

this cluster of properties tend to co-occur? Second, I want to say something about the extension of the concept; to propose a hypothesis about which cognitive systems are, in fact, modular. This second line of inquiry will provide the main structure of the discussion, the first emerging as opportunity provides targets. By the time I've finished, I shall have made the following suggestions:

(a) That the set of processors for which the modularity view currently seems most convincing is coextensive with a functionally definable subset of the cognitive systems.

(b) That there is some (more or less *a priori*) reason to believe that cognitive systems which do *not* belong to that functionally defined subset may be, in important respects, *non*-modular (e.g., mediated by horizontal faculties). And finally,

(c) I shall make some depressed remarks along the following lines: though the putatively nonmodular processes include some of the ones that we would most like to know about (thought, for example, and the fixation of belief), our cognitive science has in fact made approximately no progress in studying these processes, and this may well be *because* of their nonmodularity. It may be that, from the point of view of practicable research strategy, it is only the modular cognitive systems that we have any serious hope of understanding. In which case, convincing arguments for non-modularity should be received with considerable gloom.

## PART II A FUNCTIONAL TAXONOMY OF COGNITIVE MECHANISMS

I want to argue that the current best candidates for treatment as *modular* cognitive systems share a certain functional role in the mental life of organisms; the discussion in this section is largely devoted to saying which functional role that is. As often happens in playing cognitive science, it is helpful to characterize the functions of psychological systems by analogy to the organization of idealized computing machines. So, I commence with a brief digression in the direction of computers.

When philosophers of mind think about computers, it is often Turing machines that they are thinking about. And this is understandable. If there is an interesting analogy between minds qua

minds and computers *qua* computers, it ought to be possible to couch it as an analogy between minds and Turing machines, since a Turing machine is, in a certain sense, as general as any kind of computer can be. More precisely: if, as many of us now suppose, minds are essentially symbol-manipulating devices, it ought to be useful to think of minds on the Turing-machine model since Turing machines are (again "in a certain sense") as general as any symbol-manipulating device can be.

However, as we have already had reason to observe, Turing machines are also very simple devices; their functional architecture is exhaustively surveyed when we have mentioned a small number of interacting subsystems (tape, scanner, printer, and executive) and a small inventory of primitive machine operations (stop, start, move the tape, read the tape, change state, print). Moreover—and this is the point of present concern—Turing machines are *closed* computational systems; the sole determinants of their computations are the current machine state, the tape configuration, and the program, the rest of the world being quite irrelevant to the character of their performance; whereas, of course, organisms are forever exchanging information with their environments, and much of their psychological structure is constituted of mechanisms which function to mediate such exchanges. If, therefore, we are to start with anything like Turing machines as models in cognitive psychology, we must think of them as embedded in a matrix of subsidiary systems which affect their computations in ways that are responsive to the flow of environmental events. The function of these subsidiary systems is to provide the central machine with information about the world; information expressed by mental symbols in whatever format cognitive processes demand of the representations that they apply to.

I pause to note that the format constraint on the subsidiary systems is vital. *Any* mechanism whose states covary with environmental ones can be thought of as registering information about the world; and, given the satisfaction of certain further conditions, the output of such systems can reasonably be thought of as *representations* of the environmental states with which they covary. (See Dretske, 1981; Stampe, 1977; Fodor, forthcoming.) But if cognitive processors are *computational* systems, they have access to such information solely in virtue of the *form* of the representations in which it is

couched. Computational processes are, by definition, *syntactic*; a device which makes information available to such processes is therefore responsible for its format as well as its quality. If, for example, we think of such a device as writing on the tape of a Turing machine, then it must write *in a language that the machine can understand* (more precisely, in the language in which the machine computes). Or, to put it in a psychological-sounding way, if we think of the perceptual mechanisms as analogous to such devices, then we are saying that *what perception must do is to so represent the world as to make it accessible to thought*. The condition on appropriateness of format is by way of emphasizing that not every representation of the world will do for this purpose.

I wish that I knew what to call the "subsidiary systems" that perform this function. Here are some possibilities that I have considered and—with varying degrees of reluctance—decided to reject:

—'Perceptual systems' would be the obvious choice except that, as we shall presently see, perception is not the only psychological mechanism that functions to present the world to thought, and I would like a term broad enough to embrace them all. Moreover, as will also become apparent, there are important reasons for not viewing the subsidiary systems as effecting the fixation of belief. By contrast, perception is a mechanism of belief fixation par excellence: the normal consequence of a perceptual transaction is the acquisition of a perceptual belief. (Having entered this caveat, I shall nevertheless often speak of the subsidiary systems as mechanisms of perceptual analysis. For most purposes it is harmless to do so and it does simplify the exposition.)

—I have sometimes thought of calling these subsidiary systems 'compilers', thereby stressing that their output consists of representations that are accessible to relatively central computational processes. But that way of talking leads to difficulties too. Real compilers are functions from programs onto programs, programs themselves being (approximately) sequences of instructions. But not much of what perception makes available to thought is plausibly viewed as a program. Indeed, it is partly the attempt to force perceptual information into that mold which engenders procedural semantics, the identification of perceptual categories with action schemes, and other such aberrations of theory. (For discussion, see Fodor, 1981a, chapter 8.)

—One could try calling them 'transducers' except that, on at least one usual understanding (see Lowenstein, 1960), transducers are analog systems that take proximal stimulations onto more or less precisely covarying neural signals. Mechanisms of transduction are thus *contrasted* with computational mechanisms: whereas the latter may perform quite complicated, inference-like transformations, the former are supposed—at least ideally—to preserve the informational content of their inputs, altering *only* the format in which the information is displayed. We shall see, however, that representations at the interface between (what I have been calling) 'subsidiary' and 'central' systems exhibit levels of encoding that are quite abstractly related to the play of proximal stimulation.

Pylyshyn and I (1981) have called these subsidiary systems 'compiled transducers', using the 'compiled' part to indicate that they have an internal computational structure and the 'transducer' part to indicate that they exhibit a certain sort of informational encapsulation that will presently loom large in this discussion. I think that usage is all right given the explication, but it admittedly hasn't much to do with the conventional import of these terms and thus probably produces as much confusion as it avoids.

It is, perhaps, not surprising that computer theory provides no way of talking that does precisely the job I want to do. Computers generally interface with their environments *via some human being* (which is what makes them computers rather than robots). The programmer thus takes on the function of the subsidiary computational systems that I have been struggling to describe—viz., by providing the machine with information about the world in a form in which the machine can use it. Surprising or not, however, it is a considerable nuisance. Ingenuity having failed me completely, I propose to call them variously 'input systems', or 'input analyzers' or, sometimes, 'interface systems'. At least this terminology emphasizes that they operate relatively early on. I rely on the reader to keep it in mind, however, that input systems are *post-transductive* mechanisms according to my usage. Also that switches from one of the epithets to another usually signify no more than a yen for stylistic variation.

So, then, we are to have a trichotomous functional taxonomy of psychological processes; a taxonomy which distinguishes transducers, input systems, and central processors, with the flow of

input information becoming accessible to these mechanisms in about that order. These categories are intended to be exclusive but not, of course, to exhaust the types of psychological mechanisms that a theory of cognition might have reason to postulate. Since the trichotomy is not exhaustive, it is left wide open that there may be modular systems that do not subserve any of these functions. Among the obvious candidates would be systems involved in the motor integration of such behaviors as speech and locomotion. It would please me if the kinds of arguments that I shall give for the modularity of input systems proved to have application to motor systems as well. But I don't propose to investigate that possibility here.

Input systems function to get information into the central processors; specifically, they mediate between transducer outputs and central cognitive mechanisms by encoding the mental representations which provide domains for the operations of the latter. This does not mean, however, that input systems *translate* from the representations that transducers afford into representations in the central code. On the contrary, translation preserves informational content and, as I remarked above, the computations that input systems perform typically do not. Whereas transducer outputs are most naturally interpreted as specifying the distribution of stimulations at the 'surfaces' (as it were) of the organism, the input systems deliver representations that are most naturally interpreted as characterizing the arrangement of *things in the world*. Input analyzers are thus inference-performing systems within the usual limitations of that metaphor. Specifically, the inferences at issue have as their 'premises' transduced representations of proximal stimulus configurations, and as their 'conclusions' representations of the character and distribution of distal objects.

It is hard to see how a computer could fail to exhibit mechanisms of transduction if it is to interface with the world at all. But it is perfectly possible to imagine a machine whose computations are appropriately sensitive to environmental events but which does *not* exhibit a functional distinction between input systems and central systems. Roughly, endorsing this computational architecture is tantamount to insisting upon a perception/cognition distinction. It is tantamount to claiming that a certain class of computational problems of 'object identification' (or, more correctly, a class of

computational problems whose solutions consist in the recovery of certain proprietary descriptions of objects) has been 'detached' from the domain of cognition at large and handed over to functionally distinguishable psychological mechanisms. Perceptual analysis is, according to this model, not, strictly speaking, a species of thought. (The reader is again reminded, however, that the identification of input processing with perceptual analysis is itself only approximate. This will all presently sort itself out; I promise.)

Given the possibility in principle that the perceptual mechanisms could be continuous with the higher cognitive processes, one is tempted to ask what the point of a trichotomous functional architecture could be. What, teleologically speaking, might it buy for an organism that has transducers and central cognitive processors to have input analyzers as well? I think there probably *is* an answer to this question: Implicit in the trichotomous architecture is the isolation of perceptual analysis from certain effects of background belief and set; and, as we shall see, this has implications for both the speed and the objectivity of perceptual integration. It bears emphasis, however, that putting the teleological issues in the way I just did involves some fairly dubious evolutionary assumptions. To suppose that the issue is *Why, given that there are central processors, should there be input systems as well?* is to take for granted that the former should be viewed as phylogenetically prior to the latter. However, an equally plausible story might have it the other way 'round—viz., that input analyzers, with their (as I shall argue) relatively rigid domain specificity and automaticity of functioning, are the aboriginal prototypes of inference-making psychological systems. Cognitive evolution would thus have been in the direction of gradually freeing certain sorts of problem-solving systems from the constraints under which input analyzers labor—hence of producing, as a relatively late achievement, the comparatively domain-free inferential capacities which apparently mediate the higher flights of cognition. (See Rozen, 1976, where the plausibility of this picture of cognitive phylogeny is impressively defended.)

In any event, the justification for postulating a functionally individuated class of input analyzers distinct from central cognitive mechanisms must finally rest on two sorts of evidence: I have to show that there are interesting things that the input analyzers have in common; and I have to show that there are interesting respects

in which they differ from cognitive processes at large. The second of these burdens will be taken up in Part IV. For now, I am going to argue that the functionally specified class *input system* does pick out a "natural kind" for purposes of psychological theory construction; that there are, in fact, lots of interesting things to say about the common properties of the mechanisms that mediate input analysis.

There is, however, one more preliminary point to make before getting down to that business. To claim that the functional category *input system* picks out a natural kind is to endorse an eccentric taxonomy of cognitive processes. Eyebrows should commence to be raised starting here. For, if you ask "which *are* the psychological mechanisms that can plausibly be thought of as functioning to provide information about the distal environment in a format appropriate for central processing?" the answer would seem to be "the perceptual systems *plus language*." And this is, from the point of view of traditional ways of carving things up, an odd category.

The traditional taxonomy goes something like this: perception (vision, audition, or whatever) on the one side, and thought-and-language (the representational processes) on the other. Now, the representational character of language is self-evident, and I don't doubt the theoretical importance of the representational character of thought. (On the contrary, I think that it is *the* essential fact that an adequate theory of the propositional attitudes would have to account for. (See Fodor, 1981a, chapter 7.) ) But we're not, of course, committed to there being only one right way of assigning psychological mechanisms to functional classes. The present claim is that, for purposes of assessing the issues about modularity, a rather different taxonomy proves illuminating.

Well then, what precisely *is* the functional similarity between language mechanisms and perceptual mechanisms in virtue of which both count as 'input systems'? There is, of course, the obvious point that utterances (e.g., sentence tokens) are themselves objects to be perceptually identified, just as mountains, teacups, and four-alarm fires are. Understanding a token sentence presumably involves assigning it a structural description, this being part and parcel of computing a token-to-type relation; and that is precisely the sort of function we would expect an input system to perform. However, in stressing the functional analogy between language and percep-

tion, I have something more in mind than the fact that understanding utterances is itself a typical perceptual process.

I've said that input systems function to interpret transduced information and to make it available to central processes; and that, in the normal case, what they provide will be information about the "layout" (to borrow a term of Gibson's) of distal stimuli. How might such a system work? Heaven knows there are few harder questions; but I assume that, in the case of perception, the answer must include some such story as the following. The character of transducer outputs is determined, in some lawful way, by the character of impinging energy at the transducer surface; and the character of the energy at the transducer surface is itself lawfully determined by the character of the distal layout. Because there are regularities of this latter sort, it is possible to infer properties of the distal layout from corresponding properties of the transducer output. Input analyzers are devices which perform inferences of this sort.

A useful example is Ullman's (1979) algorithm for inferring "form from motion" in visual perception. Under assumptions (e.g., of rigidity) that distal stimuli usually satisfy, a specific sequence of transformations of the energy distributions at the retina will be reliably interpretable as having been caused by (and hence as specifying) the spatial displacement of a distal object of determinate three-dimensional shape. A device that has access to the transducer outputs can infer this shape by executing Ullman's (or some equivalent) algorithm. I assume that performing such computations is precisely the function of input systems, Ullman's case being unusual primarily in the univocality with which the premises of the perceptual inference warrant its conclusion.

Now about language: Just as patterns of visual energy arriving at the retina are correlated, in a complicated but regular way, with certain properties of distal layouts, so too are the patterns of auditory energy that excite the tympanic membrane in speech exchanges. With, of course, this vital difference: What underwrites the correlation between visual stimulations and distal layouts are (roughly) the laws of light reflectance. Whereas, what underwrites the correlation between token utterances and distal layouts is (roughly) a convention of truth-telling. In the root case, the convention is that we say of  $x$  that it is  $F$  only if  $x$  is  $F$ . Because that convention holds, it is possible to infer from what one hears said to the way that the world is.<sup>12</sup>

Of course, in neither the linguistic nor the perceptual case is the information so provided infallible. The world often isn't the way it looks to be or the way that people say it is. But, equally of course, input systems don't have to deliver apodictic truths in order to deliver quite useful information. And, anyhow, *the operation of the input systems should not be identified with the fixation of belief*. What we believe depends on the evaluation of how things look, or are said to be, *in light of background information* about (inter alia) how good the seeing is or how trustworthy the source. Fixation of belief is just the sort of thing I have in mind as a typical central process.

So much, then, for the similarity of function between the linguistic and the perceptual systems: both serve to get information about the world into a format appropriate for access by such central processes as mediate the fixation of belief. But now, is there anything to be said for exploiting this analogy? What, from the point of view of psychological theory, do we gain by postulating a functional class of perceptual-and-linguistic processes? Clearly, the proof of this pudding is *entirely* in the eating. I'm about to argue that, if we undertake to build a psychology that acknowledges this functional class as a neutral kind, we discover that the processes we have grouped together do indeed have many interesting properties in common—properties the possession of which is not entailed by their functional homogeneity. (I take it that that is what a natural kind is: a class of phenomena that have many scientifically interesting properties in common over and above whatever properties define the class.) In the present case, what the input systems have in common besides their functional similarities can be summarized in a phrase: *input systems are modules*. A fortiori, they share those properties that are characteristic of vertical faculties. Input systems are—or so I'll argue—what Gall was right about.

What follows is the elaboration of that claim, together with an occasional glimpse at the state of the evidence. I should say at the outset that not every psychologist would agree with me about what the state of the evidence is. I am arguing well in advance of (and, in some places, a little in the face of) the currently received views. So, perhaps one should take this exercise as in part a thought experiment: I'll be trying to say what you might expect the data to look like if the modularity story is true of input systems; and I'll claim that, insofar as any facts are known, they seem to be generally compatible with such expectations.

### PART III INPUT SYSTEMS AS MODULES

The modularity of the input systems consists in their possession of most or all of the properties now to be enumerated. If there are other psychological systems which possess most or all of these properties then, of course, they are modular too. It is, however, a main thesis of this work that the properties in virtue of which input systems are modular are ones which, in general, central cognitive processes do not share.

#### *III.1. Input systems are domain specific*

Let's start with this: how many input systems are there? The discussion thus far might be construed so as to suggest an answer somewhere in the vicinity of six—viz., one for each of the traditional sensory/perceptual 'modes' (hearing, sight, touch, taste, smell) and one more for language. This is *not*, however, the intended doctrine; what is proposed is something much more in the spirit of Gall's bumps. I imagine that within (and, quite possibly, across)<sup>13</sup> the traditional modes, there are highly specialized computational mechanisms in the business of generating hypotheses about the distal sources of proximal stimulations. The specialization of these mechanisms consists in constraints either on the range of information they can access in the course of projecting such hypotheses, or in the range of distal properties they can project such hypotheses about, or, most usually, on both.

Candidates might include, in the case of vision, mechanisms for color perception, for the analysis of shape, and for the analysis of three-dimensional spatial relations.<sup>14</sup> They might also include quite narrowly task-specific 'higher level' systems concerned with the visual guidance of bodily motions or with the recognition of faces of conspecifics. Candidates in audition might include computational systems that assign grammatical descriptions to token utterances; or ones that detect the melodic or rhythmic structure of acoustic arrays; or, for that matter, ones that mediate the recognition of the voices of conspecifics. There is, in fact, some evidence for the domain specificity of several of the systems just enumerated, but I suggest

the examples primarily by way of indicating the levels of grain at which input systems might be modularized.

What, then, are the arguments for the domain specificity of input systems? To begin with, there is a sense in which input systems are *ipso facto* domain specific in a way in which computational systems at large are not. This is, however, quite uninteresting, a merely semantic point. Suppose, for example, that the function of the mechanisms of visual perception is to map transduced patterns of retinal excitation onto formulas of some central computational code. Then it follows trivially that their computational domain *qua mechanisms of visual perception* is specific to the class of possible retinal outputs. Correspondingly, if what the language-processing mechanisms do is pair utterance tokens with central formulas, then their computational domains *qua mechanisms of language processing* must be whatever encodings of utterances the auditory transducers produce. In similar boring fashion, the psychological mechanisms that mediate the perception of cows are *ipso facto* domain specific *qua mechanisms of cow perception*.

