

# Neural Routing Circuits for Forming Invariant Representations of Visual Objects

Thesis by

Bruno Adolphus Olshausen

In Partial Fulfillment of the Requirements

for the Degree of

Doctor of Philosophy

California Institute of Technology

Pasadena, California

1994

(Submitted January 27, 1994)

©1994

Bruno A. Olshausen

All rights reserved

## Acknowledgements

*What a person thinks on his own, without being stimulated by the thoughts of other people, is, even in the best case, rather paltry and monotonous.*

— Albert Einstein

The work presented in this thesis is in many ways a reflection of ideas that have been suggested or provoked through discussions with others. At times it was merely a passing phrase in conversation, or a single sentence in a talk, that would crop up later in the form of an idea; and sometimes this would even lead to a major turning point. In large part, then, I owe this thesis to those who I have had the honor and pleasure of working with over the past several years.

First and foremost I thank Charlie Anderson, who has been my constant source of inspiration, providing a never ending stream of stimulating thoughts and ideas that have been the driving force behind my work. It is largely his early work on shifter circuits that forms the backbone of this thesis. David Van Essen was instrumental in guiding this work through various stages to form a detailed neurobiological model, and he has been invaluable as a mentor and source of encouragement. I am also deeply indebted to Pentti Kanerva, who provided me with my first opportunity doing

research in the Sparse Distributed Memory group at NASA Ames Research Center. It was my work there on visual coding that eventually led to my getting involved working with Charlie and David on shifter circuits.

The Caltech CNS program played a major role in making this work happen by providing its graduate students with the leeway and resources to pursue their own ideas, and by fostering an atmosphere of cooperativity among students and faculty. This provided a rich intellectual environment from which I have benefited greatly. Conversations with Mike Lewicki were especially helpful for shaping my thinking and ideas in a productive direction. Al Barr's lectures on modeling and simulation provided the impetus for developing the closed-loop control system for the autonomous routing circuit. Discussions with Bill Press spurred me to do further research on the pulvinar (presented in Appendix C), which led to the development of an anatomy database for studying pulvinar-cortical interconnectivity. Christof Koch gave many valuable suggestions. And Gilles Laurent provided advise and support in a short-term project to look for evidence of shifter circuits in the visual system of the jumping spider.

I was also fortunate to be able to spend the past year at Washington University in St. Louis. Here, Chris Lee was a primary source of inspiration and encouragement, and it was discussions with him that motivated the development of the Bayesian analogy presented in Chapter 4.

I also thank friends Ojvind Bernander and Steve Turney for their support, and for providing the much needed diversions from research in both California and Missouri, from hiking the high peaks of the Sierra Nevada to canoeing the rivers of the Ozarks.

Finally, I owe this thesis to Mom and Dad, whose love, support, and encouragement at every step along the way—from the time I was a boy building my first electric circuit up to the trials and tribulations of grad school—made it all possible.

## Abstract

This thesis presents a biologically plausible model of an attentional mechanism for forming position- and scale-invariant representations of objects in the visual world. The model relies on a set of *control neurons* to dynamically modify the synaptic strengths of intra-cortical connections so that information from a windowed region of primary visual cortex (V1) is selectively routed to higher cortical areas. Local spatial relationships (i.e., topography) within the attentional window are preserved as information is routed through the cortex, thus enabling attended objects to be represented in higher cortical areas within an object-centered reference frame that is position and scale invariant. The representation in V1 is modeled as a multiscale stack of sample nodes with progressively lower resolution at higher eccentricities. Large changes in the size of the attentional window are accomplished by switching between different levels of the multiscale stack, while positional shifts and small changes in scale are accomplished by translating and rescaling the window within a single level of the stack. The control signals for setting the position and size of the attentional window are hypothesized to originate from neurons in the pulvinar and in the deep layers of visual cortex. The dynamics of these control neurons are governed by simple differential equations that can be realized by neurobiologically plausible circuits. In pre-attentive

mode, the control neurons receive their input from a low-level “saliency map” representing potentially interesting regions of a scene. During the pattern recognition phase, control neurons are driven by the interaction between top-down (memory) and bottom-up (retinal input) sources. The model respects key neurophysiological, neuroanatomical, and psychophysical data relating to attention, and it makes a variety of experimentally testable predictions.

# Contents

<b>Acknowledgements</b>	<b>iii</b>
<b>Abstract</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 The problem: separating “what” from “where” . . . . .	2
1.2 Background . . . . .	4
Early ideas . . . . .	4
Log-polar theories . . . . .	6
Feature <i>Gemisch</i> theories . . . . .	7
Remapping . . . . .	11
1.3 Summary of the proposal . . . . .	13
1.4 Motivations for the proposal . . . . .	15
Psychology . . . . .	15
Neurobiology . . . . .	18
Computational complexity . . . . .	18
1.5 Caveats (what’s not in this thesis) . . . . .	20
1.6 What’s in this thesis . . . . .	20
<b>2 Dynamic Routing Circuits</b>	<b>22</b>

2.1	Neurobiological constraints . . . . .	22
2.2	Routing circuit architecture . . . . .	24
	The simplest possible routing circuit . . . . .	24
	Accommodating larger convergence . . . . .	30
	Multiscale input representation . . . . .	33
	Logarithmic spatial sampling . . . . .	34
2.3	Autonomous control I: single stage . . . . .	39
	System objective . . . . .	39
	Focusing attention on a blob . . . . .	42
	Recognition . . . . .	46
	Shifting attention . . . . .	53
2.4	Autonomous control II: multiple stages . . . . .	54
	Focusing on a blob . . . . .	54
	Recognition . . . . .	60
	Shifting attention . . . . .	61
	Competition among scales . . . . .	64
2.5	Summary of the model . . . . .	67
<b>3</b>	<b>Neurobiological substrates and mechanisms</b>	<b>73</b>
3.1	Routing circuit substrates . . . . .	75
	The multiscale “stack” . . . . .	77
	Intermediate stages and wiring constraints . . . . .	81
3.2	Control substrates . . . . .	84
	Control neurons . . . . .	84
	Bottom-up control sources (saliency map) . . . . .	86
	Recognition-guided control sources . . . . .	88
3.3	Gating mechanisms . . . . .	89



