the subject was asked to attend to changes in the speed of radially moving dots and when she was not. The lines correspond to the regression. The dots correspond to the observed data adjusted for confounds other than the main effects of V1 activity (dark gray dots, attention; light gray dots, no attention). Attention can be seen to augment the contribution of V1 to V5 activity.

Modulation of stimulus-specific responses



Context-specific modulation of effective connectivity

FIG. 5. Schematic illustrating the two perspectives on psychophysiological interactions. On the left the psychological or experimental factor can be thought of as modulating the contribution of one area to another in terms of neuronal activity. This can be understood as a condition or context-sensitive change in the contributory aspects of effective connectivity between the two areas. On the right the alternative interpretation is that the contribution modulates responsiveness of the target area to the psychological factor. This perspective translates into a modulation of condition-specific responses by afferents from another area.

attention (the interpretation that we have focused on) or (ii) a modulation of attention-specific responses by V1 inputs. More generally a psychophysiological interaction can be seen as (i) a context or functionally specific change in the contribution of one area to another or (ii) the modulation of responses in one area to the psychological or experimental variable by the contribution from another area. Figure 5 makes this point schematically. In the above example the contribution is from a lower area (V1) to a higher area (V5) and the first perspective seems more natural, namely that attention modulates the forward influence of V1 on V5. However, the alternative perspective is equally valid in the sense that attention-specific responses in V5 are realized only in the presence of stimulus-specific inputs from V1. In the next example we revisit this distinction when the contribution is from a higher area to a lower area.

Changes in Contribution or Modulation?

We now present a second example of psychophysiological interactions involving posterior parietal cortex. Subjects were asked to view (degraded) faces and nonface (object) controls using PET, before and after

priming with the nondegraded stimuli. Subjects were aware that they would be seeing either face or nonface stimuli irrespective of the priming. This represents a factorial design with two factors (faces vs nonfaces and priming vs nonpriming). The interaction between activity in the medial parietal region (6, -74, 36 mm) and presence of faces was most significantly expressed in the right inferotemporal region (-38, -30, -20 mm; $Z = 3.97$, $P < 0.0007$ uncorrected) in the region of the fusiform and parahippocampal gyri. Figure 6 shows the resulting SPM(t), transformed to a SPM(Z) thresholded at $P = 0.01$ (uncorrected), as a maximum intensity projection and the design matrix used to specify the statistical model. In this instance changes in medial parietal activity were introduced experimentally by preexposure of the nondegraded stimuli before some scans and not others. These results can be interpreted as a priming-dependent instantiation of attentional, memory, or learning differences in face-specific responses, in the inferotemporal region that are mediated by interactions with PP. PP has been previously implicated in face processing on the basis of an analysis of