From such truisms, it goes without saying, nothing useful follows. In particular, the modularity of a system cannot be inferred from this trivial kind of domain specificity. It is, for example, entirely compatible with the cow specificity of cow perception that the recognition of cows should be mediated by precisely the same mechanisms that effect the perception of language, or of earthquakes, or of three-masted brigantines. For example, all four could perfectly well be accomplished by one and the same set of horizontal faculties. The interesting notion of domain specificity, by contrast, is Gall's idea that there are distinct psychological mechanisms—*vertical* faculties—corresponding to distinct stimulus domains. It is this latter claim that's now at issue.

Evidence for the domain specificity of an input analyzer can be of a variety of different sorts. Just occasionally the argument is quite direct and the demonstrations correspondingly dramatic. For example, there are results owing to investigators at the Haskins Laboratories which strongly suggest the domain specificity of the perceptual systems that effect the phonetic analysis of speech. The claim is that these mechanisms are different from those which effect the perceptual analysis of auditory nonspeech, and the experiments show that how a signal sounds to the hearer does depend, in rather

startling ways, on whether the acoustic context indicates that the stimulus is an utterance. Roughly, the very same signal that is heard as the onset of a consonant when the context specifies that the stimulus is speech is heard as a "whistle" or "glide" when it is isolated from the speech stream. The rather strong implication is that the computational systems that come into play in the perceptual analysis of speech are distinctive in that they operate *only* upon acoustic signals that are taken to be utterances. (See Liberman et al., 1967; for further discussion, see Fodor, Bever, and Garrett, 1974).

The Haskins experiments demonstrate the domain specificity of an input analyzer by showing that only a relatively restricted class of stimulations can throw the switch that turns it on. There are, however, other kinds of empirical arguments that can lead to the same sort of conclusions. One that has done quite a lot of work for cognitive scientists goes like this: If you have an *eccentric* stimulus domain—one in which perceptual analysis requires a body of information whose character and content is specific to that domain—then it is plausible that psychological processes defined over that domain may be carried out by relatively special purpose computational systems. All things being equal, the plausibility of this speculation is about proportional to the eccentricity of the domain.

Comparing perceiving cows with perceiving sentences will help to show what's going on here. I really have no idea how cow perception works, but let's follow the fashions and suppose, for purposes of discussion, that we use some sort of prototype-plus-similarity-metric. That is, the perceptual recognition of cows is effected by some mechanism which provides solutions for computational problems of the form: how similar—how 'close'—is the distal stimulus to a prototypical cow? My point is that if that's the way it's done, then cow perception might be mediated by much the same mechanisms that operate in a large variety of other perceptual domains as well—in fact, in any domain that is organized around prototypes. This is because we can imagine a quite general computational system which, given a specification of a prototype and a similarity metric for an arbitrary domain of percepts, will then compute the relevant distance relations in that domain. It seems plausible, that is to say, that procedures for estimating the distance between an input and a perceptual prototype should have

pretty much the same computational structure wherever they are encountered.

It is, however, most unlikely that the perceptual recognition of sentences should be mediated by such procedures, and that is because sentence tokens constitute a set of highly eccentric stimuli: All the available evidence suggests that the computations which sentence recognizers perform must be closely tuned to a complex of stimulus properties that is quite specific to sentences. Roughly, the idea is that the structure of the sentence recognition system is responsive to universal properties of language and hence that the system works only in domains which exhibit these properties.

I take it that this story is by now pretty well known. The argument goes like this: Consider the class of *nomologically possible human languages*. There is evidence that this class constitutes quite a small subset of the logically possible linguistic systems. In particular, the nomologically possible human languages include only the ones that satisfy a set of (contingent) generalizations known as the 'linguistic universals.' One way to find out something about what linguistic universals there are is by examining and comparing *actual* human languages (French, English, Urdu, or whatever) with an eye to determining which properties they have in common. Much work in linguistics over the last twenty-five years or so has pursued this strategy, and a variety of candidate linguistic universals have been proposed, both in phonology and in syntax.

It seems quite unlikely that the existence of these universals is merely fortuitous, or that they can be explained by appeal to historical affinities among the languages that share them or by appeal to whatever pragmatic factors may operate to shape communication systems. (By pragmatic factors, I mean ones that involve general properties of communication exchanges as such, including the utilities of the partners to the exchanges. So, for example, Putnam (1961) once suggested that there are grammatical transformations because communicative efficiency is served by the deletion of redundant portions of messages, etc.) The obvious alternative to such accounts is to assume that the universals represent biases of a species-specific language-learning system, and a number of proposals have been made about how, in detail, such systems might be pretuned. It is assumed, according to all these accounts, that the language-learning mechanisms 'know about' the universals and

operate only in domains in which the universals are satisfied. (For a review, see Pinker, 1979.)

Parity of argument suggests that a similar story should hold for the mechanisms of language *perception*. In particular, the perceptual system involved is presumed to have access to information about how the universals are realized in the language it applies to. The upshot of this line of thought is that the perceptual system for a language comes to be viewed as containing quite an elaborate theory of the objects in its domain; perhaps a theory couched in the form of a grammar of the language. Correspondingly, the process of perceptual recognition is viewed as the application of that theory to the analysis of current inputs. (For some recent work on the parsing of natural language, see Marcus, 1977; Kaplan and Bresnan, in press; and Frazier and Fodor, 1978. All these otherwise quite different approaches share the methodological framework just outlined.)

To come to the moral: Since the satisfaction of the universals is supposed to be a property that distinguishes sentences from other stimulus domains, the more elaborate and complex the theory of universals comes to be the more eccentric the stimulus domain for sentence recognition. And, as we remarked above, the more eccentric a stimulus domain, the more plausible the speculation that it is computed by a special-purpose mechanism. It is, in particular, very hard to see how a device which classifies stimuli in respect of distance from a prototype could be recruited for purposes of sentence recognition. The computational question in sentence recognition seems to be not "How far to the nearest prototype?" but rather "How does the theory of the language apply to the analysis of the stimulus now at hand?"

There are probably quite a lot of kinds of relatively eccentric stimulus domains—ones whose perceptual analysis requires information that is highly specific to the domain in question. The organization of sentence perception around syntactic and phonological information does not exhaust the examples even in the case of language. So, for a further example, it is often and plausibly proposed that the processes that mediate phone recognition must have access to an internal model of the physical structure of the vocal apparatus. The argument is that a variety of constancies in speech perception seem to have precisely the effect of undoing

garble that its inertial properties produce when the vocal mechanism responds to the phonetic intentions of the speaker. If this hypothesis is correct, then phone recognition is quite closely tuned to the mechanisms of speech production (see note 13). Once again, highly tuned computations are suggestive of special-purpose processors. Analogous points could be made in other perceptual modes. Faces are favorite candidates for eccentric stimuli (see Yin, 1969, 1970; Carey, 1978); and as I mentioned above, Ullman's work has made it seem plausible that the visual recognition of three-dimensional form is accomplished by systems that are tuned to the eccentricities of special classes of rigid spatial transformations.

From our point of view, the crucial question in all such examples is: how good is the inference from the eccentricity of the stimulus domain to the specificity of the corresponding psychological mechanisms? I am, in fact, not boundlessly enthusiastic about such inferences; they are clearly a long way from apodictic. Chess playing, for example, exploits a vast amount of eccentric information, but nobody wants to postulate a chess faculty. (Well, *almost* nobody. It is of some interest that recent progress in the artificial intelligence of chess has been achieved largely by employing specialized hardware. And, for what it's worth, chess is notably one of those cognitive capacities which breeds prodigies; so it is a candidate for modularity by Gall's criteria if not by mine.) Suffice it, for the present to suggest that it is probably characteristic of many modular systems that they operate in eccentric domains, since a likely motive for modularizing a system is that the computations it performs are idiosyncratic. But the converse inference—from the eccentricity of the domain to the modularity of the system—is warranted by nothing stronger than the maxim: specialized systems for specialized tasks. The most transparent situation is thus the one where you have a mechanism that computes an eccentric domain and is also modular by independent criteria; the eccentricity of the domain rationalizes the modularity of the processor and the modularity of the processor goes some way towards explaining how the efficient computation of eccentric domains is possible.

### *III.2 The operation of input systems is mandatory*

You can't help hearing an utterance of a sentence (in a language

you know) as an utterance of a sentence, and you can't help seeing a visual array as consisting of objects distributed in three-dimensional space. Similarly, *mutatis mutandis*, for the other perceptual modes: you can't, for instance, help feeling what you run your fingers over as the surface of an object.<sup>15</sup> Marslen-Wilson and Tyler (1981), discussing word recognition, remark that "... even when subjects are asked to focus their attention on the acoustic-phonetic properties of the input, they do not seem to be able to avoid identifying the words involved. . . . This implies that the kind of processing operations observable in spoken-word recognition are mediated by automatic processes which are obligatorily applied . . . (p. 327).

The fact that input systems are apparently constrained to apply whenever they can apply is, when one thinks of it, rather remarkable. There is every reason to believe that, in the general case, the computational relations that input systems mediate—roughly, the relations between transducer outputs and percepts—are quite remote. For example, on all current theories, it requires elaborate processing to get you from the representation of a proximal stimulus that the retina provides to a representation of the distal stimuli as an array of objects in space.<sup>16</sup> Yet we apparently have no choice but to take up this computational burden whenever it is offered. In short, the operation of the input systems appears to be, in this respect, inflexibly insensitive to the character of one's utilities. You can't hear speech as noise *even if you would prefer to*.

What you can do, of course, is choose not to hear it at all—viz., not attend.<sup>17</sup> In the interesting cases—where this is achieved without deactivating a transducer (e.g., by sticking your fingers in your ears)—the strategy that works best is rather tortuous: one avoids attending to *x* by deciding to concentrate on *y*, thereby taking advantage of the difficulty of concentrating on more than one thing at a time. It may be that, when this strategy is successful, the unattended input system does indeed get selectively 'switched off', in which case there is a somewhat pickwickian sense in which voluntary control over the operation of an input system is circuitously achieved. Or it may be that the unattended input systems continue to operate but lose their access to some central processes (e.g., to those that mediate storage and report). The latter account is favored, at least for the case of language perception, in light of

a fair number of results which seem to show relatively high-level processing of the unattended channel in dichotic listening tasks (Lackner and Garrett, 1973; Corteen and Wood, 1972; Lewis, 1970). But since the experimental results in this area are not univocal, perhaps the most conservative claim is this: input analysis is mandatory in that it provides the *only* route by which transducer outputs can gain access to central processes; if transduced information is to affect thought at all, it must do so via the computations that input systems perform.

I suppose one has to enter a minor caveat. Painters, or so I'm told, learn a little to undo the perceptual constancies and thus to see the world in something like the terms that the retina must deliver—as a two-dimensional spread of color discontinuities varying over time. And it is alleged that phoneticians can be taught to hear their language as something like a sound-stream—viz., as something like what the spikes in the auditory nerves presumably encode. (Though, as a matter of fact, the empirical evidence that phoneticians are actually able to do this is equivocal; see, for example, Lieberman, 1965.) But I doubt that we should take these highly skilled phenomenological reductions very seriously as counterexamples to the generalization that input processes are mandatory. For one thing, precisely because they *are* highly skilled, they may tell us very little about the character of normal perceptual processing. Moreover, it is tendentious—and quite possibly wrong—to think of what painters and phoneticians learn to do as getting access to, as it were, raw transducer output. An at least equally plausible story is that what they learn is how to 'correct' perceptually interpreted representations in ways that compensate for constancy effects. On this latter view, "seeing the visual field" or "hearing the speech stream" are *supersophisticated* perceptual achievements. I don't know which of these stories is the right one, but the issue is clearly empirical and oughtn't to be prejudged.

Anyhow, barring the specialized achievements of painters and phoneticians, one simply cannot see the world under its retinal projection and one has practically no access to the acoustics of utterances in languages that one speaks. (You all know what Swedish and Chinese sound like; what does *English* sound like?) In this respect (and in other respects too, or so I'll presently argue) the input mechanisms approximate the condition often ascribed to re-

flexes: they are automatically triggered by the stimuli that they apply to. And this is true for both the language comprehension mechanisms and the perceptual systems traditionally so-called.

It is perhaps unnecessary to remark that it does *not* seem to be true for nonperceptual cognitive processes. We have only the narrowest of options about how the objects of perception shall be represented, but we have all the leeway in the world as to how we shall represent the objects of *thought*; outside perception, the way that one deploys one's cognitive resources, is, in general, rationally subservient to one's utilities. Here are some exercises that you can do if you choose: think of *Hamlet* as a revenge play; as a typical product of Mannerist sensibility; as a pot-boiler; as an unlikely vehicle for Greta Garbo. Think of sixteen different ways of using a brick. Think of an utterance of "All Gaul is divided into three parts" as an acoustic object. Now try *hearing* an utterance of "All Gaul is divided into three parts" as an acoustic object. Notice the difference.

No doubt there are *some* limits to the freedom that one enjoys in rationally manipulating the representational capacities of thought. If, indeed, the Freudians are right, more of the direction of thought is mandatory—not to say obsessional—than the uninitiated might suppose. But the quantitative difference surely seems to be there. There is, as the computer people would put it, "executive control" over central representational capacities; and intellectual sophistication consists, in some part, in being able to exert that control in a manner conducive to the satisfaction of one's goals—in ways, in short, that seem likely to get you somewhere. By contrast, perceptual processes apparently apply willy-nilly in disregard of one's immediate concerns. "I couldn't help hearing what you said" is one of those clichés which, often enough, expresses a literal truth; and it is what is *said* that one can't help hearing, not just what is *uttered*.

### III.3. *There is only limited central access to the mental representations that input systems compute*

It is worth distinguishing the claim that input operations are mandatory (you can't but hear an utterance of a sentence *as* an utterance of a sentence) from the claim that what might be called 'interlevels' of input representation are, typically, relatively inaccessible to con-

sciousness. Not only must you hear an utterance of a sentence as such, but, to a first approximation, you can hear it *only* that way.

What makes this consideration interesting is that, according to all standard theories, the computations that input systems perform typically proceed via the assignment of a number of intermediate analyses of the proximal stimulation. Sentence comprehension, for example, involves not only acoustic encoding but also the recovery of phonetic and lexical content and syntactic form. Apparently an analogous picture applies in the case of vision, where the recognition of a distal array as, say, a-bottle-on-a-table-in-the-corner-of-the-room proceeds via the recovery of a series of preliminary representations (in terms of visual frequencies and primal sketches *inter alia*). For a review of recent thinking about interlevels of visual representation, see Zucker, 1981).

The present point is that the subject doesn't have equal access to all of these ascending levels of representation—not at least if we take the criterion of accessibility to be the availability for explicit report of the information that these representations encode. Indeed, as I remarked above, the lowest levels (the ones that correspond most closely to transducer outputs) appear to be completely *inaccessible* for all intents and purposes. The rule seems to be that, even if perceptual processing goes from 'bottom to top' (each level of representation of a stimulus computed being more abstractly related to transducer outputs than the one that immediately preceded), still *access* goes from top down (the further you get from transducer outputs, the more accessible the representations recovered are to central cognitive systems that presumably mediate conscious report).

A plausible first approximation might be that only such representations as constitute the *final* consequences of input processing are fully and freely available to the cognitive processes that eventuate in the voluntary determination of overt behavior. This arrangement of accessibility relations is reasonable enough assuming, on the one hand, that the computational capacities of central cognitive systems are not inexhaustible in their ability to attend to impinging information and, on the other, that it is the relatively abstract products of input-processing that encode most of the news that we are likely to want to know. I said in section III.2 that the operation of input systems is relatively insensitive to the subject's

utilities. By contrast, according to this account, the architectural arrangements that govern exchanges of information between input systems and other mechanisms of cognition do reflect aspects of the organism's standing concerns.

The generalization about the relative inaccessibility of intermediate levels of input analysis is pretty rough, but all sorts of anecdotal and experimental considerations suggest that something of the sort is going on. A well known psychological party trick goes like this:

E: Please look at your watch and tell me the time.

S: (Does so.)

E: Now tell me, without looking again, what is the shape of the numerals on your watch face?

S: (Stumped, evinces bafflement and awe.) (See Morton, 1967)

The point is that visual information which specifies the shape of the numerals must be registered when one reads one's watch, but from the point of view of access to later report, that information doesn't take. One recalls, as it were, pure position with no shape in the position occupied. There are analogous anecdotes to the effect that it is often hard to remember whether somebody you have just been talking to has a beard (or a moustache, or wears glasses). Yet visual information that specifies a beard must be registered and processed whenever you recognize a bearded face. More anecdote: Almost nobody can tell you how the letters and numbers are grouped on a telephone dial, though you use this information whenever you make a phone call. And Nickerson and Adams (1979) have shown that not only are subjects unable to describe a Lincoln penny accurately, they also can't pick out an accurate drawing from ones that get it grossly wrong.

There are quite similar phenomena in the case of language, where it is easy to show that details of syntax (or of the choice of vocabulary) are lost within moments of hearing an utterance, only the gist being retained. (Which did I just say was rapidly lost? Was it the syntactic details or the details of syntax?) Yet it is inconceivable that such information is not registered somewhere in the comprehension process and, within limits, it is possible to enhance its recovery by the manipulation of instructional variables. (For edifying experiments, see Sachs, 1967; Wanner, 1968.)

These sorts of examples make it seem plausible that the relative inaccessibility of lower levels of input analysis is at least in part a matter of how priorities are allocated in the transfer of representations from relatively short- to relatively long-term memory.<sup>18</sup> The idea would be that only quite high-level representations are stored, earlier ones being discarded as soon as subsystems of the input analyzer get the goodness out of them. Or, more precisely, intermediate input representations, when not discarded, are retained only at special cost in memory or attention, the existence of such charges-for-internal-access being itself a prototypical feature of modular systems.

This is, no doubt, part of the story. Witness the fact that in tasks which minimize memory demands by requiring comparison of *simultaneously* presented stimuli, responses that are sensitive to stimulus properties specified at relatively low levels of representation are frequently *faster* than responses to properties of the sort that high-level representations mark. Here, then, the ordering of relative accessibility reverses the top-to-bottom picture proposed above. It may be worth a digression to review some relevant findings.

The classical experimental paradigm is owing to Posner (1978). S's are required to respond 'yes' to visually presented letter pairs when they are *either* font identical (*t,t; T,T*) or alphabetically identical (*t,T; T,t*). The finding is that when letters in a pair are presented *simultaneously*, response to alphabetically identical pairs that are also font identical is faster than response to pairs that are identical alphabetically but not in font. This effect diminishes asymptotically with increase in the interstimulus interval when the letters are presented *sequentially*.