<b>4 Discussion</b>	<b>92</b>
4.1 Predictions . . . . .	92
Neurophysiology . . . . .	92
Neuroanatomy . . . . .	96
Psychophysics . . . . .	97
4.2 Comparison with other models . . . . .	100
Control vs. synchronicity . . . . .	100
Control-based network models . . . . .	102
4.3 Generalizations of the model . . . . .	103
Bayesian interpretation . . . . .	103
The big picture . . . . .	106
Control as a general strategy for neural computation . . . . .	108
4.4 Unresolved issues . . . . .	110
Features instead of pixels . . . . .	110
Feedback pathways . . . . .	110
Pop-out in multiple dimensions . . . . .	111
Rotation and warp . . . . .	111
3D objects . . . . .	112
Learning . . . . .	112
<b>5 Conclusions</b>	<b>113</b>
<b>A Derivation of autonomous control dynamics</b>	<b>115</b>
A.1 Blob search . . . . .	115
A.2 Recognition . . . . .	117
<b>B Open questions about spatial frequency tuning</b>	<b>119</b>
B.1 Conflicting data . . . . .	119

B.2	How do we see so clearly? . . . . .	120
B.3	Very low-frequency cells . . . . .	123
<b>C</b>	<b>Details of pulvinar anatomy and physiology</b>	<b>125</b>
C.1	An overview of the pulvinar . . . . .	125
C.2	Doubts about pulvinar . . . . .	128
C.3	The case for pulvinar . . . . .	129
C.4	Other alternatives . . . . .	131
C.5	Which data are important? . . . . .	133
C.6	Design of experiments . . . . .	139
C.7	Conclusion . . . . .	141

## List of Figures

1.1	The model of Pitts and McCulloch (1947). . . . .	6
1.2	Fourier transform method for producing invariant representations. . . . .	9
1.3	Two objects with the same features but different spatial relationships are not equivalent. . . . .	10
1.4	Overview of the model. . . . .	14
1.5	Reference frame effects. . . . .	16
1.6	Encoding <i>what</i> and <i>where</i> with limited neural resources. . . . .	19
2.1	A simple routing circuit. . . . .	25
2.2	An illustration of “connection space.” . . . .	27
2.3	Some possible control scenarios. . . . .	29
2.4	A multistage dynamic routing circuit. . . . .	31
2.5	Connection space for the multistage dynamic routing circuit. . . . .	32
2.6	Dynamic routing circuit with a multiscale input representation. . . . .	35
2.7	The Koenderink “stack.” . . . .	36
2.8	Dynamic routing circuit with a “stack” input representation. . . . .	37
2.9	Routing circuit for a single scale within the stack circuit. . . . .	38
2.10	A simple attentional strategy for an autonomous visual system. . . . .	41

2.11	A single-stage routing circuit with a Gaussian blob presented to the input units. . . . .	42
2.12	Autonomous control. . . . .	45
2.13	Computer simulation of the autonomous routing circuit. . . . .	45
2.14	Control neuron interactions when configured into control blocks. . . . .	47
2.15	An autonomous routing circuit for recognition. . . . .	50
2.16	Computer simulation of the recognition circuit. . . . .	52
2.17	Shifting attention. . . . .	53
2.18	A two-stage routing circuit and its control. . . . .	55
2.19	Autonomous control of a multistage routing circuit. . . . .	57
2.20	Simulation of an autonomous, multistage routing circuit. . . . .	59
2.21	Multistage recognition circuit. . . . .	61
2.22	Shifting attention in the multistage routing circuit. . . . .	63
2.23	Autonomous control for the multiscale stack routing circuit. . . . .	64
2.24	Modified saliency function for scale selectivity. . . . .	66
2.25	Simulation of the stack circuit. . . . .	69
2.26	Simulation of the stack circuit. . . . .	70
2.27	Simulation of the stack circuit (local vs. global). . . . .	71
2.28	Simulation of the stack circuit (local vs. global). . . . .	72
3.1	Neurobiological substrates. . . . .	74
3.2	Schematic of the main cortical areas in the “form” pathway. . . . .	76
3.3	A six-level “stack” model for V1. . . . .	78
3.4	The stack in cortical dimensions. . . . .	80
3.5	Some possible multistage routing circuits. . . . .	82
4.1	The dynamic routing circuit interpretation of the Moran and Desimone (1985) experiment. . . . .	94

4.2	The meaning of “cycles per object.” . . . . .	99
4.3	The big picture. . . . .	107
4.4	A more general way of viewing control. . . . .	109
B.1	Relation of spatial-frequency tuning to perception. . . . .	122
C.1	A schematic of the known connections of the pulvinar. . . . .	126
C.2	The two maps of visual space within the pulvinar. . . . .	127

# Chapter 1

## Introduction

To date, the only known devices that can “see” in any meaningful sense are biological vision systems. In fact, were it not for their existence, we might be led to believe that vision as we know it is physically impossible. If we wish to understand the principles involved in vision, it behooves us to study the systems that do it well.

The goal of this thesis is to understand how a particular problem in vision is solved by the primate visual system. The approach will be to study the available neurobiological substrates and their computational properties, and then to formulate a model for solving the problem with this hardware. We shall attempt to construct the model in a detailed enough manner so that it is capable of generating useful experimental predictions.

This chapter begins with a description of the particular problem being addressed in this thesis, and the previous models that have been proposed for solving it. This is followed by a summary of the proposed model and its motivations from psychology, neurobiology, and computational complexity.

## 1.1 The problem: separating “what” from “where”

The problem I shall be concerned with is how the brain builds a neural representation of objects and their locations and sizes within a scene, given the primitive image description at the retina. That is, how do we form a representation for *what* and *where* things are in the visual world, beginning only with pixels?