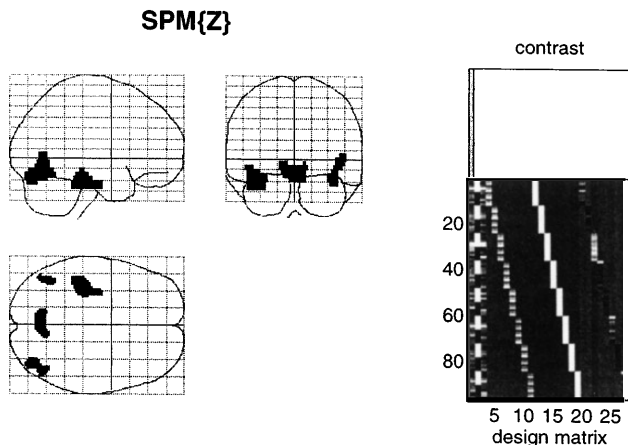


FIG. 6. (Left) Maximum intensity projection of an SPM that identifies areas whose activity can be explained on the basis of an interaction between the presence of faces in visually presented stimuli and activity in a reference location in the medial parietal cortex (6, -74, 36 mm). The largest effects were observed in the right inferotemporal and medial temporal regions. This SPM has been thresholded at $P = 0.01$ (uncorrected) and is displayed in standard format, providing three orthogonal views of the brain in the space of Talairach and Tournoux (1988). This study involved six normal subjects who had 12 PET scans while viewing degraded (binarized) faces or nonface objects presented every 3 s. Before half the scans of each stimulus type the subjects were "primed" with grayscale (nondegraded) versions of the same stimuli. This study therefore represents a 2×2 factorial design with two factors (faces vs nonfaces and primed vs not primed). The priming causes immediate and long-lasting recognition of the degraded stimuli on re-presentation. (Right) Design matrix used in this analysis. The left-hand column contains the interaction effect $\mathbf{x}_k \times \mathbf{g}_p$ in the main text, followed by the confounds. The first two confounds are the main effects of \mathbf{x}_k and \mathbf{g}_p , then subject-specific time effects, subject effects, subject-specific global effects, and a constant.

effective connectivity using structural equation modeling and PET data (McIntosh *et al.*, 1994a). Note that we could have modeled the priming effect explicitly in our design matrix (as in a conventional psychological interaction) but chose to substitute medial parietal activity in its place, enabling a more mechanistic inference, namely, not only do inferotemporal responses show a priming-dependent effect, but this effect is mediated by modulatory influences from a higher (PP) area. However, in light of the previous paragraph, is this the only interpretation?

The demonstration that right inferotemporal regions receive a significant contribution from an interaction between activity in the parietal region and the presence of faces in visual stimuli has two interpretations: either (i) parietal activity is modulating (increasing) sensitivity to afferent activity elicited by faces or (ii) inputs from parietal cortex are modulated (increased) in the context of face-specific inferotemporal processing. One way of seeing this distinction very clearly is to consider the regressions in the upper panel of Fig. 7. These present regression of activity in the inferotemporal region on (mean corrected) activity in PP for faces (dark gray dots) and nonface objects (light gray dots). The data come from the voxel identified by the cross hairs in the lower panel. In one view inferotemporal responses discriminate between face and nonface stimuli only when PP activity is high (the dark and light gray dots segregate only at extreme values of PP activity). Alternatively the slope of the regression (i.e., the contribution from PP to inferotemporal regions) is positive when faces are being processed and negative otherwise, reflecting face-specific contributory aspects of effective connectivity.

These two perspectives are reflected in the two biological interpretations of this finding presented in Dolan *et al.* (submitted for publication). First, modulatory backward projections from the parietal regions to inferotemporal cortex modulate inputs from extrastriate regions, where the intrinsic synaptic connections of the latter are responsible for configuring the face-specificity of the inferotemporal response (e.g., Perret *et al.*, 1986; Desimone *et al.*, 1984). Second, the interaction can be explained by augmented neuronal interactions between inferotemporal and cortical areas due to long-term potentiation of inferotemporal connectivity (following priming) that facilitates promulgation of activity onward, from the inferotemporal regions to the parietal cortex, that is, specific to face stimuli.

In short the activity of the inferotemporal region depends on an interaction between inputs from extrastriate areas eliciting face-specific responses and input from parietal regions (involved in more generic aspects of perceptual synthesis). These interactions can be construed as (i) a contribution from the parietal region in, and only in, the context of face processing [as