A plausible (though not mandatory) interpretation is that the representation that specifies the physical shape of the impinging stimulus is computed earlier than representations that specify its alphabetic value. (At a minimum, *some* shape information must be registered prior to alphabetic value, since alphabetic value depends upon shape.) In any event, the fact that representations of shape can drive voluntary responses suggests that they must be available to central processes at *some* point in the course of S's interaction with the stimulus. And this suggests, in turn, that the inaccessibility of font- as compared with alphabetic-information over the relatively long term must be a matter of how memory is deployed rather

than of the intrinsic opacity of low-level representations to high-level processes. It looks as though, in these cases, the relative unavailability of lower levels of input analysis is primarily a matter of the way that the subsystems of the input processors interface with memory systems. It is less a matter of information being unconscious than of its being unrecalled. (See also Crowder and Morton, 1969.)

It is unlikely, however, that this is the whole story about the inaccessibility of interlevels of input analysis. For one thing, as was remarked above, some very low levels of stimulus representation appear to be absolutely inaccessible to report. It is, to all intents and purposes (i.e., short of extensive training of the subject) impossible to elicit voluntary responses that are selectively sensitive to subphonetic linguistic distinctions (or, in the case of vision, to parameters of the retinal projection of distal objects) even though we have excellent theoretical grounds for supposing that such information must be registered somewhere in the course of linguistic (/visual) processing. And not just theoretical grounds: *we can often show that aspects of the subject's behavior are sensitive to the information that he can't report*.

For example, a famous result on the psychophysics of speech argues that utterances of syllables may be indistinguishable despite very substantial differences in their acoustic structure *so long as these differences are subphonetic*. When, however, quantitatively identical acoustic differences happen to be, as linguists say, 'contrastive'—i.e., when they mark distinctions between phones—they will be quite discriminable to the subject; as distinguishable, say, as "ba" is from "pa". It appears, in short, that there is a perceptual constancy at work which determines, in a wide range of cases, that only such acoustic differences as have linguistic value are accessible to the hearer in discrimination tasks. (See Liberman, et al., 1967.) What is equally striking, however, is that these 'inaccessible' differences *do affect reaction times*. Suppose a/a and a/b are utterance pairs such that the members of the first pair are literally acoustically identical and the members of the second differ only in *noncontrastive* acoustic properties—i.e., the acoustic distinction between a and b is subphonetic. As we have seen, it is possible to choose such properties so that the members of the a/b pair are perceptually indistinguishable (as are, of course, the members of the pair a/a).

Even so, in such cases reaction times to make the 'same' judgment for the a/a pair are reliably faster than reaction times to make the 'same' judgment for the a/b pair. (Pisoni and Tash, 1974.) The subject can't report—and presumably can't hear—the difference between signal a and signal b, but his behavior is sensitive to it all the same.

These kinds of cases are legion in studies of the constancies, and this fact bears discussion. The *typical* function of the constancies is to engender perceptual similarity in the face of the variability of proximal stimulation. Proximal variation is very often misleading; the world is, in general, considerably more stable than are its projections onto the surfaces of transducers. Constancies correct for this, so that in general percepts correspond to distal layouts *better than* proximal stimuli do. But, of course, the work of the constancies would be undone unless the central systems which run behavior were required largely to ignore the representations which encode *uncorrected* proximal information. The obvious architectural solution is to allow central systems to access information engendered by proximal stimulation only *after* it has been run through the input analyzers. Which is to say that central processes should have free access only to the *outputs* of perceptual processors, interlevels of perceptual processing being correspondingly opaque to higher cognitive systems. This, I'm claiming, is the architecture that we in fact do find.

There appears, in short, to be a generalization to state about input systems as such. Input analysis typically involves *mediated* mappings from transducer outputs onto percepts—mappings that are effected via the computation of interlevels of representation of the impinging stimulus. These intermediate representations are sometimes absolutely inaccessible to central processes, or, in many cases, they are accessible at a price: you can get at them, but only by imposing special demands upon memory or attention. Or, to put it another way: To a first approximation, input systems can be freely queried by memory and other central systems only in respect of *one* of the levels of representation that they compute; and the level that defines this interface is, in general, the one that is most abstractly related to transduced representations. This claim, if true, is substantive; and if, as I believe, it holds for input systems at large, then that is another reason to believe that the construct *input system* subsumes a natural kind.

### III.4. *Input systems are fast*

Identifying sentences and visual arrays are among the fastest of our psychological processes. It is a little hard to quantify this claim because of unclarities about the individuation of mental activities. (What precisely are the boundaries of the processes to be compared? For example, where does sentence (/scene) *recognition* stop and more central activities take over? Compare the discussion in section III.6, below.) Still, granting the imprecision, there are more than enough facts around to shape one's theoretical intuitions.

Among the simplest of voluntary responses are two-choice reactions (push the button if the *left*-hand light goes on). The demands that this task imposes upon the cognitive capacities are minimal, and a practiced subject can respond reliably at latencies on the low side of a quarter of a second. It thus bears thinking about that the recovery of semantic content from a spoken sentence can occur at speeds quite comparable to those achieved in the two-choice reaction paradigm. In particular, appreciable numbers of subjects can 'shadow' continuous speech with a quarter-second latency (shadowing is repeating what you hear as you hear it) and, contrary to some of the original reports, there is now good evidence that such 'fast shadowers' understand what they repeat. (See Marslen-Wilson, 1973.) Considering the amount of processing that must go on in sentence comprehension (unless all our current theories are *totally* wrongheaded), this finding is mind-boggling. And, mind-boggling or otherwise, it is clear that shadowing latency is an extremely conservative measure of the speed of comprehension. Since shadowing requires *repeating* what one is hearing, the 250 msec. of lag between stimulus and response includes not only the time required for the perceptual analysis of the message, but also the time required for the subject's integration of his verbalization.

In fact, it may be that the phenomenon of fast shadowing shows that the efficiency of language processing comes very close to achieving theoretical limits. Since the syllabic rate of normal speech is about 4 per second, the observed 250 msec. latency is compatible with the suggestion that fast shadowers are processing speech in syllable-length units—i.e., that the initiation of the shadower's response is commenced upon the identification of each syllable-length input. Now, work in the psychoacoustics of speech makes it look

quite likely that the syllable is the shortest linguistic unit that can be reliably identified in the speech stream (see Liberman et al., 1967). Apparently, the acoustic realizations of shorter linguistic forms (like phones) exhibit such extreme context dependence as to make them unidentifiable on a unit-by-unit basis. Only at the level of the syllable do we begin to find stretches of wave form whose acoustic properties are at all reliably related to their linguistic values. If this is so, then it suggests the following profoundly depressing possibility: the responses of fast shadowers lag a syllable behind the stimulus *not* because a quarter second is the upper bound on the speed of the mental processes that mediate language comprehension, but rather because, if the subject were to go any faster, he would overrun the ability of the speech stream to signal linguistic distinctions.<sup>19</sup>

In the attempt to estimate the speed of computation of visual processing, problems of quantification are considerably more severe. On the one hand, the stimulus is not usually spread out in time, so it's hard to determine how much of the input the subject registers before initiating his identificatory response. And, on the other hand, we don't have a taxonomy of visual stimuli comparable to the classification of utterance tokens into linguistic types. Since the question what type a linguistic token belongs to is a great deal clearer than the corresponding question for visual arrays, it is even less obvious in vision than in speech what sort of response should count as indicating that a given array has been identified.

For all of which there is good reason to believe that given a motivated decision about how to quantify the observations, the facts about visual perception would prove quite as appalling as those about language. For example, in one study by Haber (1980), subjects were exposed to 2,560 photographic slides of randomly chosen natural scenes, each slide being exposed for an interval of 10 seconds. Performance on recognition recall (ability to correctly identify a test slide as one that had been seen previously) approached 90 percent one hour after the original exposure. Haber remarks that the results "suggest that recognition of pictures is essentially perfect." Recent work by Potter (personal communication) indicates that 10 seconds of exposure is actually a great deal more than subjects need to effect a perceptual encoding of the stimulus adequate to mediate this near-perfect performance. According to Pot-

ter, S's performance in the Haber paradigm asymptotes at an exposure interval of about 2 seconds per slide.

There are some other results of Potter's (1975) that make the point still more graphically. S is shown a sequence of slides of magazine photographs, the rate of presentation of the slides being the experimentally manipulated variable. Prior to each sequence, S is provided with a brief description of an object or event that may appear in one or another slide—e.g., a boat, two men drinking beer, etc. S is to attend to the slides, responding when he sees one that satisfies the description. Under these conditions, S's respond with better than 70 percent accuracy when each slide is exposed for 125 msec. Accuracy asymptotes (at around 96 percent) at exposure times of 167 msec. per slide. It is of some interest that S's are as good at this task as they are at recognition recall (i.e., at making the global judgment that a given slide is one that they have seen before).

Two first-blush morals should be drawn from such findings about the computational efficiency of input processes. First, it contrasts with the relative slowness of paradigmatic central processes like problem-solving; and, second, it is presumably no accident that these very fast psychological processes are mandatory.

The first point is, I suppose, intuitively obvious: one can, and often does, spend hours thinking about a problem in philosophy or chess, though there is no reason to suppose that the computational complexity of these problems is greater than that of the ones that are routinely solved effortlessly in the course of perceptual processing. Indeed, the puzzle about input analysis is precisely that the computational complexity of the problem to be solved doesn't seem to predict the difficulty of solving it; or, rather, if it does, the difference between a 'hard' problem and an 'easy' one is measured not in months but in milliseconds. This dissimilarity between perception and thought is surely so adequately robust that it is unlikely to be an artifact of the way that we individuate cognitive achievements. It is only in trick cases, of the sorts that psychologists devise in experimental laboratories, that the perceptual analysis of an utterance or a visual scene is other than effectively instantaneous. What goes on when you parse a standard psycholinguistic poser like "the horse raced past the barn fell" is, almost certainly, *not* the same sort of processing that mediates sentence recognition in the normal case. They even *feel* different.

Second, it may well be that processes of input analysis are fast because they are mandatory. Because these processes are automatic, you save computation (hence time) that would otherwise have to be devoted to deciding whether, and how, they ought to be performed. Compare: eyeblink is a fast response because it is a reflex—i.e., because you don't have to decide whether to blink your eye when someone jabs a finger at it. Automatic responses are, in a certain sense, deeply unintelligent; of the whole range of computational (and, eventually, behavioral) options available to the organism, only a stereotyped subset is brought into play. But what you save by indulging in this sort of stupidity is *not having to make up your mind*, and making your mind up takes time. Reflexes, whatever their limitations, are not in jeopardy of being sickled o'er with the pale cast of thought. Nor are input processes, according to the present analysis.

There is, however, more than this to be said about the speed of input processes. We'll return to the matter shortly.

### *III.5. Input systems are informationally encapsulated*

Some of the claims that I'm now about to make are in dispute among psychologists, but I shall make them anyway because I think that they are true. I shall run the discussion in this section largely in terms of language, though, as usual, it is intended that the morals should hold for input systems at large.

I remarked above that, almost certainly, understanding an utterance involves establishing its analysis at several different levels of representation: phonetic, phonological, lexical, syntactic, and so forth. Now, in principle, information about the probable structure of the stimulus at any of these levels could be brought to bear upon the recovery of its analysis at any of the others. Indeed, in principle *any* information available to the hearer, including meteorological information, astrological information, or—rather more plausibly—information about the speaker's probable communicative intentions could be brought to bear at any point in the comprehension process. In particular, it is entirely possible that, in the course of computing a structural description, information that is specified only at relatively high levels of representation should be 'fed back' to determine analyses at relatively lower levels.<sup>20</sup> But

though this is possible in principle, the burden of my argument is going to be that the operations of input systems are in certain respects unaffected by such feedback.

I want to emphasize the 'in certain respects'. For there exist, in the psychological literature, dramatic illustrations of the effects of information feedback upon some input operations. Consider, for example, the 'phoneme restoration effect' (Warren, 1970). You make a tape recording of a word (as it might be, the word "legislature") and you splice out one of the speech sounds (as it might be, the 's'), which you then replace with a tape recording of a cough. The acoustic structure of the resultant signal is thus /legi(cough)lature/. But what a subject will *hear* when you play the tape to him is an utterance of /legislature/ with a cough 'in the background'. It surely seems that what is going on here is that the perceived phonetic constituency of the utterance is determined not just by the transduced information (not just by information specified at subphonetic levels of analysis) but also by higher-level information about the probable lexical representation of the utterance (i.e., by the subject's guess that the intended utterance was probably /legislature/).

It is not difficult to imagine how this sort of feedback might be achieved. Perhaps, when the stimulus is noisy, the subject's mental lexicon is searched for a 'best match' to however much of the phonetic content of the utterance has been securely identified. In effect, the lexicon is queried by the instruction 'Find an entry some ten phones long, of which the initial phone sequence is /legi/ and the terminal sequence is /lature/.' The reply to this query constitutes the lexical analysis under which the input is heard.

Apparently rather similar phenomena occur in the case of visual scotoma (where neurological disorders produce a 'hole' in the subject's visual field). The evidence is that scotoma can mask quite a lot of the visual input without creating a phenomenal blind spot for the subject. What happens is presumably that information about higher-level redundancies is fed back to 'fill in' the missing sensory information. Some such process also presumably accounts for one's inability to 'see' one's retinal blind spot.

These sorts of considerations have led to some psychologists (and many theorists in AI) to propose relentlessly top-down models

of input analysis, in which the perceptual encoding of a stimulus is determined largely by the subject's (conscious or unconscious) beliefs and expectations, and hardly at all by the stimulus information that transducers provide. Extreme examples of such feedback-oriented approaches can be found in Schank's account of language comprehension, in Neisser's early theorizing about vision, and in 'analysis by synthesis' approaches to sentence parsing. Indeed, a sentimental attachment to what are known generically as 'New Look' accounts of perception (Bruner, 1973) is pervasive in the cognitive science community. It will, however, be a main moral of this discussion that the involvement of certain sorts of feedback in the operation of input systems would be incompatible with their modularity, at least as I propose to construe the modularity thesis. One or other of these doctrines will have to go.

In the long run, which one goes will be a question of how the data turn out. Indeed, a great deal of the empirical interest of the modularity thesis lies in the fact that the experimental predictions it makes tend to be diametrically opposed to the ones that New Look approaches license. But experiments to one side, there are some *prima facie* reasons for doubting that the computations that input systems perform could have anything like unlimited access to high-level expectations or beliefs. These considerations suggest that even if there are *some* perceptual mechanisms whose operations are extensively subject to feedback, there must be others that compute the structure of a percept largely, perhaps solely, in isolation from background information.

For one thing, there is the widely noted persistence of many perceptual illusions (e.g., the Ames room, the phi phenomenon, the Muller-Lyre illusion in vision; the phoneme restoration and click displacement effects in speech) even in defiance of the subject's explicit knowledge that the percept is illusory. The very same subject who can tell you that the Muller-Lyre arrows are identical in length, who indeed has seen them measured, still finds one looking longer than the other. In such cases it is hard to see an alternative to the view that at least *some* of the background information at the subject's disposal is inaccessible to at least some of his perceptual mechanisms.

An old psychological puzzle provides a further example of this kind. When you move your head, or your eyes, the flow of images

across the retina may be identical to what it would be were the head and eyes to remain stationary while the scene moves. So: why don't we experience apparent motion when we move our eyes? Most psychologists now accept one or other version of the "corollary discharge" answer to this problem. According to this story, the neural centers which initiate head and eye motions communicate with the input analyzer in charge of interpreting visual stimulations (See Buzzi, 1968). Because the latter system knows what the former is up to, it is able to discount alterations in the retinal flow that are due to the motions of the receptive organs.

Well, the point of interest for us is that this visual-motor system is informationally encapsulated. Witness the fact that, if you (gently) push your eyeball with your finger (as opposed to moving it in the usual way: by an exercise of the will), you *do* get apparent motion. Consider the moral: when you voluntarily move your eyeball with your finger, you certainly are possessed of the information that it's your eye (and not the visual scene) that is moving. This knowledge is absolutely explicit; if I ask you, you can *say* what's going on. But this explicit information, available to you for (e.g.) report, is *not* available to the analyzer in charge of the perceptual integration of your retinal stimulations. That system has access to corollary discharges from the motor center *and to no other information that you possess*. Modularity with a vengeance.

We've been surveying first blush considerations which suggest that at least some input analyzers are encapsulated with respect to at least some sorts of feedback. The next of these is a point of principle: feedback works only to the extent that the information which perception supplies is redundant; and it *is* possible to perceptually analyze arbitrarily unredudant stimulus arrays. This point is spectacularly obvious in the case of language. If I write "I keep a giraffe in my pocket," you are able to understand me despite the fact that, on even the most inflationary construal of the notion of context, there is nothing in the context of the inscription that would have enabled you to predict either its form or its content. In short, feedback is effective only to the extent that, *prior* to the analysis of the stimulus, the perceiver knows quite a lot about what the stimulus is going to be like. Whereas, the *point* of perception is, surely, that it lets us find out how the world is even when the world is some way that we *don't* expect it to be. The teleology of

perceptual capacities presupposes a considerably-less-than-omniscient-organism; they'd be no use to God. If you already know how things are, why look to *see* how things are?<sup>21</sup>

So: The perceptual analysis of *unanticipated* stimulus layouts (in language and elsewhere) is possible only to the extent that (a) the output of the transducer is insensitive to the beliefs/expectations of the organism; and (b) the input analyzers are adequate to compute a representation of the stimulus from the information that the transducers supply. This is to say that the perception of novelty depends on bottom-to-top perceptual mechanisms.

There is a variety of ways of putting this point, which is, I think, among the most important for understanding the character of the input systems. Pylyshyn (1980) speaks of the "cognitive impenetrability" of perception, meaning that the output of the perceptual systems is largely insensitive to what the perceiver presumes or desires. Pylyshyn's point is that a condition for the *reliability* of perception, at least for a fallible organism, is that it generally sees what's there, not what it wants or expects to be there. Organisms that don't do so become deceased.

Here is another terminology for framing these issues about the direction of information flow in perceptual analysis: Suppose that the organism is given the problem of determining the analysis of a stimulus at a certain level of representation—e.g., the problem of determining which sequence of words a given utterance encodes. Since, in the general case, transducer outputs underdetermine perceptual analyses,<sup>22</sup> we can think of the solution of such problems as involving processes of nondemonstrative inference. In particular, we can think of each input system as a computational mechanism which projects and confirms a certain class of hypotheses on the basis of a certain body of data. In the present example, the available hypotheses are the word sequences that can be constructed from entries in the subject's mental lexicon, and the perceptual problem is to determine which of these sequences provides the right analysis of the currently impinging utterance token. The mechanism which solves the problem is, in effect, the realization of a *confirmation function*: it's a mapping which associates with each pair of a lexical hypothesis and some acoustic datum a value which expresses the degree of confirmation that the latter bestows upon the former. (And similarly, *mutatis mutandis*, for the nondemonstrative infer-

ences that the other input analyzers effect.) I emphasize that construing the situation this way involves no commitment to a detailed theory of the operation of perceptual systems. Any nondemonstrative inference can be viewed as the projection and confirmation of a hypothesis, and I take it that perceptual inferences must in general be nondemonstrative, since their underdetermination by sensory data is not in serious dispute.