Our current, limited understanding of how this is done is that the brain extracts progressively more complex forms of structure at each stage of visual processing. For example, in the retina, ganglion cells appear to code for contrast, and their physiological responses can be understood in terms of a process that reduces the “redundancy” (or structure) present in natural images (Atick and Redlich, 1990). This principle can also be extended to understand the properties of orientation selective cells in primary visual cortex (Field, 1987; Li and Atick, 1994). At progressively higher levels of cortical processing (V2, V4), one finds that cells can respond quite selectively to more complex forms of structure, such as spiral or hyperbolic patterns (Gallant et al., 1993). It is far less understood what forms of structure are being represented by cells at these stages, though, or what principle of information processing is being followed. Finally, at the highest levels of form processing (IT, STS), cells appear to code for the presence of specific complex objects, such as hands and faces, without regard for their position or size (Gross et al., 1972; Rolls and Baylis, 1986).

At the lower stages of visual processing, it seems feasible to extract local structure in parallel at different scales and positions across the retina. For example, millions of simple cells in area V1 code for the local orientation of contrast, each at a different position and scale. If this massively parallel coding scheme were to continue for progressively more numerous and complex features, though, the number of cells required to code for every position and scale would grow rapidly. Sooner or later, combinatorial explosion catches up with you, and the number of neural resources required will

exceed what is available in the brain.

Consider the problem of learning a particular person's face. What distinguishes one face from another face are the particular features—eye, nose, mouth, etc.—and their particular spatial relationships to each other. For example, a distinguishing characteristic may be a small upper lip that separates the nose from the mouth. If this arrangement is learned by a certain cell, or population of cells, at one position and size on the retina, then how can this knowledge be retained when the face is presented at a different position and size? Encoding the face at each position and size—or even a coarse-coding—is certainly not an efficient solution, since it would quickly deplete our neural resources to represent the countless number of objects we can readily recognize. In addition, this scheme would require that objects be presented in all possible configurations before they could be recognized at any arbitrary size or position on the retina.

There seem to be two major factors at work here: *what* the object is (i.e., the identity of the face) is determined by the particular features and spatial arrangements that make it up, independent of *where* it is (i.e., its location and size) in the image. An efficient solution would thus be to encode *what* separately from *where* so they are represented independently. The number of neural resources required would then be dramatically reduced, since it would no longer be necessary to represent the conjoint space of *what* and *where* (see Fig. 1.6). In addition, the learning of any *what* could naturally be generalized to any *where*.

Despite years of neurobiological research and the many models of recognition that have been proposed, there still exists no coherent neural model to explain how we could form independent representations of *what* and *where*. The model proposed in this thesis is an attempt to understand how this may be accomplished—albeit in a simplified form—by the neural machinery of the primate brain.



## 1.2 Background

We will review the previous work on this problem by first describing some of the early ideas, since they are so very different in the form of their approach. The more modern approaches can essentially be subdivided into three major classes: those utilizing log-polar transforms, those based on what I shall call “feature *Gemisch*” theories, and those that remap visual information from one reference frame to another. We consider each of these in turn.

### Early ideas

The Greeks were the first to ponder the general question of how, by seeing different examples of something, we learn that all the examples are instances of the same thing (Anderson and Rosenfeld, 1988). In more recent history, this question was picked up by the Gestaltists, who pondered how it is that we can recognize a square as a square no matter where it appears in the visual field. They made no real progress, though, in thinking about the problem in neural terms.

Probably the first to ponder a neural theory for how objects are represented independent of position and size was the early psychobiologist K.S. Lashley. From his own description of the problem (Lashley, 1942), one gets the feeling that he was completely perplexed by it:

Visual fixation can be held accurately for only a moment, yet, in spite of changes in direction of gaze, an object remains the same object. An indefinite number of combinations of retinal cells and afferent paths are equivalent in perception and in the reactions they produce. This is the most elementary problem of cerebral function and I have come to doubt that any progress will be made toward a genuine understanding of nervous integration until the problem of equivalent nervous connections, or as it is more generally termed, of stimulus equivalence, is solved. (p. 304)

The solution that he proposed was that the presentation of an object on the retina

would form waves of activity in the cortex that emanate from the object via the horizontal connections. These waves would then form characteristic interference patterns for each object, without regard to its position or size. The model did not get very specific beyond this rather general notion, though, and it is certainly at odds with our modern understanding of the brain.

The first to actually propose a neurobiologically detailed model were Pitts and McCulloch (1947). As they stated it, “We seek general methods for designing nervous nets which recognize figures in such a way as to produce the same output for every input belonging to the figure. We endeavor particularly to find those which fit the histology and physiology of the actual structure.” Their model attempted not only to account for invariance to visual forms, but for invariance to auditory forms as well—e.g., recognizing a chord regardless of pitch. It is worthwhile to describe their model in some detail, because it has so much in common with the model proposed in this thesis. It is best explained by examining the circuit they proposed for the auditory system, illustrated in Figure 1.1. The distribution of excitation along Heschel’s gyrus (the area of cortex which is organized tonotopically) is represented by the function  $\phi(x)$ . This distribution is transmitted slantwise through a translator which reads the distribution shifted by an amount  $a$ , depending on the level of the translator. When level  $M_a$  is turned on, it sends its distribution,  $\phi(x + a)$ , to a common layer in the “depths.” A sweep control turns on one level at a time sequentially (thought to be set to the alpha rhythm). Thus, the distribution of excitation produced on Heschel’s gyrus by a particular chord will move uniformly back and forth in the depths, preserving intervals. Some appropriate mechanism (they don’t exactly specify) is then used to average over this group to recognize the sweeping distribution as a whole. A similar mechanism was proposed to explain how an object can be recognized independent of size. In this case, a two-dimensional pattern in V1 would be sequentially zoomed in and out, preserving form. Although this model did not turn out to be