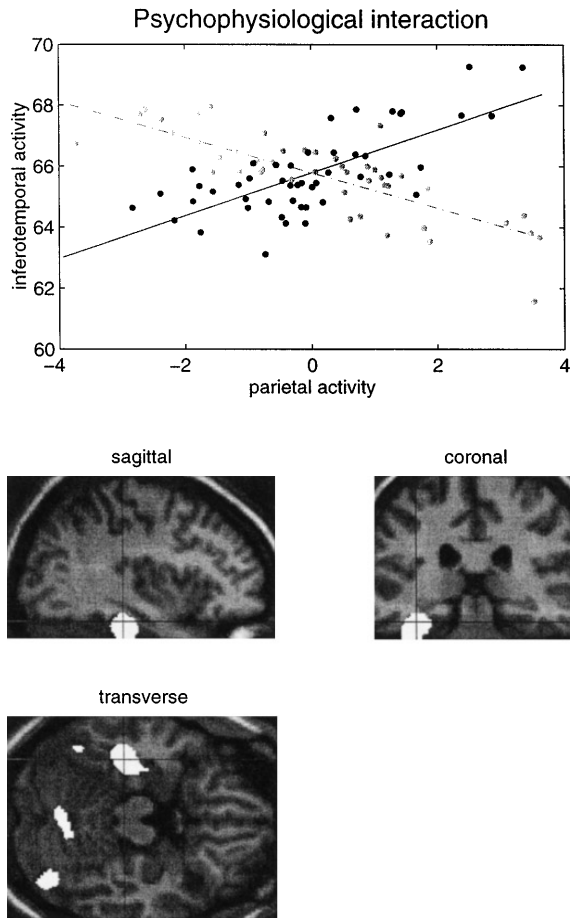


FIG. 7. (Top) The activity in the most significant voxel ($Z = 3.97$, $P < 0.0007$ uncorrected at $-38, -30, -20$ mm) is displayed as a function of (mean corrected) activity in the medial parietal voxel used to specify the psychophysiological interaction (at $6, -74, 36$ mm). The dark gray dots correspond to (adjusted for confounds only) activity while viewing faces and the light gray dots correspond to nonface stimuli. The essence of this effect can be seen by (i) noting that this region differentiates between faces and nonfaces when, and only when, medial parietal activity is high, or alternatively (ii) noting that the contribution of parietal activity to inferotemporal activity is greater in the context of face-processing. (Bottom) Suprathreshold regions from the previous figure superimposed on a T_1 -weighted structural image. The cross-hairs localize the voxel from which the above data were taken.

a result of early perceptual learning in the inferotemporal regions that facilitate reciprocal interactions (Roelfsema *et al.*, 1997) between inferotemporal and parietal areas when primed faces are seen] or (ii) a modulation of face-specific responses by parietal afferents where, in this example, the modulatory input is dependent on the reciprocal interactions above. The interpretation of this interaction is more complicated than some because PP activity must itself be a function of processing in the face-specific regions to which it contributes. In other instances the genesis of top-down modulation may be more straightforward than in the example considered here, for example, stimulus presentation during inde-

pendent manipulations of attentional set as in the previous example using radial motion.

DISCUSSION

In this paper we have introduced the notion of testing for interactions between brain responses and some experimental (sensory or task-related) parameter in neuroimaging. We have referred to these interaction effects as psychophysiological interactions to distinguish them from interactions based solely on designed effects in conventional factorial designs and interactions among neurophysiological measurements (i.e., physiological interactions). We have framed psychophysiological interactions in terms of functional integration by noting that the degree to which activity in one area can be predicted on the basis of activity in another corresponds to the contribution of the second to the first. The potential importance of these approaches lies in (i) conferring a degree of functional specificity on the contribution aspect of effective connectivity and (ii) providing a model of modulation, where the contribution can be considered to modulate responses to the psychological or stimulus-specific parameter used to define the interaction. Although distinct in neurobiological terms, these are equivalent perspectives on the same underlying interaction.

Psychophysiological interactions are interesting from two points of view. First, the explanatory variables used to predict activity (i.e., the response variable) in any brain region comprise a standard predictor variable based on the experimental design (e.g., the presence or absence of a particular stimulus attribute) and a response variable from another part of the brain. The second reason that this analysis is interesting is that it uses techniques usually used to make inferences about functional specialization to infer something about functional integration, in this instance, context-sensitive contributions.

Psychophysiological Interactions and Changes in Effective Connectivity

We have interpreted the interactions between context and activity in a remote area from two perspectives. First, this interaction could represent a context-sensitive change in the contribution of the remote area to the area in question and second, it could represent enhanced responses in the reference area, to the context, that are mediated by the (fixed) contribution from the distal area. The first interpretation is related conceptually to approaches that look for context-sensitive changes in functional or effective connectivity. The most closely related work in this viewpoint is the regression analysis of changes in effective connectivity presented in Friston *et al.* (1993b). In this work time-dependent changes in effective connectivity to the