Looked at this way, the claim that input systems are informationally encapsulated is equivalent to the claim that the data that can bear on the confirmation of perceptual hypotheses includes, in the general case, considerably less than the organism may know. That is, the confirmation function for input systems does not have access to all of the information that the organism internally represents; there are restrictions upon the allocation of internally represented information to input processes.

Talking about the direction of information flow in psychological processes and talking about restrictions upon the allocation of information to such processes are thus two ways of talking about the same thing. If, for example, we say that the flow of information in language comprehension runs directly from the determination of the phonetic structure of an utterance to the determination of its lexical content, then we are saying that only phonetic information is available to whatever mechanism decides the level of confirmation of perceptual hypotheses about lexical structure. On that account, such mechanisms are encapsulated with respect to *nonphonetic* information; they have no access to such information; not even if it is *internally represented, accessible to other cognitive processes* (i.e., to cognitive processes other than the assignment of lexical analyses to phone sequences) and *germane* in the sense that if it were brought to bear in lexical analysis, it would affect the confirmation levels of perceptual hypotheses about lexical structure.

I put the issue of informational encapsulation in terms of constraints on the data available for hypothesis confirmation because doing so will help us later, when we come to compare input systems with central cognitive processes. Suffice it to say, for the moment, that this formulation suggests another possible reason why input systems are so fast. We remarked above that the computations that input systems perform are mandatory, and that their being so saves time that would otherwise have to be used in executive decision-

making. We now add that input systems are bull-headed and that this, too, makes for speed. The point is this: to the extent that input systems are informationally encapsulated, of all the information that might *in principle* bear upon a problem of perceptual analysis only a portion (perhaps only quite a small and stereotyped portion) is actually admitted for consideration. This is to say that speed is purchased for input systems by permitting them to ignore lots of the facts. Ignoring the facts is not, of course, a good recipe for problem-solving in the general case. But then, as we have seen, input systems don't function in the *general* case. Rather, they function to provide very special kinds of representations of very specialized inputs (to pair transduced representations with formulas in the domains of central processes). What operates in the general case, and what is sensitive, at least in principle, to *everything* that the organism knows, are the central processes themselves. Of which more later.

I should add that these reflections upon the value of bull-headedness do not, as one might suppose, entirely depend upon assumptions about the speed of memory search. Consider an example. Ogden Nash once offered the following splendidly sane advice: "If you're called by a panther/don't anther." Roughly, we want the perceptual identification of panthers to be very fast and to err, if at all, only on the side of false positives. If there is a body of information that must be deployed in such perceptual identifications, then we would prefer not to have to recover that information from a large memory, assuming that the speed of access varies inversely with the amount of information that the memory contains. This is a way of saying that we do not, on that assumption, want to have to access panther-identification information from the (presumably *very large*) central storage in which representations of background-information-at-large are generally supposed to live. Which is in turn to say that we don't want the input analyzer that mediates panther identification to communicate with the central store on the assumption that large memories are searched slowly.

Suppose, however, that random access to a memory is *insensitive* to its size. Even so panther-identification (and, mutatis mutandis, other processes of input analysis) had better be insensitive to much of what one knows. Suppose that we can get at *everything* we know about panthers *very fast*. We still have the problem of deciding,

for each such piece of information retrieved from memory, how much inductive confirmation it bestows upon the hypothesis that the presently observed black-splotch-in-the-visual-field is a panther. The point is that in the rush and scramble of panther identification, there are many things I know about panthers whose bearing on the likely pantherhood of the present stimulus *I do not wish to have to consider*. As, for example, that my grandmother abhors panthers; that every panther bears some distant relation to my Siamese cat Jerrold J.; that there are no panthers on Mars; that there is an Ogden Nash poem about panthers . . . etc. Nor is this all; for, in fact, the property of being 'about panthers' is not one that can be surefootedly relied upon. Given enough context, practically everything I know can be construed as panther related; and, *I do not want to have to consider everything I know* in the course of perceptual panther identification. In short, the point of the informational encapsulation of input processes is not—or not solely—to reduce the memory space that must be searched to find information that is perceptually relevant. The primary point is to so restrict the number of *confirmation relations* that need to be estimated as to make perceptual identifications fast. (I am indebted to Scott Fahlman for raising questions that provoked the last two paragraphs.)<sup>23</sup>

The informational encapsulation of the input systems is, or so I shall argue, the essence of their modularity. It's also the essence of the analogy between the input systems and reflexes; reflexes are informationally encapsulated with bells on.

Suppose that you and I have known each other for many a long year (we were boys together, say) and you have come fully to appreciate the excellence of my character. In particular, you have come to know perfectly well that under no conceivable circumstances would I stick my finger in your eye. Suppose that this belief of yours is both explicit and deeply felt. You would, in fact, go to the wall for it. Still, if I jab my finger near enough to your eyes, and fast enough, you'll blink. To say, as we did above, that the blink reflex is mandatory is to say, *inter alia*, that it has no access to what you know about my character or, for that matter, to any other of your beliefs, utilities and expectations. For this reason the blink reflex is often produced when sober reflection would show it to be uncalled for; like panther-spotting, it is prepared to trade false positives for speed.

That is what it is like for a psychological system to be informationally encapsulated. If you now imagine a system that is encapsulated in the way that reflexes are, but also computational in a way that reflexes are not, you will have some idea of what I'm proposing that input systems are like.

It is worth emphasizing that being modular in this sense is not quite the same thing as being autonomous in the sense that Gall had in mind. For Gall, if I read him right, the claim that the vertical faculties are autonomous was practically equivalent to the claim that there are no horizontal faculties for them to share. Musical aptitude, for example, is autonomous *in that* judging musical ideas shares no cognitive mechanisms with judging mathematical ideas; remembering music shares no cognitive mechanisms with remembering faces; perceiving music shares no cognitive mechanisms with perceiving speech; and so forth.

Now, it is unclear to what extent the input systems *are* autonomous *in that* sense. We do know, for example, that there are systematic relations between the amount of computational strain that decoding a sentence places on the language handling systems and the subject's ability to perform simultaneous nonlinguistic tasks quickly and accurately. 'Phoneme monitor' (Foss, 1970) techniques, and others, can be used to measure such interactions, and the results suggest a picture that is now widely accepted among cognitive psychologists: Mental processes often compete for access to resources variously characterized as attention, short-term memory, or work space; and the result of allocating such resources to one of the competing processes is a decrement in the performance of the others. How general this sort of interaction is is unclear in the present state of the art (for contrary cases, suggesting isolated work spaces for visual imagery on the one hand and verbal recall on the other, see Brooks, 1968). In any event, where such competition does obtain, it is a counterexample to autonomy in what I am taking to be Gall's understanding of that notion.<sup>24</sup>

On the other hand, we can think of autonomy in a rather different way from Gall's—viz., in terms of informational encapsulation. So, instead of asking what access language processes (e.g.) have to computational resources that other systems also share, we can ask what access they have to the *information* that is available to other systems. If we do look at things this way, then the question "how

much autonomy?" is the same question as "how much constraint on information flow?" In a nutshell: one way that a system can be autonomous is by being encapsulated, by not having access to facts that other systems know about. I am claiming that, whether or not the input systems are autonomous in Gall's sense, they are, to an interesting degree, autonomous in this informational sense.

However, I have not yet given any arguments (except some impressionistic ones) to show that the input systems actually are informationally encapsulated. In fact, I propose to do something considerably more modest: I want to suggest some caveats that ought to be, but frequently aren't, observed in interpreting the sorts of data that have usually been alleged in support of the contrary view. I think that many of the considerations that have seemed to suggest that input processes are cognitively penetrable—that they are importantly affected by the subject's belief about context, or his background information, or his utilities—are, in fact, equivocal or downright misleading. I shall therefore propose several ground rules for evaluating claims about the cognitive penetrability of input systems; and I'll suggest that, when these rules are enforced, the evidence for 'New Look' approaches to perception begins to seem not impressive. My impulse in all this is precisely analogous to what Marr and Poggio say motivates their work on vision: "... to examine ways of squeezing the last ounce of information from an image before taking recourse to the descending influence of high-level interpretation on early processing" (1977, pp. 475–476).

(a) Nobody doubts that the information that input systems provide must somehow be reconciled with the subject's background knowledge. We sometimes know that the world can't really be the way that it looks, and such cases may legitimately be described as the correction of input analyses by top-down information flow. (This, ultimately, is the reason for refusing to identify input analysis with perception. The point of perception is the fixation of belief, and the fixation of belief is a *conservative* process—one that is sensitive, in a variety of ways, to what the perceiver already knows. Input analysis may be informationally encapsulated, but perception surely is not.) However, to demonstrate *that* sort of interaction between input analyses and background knowledge is not, in and of itself, tantamount to demonstrating the cognitive penetrability of the former; you need also to show that the locus of the top-down effect

is *internal* to the input system. That is, you need to show that the information fed back interacts with interlevels of input-processing and not merely with the final results of such processing. The penetrability of a system is, by definition, its susceptibility to top-down effects at stages *prior* to its production of output.

I stress this point because it seems quite possible that input systems specify only relatively shallow levels of representation (see the next section). For example, it is quite possible that the perceptual representation delivered for a token sentence specifies little more than the type to which the token belongs (and hence does *not* specify such information as the speech act potential of the token, still less the speech act performed by the tokening). If this is so, then data showing effects of the hearer's background information on, e.g., his estimates of the speaker's communicative intentions would *not* constitute evidence for the cognitive penetration of the presumptive language-comprehension module; by hypothesis, the computations involved in making such estimates would not be among those that the language-comprehension module *per se* performs. Similarly, *mutatis mutandis*, in the case of vision. There is a great deal of evidence for context effects upon certain aspects of visual object recognition. But such evidence counts for nothing in the present discussion unless there is independent reason to believe that these aspects of object recognition are part of visual input analysis. Perhaps the input system for vision specifies the stimulus only in terms of "primal sketches" (for whose cognitive impenetrability there is, by the way, some nontrivial evidence. See Marr and Nishihara (1978).) The problem of assessing the degree of informational encapsulation of input systems is thus not independent of the problem of determining how such systems are individuated and what sorts of representations constitute their outputs. I shall return to the latter issue presently; for the moment, I'm just issuing caveats.

(b) Evidence for the cognitive penetrability of some computational mechanism that does what input systems do is not, in and of itself, evidence for the cognitive penetrability of input systems.

To see what is at issue here, consider some of the kinds of findings that have been taken as decisively exhibiting the effects of background expectations upon language perception. A well known way of estimating such expectations is the use of the so-called Cloze

procedure. Roughly, S is presented with the first  $n$  words of a sentence and is asked to complete the fragment. Favored completions (as, for example, "salt" in the case of the fragment "I have the pepper, but would you please pass the \_\_\_\_") are said to be "high Cloze" and are assumed to indicate what the subject would expect a speaker to say next if he had just uttered a token of the fragment. An obvious generalization allows the estimation of the Cloze value at each point in a sentence, thereby permitting experiments in which the average Cloze value of the stimulus sentences is a manipulated variable.

It is quite easy to show that relative Cloze value affects S's performance on a number of experimental tasks, and it is reasonable to infer from such demonstrations that whatever mechanisms mediate the performance of these tasks must have access to S's expectations about what speakers are likely to say, hence not just to the 'stimulus' (e.g., acoustic) properties of the linguistic token under analysis. (For an early review of the literature on redundancy effects in sentence processing, see Miller and Isard, 1963.) So, for example, it can be shown that the accuracy of S's perception of sentences heard under masking noise is intimately related to the average Cloze value of the sentences: high Cloze sentences can be understood under conditions of greater distortion than the perception of low Cloze sentences tolerates. (Similarly, high Cloze sentences are, in general, more easily remembered than low Cloze sentences; recognition thresholds for words that are high Cloze in a context are lower than those for words that are low Cloze in that context; and so forth.)

The trouble with such demonstrations, however, is that although they show that there exist *some* language-handling processes that have access to the hearer's expectations about what is likely to be said, they do *not* show that the input systems enjoy such access. For example, it might be argued that, in situations where the stimulus is acoustically degraded, the subject is, in effect, encouraged to guess the identity of the material that he can't hear. (Similarly, *mutatis mutandis*, in memory experiments where a reasonable strategy for the subject is to guess at such of the material as he can't recall.) Not surprisingly, in such circumstances, the subject's background information comes into play with measurable effect. The question, however, is whether the psychological mechanisms

deployed in the slow, relatively painful, highly attentional process of reconstructing noisy or otherwise degraded linguistic stimuli are the same mechanisms which mediate the automatic and fluent processes of normal speech perception.

That this question is not merely frivolous is manifested by results such as those of Fishler and Bloom (1980). Using a task in which sentences are presented in clear, they found only a marginal effect of high Cloze on the recognition of test words, and such effects vanished entirely when the stimuli were presented at high rates. (High presentation rates presumably discourage guessing; guessing takes time.) By contrast, words that are 'semantically anomalous' in context showed considerable inhibition in comparison with neutral controls. This last finding is of interest because it suggests that at least some of the effects of sentence context in speech recognition must be, as psychologists sometimes put it, 'post-perceptual'. In our terminology, these processes must operate *after* the input system has provided a (tentative) analysis of the lexical content of the stimulus. The point is that even if the facilitation of redundant items is mediated by predictive, expectation-driven mechanisms, the inhibition of contextually anomalous items cannot be. It is arguable that, in the course of speech perception, one is forever making such predictions as that 'pepper' will occur in 'salt and ----'; but surely one can't also be forever predicting that 'dog', 'tomorrow', and all the other anomalous expressions will *not* occur there.<sup>25</sup> The moral is: some processes which eventuate in perceptual identifications are, doubtless, cognitively penetrated. But this is compatible with the informational encapsulation of the input systems themselves. Some traditional enthusiasm for context-driven perceptual models may have been prompted by confusion on this point.

(c) The claim that input systems are informationally encapsulated must be very carefully distinguished from the claim that there is top-down information flow *within* these systems. These issues are very often run together, with consequent exaggeration of the well-groundedness of the case against encapsulation.

Consider, once again, the phoneme restoration effect. Setting aside the general caution that experiments with distorted stimuli provide dubious grounds for inferences about speech perception in clear, phoneme restoration provides considerable *prima facie*

evidence that phone identification has access to what the subject knows about the lexical inventory of his language. If this interpretation is correct, then phoneme restoration illustrates top-down information flow in speech perception. It does *not*, however, illustrate the cognitive penetrability of the language input system. To show that that system is penetrable (hence informationally unencapsulated), you would have to show that its processes have access to information that is not specified at any of the levels of representation that the language input system computes; for example, that it has generalized access to what the hearer knows about the probable beliefs and intentions of his interlocutors. If, by contrast, the 'background information' deployed in phoneme restoration is simply the hearer's knowledge of the words in his language, then that counts as top-down flow within the language module; on any remotely plausible account, the knowledge of a language includes knowledge of its lexicon.

The most recent work in phoneme restoration makes this point with considerable force. Samuel (1981) has shown that both information about the lexical inventory and 'semantic' information supplied by sentential context affect the magnitude of the phoneme restoration effect. Specifically, you get more restoration in words than in (phonologically possible) nonwords, and you get more restoration when a word is predictable in sentence context than when the context is neutral. This looks like the penetration of phone recognition by both lexical and 'background' information, but the appearance is misleading. In fact, Samuel's data suggest that, of the two effects, *only the former* is strictly perceptual, the latter operating in consequence of a response bias to report predictable words as intact. (Detection theoretically: the word/nonword difference affects  $\delta'$ , whereas the neutral context/predictive context difference affects  $\beta$ .) As Samuel points out, the amount of restoration is inversely proportional to S's ability to distinguish the stimulus word with a phone missing from an undistorted token of the same type; and, on Samuel's data, this discrimination is actually *better* for items that are highly predictable in context than for items that aren't. Another case, in short, where what had been taken to be an example of context-driven prediction in perception is, in fact, an effect of the biasing of post-perceptual decision processes.

The importance of distinguishing cognitive penetration from in-

transmodular effects can be seen in many other cases where predictive analysis in perception is demonstrable. It is, for example, probable (though harder to show than one might have supposed) that top-down processes are involved in the identification of the surface constituent structure of sentences (see Wright, 1982). For example, it appears that the identification of nouns is selectively facilitated in contexts like T A ———, the identification of verbs is selectively facilitated in contexts like T N ———, and so forth. Such facilitation indicates that the procedures for assigning lexical items to form classes have access to information about the general conditions upon the well-formedness of constituent structure trees.

Now, it is a question of considerable theoretical interest whether, and to what extent, predictive analysis plays a role in parsing; but this issue must be sharply distinguished from the question whether the parser is informationally encapsulated. Counterexamples to encapsulation must exhibit the sensitivity of the parser to information that is not specified internal to the language-recognition module, and constraints on syntactic well-formedness are paradigms of information that does *not* satisfy this condition. The issue is currently a topic of intensive experimental and theoretical inquiry; but as things stand I know of no convincing evidence that syntactic parsing is ever guided by the subject's appreciation of semantic context or of 'real world' background. Perhaps this is not surprising; there are, in general, so many syntactically different ways of saying the same thing that even if context allowed you to estimate the *content* of what is about to be said, that information wouldn't much increase your ability to predict its *form*.<sup>26</sup>

These questions about where the interacting information comes from (whether it comes from inside or outside the input system) take on a special salience in light of the following consideration: it is possible to imagine ways in which mechanisms *internal* to a module might contrive to, as it were, mimic effects of cognitive penetration. The operation of such mechanisms might thus invite overestimations of the extent to which the module has access to the organism's general informational resources. To see how this might occur, let's return to the question of contextual facilitation of word recognition; traditionally a parade case for New Look theorizing, but increasingly an area in which the data are coming to seem equivocal.

Here are the bare bones of an ingenious experiment of David Swinney's (1979; for further, quite similar, results, see Tannenhaus, Leirnau, and Seidenberg, 1979). The subject listens to a stimulus sentence along the lines of "Because he was afraid of electronic surveillance, the spy carefully searched the room for bugs." Now, we know from previous research that the response latencies for 'bugs' (say, in a word/nonword decision task) will be faster in this context, where it is relatively predictable, than in a neutral context where it is acceptable but relatively low Cloze. This seems to be—and is traditionally taken to be—the sort of result which demonstrates how expectations based upon an intelligent appreciation of sentential context can guide lexical access; the subject predicts 'bugs' before he hears the word. His responses are correspondingly accelerated whenever his prediction proves true. Hence, cognitive penetration of lexical access.