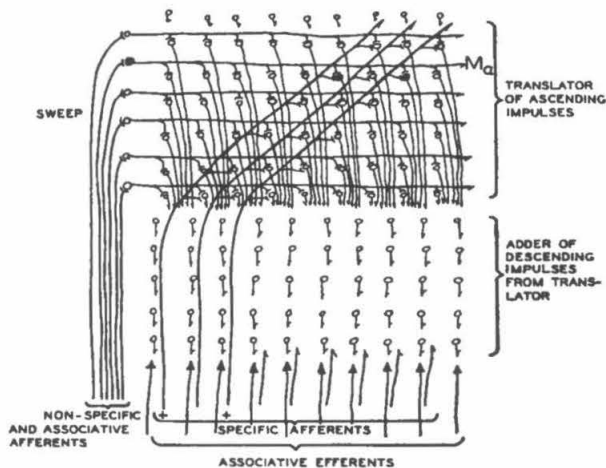


Figure 1.1: **The model of Pitts and McCulloch (1947).**

The function  $\phi(x)$  is represented on the “specific afferents” and is transmitted slantwise through the translator. Here, level  $M_a$  is turned on by the sweep control (*left*), and the shifted function  $\phi(x + a)$  is sent into the “adder” circuit below.

true, it is admirable for its detailed neuroanatomical correlates and its attempt to be consistent with the known physiology of the day (alpha rhythms, etc.). In addition, it is probably the first explicit use of a gating mechanism in a model for remapping sensory information. As we shall soon see, this basic idea forms the underpinning of the model proposed in this thesis.

## Log-polar theories

It was observed early on by Fischer (1973) and Schwartz (1977) that the transformation from retina to cortex is approximately a log-polar transform. Schwartz proposed that the transform could be described by the function  $\log(z + a)$ , where  $z$  is the complex variable  $x + iy$  ( $x$  and  $y$  being the Cartesian coordinates in the image), and  $a$  is a scalar offset. Under this transformation, changes in the scale and orientation of an

object centered in the Cartesian coordinate system are converted into approximate horizontal and vertical shifts in the log-polar, or cortical, coordinate system. While this is indeed an interesting property that naturally falls out of a log-polar transform, it seems doubtful that it could serve as the principal means by which size and orientation invariance are achieved. One problem is that an object would have to be centered fairly precisely on the fovea in order for this property to hold. If the object happened to be offset from the fovea, it would undergo a very strange transformation when rescaled or rotated in the retina. Moreover, even if the object were centered on the fovea, some appropriate shifting mechanism would be required to remove the translation in log-polar space to achieve an invariant representation, but none was proposed. Baron (1987) has made a similar proposal that elaborates on this scheme for shifting, rescaling, and rotating objects into a canonical representation, but again without any specific mechanisms for performing the requisite shift operation.

This technique has been applied in conjunction with a Fourier transform in a number of machine vision systems (Wechsler and Zimmerman 1988; Carpenter and Grossberg, 1987b). However, these systems are highly constrained in the types of images and visual environments they can handle (see below).

## **Feature *Gemisch* theories**

A number of theories have been built on the notion that the brain removes translation, scale, and other variations by transforming the representation of an object in such a way that the presence of certain features of the object are preserved, but information about the spatial relationships among the features is lost. Objects are then identified by looking for the right mixture, or *Gemisch*, of features without regard for their specific spatial relationships to each other. These models fall into two major classes: those based on Fourier transforms, and those based on hierarchies of feature selective,

position insensitive cells.

The basic idea behind using a Fourier transform is that the amplitude spectrum of the transform is invariant to shifts in the image. That is, for the Fourier transform,  $\hat{I}(u, v)$ , of an image,  $I(x, y)$ , the amplitude spectrum,

$$|\hat{I}(u, v)| = \sqrt{\text{Re}\{\hat{I}\}^2 + \text{Im}\{\hat{I}\}^2} \quad (1.1)$$

will remain constant (barring edge effects) no matter how the pattern is positioned in the input image  $I(x, y)$ . This is illustrated in Figure 1.2. Pollen et al. (1971) hypothesized that this transform is computed on a local scale by the complex cells in V1, and that it constitutes the beginning stage of translation invariant perception. Cavanagh (1978, 1985) expanded on this basic idea by proposing that local log-polar frequency transforms are computed in V1 and then summed globally to form a global log-polar frequency transform. Taking a Fourier transform of this new representation (in IT) would then result in a scale and rotation invariant transformation.

Although the invariance properties of these transforms are interesting, it is difficult to see how they could account for translation and scale invariant perception. Perhaps the most serious difficulty is that these transforms will be highly sensitive to spurious patterns in the input, such as occlusions or shadows. In addition, because the phase is ignored, there are a multitude of nonsensical images that will result in the same amplitude spectrum for any given object (Fig. 1.2c). If this method were in fact used for achieving invariant perception, then one would expect these images to be perceived equivalently.

The other class of feature *Gemisch* methods, based on hierarchies of feature selective, position insensitive cells, was first introduced by Fukushima (1980) in his “Neocognitron” model. This model uses successive stages of feature extraction and position invariance to build an invariant representation of objects at the highest stage,

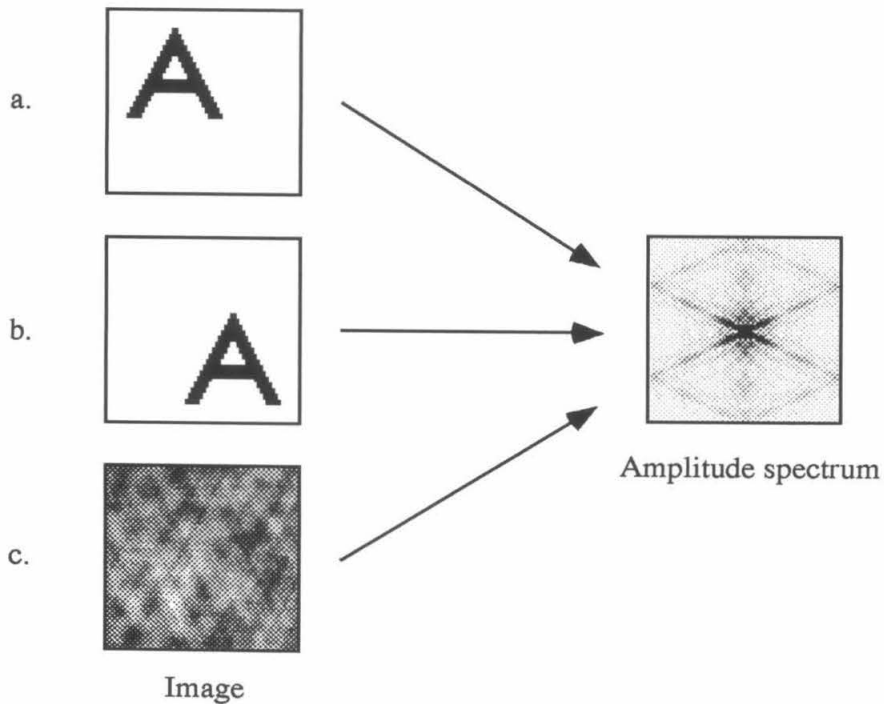


Figure 1.2: **Fourier transform method for producing invariant representations.**

The amplitude spectrum of the Fourier transform is invariant to shifts in an image. Thus, images *a* and *b* have equivalent amplitude spectra. A negative side-effect, though, is that non-sensical images such as *c* (produced by randomizing the phases of image *a*) also have the same amplitude spectrum.