prefrontal cortex were estimated by solving Eq. (1) using data from the beginning and end of an experiment and characterizing the differences. A similar approach was used to assess modulatory connections between V2 and V1 (Friston *et al.*, 1995a). In the latter study we used multiple regression to analyze the effective connectivity from all voxels in extrastriate cortex to V1 under two different contexts, when V1 activity was high and when V1 activity was low. The only region that showed a change in effective connectivity was V2, suggesting that V1 intrinsically modulated its sensitivity to V2 inputs in a modulatory fashion. Note that both these analyses included the effects from many possible sources of input and were framed in terms of *effective connectivity*. In this paper the effects modeled derive from only one region and are framed in terms of *contribution*. The analysis presented in this paper does not constitute an analysis of effective connectivity for this reason (but see below). Conceptually related approaches include the comparison of connection strengths in different contexts using interregional correlations or structural equation modeling: Horwitz *et al.* (1992) have shown profound changes in correlations among visual processing areas in face-matching, relative to dot (position)-matching tasks. However, in relation to psychophysiological interactions, it should be noted that a test for a change in regression slope does not constitute a test for a change in correlation, even in the context of a single explanatory region.

MacIntosh *et al.* (1994b) have used structural equation modeling to show dramatic and opposite changes in effective connectivity, in the dorsal and ventral visual pathways, with the delay period in a working memory task. Psychophysiological interactions can be assessed using these sorts of approaches: Consider the visual attention to motion study presented above. A structural equation model of the effective connectivity might show an enhanced connection strength from V1 to V5 under attention to motion, in relation to the no-attention context. For this to work, however, one has to specify V5 in the model, whereas in the current approach this region emerged spontaneously from the analysis. Generally, however, it would be difficult to test for psychophysiological interactions with differences in functional or effective connectivity because the psychological factor will not necessarily have two discrete levels (for example, it may be a parametric variable). In this more general case a structural equation model would have to involve the explicit introduction of interaction terms as "moderator variables" (Kenny and Judd, 1984).

The qualifications usually associated with estimates of effective connectivity should also be borne in mind when interpreting psychophysiological interactions: Namely (i) one cannot guarantee that the effective connections are direct (i.e., they could be mediated

through other areas) and (ii) there is always the problem of common input. In the context of psychophysiological interactions this means that a third area, which shows context-sensitive responses, may be providing input to the two areas implicated in the psychophysiological interaction. In a sense, however, this common input to both areas is itself context-sensitive and could be identified using psychophysiological interactions with the third area as the source region.

Psychophysiological Interactions and Factorial Designs

Psychophysiological interactions generally depend on factorial experimental designs, wherein one can introduce neuropsychological changes in one brain system that are uncorrelated with the stimulus or cognitive context one hopes to see an interaction with. We make this point explicit, suggesting that this is another example of the usefulness of factorial experiments: Although it is possible to test for psychophysiological interactions in almost any experimental design, the use of factorial designs ensures that any psychophysiological interactions will be detected with a fair degree of sensitivity. This is because the activities in the source area, the psychological context, and the interaction between them, will be roughly orthogonal and therefore one can use the first two as confounds with impunity. The converse situation, in which only one stimulus or cognitive factor has been changed, may render the activity in the source area and changes in the factor correlated. If this is the case, there is no guarantee that the interaction will be independent of either and its effect may be difficult to detect in the presence of the "main effects."

Neurobiological Mechanisms

To illustrate some of the finer aspects of interpreting psychophysiological interactions, and to illustrate the importance of neurobiological constraints on these interpretations, we will consider the possible mechanisms behind the psychophysiological interaction in the perceptual priming example presented above in some detail. The two interpretations (see Fig. 5) of this interaction (between parietal activity and face stimuli) can be construed as different perspectives on the same phenomenon. However, any mechanistic explanation must accommodate several facts. First, the priming-dependent effect (increased inferotemporal responses following priming by the nondegraded and immediately recognizable faces) must be due either to delay period activity that endures after priming, in some neuronal system, or to long-term changes in synaptic efficacy. The fact that priming effects persist for several days and the delay period employed in this study lasts for many minutes suggests that the interaction is medi-