You can, or so it seems, gild this lily. Suppose that, instead of measuring reaction time for word/nonword decisions on 'bugs', you simultaneously present (flashed on a screen that the subject can see) a different word belonging to the same (as one used to say) 'semantic field' (e.g., 'microphones'). If the top-down story is right in supposing that the subject is using semantic/background information to predict lexical content, then 'microphones' is as good a prediction in context as 'bugs' is, so you might expect that 'microphones', too, will exhibit facilitation as compared with a neutral context. And so it proves to do. Cognitive penetration of lexical access with bells on, or so it would appear.

But the appearance is misleading. For Swinney's data show that if you test with 'insects' instead of 'microphones', you get the same result: facilitation as compared with a neutral context. Consider what this means. 'Bugs' has two paraphrases: 'microphones' and 'insects'. But though only one of these is contextually relevant, *both* are contextually facilitated. This looks a lot less like the intelligent use of contextual/background information to guide lexical access. What it looks like instead is some sort of associative relation among lexical forms (between, say, 'spy' and 'bug'); a relation pitched at a level of representation sufficiently superficial to be *insensitive* to the semantic content of the items involved. This possibility is important for the following reason: If facilitation is mediated by merely interlexical relations (and not by the interaction of background

information with the semantic content of the item and its context), then the information that is exploited to produce the facilitation can be represented *in the lexicon*; hence *internal to the language recognition module*. And if that is right, then contextual facilitation of lexical access is *not* an argument for the cognitive penetration of the module. It makes a difference, as I remarked above, where the penetrating information comes from.

Let's follow this just a little further. Suppose the mental lexicon is a sort of connected graph, with lexical items at the nodes and with paths from each item to several others. We can think of accessing an item in the lexicon as, in effect, exciting the corresponding node; and we can assume that one of the consequences of accessing a node is that excitation spreads along the pathways that lead from it. Assume, finally, that when excitation spreads through a portion of the lexical network, response thresholds for the excited nodes are correspondingly lowered. Accessing a given lexical item will thus decrease the response times for items to which it is connected. (This picture is familiar from the work of, among others, Morton, 1969, and Collins and Loftus, 1975; for relevant experimental evidence, see Meyer and Schvaneveldt, 1971.)

The point of the model-building is to suggest how mechanisms internal to the language processor could mimic the effects that cognitive penetration would produce if the latter indeed occurred. In the present example, what mimics the background knowledge that (roughly) spies have to do with bugs is the existence of a connection between the node assigned to the word 'spy' and the node assigned to the word 'bug'. Facilitation of 'bug' in spy contexts is affected by the excitation of such intralexical connections.

Why should these intralexical connections exist? Surely not just in order to lead psychologists to overestimate the cognitive penetrability of language-processing. In fact, if one works the other way 'round and assumes that the input systems are encapsulated, one might think of the mimicry of penetration as a way that the input processors contrive to make the best of their informational isolation. Presumably, what encapsulation buys is speed; and, as we remarked above, it buys speed at the price of unintelligence. It would, one supposes, take a lot of time to make reliable decisions about whether there is the kind of relation between spies and bugs that makes it on balance likely that the current token of 'spy' will

be followed by a token of 'bug'. But that is precisely the kind of decision that the subject would have to make if the contextual facilitation of lexical access were indeed an effect of background knowledge interacting with the semantic content of the context. The present suggestion is that no such intelligent evaluation of the options takes place; there is merely a brute facilitation of the recognition of 'bug' consequent upon the recognition of 'spy'. The condition of this brute facilitation buying anything is that it should be possible, with reasonable accuracy, to mimic what one knows about connectedness *in the world* by establishing corresponding connections among entries in the mental lexicon. In effect, the strategy is to use the structure of interlexical connections to mimic the structure of knowledge. The mimicry won't be precise (a route from 'spy' to 'insect' will be generated as a by-product of the route from 'spy' to 'bug'). But there's no reason to doubt that it may produce savings over all.

Since we are indulging speculations, we might as well indulge this one: It is a standing mystery in psychology why there should be interlexical associations at all; why subjects should exhibit a reliable and robust disposition to associate 'salt' with 'pepper', 'cat' with 'dog', 'mother' with 'father', and so forth. In the heyday of associationism, of course, such facts seemed quite *unmysterious*; they were, indeed, the stuff of which the mental life was supposed to be made. On one account the utterance of a sentence was taken to be a chained response, and associations among lexical items were what held the links together. According to still earlier tradition, the postulation of associative connections between Ideas was to be the mechanism for reconstructing the notion of degree of belief. None of this seems plausible now, however. Belief is a matter (not of association but) of *judgment*; sentence production is a matter (not of association but) of *planning*. So, what on earth are associations for?

The present suggestion is that associations are the means whereby stupid processing systems manage to behave as though they were smart ones. In particular, interlexical associations are the means whereby the language processor is enabled to act as though it knows that spies have to do with bugs (whereas, in fact, it knows no such thing). The idea is that, just as the tradition supposed, terms for things frequently connected in experience become them-

selves connected in the lexicon. Such connection is *not* knowledge; it is not even judgment. It is simply the mechanism of the contextual adjustment of response thresholds. Or, to put the matter somewhat metaphysically, the formation of interlexical connections buys the synchronic encapsulation of the language processor at the price of its cognitive penetrability *across time*. The information one has about how things are related in the world is inaccessible to modulate lexical access; that is what the encapsulation of the language processor implies. But one's experience of the relations of things in the world does affect the structure of the lexical network—viz., by instituting connections among lexical nodes. If the present line of speculation is correct, these connections have a real, if modest, role to play in the facilitation of the perceptual analysis of speech. The traditional, fundamental, and decisive objection to association is that it is too stupid a relation to form the basis of a mental life. But stupidity, when not indulged in to excess, is a virtue in fast, peripheral processes; which is exactly what I have been supposing input processes to be.

I am not quite claiming that all the putative effects of information about background (context, etc.) on sentence recognition are artifacts of connections in the lexical network (though, as a matter of fact, such experimental attempts as I've seen to demonstrate a residual effect of context after interlexical/associative factors are controlled for strike me as not persuasive). I am claiming only that the possibility of such artifacts contaminates quite a lot of the evidence that is standardly alleged. The undoubted fact that "semantically" coherent text is relatively easy to process does not, in and of itself, demonstrate that the input system for language has access to what the organism knows about how the world coheres. Such experimental evidence as supported early enthusiasms for massively top-down perceptual models was, I think, sexy but inconclusive; and the possibility of a modular treatment of input processes provides motivation for its reconsideration. The situation would seem to be paradigmatically Kuhnian: the data look different to a jaundiced eye.

Consider the provenance of New Look theorizing. Cognitive psychologists in the '40s and '50s were faced with the proposal that perception is *literally* reflexive; for example, that the theory of perception is reducible without residue to the theory of discrim-

inative operant response. It was natural and admirable in such circumstances to stress the 'intelligence' of perceptual integration. However, in retrospect it seems that the intelligence of perceptual integration may have been seriously misconstrued by those who were most its partisans.

In the ideal condition—one approached more frequently in the textbooks than *in rerum naturae*, to be sure—reflexes have two salient properties. They are computationally simple (the stimulus is "directly connected" to the response), and they are informationally encapsulated (see above). I'm suggesting that New Look theories failed to distinguish these properties. They thus assumed, wrongly, that the disanalogy between perceptual and reflexive processes consisted in the capacity of the former to access and exploit background information. From the point of view of the modularity thesis, this is a case of the right intuition leading to the wrong claim. Input systems *are* computationally elaborated. Their typical function is to perform inference-like operations on representations of impinging stimuli. Processes of input analysis are thus unlike reflexes in respect of the character and complexity of the operations that they perform. But this is quite compatible with reflexes and input processes being similar in respect of their informational encapsulation; in this latter respect, both of them contrast with "central processes"—problem-solving and the like—of which cognitive penetrability is perhaps the most salient feature, or so I shall argue below. To see that informational encapsulation and computational elaboration are compatible properties, it is only necessary to bear in mind that unencapsulation is the exploitation of information from *outside* a system; a computationally elaborated system can thus be encapsulated if it stores the information that its computations exploit. Encapsulation is a matter of foreign affairs; computational elaboration begins at home.

It may be useful to summarize this discussion of the informational encapsulation of input systems by comparing it with some recent, and very interesting, suggestions owing to the philosopher Steven Stich (1978). Stich's discussion explores the difference between belief and the epistemic relation that is alleged to hold between, for example, speaker/hearers and the grammar of their native language (the relation that Chomsky calls 'cognizing'). Stich supposes, for purposes of argument, that the empirical evidence shows that

speakers *in some sense* 'know' the grammar of their native language; his goal is to say something about what that sense is.

Let us call the epistemic relation that a native speaker has to the grammar of his language *subdoxastic belief*.<sup>27</sup> Stich suggests that there are two respects in which subdoxastic beliefs differ from beliefs strictly so-called. In the first place, as practically everybody has emphasized, subdoxastic beliefs are *unconscious*. But, Stich adds, subdoxastic beliefs are also typically "inferentially unintegrated." The easiest way to understand what Stich means by this is to consider one of his examples.

If a linguist believes a certain generalization to the effect that no transformation rule exhibits a certain characteristic, and if he comes to (nonsubdoxastically) believe a given transformation which violates the generalization, he may well infer that the generalization is false. But merely having the rule stored (in the way that we are assuming all speakers of the language do) does not enable the linguist to draw the inference. . . . Suppose that for some putative rule, you have come to believe that if *r* then Chomsky is seriously mistaken. Suppose further that, as it happens, *r* is in fact among the rules stored by your language processing mechanism. The belief along with the subdoxastic state will not lead to the belief that Chomsky is seriously mistaken. By contrast, if you believe (perhaps even mistakenly) that *r*, then the belief that Chomsky is seriously mistaken is likely to be inferred. [pp. 508-509]

Or, as Stich puts the argument at another point, "It is characteristic of beliefs that they generate further beliefs via inference. What is more, beliefs are inferentially promiscuous. Provided with a suitable set of supplementary beliefs, almost any belief can play a role in the inference to any other. . . . (However) subdoxastic states, as contrasted with beliefs, are largely inferentially isolated from the large body of inferentially integrated beliefs to which a subject has (conscious) access."

Now, as Stich clearly sees, the proposal that subdoxastic states are typically both unconscious and inferentially unintegrated raises a question—viz., *Why should these two properties co-occur?* Why should it be, to put it in my terminology, that subdoxastic states

are typically encapsulated with respect to the processes which affect the inferential integration of beliefs?

Notice that there is a kind of encapsulation that follows from unconsciousness: an unconscious belief cannot play a role as a premise in the sort of reasoning that goes on in the conscious drawing of inferences. Stich is, however, urging something more interesting than this trivial truth. Stich's claim is that subdoxastic beliefs are largely inaccessible even to *unconscious* mental processes of belief fixation. If this claim is true, the question does indeed arise why it should be so.

I want to suggest, however, that the question doesn't arise because, as a matter of act, subdoxastic beliefs are not in general encapsulated; or, to put it more precisely, they are not in general encapsulated *qua* subdoxastic. Consider, as counterexamples, one's subdoxastic views about inductive and deductive warrant; for example, one's subdoxastic acquiescence in the rule of *modus ponens*. On the sort of psychological theory that Stich has in mind, subdoxastic knowledge of such principles must be accessible to practically all mental processes, since practically all inferential processes exploit them in one way or another. One's subdoxastic beliefs about validity and confirmation are thus quite unlike one's subdoxastic beliefs about the rules of grammar; though both are unconscious, the former are paradigms of promiscuous and unencapsulated mental states. So the connection between unconsciousness and encapsulation cannot be *intrinsic*.

Nevertheless, I think that Stich is onto something important. For, though much unconscious information must be widely accessible to processes of fixation of belief, it is quite true that very many of the examples of unconscious beliefs for which there is currently good empirical evidence are encapsulated. This is because most of our current cognitive science is the science of input systems, and, as we have seen, *informational encapsulation is arguably a pervasive feature of such systems*. Input systems typically do not exchange subdoxastic information with central processes or with one another.

Stich almost sees this point. He says that "subdoxastic states occur in a variety of separate, special purpose cognitive systems" (p. 508). True enough; but they must also occur in integrated, general purpose systems (in what I'm calling "central" systems), assuming that much of the fixation of belief is both unconscious

and subserved by inferential mechanisms of that kind. The point is: subdoxastic states are informationally encapsulated *only* insofar as they are states of special purpose systems (e.g., states of input analyzers). Practically all psychologically interesting cognitive states are unconscious; but it is only the beliefs accessible to modules that are subdoxastic by the second of Stich's criteria as well.

### *III.6. Input analyzers have 'shallow' outputs.*

The question where to draw the line between observation and inference (in the psychological version, between perception and cognition) is one of the most vexed, and most pregnant, in the philosophy of science. One finds every opinion from the extreme 'foundationalist' view, which restricts observation to processes that issue in infallible introspective reports, to the recent revisionism which denies that the distinction is in any respect principled. (Hanson, 1958, for example, holds that a physicist can *see* that the cloud chamber contains a proton track in the same sense of 'see' that is operative when Smith sees that there's a spot on Jones' tie.) Sometimes the argument for this sort of view is based explicitly on accounts of perception borrowed from New Look psychology, which suggests that *all* perception is ineliminably and boundlessly theory laden; see Goodman (1978).

Philosophers have cared about the observation/inference distinction largely for epistemological reasons; what is (nondemonstratively) inferred is supposed to run an inductive risk from which what is observed is supposed to be free. And it has seemed important to some epistemologists that whatever count as the data statements of a science should be isolated from such risk, the idea being that unless some contingent truths are certain, no empirical theory can compel rational belief.

I am not myself much moved by the idea that inductive warrant is inherited upward in science from a base level of indubitable truths; and barring some such assumption, the philosophical problem of making the observation/theory distinction rigorous seems less consequent than was once supposed. However, the corresponding psychological problem of saying where perceptual processes interface with cognitive ones must be addressed by anyone who takes the postulation of modular input systems seriously. For

one thing, it is a point of definition that distinct functional components cannot interface *everywhere* on pain of their ceasing to be distinct. It is this consideration that flow-chart notation captures by drawing boxes around the processing systems it postulates. That only the inputs and outputs of functionally individuated systems can mediate their information exchanges is tautological.

Moreover, we have seen that the plausibility of claims for the informational encapsulation of an input system depends very much on how one draws the distinction between its *outputs* and its *internal levels* of representation. Since it is common ground that there must be *some* mental processes in which perception interacts with background knowledge and with utilities, the issue about informational encapsulation is whether such interactions take place *internal* to the input systems. But the question what is internal to a system, and the question what is to count as the output of the system, are patently two ways of asking the same thing.

In general, the more constrained the information that the outputs of perceptual systems are assumed to encode—the shallower their outputs, the more plausible it is that the computations that effect the encoding are encapsulated. If, for example, the visual analysis system can report only upon the shapes and colors of things (all higher-level integrations being post-perceptual) it is correspondingly plausible that all the information that system exploits may be represented internal to it. By contrast, if the visual system can deliver news about protons (as a psychologized version of the Hanson story would suggest), then the likelihood that visual analysis is informationally encapsulated is negligible. Chat about protons surely implies free access to quite a lot of what I have been calling 'background knowledge'.

In this section I want to make a few, highly speculative suggestions about how the outputs of the language and visual processors might be characterized—that is, about the level of representation at which these systems interface with central processes. I shall rely heavily on the assumptions that input computations are very fast, and that their outputs are typically phenomenologically salient (see above). Consonant with these assumptions, I shall argue that there are some reasonable proposals to make about how to distinguish visual and linguistic perception from the cognitive processes with which they interface. It turns out, however, that there is nothing *episte-*

mologically special about the levels of representation which constitute the outputs of the visual (/linguistic) processing mechanisms. So if, in the spirit of epistemology naturalized, one leaves it to psychologists to draw the observation/theory distinction, then, according to these proposals, there is nothing epistemologically interesting about that distinction. For example, it does *not* correspond to the distinction between what we infallibly know and what we merely justifiably surmise. This seems to me, if anything, to argue in favor of drawing the line where I propose to draw it; still this version of naturalized epistemology may strike some epistemologists as far too deflationary.

What representation of an utterance does the language input processor compute? Or, to put the question in the context of the preceding discussion, which phenomenologically accessible properties of an utterance are such that, on the one hand, their recovery is mandatory, fast, and relevant to the perceptual encoding of the utterance and, on the other, such that their recovery might be achieved by an informationally encapsulated computational mechanism? Clearly, there is a wide choice of properties of utterances that *could* be computed by computational systems whose access to background information is, in one way or another, interestingly constrained—the duration of the utterance, e.g. For all that, there is, in the case of language, a glaringly obvious galaxy of candidates for modular treatment—viz., those properties that utterances have in virtue of some or other aspects of their linguistic structure (where this means, mostly, grammatical and/or logical form). Making these notions clear is notoriously hard; but the relevant intuitions are easy enough to grasp.

Whether John's utterance of "Mary might do it, but Joan is above that sort of thing" is ironical, say, is a question that can't be answered short of using a lot of what you know about John, Mary, and Joan. Worse yet, there doesn't seem to be any way to say, in the general case, how much, or precisely what, of what you know about them might need to be accessed in making such determinations. *Maybe* an interestingly encapsulated system could reliably recognize the irony (sincerity, metaphoricalness, rhetoricalness, etc.) of utterances, but there are certainly no plausible proposals about how this might be so. It looks as though recognizing such properties of utterances is typically an exercise in "inference to the best explanation": given

what I know about John, and about what John thinks about Mary and Joan, he *couldn't* have meant that literally . . . etc. These are, of course, precisely the sorts of inferences that you would *not* expect encapsulated systems to perform. The "best" explanation is the one you want to accept *all things considered*, and encapsulated systems are prohibited by definition from considering all things.

Compare the computational problems involved in the recognition of linguistic form. The idea here is that the grammatical and logical structure of an utterance is uniquely determined (or, more precisely, uniquely determined up to ambiguity) by its phonetic constituency; and its phonetic constituency is uniquely determined in turn by certain of its acoustic properties (mutatis mutandis, the linguistic properties of written tokens are uniquely determined by certain properties of their shapes). "Acoustic" properties, according to this usage, are ipso facto transducer-detectable; so an input system that has access to the appropriate transduced representations of an utterance knows everything about the utterance that it needs to know to determine which sentential type it is a token of and, probably, what the logical form of the utterance is.<sup>28</sup> In short, if you are looking for an interesting property of utterances that might be computed by rigidly encapsulated systems—indeed, a property that might even be computed by largely bottom-to-top processors—then the type-identity of the utterance, together, perhaps, with its logical form would seem to be a natural candidate.