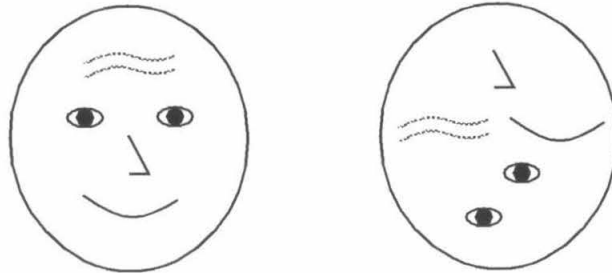


Figure 1.3: **Two objects with the same features but different spatial relationships are not equivalent.**

corresponding to inferotemporal cortex. Position invariance is achieved progressively at each stage by cells that simply summate over a local group of features within the layer below. Although this method does seem to meet with limited success in recognizing a variety of digits, it is unsatisfactory as an explanation for how we achieve position and size invariance. As with the Fourier transform, occlusions or other spurious detail will tend to confuse the system, since these things will be registered as features at the low-levels. Also, in this model all information about the position and size of an object are lost, whereas in human perception this information is readily retained. In addition, this network has been applied only to relatively small images where the object has already been centered and scaled more or less correctly. It is difficult to see how such a system would scale up for larger input arrays with multiple objects simultaneously present.<sup>1</sup> It also becomes problematic when one considers how complex objects, such as faces, would be distinguished from other objects containing the same features but with different spatial relationships, as illustrated in Figure 1.3.

LeCun et al. (1990) have successfully used a similar method to train an artificial

---

<sup>1</sup>In later work (Fukushima, 1987), an attentional model was proposed for dealing with an input array containing multiple or even overlapping objects. However, since there is no spatial coherence within the attentional beam, the system could just as easily combine features from different objects in vastly disparate parts if the visual field in order to construct any given object.

neural network to recognize handwritten digits in zipcodes. It should be noted, however, that although this network performs well on this specific task, it loses important information about the shape and style of the characters—information that we readily retain and use to advantage in interpreting other digits within the zip code.

## Remapping

A third major class of theories has been built upon the idea of remapping the representation of an object from one reference frame to another. This general idea seems to have been first proposed by Attneave (1954), who suggested that an efficient method for encoding a shape would be to first establish a center-point for the shape and then encode the spatial relationships among the features with respect to this point. He did not propose an explicit neural mechanism for accomplishing this, though.

Hinton (1981a) described a network model for transforming reference frames that utilized a set of *mapping units* to appropriately gate the connections between a set of input and output units depending on the shift, scale, and rotation of the transformation. A set of object units (grandmother cells) coded for spatial relationships among features represented on the output units. Given an image of an object at a particular position, size, and orientation, the network would then relax via a collective computation to represent the identity of the object and the reference frame transformation on separate sets of units. In later work where this network was actually simulated (Hinton and Lang, 1985), it was shown that the network could account for the illusory conjunction effects demonstrated by Triesman and Schmidt (1982). This effect resulted from the relaxation process getting confused with very short presentation times. Although this network seemed promising as a neural model for transforming reference frames, it was not developed beyond a simple connectionist network, and so it has little predictive value in neurobiology as it stands.



A somewhat different network for remapping reference frames has been described by Von der Malsburg and Bienenstock (1986). In contrast to Hinton’s network, this network does not use explicit gating units to change the connection strengths between input and output units. Instead, the strengths of the connections are based upon the *synchronization* between units. Input and output units fire in bursts, and the connections between those units with synchronized bursting are strengthened on a very short time scale. Although this network has been proposed as a neurobiological model for how connections may be changed dynamically, it is somewhat lacking in specific neurobiological substrates in the visual system. In addition, it is difficult to see how synchronicity alone could change the connections in the very specific and coordinated point-to-point fashion required to perform a reference frame transformation. The way that one input-output connection changes will somehow need to influence how other connections are changed, and synchronicity alone provides no natural avenue for doing this. (A related method has been successively employed in the “dynamic link architecture” of Buhmann et al. (1990) for face recognition, although without regard for the implementation details for changing connection strengths.)

A more neurobiologically detailed mechanism for transforming reference frames has been proposed by Anderson and Van Essen (1987). In their “shifter circuit” model, it was hypothesized that a set of control neurons (analogous to Hinton’s gating units) would dynamically shift the alignment of neural input and output arrays—without loss of spatial relationships—by multiplicative gating on dendrites. It was suggested that such a network could account for image stabilization in V1, and also that it could serve as an attentional mechanism for routing a region of interest in V1 onto a set of output nodes in a high-level area. However, it was not shown how such a circuit could rescale information, or more importantly, how the control neurons of the circuit could be automatically driven to position and scale the attentional window in *the image*.

### 1.3 Summary of the proposal

The proposal advanced in this thesis falls in the third class of models mentioned previously, based on remapping between reference frames, and builds upon the “shifter circuit” model of Anderson and Van Essen. In this thesis, I propose that the brain forms position- and scale-invariant representations of objects by an attentional process that selectively routes information from a region of interest in V1 into higher cortical areas.