ated by changes in synaptic efficacy. Second, inferotemporal responses are face-specific whereas the parietal responses are not (they are recognition-specific, also responding to the recognition of nonface objects). The first explanation (top-down modulation of face-specific responses) would require recognition-dependent responses in parietal regions to occur before the face-selective inferotemporal responses they were consolidating. Assuming that parietal unit responses are not in themselves face-selective, this would violate the laws of causality. The second explanation (recognition-dependent responses that are facilitated by rapid perceptual learning in, and only in, the feed-forward pathways through inferotemporal regions) has similar shortcomings because it posits recognition-dependent changes in synaptic efficacy that precede recognition. In other words, those components of associative plasticity that depend on a perceptual synthesis involving higher order areas cannot occur unless the synthetic neuronal systems "reach back" and influence plasticity in earlier feed-forward pathways. A mechanism that resolves these difficulties is rapid perceptual learning, mediated by potentiation of synaptic connections intrinsic to the inferotemporal region that involves, or depends on, parietal afferents. Simple associative plasticity, in the context of nonlinear unit responses, would be a sufficient mechanism to implement this. In this construct perceptual learning is facilitated by recognition through consolidation of feed-forward synaptic inputs to inferotemporal units that depends on conjoint activation of parietal afferents. Recognition-dependent parietal responses, implicated in perceptual integration and synthesis, are themselves dependent on inferotemporal inputs during presentation of easily recognized priming stimuli. This (somewhat heuristic) model, which depends on the interaction between extrastriate and parietal inputs to inferotemporal regions, reconciles both of the above views in the sense that: (i) face-specific responses are augmented by parietal afferents though selective potentiation of the presynaptic inputs to the inferotemporal region and (ii) parietal activity is vicariously modulated by potentiation of feed-forward connections intrinsic to the inferotemporal region. In other words, the parietal contribution can be seen as a "self-connection" to the inferotemporal region, through PP, that is facilitated by face-specific perceptual learning. The question that remains is: "Is the perceptual learning during priming sufficient to explain priming-dependent augmentation of inferotemporal responses to degraded stimuli or are the recognition-dependent responses in parietal regions that ensue required to maintain enhanced responses during presentation of the degraded stimuli?" In other words, although the parietal responses may be crucial during the presentation of the priming stimuli, they may be incidental once changes in synaptic efficacy have occurred. Unfortun-

nately we cannot answer this question using neuroimaging alone (although some interesting experiments using magneto-stimulation suggest themselves).

Extensions

In this paper we have restricted the analysis of psychological interactions to the simplest formulation. It is, of course, possible to include the activity of other areas as in Eq. (2a). This would make the models more comprehensive and closer to multiple regression approaches to effective connectivity (see Physiological Interaction and Nonlinear Models of Effective Connectivity). Another extension would be to use polynomial expansions (Buechel *et al.*, 1996) to test for interactions that are expressed in terms of nonlinear functions of activity in distal areas.

ACKNOWLEDGMENTS

The work was supported by the Wellcome Trust. We thank Chris Frith, Cathy Price, Richard Frackowiak, and the two anonymous reviewers for invaluable comments.

REFERENCES

- Aertsen, A., and Preissl, H. 1991. Dynamics of activity and connectivity in physiological neuronal networks. In *Non Linear Dynamics and Neuronal Networks* (H. G. Schuster, Ed.), pp. 281–302. VCH, New York.
- Buechel, C., and Friston, K. J. Assessing modulatory interactions among visual areas using nonlinear structural equation modelling and fMRI. Submitted for publication.
- Buechel, C., Wise, R. S. J., Mummary, C., and Friston, K. J. 1996. Nonlinear regression in parametric activation studies. *NeuroImage* **4**:60–66.
- Desimone, R., Albright, T. D., Gross, C. G., and Bruce, C. 1984. Stimulus selective properties of inferior temporal neurons in the macaque. *J. Neurosci.* **4**:2051–2062.
- Dolan, R. J., Fink, G. R., Rolls, E., Booth, M., Frackowiak, R. S. J., and Friston, K. J. How the brain learns to see in an impoverished context. Submitted for publication.
- Fletcher, P. C., Frith, C. D., Grasby, P. M., Shallice, T., Frackowiak, R. S. J., and Dolan, R. J. 1995. Brain systems for encoding and retrieval of auditory-verbal memory. *Brain* **118**:401–416.
- Friston, K. J., Frith, C., Passingham, R. E., Liddle, P., and Frackowiak, R. S. J. 1992. Motor practice and neurophysiological adaptation in the cerebellum: A positron tomography study. *Proc. R. Soc. London Ser. B* **248**:223–228.
- Friston, K. J., Frith, C. D., Liddle, P. F., and Frackowiak, R. S. J. 1993a. Functional connectivity: The principal component analysis of large (PET) data sets. *J. Cereb. Blood Flow Metab.* **13**:5–14.
- Friston, K. J., Frith, C. D., and Frackowiak, R. S. J. 1993b. Time-dependent changes in effective connectivity measured with PET. *Hum. Brain Mapping* **1**:69–80.
- Friston, K. J., Ungerleider, L. G., Jezzard, P., and Turner, R. 1995a. Characterizing modulatory interactions between areas V1 and V1 in human cortex: A new treatment of functional MRI data. *Hum. Brain Mapping* **2**:211–224.
- Friston, K. J., Holmes, A. P., Worsley, K. J., Poline, J.-B., Frith, C. D., and Frackowiak, R. S. J. 1995b. Statistical parametric maps in