It is thus worth stressing that type-identity and at least some aspects of logical form are phenomenologically salient and are patently recognized 'on line'; moreover, the computation of type-identity is clearly an essential part of the overall process of language comprehension. In the general case, you can't understand what the speaker has said unless you can at least figure out which sentence he has uttered.

Is there, then, an encapsulated analyzer for logical and grammatical form? All the arguments are indirect; but, for what it's worth, it's rather hard to see how some of the processes that recognize logical and grammatical form could be anything but encapsulated. Background information can be brought to bear in perceptual analysis only where the property that is recognized is, to some significant extent, redundant in the context of recognition. But, as we remarked above, there doesn't seem to be much re-

dundancy between context variables and the *form* of an utterance, however much context may predict its *content*. Even if you know precisely what someone is going to say—in the sense of knowing precisely which proposition he is going to assert—the knowledge buys you very little in predicting the type/token relation for his utterance; there are simply too many linguistically different ways of saying the same thing.

It is not, therefore, surprising that the more extreme proposals for context-driven language recognizers do *not* generally proceed by using contextual information to identify grammatical relations. Instead, they proceed whenever possible directly from a lexical analysis to a "conceptual" analysis—one which, in effect, collapses across synonymous tokens regardless of their linguistic type. It is unclear to me whether such models are proposed as serious candidates for the explanation of human communicative capacities, though sometimes I fear that they may be. (See, e.g., Schank and Abelson, 1975; for experimental evidence that linguistic form continues to have its effect as semantic integration increases, precisely as one would expect if the recovery of logical syntactic form is mandatory, see Forster and Olberi, 1973.) To put the point in a nutshell: linguistic form recognition can't be context-driven because context doesn't determine form; if linguistic form is recognized at all, it must be by largely encapsulated processes.

So the present proposal is that the language-input system specifies, for any utterance in its domain, its linguistic and maybe its logical form. It is implicit in this proposal that it does no more than that<sup>29</sup>—e.g., that it doesn't recover speech-act potential (except, perhaps, insofar as speech-act potential may be correlated with properties of form, as in English interrogative word order). As I suggested, the main argument for this proposal is that, on the one hand, type/token relations surely must be computed in the course of sentence comprehension and, on the other, it is hard to see how anything much richer than type/token relations could be computed by an informationally encapsulated processor. All this comports with the strong intuition that while there could perhaps be an algorithm for parsing, there surely could not be an algorithm for estimating communicative intentions in anything like their full diversity. Arguments about what an author *meant* are thus able to be interminable in ways in which arguments about what he *said* are not.

This is all pretty loose. Most discussions in linguistics and psycholinguistics have been primarily interested in establishing *minimal* conditions on the output of the sentence processor, e.g., by demonstrating that one or another level of linguistic representation is "psychologically real" and recovered on line. By contrast, the problem that arises in discussions of modularity is typically of the form: What is the *most* that an encapsulated processor should be supposed to compute? Which aspects of the input can plausibly be recognized without generalized appeal to background data? There is, however, one area of language research in which issues of this latter sort have been extensively discussed. It may be worth a brief recapitulation here, since it provides quite a clear illustration of what problems about determining the level of the perception/cognition interface are like.

Consider again the question of the vocabulary of an utterance (as opposed to its logicosyntactic form on the one hand and its propositional content on the other). Since I have assumed that input-processing yields type identifications, I am committed to the claim that the language processor delivers, for each input utterance, a representation which specifies its lexical constituents *inter alia*. (Utterances which differ in their lexical constituents are, of course, *ipso facto* distinct in type.) The present question is whether it is plausible to suppose that the language-input system provides still deeper representations at the lexical level.

A view that has been influential in both linguistics and psychology suggests that it does. According to this view, understanding an utterance involves recovering the *definitions* of such definable lexical items as it may contain. So, for example, understanding a token of "John is a bachelor" involves representing the utterance as containing a word that means *unmarried man*. Note that this is a claim about processes of *comprehension* and not, e.g., about inferential operations which may be applied to the internal representation of the utterance *after* it has been understood. It is thus natural to interpret the claim as implying that the recovery of definitions of lexical items takes place during input processing (viz., *interior* to the putative language module). We would thus expect, if the claim is true, that the recovery of definitional information should exhibit the typical properties of input processes: it should happen fast, it should be mandatory (insensitive to task demands), etc.

The alternative view is that the "surface" vocabulary of an utterance is preserved at the level of representation where the language processor interfaces with cognitive processes at large. There should thus be no level of analysis specified by the language-input system at which "... bachelor ..." and "... unmarried man ..." receive identical representations (though, of course, *postcomprehension* inferential processes may indeed identify them as synonymous. One could imagine that such *postcomprehension* inferences might be mediated by the application of "meaning postulates" in something like the sense of Carnap (1960); for discussion, see Kintsch (1974), Fodor, Fodor and Garrett (1975).)

The currently available experimental evidence supports the latter view. (See Fodor et al., 1980.) In fact, so far as I know, there have been *no* convincing data in favor of the claim that representations of definitional content engage *any* sentence-comprehension process. The importance of imposing appropriate task demands in experimental tests of this claim can, however, hardly be overemphasized. There is, e.g., no doubt at all that definitionally related sentences tend to be conflated in experiments that require not just comprehension but recall as well. This is quite consonant with the view that memory is an inferential process par excellence (see Bartlett, 1932).

If these observations are correct, they strongly suggest that input-processing for language provides no semantic analysis "inside" lexical items. Or, to put it another way, the functionally defined level *output of the language processing module* respects such *structurally defined* notions as *item in the morphemic inventory of the language*. It is of primary importance to see that there is no *a priori* reason why this should be true.<sup>30</sup> That is, there is no *a priori* reason why the representations of utterances that are computed by fast, mandatory, informationally encapsulated, etc., etc., processes should constitute a representational level by *any* independent criteria. But, in the case of language at least, there is some *a posteriori* reason to believe that they do: on the one hand, there is strong evidence that such notions as *morphemic level* and *syntactic level* pick out coherent classes of representations; and, on the other, there are at least reasonable grounds for supposing that it is representations at these sorts of levels that the input system delivers.

By the way, the (presumptive) fact that the representations which

input systems recover constitute linguistic natural kinds is a strong argument that the concept *input process* itself picks out a natural kind. Suppose that the representations of utterances that are recovered by fast, informationally encapsulated, mandatory, etc. processes turned out to specify, e.g., the second phoneme of the third word of each utterance, the intonation contour of its last five syllables, and the definitions of all the words that it contains which begin with 'u'. Since this collection of properties has no theoretical interest whatever, we would be inclined to infer that there is, to that extent, nothing interesting about the class of psycholinguistic processes that are fast, mandatory, and informationally encapsulated. But, apparently, that is *not* the sort of thing that we find. What we find instead is that the fast, mandatory ... etc. processes deliver representations of utterances which make perfectly good sense considered as representations of utterances; representations which specify, for example, morphemic constituency, syntactic structure, and logical form. This is just the sort of thing you would expect if the fast, mandatory ... etc. processes form a system that is functionally relevant to language comprehension. In particular, it is just what you would expect if language comprehension is effected by the sort of system that I am calling a module.

If I am inclined to harp on these points, it is because the opposed view—that sentence-processing grades off insensibly into inference and the appreciation of context; into general cognition in short—is actually predominant in the field. (Especially on the West Coast, where gurus teach that the All is One.) Suffice it to say that the choice between these pictures is empirical—not a matter of taste—and that such evidence as is actually germane seems not unfavorable to the modularity view.

The preceding discussion provides a context for raising analogous issues about vision. If the modularity story is to be plausible here, the output of the visual processor must be reasonably shallow (it should not categorize visual stimuli in such terms as *proton trace*), and it must form a level of representation by some independent criterion—i.e., there should be interesting things to say about the output representations other than that they are, *de facto*, the kinds of representations that the visual processor puts out.

Moreover, various candidates that satisfy the shallowness test and the levels test must nevertheless be rejected on grounds of

phenomenological inaccessibility.<sup>31</sup> I am thinking of such representations as Marr's 'primal', '2.5 D', and '3 D' sketch (Marr and Nishihara, 1978). Such representations are certainly shallow enough. Indeed, they would seem to be too shallow. If we accept them as defining visual processor outputs, we shall have to say that even object recognition is not, strictly speaking, a phenomenon of visual perception, since, at these levels of representation, only certain geometric properties of the stimulus are specified. But, surely, from the point of view of phenomenological accessibility, perception is above all the recognition of objects and events. Shallower systems of representation can therefore constitute only interlevels of input analysis. What, then, is its output?

One of the most interesting ideas in recent cognitive theorizing is that there is a level of 'basic' perceptual objects (or, to use a slightly less misleading terminology, of basic perceptual categories). This notion is explored extensively in Brown (1958) and in Rosch et al. (1976), but a quick presentation may make the point. Consider a category hierarchy like *poodle*, *dog*, *mammal*, *animal*, *physical object*, *thing*. Roughly, the following seems to be true of such sets of categories: they effect a taxonomy of objects at increasing levels of abstractness, such that a given entity may belong to any or all of them, and such that the potential extensions of the categories increase as you go up the hierarchy (there are, as it were, more possible dogs than possible poodles; more possible animals than possible dogs; and so forth). Moreover, this is an *implicational* hierarchy in the sense that it is somehow *necessary* that whatever satisfies a category at the *n*th level of abstraction must always satisfy every category at higher-than-*n* levels of abstraction. (I don't care, for present purposes [actually, I don't think I care at all] whether this necessity is analytic or even whether it is linguistic. Suffice it that it is no accident that every poodle is a dog.)

The idea of *basic* categories is that some of the levels of abstraction in such implicational hierarchies have peculiar psychological salience. Intuitively, salience clusters at the "middle" levels of abstraction (in the present case, *dog* rather than *poodle* or *thing*). There is, alas, no independent definition of "middle," and it is quite conceivable that intuitions about which levels are in the middle just *are* intuitions of relative salience. Still, the fact seems to be that the following cluster of psychological properties tend to con-

verge on the same member (or members) of each implicational hierarchy; that is, whatever member(s) of a hierarchy has one of them is also quite likely to have the rest. A category that has them all is paradigmatically basic.

(a) The basic category of a hierarchy often turns out to correspond to the high-frequency item in vocabulary counts; "dog" is thus a higher-frequency lexical item than either "animal" or "poodle."

(b) The word for the basic category of a hierarchy tends to be learned earlier than words that express other levels in the hierarchy (Anglin, 1979).

(c) The basic category is often the least abstract member of its hierarchy that is monomorphemically lexicalized. Compare "Sheraton wing-back armchair"; "armchair"; "chair"; "furniture"; "artifact"; "physical object . . ." In some domains there is evidence that the monomorphemic lexicalization of the basic category is universal—for example, there are few or no languages that have a single word for what we would call "a washed-out pinkish red" while coding what we would call plain "red" polymorphemically. (See Berlin and Kay, 1969.) As with (a) and (b), it seems natural to interpret (c) as a linguistic reflex of the relative psychological salience of the basic category as compared with other members of its hierarchy.

(d) Basic categories are natural candidates for ostensive introduction. "Dog" is ostensively definable for a child who hasn't learned "poodle," but it is probably not possible to teach "poodle" ostensively to a child who hasn't got "dog"; and it probably is not possible to teach "animal" ostensively to a child who hasn't got at least some animal words at the same level as "dog." This becomes glaringly obvious if one thinks about the relative ostensive definability of, e.g., "pale red," "red," and "color." Once again, it seems plausible to connect the relative ostensive definability of a word with the relative psychological salience of the property that the word expresses. (For a discussion of the implications of the correlation between basicness and ostensive definability, see Fodor, 1981a, chap. 10.)

(e) Basic categorizations yield 'information peaks' in the following sense. Ask a subject to list all the properties that come to mind when he thinks of *animals*; then ask him to list all the properties that come to mind when he thinks of *dogs*; and then ask him to

list all the properties that come to mind when he thinks of *poodles*. One finds that one gets quite a lot more properties for *dog* than for *animal*, whereas the properties listed for *poodles* include very few more than one got for *dog*.<sup>32</sup> (See Rosch, et al., 1976.) It seems that—in some sense that is admittedly not very clear—basic categorizations are the ones that encode the most information per unit judgment. Taken together with Paul Grice's "maxim of quantity" (be informative) and his "maxim of manner" (be succinct), this observation predicts the following bit of pragmatics:

(f) Basic categories are the natural ones to use for describing things, *ceteris paribus*. "*Ceteris paribus*" means something like 'assuming that there are no special task demands in play'. You say to me, 'What do you see out the window?' I reply, 'A lady walking a dog', (rather than, e.g., 'A lady walking an animal' on the one hand, or 'A lady walking a silver-grey, miniature, poodle bitch', on the other. The point to notice here is that, all things being equal, the first is the preferred level of description even where I may happen to know enough to provide the third.

I assume that these linguistic facts are surface reflections of a deeper psychological reality, to wit:

(g) Basic categorizations are phenomenologically *given*; they provide, as it were, the natural level for describing things to *oneself*. A glance out the window thus reveals: a lady walking a *dog*, rather than a lady walking a silver-grey, miniature... etc. (Of course, sustained inspection alters all this. But phenomenological salience is accessibility *without* sustained inspection.) You might predict from these intuitions that perceptual identifications which involve the application of basic categories ought to be fast as compared to applications of either more or less abstract members of their implication hierarchies. There is, in fact, experimental evidence that this is true. (See Intraub, 1981.)

(h) Basic categories are typically the most abstract members of their implication hierarchies which subtend individuals of approximately similar appearance (Rosch, et al., 1976). So, roughly, you can draw something that is just a dog, but you can't draw something that is just an animal; you can draw something that is just a chair, but you can't draw something that is just furniture.

This observation suggests that, to a first approximation, basic categorizations (unlike categorizations that are more abstract) can

be made, with reasonable reliability, on the basis of the visual properties of objects. It thus returns us to the issue of perception. Since input systems are, by assumption, informationally encapsulated (no generalized top-down access to background information), the categorizations such systems effect must be comprehensively determined by properties that the visual transducers can detect: shape, color, local motion, or whatever. Input systems aren't, of course, confined to encoding properties like shape and color, but they *are* confined—in virtue of their informational encapsulation—to categorizations which can be inferred, with reasonable accuracy, from such "purely visual" properties of the stimulus.<sup>33</sup> (Compare: the language processor is confined to recovering properties of the input token that can be inferred, with reasonable accuracy, from its acoustic properties—hence to recovering linguistic form rather than, say, the speaker's metaphorical intent.)

Putting it all together, then: basic categorizations are typically the most abstract members of their inferential hierarchies that *could* be assigned by an informationally encapsulated visual-input analyzer; more abstract categorizations are not reliably predicted by *visual* properties of the distal stimulus. And basic categorizations are the ones that you would want the input systems to deliver assuming that you are interested in maximizing the information per unit of perceptual integration (as, presumably, you are). So, the suggestion is that the visual-input system delivers basic categorizations.<sup>34</sup>

A lot follows from this suggestion: for example, that in one useful sense of the observation/theory distinction, dogs but not protons count as observed; that the outputs of the visual processor—like the outputs of the language processor—constitute a level of representation on grounds independent of the fact that they happen to be the set of representations that some input system delivers; that it is no accident that the phenomenologically accessible categorizations are expressed by ostensively definable words. And so forth. I leave it to the reader to draw the morals. Suffice it that the notion that visual analyses are computed by an informationally encapsulated system leads to the prediction that there should be some set of representations which are (roughly) shape-assignable on the one hand, and which, on the other hand, play a specially central role in the mental life of the organism. The pregnancy of the basic category construct suggests that this prediction is true.

### III.7. Input systems are associated with fixed neural architecture

Martin Gardner has a brief discussion of Gall in his *In the Name of Science* (1952). Gardner remarks that "Modern research on the brain has, as most everyone knows, completely demolished the old 'faculty psychology'. Only sensory centers are localized" (p. 293). The argument moves breathtakingly fast. Is faculty psychology literally incompatible with, say, an equipotential brain? Remember that faculties are, in the first instance, functionally rather than physiologically individuated. And perhaps *localization* isn't precisely the notion that Gardner wants, since, after all, there might be neural specificity of some functions that aren't localized in the sense of being associated with large, morphologically characterizable brain regions. Still, if you read "perceptual" for "sensory", and if you add language, and if you don't worry about the localization of motor and other noncognitive functions, there is something to what Gardner says. In particular, it seems that there is characteristic neural architecture associated with each of what I have been calling the input systems. Indeed, the following, stronger, claim seems to be approximately true: *all* the cases of massive neural structuring to which a content-specific cognitive function can confidently be assigned appear to be associated with input analysis, either with language or with perception. There is, to put it crudely, no known brain center for *modus ponens*.

I shall return presently to consider the implications of this observation. Suffice it, for the moment, that the intimate association of modular systems with neural hardwiring is pretty much what you would expect given the assumption that the key to modularity is informational encapsulation. Presumably, hardwired connections indicate privileged paths of informational access; the effect of hardwiring is thus to facilitate the flow of information from one neural structure to another. But, of course, what counts as relative *facilitation* when viewed one way counts as relative *encapsulation* when viewed the other way. If you facilitate the flow of information from A to B by hardwiring a connection between them, then you provide B with a kind of access to A that it doesn't have to locations C, D, E, . . . This sort of differential accessibility makes sense for a system only under the condition that it wants faster (easier, more contin-

uous, anyhow cheaper) access to A than it does to C, D, E, and the rest. That is, it makes sense only for a system whose informational demands are relatively skewed. There is, in particular, no point in hardwiring the connections of paradigmatic *unencapsulated* systems—ones whose informational demands may be imposed anywhere at any time. Neural architecture, I'm suggesting, is the natural concomitant of informational encapsulation.

Anyhow, we do find neurological structure associated with the perceptual systems and with language. Whatever the right interpretation of this finding may be, it provides yet another reason to believe that the input systems constitute a natural kind.

### III.8. Input systems exhibit characteristic and specific breakdown patterns

The existence of—and analogies between—relatively well defined pathological syndromes in the perceptual systems on the one hand and the language-processing mechanisms on the other has been too frequently noted to require much discussion here. There seems to be general agreement that the agnosias and aphasias constitute patterned failures of functioning—i.e., they cannot be explained by mere quantitative decrements in global, horizontal capacities like memory, attention, or problem-solving. This is hardly surprising if, on the one hand, input analysis is largely effected by specific, hardwired neural circuitry and, on the other, the pathologies of the input systems are caused by insult to these specialized circuits.