The overall scheme is illustrated in Figure 1.4. It is assumed that relatively simple, local image features, such as orientation, texture, motion, etc., are extracted within the lower cortical areas at different scales and positions in parallel. The collection of features falling within the window of attention are then brought into a higher cortical area, with spatial relationships intact, for further analysis. Complex spatial relationships are then coded for only within the window of attention, with variations in position and scale removed. Information within the object-centered reference frame is represented with a fixed number of “sample nodes,” the consequence of which is that a small window of attention will capture information in the retina with higher resolution than will a large window.

It is proposed that the control neurons dynamically modify intracortical connection strengths via multiplicative couplings with the inputs at each stage of cortical processing. The control neurons themselves are driven in one of two modes: a pre-attentive mode, or a recognition mode. In the pre-attentive mode, the control neurons receive their input from a low-level “saliency map” representing the position and size of potentially interesting regions (i.e., potential objects) in a scene. In recognition mode, control neurons are driven by the interaction between top-down (memory) and bottom-up (retinal input) sources. The circuit is thus capable of running as an autonomous, closed-loop system, without the need for external commands.

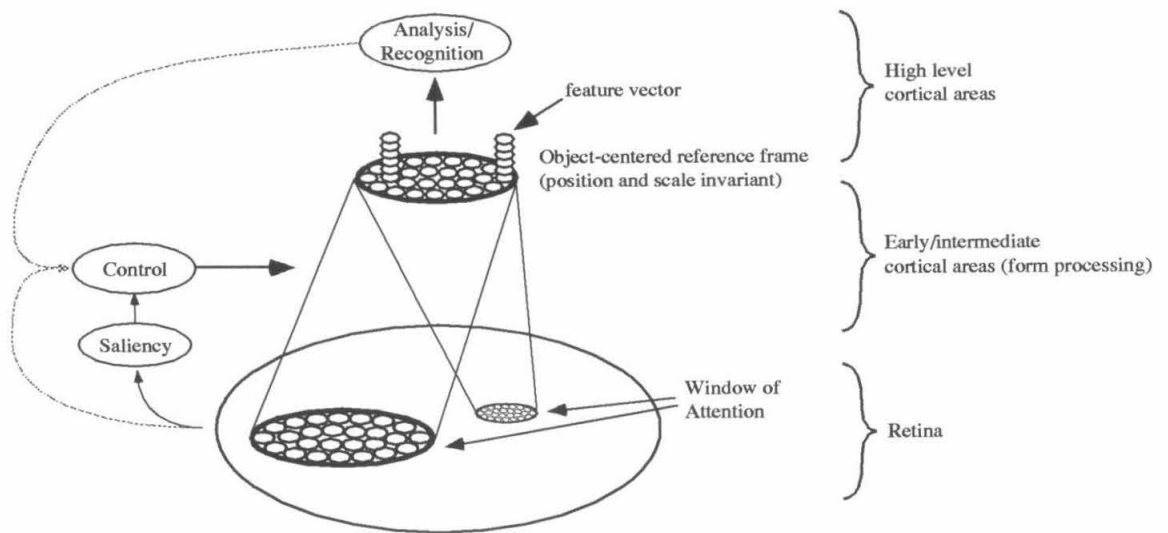


Figure 1.4: **Overview of the model.**

An attentional mechanism selectively routes information from a region of interest within the retinal image into higher cortical areas, with spatial relationships intact. Complex spatial relationships are then coded for only within the window of attention, with variations in position and scale removed. Control neurons dynamically modify intra-cortical connection strengths to set the position and size of the attentional window, and are driven bottom-up by a “saliency map” indicating interesting regions of the input to attend to.

The system as a whole represents the aspects of *what* and *where* independently, on different sets of neurons. *What* is represented by the contents of the window of attention (or some label attached to it in a higher area), and *where* is represented by the activities of the control neurons, which code for the position and size of the window of attention.

## 1.4 Motivations for the proposal

Two features of the above proposal that distinguish it from most previous models for forming invariant representations are 1) the preservation of spatial relationships within the window of attention all the way to high level areas, and 2) the explicit use of control neurons and switches for gating information flow through the visual cortex. It is worth stating from the outset a few of the insights from psychology, neurobiology, and computational complexity that have motivated these choices.

### Psychology

One of the more powerful psychological effects that support the idea of a high-level, spatial representation of an object are the so-called “reference frame effects.” Attneave (1965), Rock (1973), Hinton (1979), and Palmer (1983), among others, have shown that the reference frame that is imposed in viewing or imagining an object has a profound effect on the interpretation of its shape. For example, an isolated square rotated by 45 degrees may be interpreted equally well as a diamond shape or as a square standing on its corner. Which of these interpretations is chosen is affected drastically by the context in which the object is placed, as illustrated in Figure 1.5. Placing other items above and below the object imposes a vertically oriented reference frame, which leads to the interpretation of a diamond (Fig. 1.5a). On the other

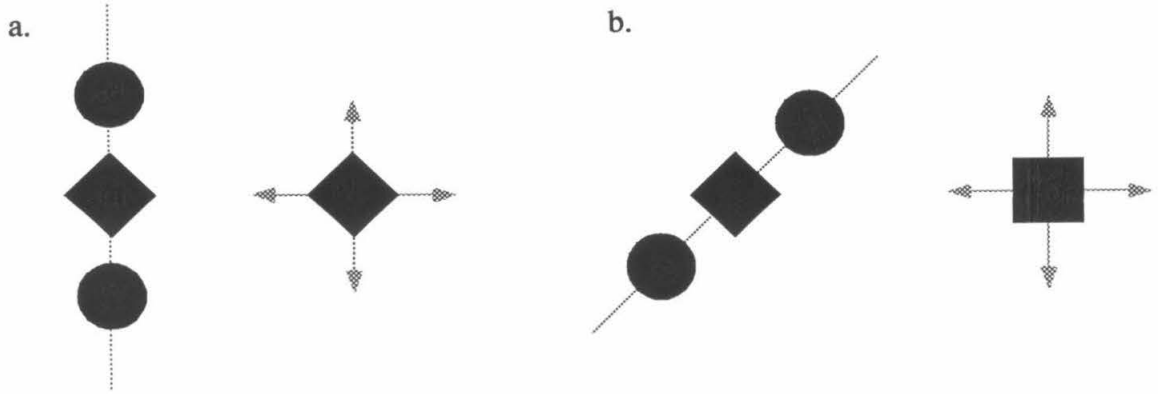


Figure 1.5: **Reference frame effects.**

A square rotated by  $45^\circ$  may be interpreted equally well as a diamond shape or as a square standing on its corner. *a*, Imposing a vertical reference frame results in the perception of a diamond, which can be explained if the diamond is stored internally within this frame of reference. *b*, Imposing a diagonal reference frame results in the perception of a square, which agrees with the internal reference frame shown.