- functional imaging: A general linear approach. *Hum. Brain Mapping* **2**:189–210.
- Friston, K. J., Holmes, A. P., Poline, J.-B., Grasby, P. J., Williams, S. C. R., Frackowiak, R. S. J., and Turner, R. 1995c. Analysis of fMRI time-series revisited. *NeuroImage* **2**:45–53.
- Friston, K. J., Ashburner, J., Frith, C. D., Poline, J.-B., Heather, J. D., and Frackowiak, R. S. J. 1995d. Spatial registration and normalisation of images. *Hum. Brain Mapping* **2**:165–189.
- Friston, K. J., Williams, S., Howard, R., Frackowiak, R. S. J., and Turner, R. 1996. Movement related effects in fMRI time series. *Mag. Res. Med.* **35**:346–355.
- Frith, C. D., and Friston, K. J. 1997. The role of the thalamus in 'top down' modulation of attention to sound. *NeuroImage* **0**:00–00.
- Gerstein, G. L., and Perkel, D. H. 1969. Simultaneously recorded trains of action potentials: Analysis and functional interpretation. *Science* **164**:828–830.
- Horwitz, B., Grady, C. L., Haxby, J. V., Ungerleider, L. G., Schapiro, M. B., Mishkin, M., and Rapoport, S. I. 1992. Functional associations among human posterior extrastriate brain regions during object and spatial vision. *J. Cog. Neurosci.* **4**:311–322.
- Kenny, D. A., and Judd, C. M. 1984. Estimating nonlinear and interactive effects of latent variables. *Psychol. Bull.* **96**: 201–210.
- McIntosh, A. R., Grady, C. L., Ungerleider, L. G., Haxby, J. V., Rapoport, S. I., and Horwitz, B. 1994a. Network analysis of cortical visual pathways mapped with PET. *J. Neurosci.* **14**:655–666.
- McIntosh, A. R., Horwitz, B., Haxby, J. V., Ungerleider, L. G., and Grady, C. L. 1994b. Functional cortical networks during short-term recognition memory for faces. *Abstr. Soc. Neurosci.* **20**:362.
- Paus, T., Marrett, S., Worsley, K. J., and Evans, A. 1996. Imaging motor-to-sensory discharges in the human brain: An experimental toll for the assessment of functional connectivity. *NeuroImage* **4**:78–86.
- Perret, D. I., Mistlin, A. J., Potter, D. D., Smith, P. A. J., Head, A. S., Chitty, A. J., Broenimann, R., Milner, A. D., and Jeeves, M. A. 1986. Functional organization of visual neurones processing face identity. In *Aspects of Face Processing* (H. Ellis, M. A. Jeeves, F. Newcombe, and A. W. Young, Eds.), pp. 187–198. Nijhoff, Dordrecht.
- Roelfsema, P. R., Engel, A. K., Konig, P., and Singer, W. 1997. Visuomotor integration is associated with zero time-lag synchronization among cortical areas. *Nature* **385**:157–161.
- Talairach, J., and Tournoux, P. 1988. *A Co-planar Stereotaxic Atlas of a Human Brain*. Thieme, Stuttgart.
- Zeki, S., Watson, J. D. G., Lueck, C. J., Friston, K. J., Kennard, C., and Frackowiak, R. S. J. 1991. A direct demonstration of functional specialisation in human visual cortex. *J. Neurosci.* **11**:641–649.