Contrast the central processes, which do not appear to be intimately associated with specific neural architecture and also do not appear to be prone to well defined breakdown syndromes. (It used to be thought that schizophrenia is a "pathology of thought," but I gather this view is no longer very popular.)

I don't, however, wish to overplay this point. Any psychological mechanism which is functionally distinct may presumably be selectively impaired, horizontal faculties included. There may thus quite possibly be pathologies of, say, memory or attention that are not domain specific in the way that the aphasias and agnosias are supposed to be; see, e.g., Milner, Corbin, and Teuber (1968). If so, then that is evidence (*contra* Gall) that such capacities are mediated by bona fide faculties and that they are horizontally organized. As

previously remarked, the possibility of advancing mixed models in this area ought not to be ignored.

### *III.9. The ontogeny of input systems exhibits a characteristic pace and sequencing*

The issues here are so very moot, and the available information is so fragmentary, that I offer this point more as a hypothesis than a datum. There are, however, straws in the wind. There is now a considerable body of findings about the ontogenetic sequencing of language acquisition, and there are some data on the very early visual capacities of infants. These results are compatible, so far, with the view that a great deal of the developmental course of the input systems is endogenously determined. On the one hand, the capacity of infants for visual categorization appears to have been very seriously underestimated by empiricist theorizing (see the recent work of Spelke, 1982; Meltzoff, 1979; Bower, 1974; and others). And, on the other hand, linguistic performance—though obviously not present in the neonate—appears to develop in an orderly way that is highly sensitive to the maturational state of the organism, and surprisingly insensitive to deprivation of environmental information. (Goldin-Meadow and Feldman, 1977; Gleitman, 1981.) Moreover, language development appears to respect many of the universals of adult grammatical organization even at quite early stages (see Brown, 1973, and, papers in Takavolian, 1981). There have been occasional attempts to account for such apparently domain-specific features of ontogeny by appeal to the developing structure of ‘problem-solving heuristics’ or of ‘general intelligence’; but they have been half-hearted and, in my view, quite unsuccessful when contemplated in detail. (For extensive discussion of these issues, see Piatelli-Palmarini, 1980, and the reviews by Marshall, 1981, and by Pylyshyn, 1981.) For what it’s worth, then, no facts now available contradict the claim that the neural mechanisms subserving input analysis develop according to specific, endogenously determined patterns under the impact of environmental releasers. This picture is, of course, quite compatible with the view that these mechanisms are instantiated in correspondingly specific, hardwired neural structures. It is also compatible with the suggestion that much of the information at the disposal of such systems is

innately specified; as, indeed, vertical faculty theorists from Gall to Chomsky have been wont to claim.

I have been arguing that the psychological systems whose operations “present the world to thought” constitute a natural kind by criteria independent of their similarity of function; there appears to be a cluster of properties that they have in common but which, *qua* input analyzers, they might perfectly well not have shared.<sup>35</sup> We can abbreviate all this by the claim that the input systems constitute a family of modules: domain-specific computational systems characterized by informational encapsulation, high-speed, restricted access, neural specificity, and the rest.

Let’s suppose, probably contrary to fact, that you have found this story convincing. So, you are pretending to believe, for purposes of the following discussion, that the input systems are modular. If you actually did believe this, you would surely be led to pose the following question: are cognitive mechanisms *other than* input systems also modular? Or are the properties of being modular and being an input system coextensive? We are thus, finally, about to raise what I regard as the main issue: whether modularity is (as Gall, for example, thought it was) the *general* fact about the organization of the mind. I am going to suggest that at least some cognitive systems are nonmodular, and then I’m going to explore a variety of consequences of their (putative) nonmodularity.

## PART IV CENTRAL SYSTEMS

Vertical faculties are domain specific (by definition) and modular (by hypothesis). So the questions we now want to ask can be put like this: Are there psychological processes that can plausibly be assumed to cut across cognitive domains? And, if there are, is there reason to suppose that such processes are subserved by nonmodular (e.g., informationally unencapsulated) mechanisms?

The answer to the first of these questions is, I suppose, reasonably clear. Even if input systems are domain specific, there must be some cognitive mechanisms that are not. The general form of the argument goes back at least to Aristotle: the representations that

input systems deliver have to interface somewhere, and the computational mechanisms that effect the interface must ipso facto have access to information from more than one cognitive domain. Consider:

(a) We have repeatedly distinguished between what the input systems compute and what the organism (consciously or subdoxastically) *believes*. Part of the point of this distinction is that input systems, being informationally encapsulated, typically compute representations of the distal layout on the basis of less information about the distal layout than the organism has available. Such representations want correction in light of background knowledge (e.g., information in memory) and of the simultaneous results of input analysis in other domains (see Aristotle on the 'common sense'). Call the process of arriving at such corrected representations "the fixation of perceptual belief." To a first approximation, we can assume that the mechanisms that effect this process work like this: they look simultaneously at the representations delivered by the various input systems and at the information currently in memory, and they arrive at a best (i.e., best available) hypothesis about how the world must be, given these various sorts of data.<sup>36</sup> But if there are mechanisms that fix perceptual belief, and if they work in anything like this way, then these mechanisms are not domain specific. Indeed, the point of having them is precisely to ensure that, wherever possible, what the organism believes is determined by all the information it has access to, regardless of which cognitive domains this information is drawn from.

(b) We use language (inter alia) to communicate our views on how the world is. But this use of language is possible only if the mechanisms that mediate the production of speech have access to what we see (or hear, or remember, or think) that the world is like. Since, by assumption, such mechanisms effect an interface among vertical faculties, they cannot themselves be domain specific. More precisely, they must at least be *less* domain specific than the vertical faculties are.<sup>37</sup>

(c) One aspect of the 'impenetrability' of the input systems is, we assumed, their insensitivity to the utilities of the organism. This assumption was required in part to explain the *veridicality* of perception given that the world doesn't always prove to be the way that we would prefer it to be. However, an interface between per-

ception and utilities must take place *somewhere* if we are to use the information that input systems deliver in order to determine how we ought to act. (Decision theories are, to all intents and purposes, models of the structure of this interface. The point is, roughly, that wishful seeing is avoided by requiring interactions with utilities to occur *after*—not *during*—perceptual integration.) So, again, the moral seems to be that there must be some mechanisms which cross the domains that input systems establish.

For these and other similar reasons, I assume that there must be relatively nondenominational (i.e., domain-inspecific) psychological systems which operate, *inter alia*, to exploit the information that input systems provide. Following the tradition, I shall call these "central" systems, and I will assume that it is the operation of these sorts of systems that people have in mind when they talk, pretheoretically, of such mental processes as thought and problem-solving. Central systems may be domain specific in *some* sense—we will consider this when we get to the issues about 'epistemic boundedness'—but at least they aren't domain specific in the way that input systems are. The interesting question about the central systems is whether, being nondenominational, they are also non-modular in other respects as well. That is, whether the central systems fail to exhibit the galaxy of properties that lead us to think of the input systems as a natural kind—the properties enumerated in Part III.

Briefly, my argument is going to be this: we have seen that much of what is typical of the input systems is more or less directly a product of their informational encapsulation. By contrast, I'll claim that central systems are, in important respects, *unencapsulated*, and that it is primarily for this reason that they are not plausibly viewed as modular. Notice that I am not going to be arguing for a tautology. It is perfectly possible, in point of logic, that a system which is *not* domain specific might nevertheless be encapsulated. Roughly, domain specificity has to do with the range of questions for which a device provides answers (the range of inputs for which it computes analyses); whereas encapsulation has to do with the range of information that the device consults in deciding what answers to provide. A system could thus be domain specific but unencapsulated (it answers a relatively narrow range of questions but in doing so it uses whatever it knows); and a system could be nondenomi-

national but encapsulated (it will give some answer to any question; but it gives its answers off the top of its head—i.e., by reference to less than all the relevant information). If, in short, it is true that only domain-specific systems are encapsulated, then that truth is interesting. Perhaps it goes without saying that I am not about to demonstrate this putative truth. I am, however, about to explore it.

So much for what I'm going to be arguing for. Now a little about the strategy of the argument. The fact is that there is practically no direct evidence, pro or con, on the question whether central systems are modular. No doubt it is possible to achieve some gross factoring of "intelligence" into "verbal" versus "mathematical/spatial" capacities; and no doubt there is something to the idea of a corresponding hemispheric specialization. But such dichotomies are *very* gross and may themselves be confounded with the modularity of the input systems—that is to say, they give very little evidence for the existence of domain-specific (to say nothing of modular) systems other than the ones that subserve the functions of perceptual and linguistic analysis.

When you run out of direct evidence, you might just as well try arguing from analogies, and that is what I propose to do. I have been assuming that the typical function of central systems is the fixation of belief (perceptual or otherwise) by nondemonstrative inference. Central systems look at what the input systems deliver, and they look at what is in memory, and they use this information to constrain the computation of 'best hypotheses' about what the world is like. These processes are, of course, largely unconscious, and very little is known about their operation. However, it seems reasonable enough that something can be inferred about them from what we know about *explicit* processes of nondemonstrative inference—viz., from what we know about empirical inference in science. So, here is how I am going to proceed. First, I'll suggest that scientific confirmation—the nondemonstrative fixation of belief in science—is typically unencapsulated. I'll then argue that if, pursuing the analogy, we assume that the central psychological systems are also unencapsulated, we get a picture of those systems that is, anyhow, not radically implausible given such information about them as is currently available.

The nondemonstrative fixation of belief in science has two prop-

erties which, though widely acknowledged, have not (so far as I know) yet been named. I shall name them: confirmation in science is *isotropic* and it is *Quinean*. It is notoriously hard to give anything approaching a rigorous account of what being isotropic and Quinean amounts to, but it is easy enough to convey the intuitions.

By saying that confirmation is isotropic, I mean that the facts relevant to the confirmation of a scientific hypothesis may be drawn from anywhere in the field of previously established empirical (or, of course, demonstrative) truths. Crudely: everything that the scientist knows is, in principle, relevant to determining what else he ought to believe. In principle, our botany constrains our astronomy, if only we could think of ways to make them connect.

As is usual in a methodological inquiry, it is possible to consider the isotropy of confirmation either normatively (as a principle to which we believe that rational inductive practice *ought* to conform) or sociologically (as a principle which working scientists actually adhere to in assessing the degree of confirmation of their theories). In neither case, however, should we view the isotropy of confirmation as merely gratuitous—or, to use a term of Rorty's (1979) as merely "optional." If isotropic confirmation 'partially defines the language game that scientists play' (remember when we used to talk that way?), that is because of a profound conviction—partly metaphysical and partly epistemological—to which scientists implicitly subscribe: the world is a connected causal system *and we don't know how the connections are arranged*. Because we don't, we must be prepared to abandon previous estimates of confirmational relevance as our scientific theories change. The point of all this is: confirmational isotropy is a reasonable property for nondemonstrative inference to have because the goal of nondemonstrative inference is to determine the truth about a causal mechanism—the world—of whose workings we are arbitrarily ignorant. That is why our institution of scientific confirmation is isotropic, and it is why it is plausible to suppose that what psychologists call "problem-solving" (i.e., nondemonstrative inference in the service of individual fixation of belief) is probably isotropic too.

The isotropy of scientific confirmation has sometimes been denied, but never, I think, very convincingly. For example, according to some historians it was part of the Aristotelian strategy against Galileo to claim that no data other than observations of the movements

of astronomical objects could, in principle, be relevant to the (dis)confirmation of the geocentric theory. Telescopic observations of the phases of Venus were thus ruled irrelevant *a priori*. In notably similar spirit, some linguists have recently claimed that no data except certain specified kinds of facts about the intuitions of native speakers could, in principle, be relevant to the (dis)confirmation of grammatical theories. Experimental observations from psycholinguistics are thus ruled irrelevant *a priori*. However, this sort of methodology seems a lot like special pleading: you tend to get it precisely when cherished theories are in trouble from *prima facie* disconfirming data. Moreover, it often comports with Conventionalist construals of the theories so defended. That is, theories for which nonisotropic confirmation is claimed are often viewed, even by their proponents, as merely mechanisms for making predictions; what is alleged in their favor is predictive adequacy rather than correspondence to the world. (Viewed from our perspective, nonisotropic confirmation is, to that extent, not a procedure for fixation of belief, since, on the Conventionalist construal, the predictive adequacy of a theory is *not* a reason for believing that the theory is *true*.)

One final thought on the isotropy issue. We are interested in isotropic systems because such systems are *ipso facto* unencapsulated. We are interested in scientific confirmation because (a) there is every reason to suppose that it is isotropic; (b) there is every reason to suppose that it is a process fundamentally similar to the fixation of belief; and (c) it is perhaps the only "global", unencapsulated, wholistic cognitive process about which anything is known that's worth reporting. For all that, scientific *confirmation* is probably not the best place to look if you want to see cognitive isotropy writ large. The best place to look, at least if one is willing to trust the anecdotes, is scientific *discovery*.

What the anecdotes say about scientific discovery—and they say it with a considerable show of univocality (see, e.g., papers in Ortony, 1979)—is that some sort of 'analogical reasoning' often plays a central role. It seems to me that we are thoroughly in the dark here, so I don't propose to push this point very hard. But it really does look as though there have been frequent examples in the history of science where the structure of theories in a new subject area has been borrowed from, or at least suggested by,

theories *in situ* in some quite different domain: what's known about the flow of water gets borrowed to model the flow of electricity; what's known about the structure of the solar system gets borrowed to model the structure of the atom; what's known about the behavior of the market gets borrowed to model the process of natural selection, which in turn gets borrowed to model the shaping of operant responses. And so forth. The point about all this is that "analogical reasoning" would seem to be isotropy in the purest form: a process which depends precisely upon the transfer of information among cognitive domains previously assumed to be mutually irreverent. By definition, encapsulated systems do not reason analogically.

I want to suggest two morals before I leave this point. The first is that the closer we get to what we are pretheoretically inclined to think of as the 'higher,' 'more intelligent', less reflexive, less routine exercises of cognitive capacities, the more such global properties as isotropy tend to show up. I doubt that this is an accident. I suspect that it is precisely its possession of such global properties that we have in mind when we think of a cognitive process as paradigmatically intelligent. The second moral preshadows a point that I shall jump up and down about further on. It is striking that, while everybody thinks that analogical reasoning is an important ingredient in all sorts of cognitive achievements that we prize, nobody knows anything about how it works; not even in the dim, in-a-glass-darkly sort of way in which there are some ideas about how confirmation works. I don't think that this is an accident either. In fact, I should like to propose a generalization; one which I fondly hope will some day come to be known as 'Fodor's First Law of the Nonexistence of Cognitive Science'. It goes like this: the more global (e.g., the more isotropic) a cognitive process is, the less anybody understands it. Very global processes, like analogical reasoning, aren't understood at all. More about such matters in the last part of this discussion.

By saying that scientific confirmation is Quinean, I mean that the degree of confirmation assigned to any given hypothesis is sensitive to properties of the entire belief system; as it were, the shape of our whole science bears on the epistemic status of each scientific hypothesis. Notice that being Quinean and being isotropic are not the same properties, though they are intimately related. For example, if scientific confirmation is isotropic, it is quite possible

that some fact about photosynthesis in algae should be relevant to the confirmation of some hypothesis in astrophysics ("the universe in a grain of sand" and all that). But the point about being Quinean is that we might have two astrophysical theories, both of which make the same predictions about algae and about everything else that we can think of to test, but such that one of the theories is better confirmed than the other—e.g., on grounds of such considerations as simplicity, plausibility, or conservatism. The point is that simplicity, plausibility, and conservatism are properties that theories have in virtue of their relation to the whole structure of scientific beliefs *taken collectively*. A measure of conservatism or simplicity would be a metric over *global* properties of belief systems.

Consider, by way of a simple example, Goodman's original (1954) treatment of the notion of projectability. We know that two hypotheses that are equivalent in respect of all the available data may nevertheless differ in their level of confirmation depending on which is the more projectable. Now, according to Goodman's treatment, the projectability of a hypothesis is inherited (at least in part) from the projectability of its vocabulary, and the projectability of an item of scientific vocabulary is determined by the (weighted?) frequency with which that item *has been projected* in previously successful scientific theories. So, the whole history of past projections contributes to determining the projectability of any given hypothesis on Goodman's account, and the projectability of a hypothesis (partially) determines its level of confirmation. Similarly with such notions as simplicity, conservatism, and the rest if only we knew how to measure them.

The idea that scientific confirmation is Quinean is by no means untendentious. On the contrary, it was a legacy of traditional philosophy of science—one of the "dogmas of Empiricism" (Quine, 1953) that there must be *semantic* connections between each theory statement and some data statements. That is, each hypothesis about "unobservables" must *entail* some predictions about observables, such entailments holding in virtue of the meanings of the theoretical terms that the hypotheses contain.<sup>38</sup> The effect of postulating such connections would be to determine a priori that certain data would disconfirm certain hypotheses, *whatever the shape of the rest of one's science might be*. For, of course, if H entails O, the discovery that  $\neg O$  would entail that  $\neg H$ . To that extent, the (dis)confirmation of

H by  $\neg O$  is independent of global features of the belief system that H and O belong to. To postulate meaning relations between data statements and theory statements is thus to treat confirmation as a *local* phenomenon rather than a global one.

I emphasize this consideration because analogous semantic proposals can readily be found in the psychological literature. For example, in the sorts of cognitive theories espoused by, say, Bruner or Vygotsky (and, more recently, in the work of the "procedural" semanticists), it is taken for granted that there must be connections of meaning between 'concepts' and 'percepts'. Basically, according to such theories, concepts are recipes for sorting stimuli into categories. Each recipe specifies a (more or less determinate) galaxy of tests that one can perform to effect a sorting, and each stimulus category is identified with a (more or less determinate) set of outcomes of the tests. To put the idea crudely but near enough for present purposes, there's a rule that you can test for *dog* by finding out if a thing barks, and the claim is that this rule is constitutive (though not, of course, exhaustive) of the concept *dog*. Since it is alleged to be a *conceptual* truth that whether it barks is relevant to whether it's a dog, it follows that the confirmation relation between "a thing is a dog" and "it barks" is insensitive to global properties of one's belief system. So considerations of theoretical simplicity etc. *could* not, even in principle, lead to the conclusion that whether it barks is *irrelevant* to whether it's a dog. To embrace that conclusion would be to change the concept.

This sort of example makes it clear how closely related being Quinean and being isotropic are. Since, on the view just scouted, it is a matter of *meaning* that barking is relevant to dogness, it is not possible to discover on empirical grounds that one was wrong about that relevancy relation. But isotropy is the principle that *any* fact may turn out to be (ir)relevant to the confirmation of any other. The Bruner-Vygotsky-procedural semantics line is thus incompatible with the isotropy of confirmation as well as with its Quineanness.