hand, placing other items along the diagonal imposes a diagonally oriented reference frame, which yields the interpretation of a square (Fig. 1.5*b*). This effect can be explained if the diamond were to be stored within an internal reference frame that has its principal axes running through the corners, and the square were to be stored in an internal reference frame that has its principal axes running perpendicular to the edges of the square. The effect is not so easily explained, however, by theories based on a feature *Gemisch*.

Another insight from psychology that indicates that we may work with spatial image representations rather than a feature *Gemisch* are the results from mental imagery and shape comparison studies. For example, the experiments of Shepard and Metzler (1971) and Larsen and Bundesen (1978) have shown that the time required to compare two shapes differing in orientation and size is a linearly increasing function of the rotation angle or scaling factor required to transform one shape into the other.

These results would be consistent with a model in which we internally rotate and scale the representations of objects in order to match their shapes. In addition, Kosslyn and Schwartz (1978) have provided several telling experiments which indicate that mental images preserve the metric, spatial properties of objects as they are perceived. Thus, any model of perceptual processing would presumably need to preserve this information in higher-level processing in order to be truthful to human perception.

A number of psychophysical studies suggest not only that we work with high-level spatial representations, but also that we work with an attentional window that contains a fixed number of spatial elements, or “pixels,” for representing spatial information. For example, studies of spatial acuity (Toet et al., 1987) and recognition (Sperling et al., 1985; Campbell, 1985), suggest that once our window is set to a particular size, adding information beyond a certain critical resolution relative to the window size yields little or no incremental improvement in performance. In addition, the general Weber law effects observed in spatial frequency, or spatial interval, discrimination are consistent with the notion that we spread a spatial grid with a fixed number of divisions over the stimulus of interest, thereby limiting the accuracy of our spatial judgements to a proportion of the grid element spacing. There are also some interesting introspective observations one can make: For example, try forming a mental image of an object and then zoom in on a specific aspect of it. One will notice that the resolution with which this specific aspect is imagined is much greater than when the entire object is imagined as a whole; it seems nearly impossible to hold both the view of the entire object and a high-resolution view of a specific aspect simultaneously.

## Neurobiology

If one takes seriously the evidence from psychology, there needs to exist some way of preserving information about spatial relationships at higher levels of visual processing. A highly efficient and natural way to do this would be to encode the spatial relationships explicitly within a neural map. This way, local spatial relationships are encoded simply as local neural relationships: encoding the fact that one feature is to the right of another requires simply that the neuron representing that feature be situated “to the right” of the neuron representing the other feature within the neural substrate. Indeed, it is difficult to conceive of a scheme in which spatial relationships are not encoded in this way. Presumably, a tag of some sort would need to be attached to each feature stating how it is related to the others, and as yet there are no concrete neural models for how this might be accomplished.

If one were quick to judge the immediate neurobiological evidence, one might be inclined to believe that there is no retinotopic order in higher cortical areas because of the observations made to date in anesthetized animals. However, as we shall see later (Discussion, Section 4.1), interpreting the response of cortical cells in these high level areas will depend critically on the attentional state of the animal. Consequently then, there is insufficient evidence to paint a clear picture one way or the other as yet.

## Computational complexity

From a computational viewpoint, a prime motivation for this proposal is that *what* and *where* are represented independently, resulting in an efficient use of computational resources. This is illustrated in Figure 1.6. If a fixed number of neurons are available to represent *what* and *where*, then a system that did not separate these variables would have these resources distributed over the joint space of *what* and *where* (Fig. 1.6a).

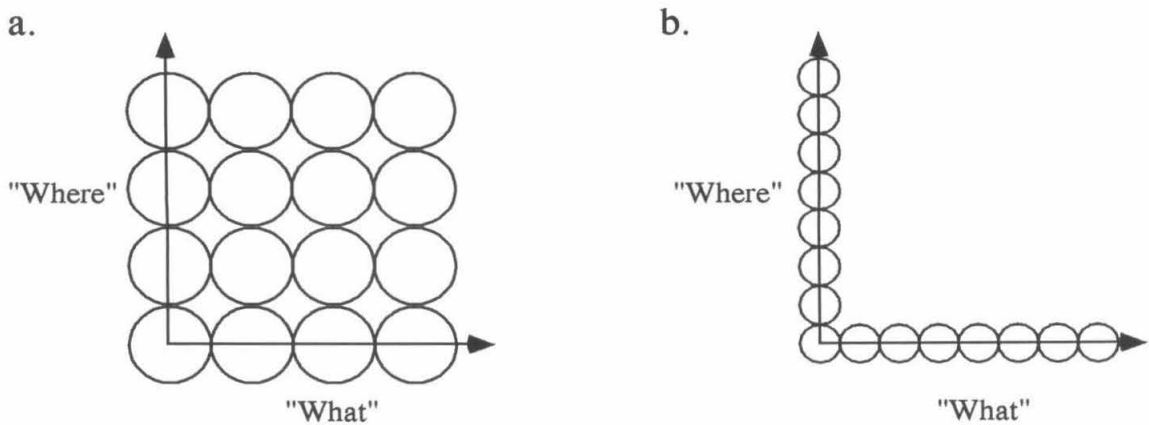


Figure 1.6: **Encoding *what* and *where* with limited neural resources.** *a*, 16 neurons code for the joint space of *what* and *where* by spreading their receptive field's over this space. *b*, The same number of neurons are spread along the independent dimensions of *what* and *where*, resulting in a higher resolution representation of these dimensions.