In saying that confirmation is isotropic and Quinean, I am thus consciously disagreeing with major traditions in the philosophy of science and in cognitive psychology. Nevertheless, I shall take it for granted that scientific confirmation is Quinean and isotropic. (Those who wish to see the arguments should refer to such classic papers in the modern philosophy of science as Quine, 1953, and

Putnam, 1962.) Moreover, since I am committed to relying upon the analogy between scientific confirmation and psychological fixation of belief, I shall take it for granted that the latter must be Quinean and isotropic too, hence that the Bruner-Vygotsky-procedural semantics tradition in cognitive psychology must be mistaken. I propose, at this point, to be both explicit and emphatic. The argument is that the central processes which mediate the fixation of belief are typically processes of rational nondemonstrative inference and that, since processes of rational nondemonstrative inference are Quinean and isotropic, so too are central processes. In particular, the theory of such processes must be consonant with the principle that the level of acceptance of any belief is sensitive to the level of acceptance of any other and to global properties of the field of beliefs taken collectively.

Given these assumptions, I have now got two things to do: I need to show that this picture of the central processes is broadly incompatible with the assumption that they are modular, and I need to show that it is a picture that has some plausibility independent of the putative analogy between cognitive psychology and the philosophy of science.

I take it that the first of these claims is relatively uncontroversial. We argued that modularity is fundamentally a matter of informational encapsulation and, of course, informationally encapsulated is precisely what Quinean/isotropic systems are not. When we discussed input systems, we thought of them as mechanisms for projecting and confirming hypotheses. And we remarked that, viewed that way, the informational encapsulation of such systems is tantamount to a constraint on the confirmation metrics that they employ; the confirmation metric of an encapsulated system is allowed to 'look at' only a certain restricted class of data in determining which hypothesis to accept. If, in particular, the flow of information through such a system is literally bottom-to-top, then its informational encapsulation consists in the fact that the *i*th-level hypotheses are (dis)confirmed solely by reference to lower-than-*i*th level representations. And even if the flow of data is unconstrained *within* a module, encapsulation implies constraints upon the access of intramodular processes to extramodular information sources. Whereas, by contrast, isotropy is by definition the property that a system has when it can look at anything it knows about in the

course of determining the confirmation levels of hypotheses. So, in general, the more isotropic a confirmation metric is, the more heterogeneous the provenance of the data that it accepts as relevant to constraining its decisions. Scientific confirmation is isotropic in the limit in this respect; it provides a model of what the *nonmodular* fixation of belief is like.

Similarly with being Quinean. Quinean confirmation metrics are *ipso facto* sensitive to global properties of belief systems. Now, an informationally encapsulated system *could*, strictly speaking, nevertheless be Quinean. Simplicity, for example, could constrain confirmation even in a system which computes its simplicity scores over some arbitrarily selected subset of beliefs. But this is mere niggling about the letter. In spirit, global criteria for the evaluation of hypotheses comport most naturally with isotropic principles for the relevance of evidence. Indeed, it is only on the assumption that the selection of evidence is isotropic that considerations of simplicity (and other such global properties of hypotheses) are *rational* determinants of belief. It is epistemically interesting that H & T is a simpler theory than -H & T where H is a hypothesis to be evaluated and T is the rest of what one believes. But there is no interest in the analogous consideration where T is some *arbitrarily delimited* subset of one's beliefs. Where relevance is non-isotropic, assessments of relative simplicity can be gerrymandered to favor any hypothesis one likes. This is one of the reasons why the operation of (by assumption informationally encapsulated) input systems should not be identified with the fixation of perceptual belief; not, at least, by those who wish to view the fixation of perceptual belief as by and large a rational process.

So it seems clear that isotropic/Quinean systems are *ipso facto* unencapsulated; and if unencapsulated, then presumably non-modular. Or rather, since this is all a matter of degree, we had best say that *to the extent that* a system is Quinean and isotropic, it is also nonmodular. If, in short, isotropic and Quinean considerations are especially pressing in determining the course of the computations that central systems perform, it should follow that these systems differ in their computational character from the vertical faculties.

We are coming close to what we started out to find: an overall taxonomy of cognitive systems. According to the present proposal,

there are, at a minimum, two families of such systems: modules (which are, relatively, domain specific and encapsulated) and central processes (which are, relatively, domain neutral and isotropic/Quinean). We have suggested that the characteristic function of modular cognitive systems is input analysis and that the characteristic function of central processes is the fixation of belief. If this is right, then we have three ways of taxonomizing cognitive processes which prove to be coextensive:

FUNCTIONAL TAXONOMY: input analysis versus fixation of belief

TAXONOMY BY SUBJECT MATTER: domain specific versus domain neutral

TAXONOMY BY COMPUTATIONAL CHARACTER: encapsulated versus Quinean/isotropic

I repeat that this coextension, if it holds at all, holds contingently. Nothing in point of logic stops one from imagining that these categories cross-classify the cognitive systems. If they do not, then that is a fact about the structure of the mind. Indeed, it is a *deep* fact about the structure of the mind.

All of which would be considerably more impressive if there were better evidence for the view of central processes that I have been proposing. Thus far, that account rests entirely on the analogy between psychological processes of belief fixation and a certain story about the character of scientific confirmation. There is very little that I can do about this, given the current underdeveloped state of psychological theories of thought and problem-solving. For what it's worth, however, I want to suggest two considerations that seem relevant and promising.

The first is that the difficulties we encounter when we try to construct theories of central processes are just the sort we would expect to encounter if such processes are, in essential respects, Quinean/isotropic rather than encapsulated. The crux in the construction of such theories is that there seems to be no way to delimit the sorts of informational resources which may affect, or be affected by, central processes of problem-solving. We can't, that is to say, plausibly view the fixation of belief as effected by computations over bounded, local information structures. A graphic example of this sort of difficulty arises in AI, where it has come to be known as the "frame problem" (i.e., the problem of putting a "frame"

around the set of beliefs that may need to be revised in light of specified newly available information. Cf. the discussion in McCarthy and Hayes, 1969, from which the following example is drawn).

To see what's going on, suppose you were interested in constructing a robot capable of coping with routine tasks in familiar human environments. In particular, the robot is presented with the job of phoning Mary and finding out whether she will be late for dinner. Let's assume that the robot 'knows' it can get Mary's number by consulting the directory. So it looks up Mary's number and proceeds to dial. So far, so good. But now, notice that commencing to dial has all sorts of direct and indirect effects on the state of the world (including, of course, the internal state of the robot), and some of these effects are ones that the device needs to keep in mind for the guidance of its future actions and expectations. For example, when the dialing commences, the phone ceases to be free to outside calls; the robot's fingers (or whatever) undergo appropriate alterations of spatial location; the dial tone cuts off and gets replaced by beeps; something happens in a computer at Murray Hill; and so forth. Some (but, in principle, not all) such consequences are ones that the robot must be designed to monitor since they are relevant to "updating" beliefs upon which it may eventually come to act. Well, *which* consequences? The problem has at least the following components. The robot must be able to identify, with reasonable accuracy, those of its previous beliefs whose truth values may be expected to alter as a result of its current activities; and it must have access to systems that do whatever computing is involved in effecting the alterations.

Notice that, unless these circuits are arranged correctly, things can go absurdly wrong. Suppose that, having consulted the directory, the robot has determined that Mary's number is 222-2222, which number it commences to dial, pursuant to instructions previously received. But now it occurs to the machine that *one of the beliefs that may need updating in consequence of its having commenced dialing is its (recently acquired) belief about Mary's telephone number*. So, of course, it stops dialing and goes and looks up Mary's telephone number (again). Repeat, *da capo*, as many times as may amuse you. Clearly, we have here all the makings of a computational trap. Unless the robot can be assured that some of its beliefs are invariant under some of its actions, it will never get to *do* anything.

How, then, does the machine's program determine which beliefs the robot ought to reevaluate given that it has embarked upon some or other course of action? What makes this problem so hard is precisely that it seems unlikely that any *local* solution will do the job. For example, the following truths appear to be self-evident: First, that there is no fixed set of beliefs such that, for any action, those and only those beliefs are the ones that require reconsideration. (That is, which beliefs are up for grabs depends intimately upon which actions are performed and upon the context of the performances. There are *some*—indeed, indefinitely many—actions which, if performed, *should* lead one to consider the possibility that Mary's telephone number has changed in consequence.) Second, new beliefs don't come docketed with information about which old beliefs they ought to affect. On the contrary, we are forever being surprised by the implications of what we know, including, of course, what we know about the actions we perform. Third, the set of beliefs apt for reconsideration cannot be determined by reference to the recency of their acquisition, or by reference to their generality, or by reference to merely semantic relations between the contents of the beliefs and the description under which the action is performed . . . etc. Should any of these propositions seem less than self-evident, consider the special case of the frame problem where the robot is a mechanical scientist and the action performed is an experiment. Here the question 'which of my beliefs ought I to reconsider given the possible consequences of my action' is transparently equivalent to the question "What, in general, is the optimal adjustment of my beliefs to my experiences?" This is, of course, exactly the question that a theory of confirmation is supposed to answer; and, as we have been at pains to notice, confirmation is not a relation reconstructible by reference to local properties of hypotheses or of the data that bear upon them.

I am suggesting that, as soon as we begin to look at cognitive processes other than input analysis—in particular, at central processes of nondemonstrative fixation of belief—we run into problems that have a quite characteristic property. They seem to involve isotropic and Quinean computations; computations that are, in one or other respect, sensitive to the whole belief system. This is exactly what one would expect on the assumption that nondemonstrative fixation of belief really is quite like scientific confirmation, and that

scientific confirmation is itself characteristically Quinean and isotropic. In this respect, it seems to me, the frame problem is paradigmatic, and in this respect the seriousness of the frame problem has not been adequately appreciated.

For example, Raphael (1971) comments as follows: "(An intelligent robot) will have to be able to carry out tasks. Since a task generally involves some change in the world, it must be able to update its model (of the world) so it remains as accurate during and after the performance of a task as it was before. Moreover, it must be able to *plan* how to carry out a task, and this planning process usually requires keeping 'in mind' simultaneously a variety of possible actions and corresponding models of hypothetical worlds that would result from those actions. The bookkeeping problems involved with keeping track of these hypothetical worlds account for much of the difficulty of the frame problem" (p. 159). This makes it look as though the problem is primarily (a) how to notate the possible worlds and (b) how to keep track of the *demonstrative* consequences of changing state descriptions. But the deeper problem, surely, is to keep track of the *nondemonstrative* consequences. Slightly more precisely, the problem is, given an arbitrary belief world W and a new state description 'a is F', what is the appropriate successor belief world W'? What ought the device to believe, given that it used to believe W and now believes that a is F? But this isn't just a bookkeeping problem; it is the general problem of inductive confirmation.<sup>39</sup>

So far as I can tell, the usual assumption about the frame problem in AI is that it is somehow to be solved 'heuristically'. The idea is that, while nondemonstrative confirmation (and hence, presumably, the psychology of belief fixation) is isotropic and Quinean *in principle*, still, given a particular hypothesis, there are, in practice, heuristic procedures for determining the range of effects its acceptance can have on the rest of one's beliefs. Since these procedures are by assumption merely heuristic, they may be assumed to be local—i.e., to be sensitive to less than the whole of the belief systems to which they apply. Something like this may indeed be true; there is certainly considerable evidence for heuristic short-cutting in belief fixation, deriving both from studies of the psychology of problem-solving (for a recent review, see Nisbett and Ross, 1980) and from the sociology of science (Kuhn, 1970). In such cases, it is possible

to show how potentially relevant considerations are often systematically ignored, or distorted, or misconstrued in favor of relatively local (and, of course, highly fallible) problem-solving strategies. Perhaps a bundle of such heuristics, properly coordinated and rapidly deployed, would suffice to make the central processes of a robot as Quinean and isotropic as yours, or mine, or the practicing scientist's ever actually succeed in being. Since there are, at present, no serious proposals about what heuristics might belong to such a bundle, it seems hardly worth arguing the point.

Still, I am going to argue it a little.

There are those who hold that ideas recently evolved in AI—such notion as, e.g., those of 'frame' (see Minsky, 1975)<sup>40</sup> or 'script' (see Schank and Abelson, 1975)—will illuminate the problems about the globality of belief fixation since they do, in a certain sense, provide for placing a frame around the body of information that gets called when a given sort of problem is encountered. (For a discussion that runs along these optimistic lines, see Thagard, 1980.) It seems to me, however, that the appearance of progress here is entirely illusory—a prime case of confusing a notation with a theory.

If there were a principled solution to the frame problem, then no doubt that solution could be expressed as a constraint on the scripts, or frames, to which a given process of induction has access. But, lacking such a solution, there is simply no content to the idea that only the information represented in the frame (/script) that a problem elicits is computationally available for solving the problem. For one thing, since there are precisely no constraints on the individuation of frames (/scripts), *any* two pieces of information can belong to the same frame (/script) at the discretion of the programmer. This is just a way of saying that the solution of the frame problem can be accommodated to the frame (/script) notation *whatever that solution turns out to be*. Which is just another way of saying that the notation does not constrain the solution. Second, it is a widely advertised property of frames (/scripts) that they can cross-reference to one another. The frame for Socrates says, among other things, 'see Plato' . . . and so forth. There is no reason to doubt that, in any developed model, the system of cross-referencing would imply a graph in which there is a route (of greater or lesser length) from each point to any other. But now we have the frame problem all over again, in the form: Which such paths should

actually be traversed in a given case of problem-solving, and what should bound the length of the trip? All that has happened is that, instead of thinking of the frame problem as an issue in the logic of confirmation, we are now invited to think of it as an issue in the theory of executive control (a change which there is, by the way, no reason to assume is for the better). More of this presently.

For now, let's summarize the major line of argument. If we assume that central processes are Quinean and isotropic, then we ought to predict that certain kinds of problems will emerge when we try to construct psychological theories which simulate such processes or otherwise explain them; specifically, we should predict problems that involve the characterization of nonlocal computational mechanisms. By contrast, such problems should not loom large for theories of psychological modules. Since, by assumption, modular systems are informationally encapsulated, it follows that the computations they perform are relatively local. It seems to me that these predictions are in reasonably good accord with the way that the problems of cognitive science have in fact matured: the input systems appear to be primarily stimulus driven, hence to exploit computational processes that are relatively insensitive to the general structure of the organism's belief system. Whereas, when we turn to the fixation of belief, we get a complex of problems that appear to be intractable precisely because they concern mental processes that aren't local. Of these, the frame problem is, as we have seen, a microcosm.

I have been marshaling considerations in favor of the view that central processes are Quinean/isotropic. That is what the analogy to scientific confirmation suggests that they ought to be, and the structure of the problems that arise in attempts to model central processes is quite compatible with that view of them. I now add that the view of central processes as computationally global can perhaps claim some degree of neurological plausibility. The picture of the brain that it suggests is a reasonably decent first approximation to the kind of brain that it appears we actually have.

When we discussed input analyzers, I commented on the natural connection between informational encapsulation and fixed neural architecture. Roughly, standing restrictions on information flow imply the option of hardwiring. If, in the extreme case, system B is required to take note of information from system A and is allowed

to take note of information from nowhere else, you might as well build your brain with a permanent neuroanatomical connection from A to B. It is, in short, reasonable to expect biases in the distribution of information to mental processes to show up as structural biases in neural architecture.

Consider, by contrast, Quinean/isotropic systems, where more or less any subsystem may want to talk to any other at more or less any time. In this case, you'd expect the corresponding neuroanatomy to be relatively diffuse. At the limit, you might as well have a random net, with each computational subsystem connected, directly or indirectly, with every other; a kind of wiring in which you get a minimum of stable correspondence between neuroanatomical form and psychological function. The point is that in Quinean/isotropic systems, it may be *unstable, instantaneous* connectivity that counts. Instead of hardwiring, you get a connectivity that changes from moment to moment as dictated by the interaction between the program that is being executed and the structure of the task in hand. The moral would seem to be that computational isotropy comports naturally with neural isotropy (with what Lashley called "equipotentiality" of neural structure) in much the same way that informational encapsulation comports naturally with the elaboration of neural hardwiring.

So, if input analysis is modular and thought is Quinean/isotropic, you might expect a kind of brain in which there is stable neural architecture associated with perception-and-language but not with thought. And, I suggest, this seems to be pretty much what we in fact find. There is, as I remarked above, quite a lot that can be said about the neural specificity of the perceptual and linguistic mechanisms: at worst we can enumerate in some detail the parts of the brain that handle them; and at best we can exhibit characteristic neural architecture in the areas where these functions are performed. And then there are the rest of the higher brain systems (cf. what used to be called "association cortex"), in which neural connectivity appears to go every which way and the form/function correspondence appears to be minimal. There is some historical irony in all this. Gall argued from a (vertical) faculty psychology to the macroscopic differentiation of the brain. Flourens, his archantagonist, argued from the unity of the Cartesian ego to the brain's equipotentiality (see Bynum, op. cit.). The present suggestion is that they were *both* right.<sup>41</sup>

I am, heaven knows, not about to set up as an expert on neuropsychology, and I am painfully aware how impressionistic this all is. But while we're collecting impressions, I think the following one is striking. A recent issue of *Scientific American* (September, 1979) was devoted to the brain. Its table of contents is quite as interesting as the papers it contains. There are, as you might expect, articles that cover the neuropsychology of language and of the perceptual mechanisms. But there is nothing on the neuropsychology of thought—presumably because nothing is known about the neuropsychology of thought. I am suggesting that there is a good reason why nothing is known about it—namely, that there is nothing to know about it. You get form/function correspondence for the modular processes (specifically, for the input systems); but, in the case of central processes, you get an approximation to universal connectivity, hence no stable neural architecture to write *Scientific American* articles about.

To put these claims in a nutshell; there are *no* content-specific central processes for the performance of which correspondingly specific neural structures have been identified. Everything we now know is compatible with the claim that central problem-solving is subserved by equipotential neural mechanisms. This is precisely what you would expect if you assume that the central cognitive processes are largely Quinean and isotropic.

#### PART V

#### CAVEATS AND CONCLUSIONS

We now have before us what might be called a 'modified' modularity theory of cognitive processes. According to this theory, Gall was right in claiming that there are vertical faculties (domain specific computational mechanisms). Indeed, a still stronger claim is plausible: that the vertical faculties are modules (informationally encapsulated, neurologically hardwired, innately specified and so forth). But nonmodular cognitive systems are also acknowledged, and it is left open that these latter may exhibit features of horizontal organization. Roughly speaking, on this account, the distinction between vertical and horizontal modes of computational organization is taken to be coextensive with the functional distinction