By contrast, a system that represents *what* and *where* independently allows the same number of resources to code for a greater number of objects and their positions and sizes (Fig. 1.6*b*). The price we pay for this savings, however, is that an attentional mechanism is now required in order to ensure that for an image containing multiple objects, a unique *what* is paired with a unique *where* at any instant in time. It is thus necessary to introduce switching elements for gating information from one locus at a time.

This model also makes computational sense in that it brings visual information into a standard format that is more amenable to a general purpose analysis. Consider for example the problem of examining random objects, such as a rock, a crumpled piece of newspaper, or an alphanumeric character in a very strange font. It is hard to imagine how any of the feature *Gemisch* theories described above could represent such objects in a meaningful way that would, for example, enable one to draw the object, or to describe general metric properties of its shape.



## 1.5 Caveats (what’s not in this thesis)

No model can address all issues as once. Out of necessity, we will be ignoring some issues in order to focus on the problem of interest. For example, while there are many aspects to visual attention—color, motion, etc.—we will concentrate here on the spatial aspect of attention. Furthermore, we will be considering mainly the automatic (involuntary) component of attention, as opposed to the consciously controlled (voluntary) form of attention.

We will also not be considering the nature of the representation of visual information, such as the orientation selective cells and other feature processing known to exist in the visual cortex. It will be assumed that routing can be considered somewhat independently of these issues, and that the features can be bundled together into what we will denote as a “sample node.”

Also, the main emphasis of this model is to address how invariant representations of objects are formed, not how objects are recognized beyond this point. The problem of visual recognition per se is an extremely difficult and involved problem beyond the scope of this thesis.

## 1.6 What’s in this thesis

We begin in Chapter 2 with a discussion of the basic principles of routing, and we develop a model routing circuit that is capable of autonomously attending and recognizing objects over a wide range of positions and sizes in its input array. In Chapter 3, we discuss the proposed neurobiological substrates for routing and neural mechanisms for gating information flow in the cortex. Chapter 4 then discusses the predictions of the model, comparisons to other models of attention and invariant object representation, generalizations of the model, and unresolved issues. Conclusions are presented

in Chapter 5.

## Chapter 2

# Dynamic Routing Circuits

In this chapter we shall derive a model routing circuit that autonomously forms position- and size-invariant representations of objects in an image. The neurobiological substrates for this model will be discussed in detail in the next chapter, so we will consider here only the major neurobiological factors that significantly constrain our design options. Initially, we will put aside the issue of “what controls the control neurons” and concentrate on developing a routing circuit that satisfies the major neurobiological constraints. Then, in Sections 2.3 and 2.4, we describe how the control neurons may be driven autonomously to set the position and size of the window of attention, and how the routing circuit may be coupled to an associative memory in order to further guide the attentional window during recognition.

### 2.1 Neurobiological constraints

There are basically four key observations from neurobiology that will influence our design of a neural routing circuit.

1. Limited fan-in: Cortical neurons are typically limited to about  $10^3$ - $10^4$  total inputs (Cherniak, 1990; Douglas and Martin, 1990a). Thus, if the total input-output convergence of the routing circuit exceeds this amount, the circuit must be broken into multiple stages. Control must then be coordinated among these stages.
2. Limited fan-out: We assume that neural fan-out is also limited to about  $10^3$ - $10^4$  (on average), and so each control neuron can modify at most this many synapses. Since there will likely be many more than 1000 synapses to modify in order to realize a given position and size of the window of attention, control will need to be modularized so that each control neuron modifies a local group of synapses. Multiple control neurons must then cooperate and act together in order to establish a global position and size of the window of attention.
3. Multiscale input representation: Cells in visual cortex are tuned to different spatial-frequencies, with typical bandwidth hovering around 1-1.5 octaves (De Valois et al., 1982). High frequency cells will integrate information over a small region of visual space, while low frequency cells will integrate information over a large region. This type of representation can be incorporated into the routing circuit advantageously by selectively routing from high or low frequency cells depending on whether the window is small or large, respectively. This way, much of the image blurring required for rescaling can be accomplished by switching between filters, rather than requiring the routing circuit to blur over a wide dynamic range.
4. Logarithmic spatial sampling: The spacing between retinal ganglion cells increases linearly with eccentricity, resulting in an essentially logarithmic transformation of visual space within the cortex. This will significantly affect the

routing and control architecture, because the size of the attentional window will depend on its eccentricity.

In this chapter, we will be working with a scaled-down routing circuit, so many of the constraints will be scaled-down proportionally. Then, in Chapter 3, we will take the concepts developed here and extend them to very large routing circuits on the scale of those proposed to exist in the brain.

## 2.2 Routing circuit architecture

In this section, we shall build up in steps a routing circuit architecture that meets the above constraints. We first introduce some of the basic issues of routing using the smallest, simplest circuit possible. This circuit will then be progressively modified in order to accommodate greater input-output convergence, a multi-resolution input representation, and an approximately logarithmic input sampling lattice.

### The simplest possible routing circuit

Figure 2.1 illustrates a simple, one-dimensional routing circuit composed of nine input nodes, five output nodes, and set of control neurons. Each output node receives five inputs via dynamically modifiable links to the input layer. The strength of these links are determined by the control neurons, which make multiplicative couplings with the inputs.

More precisely, the activities of the output nodes,  $I_i^{out}$ , are computed from the activities of the input nodes,  $I_j^{in}$  and control neurons,  $c_k$ , according to

$$I_i^{out} = \sum_j \sum_k c_k \Gamma_{ijk} I_j^{in} \quad (2.1)$$

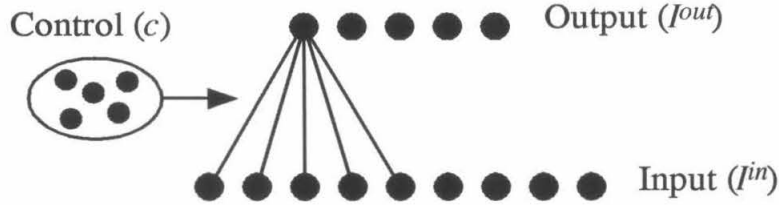


Figure 2.1: **A simple routing circuit.**

Each output node receives five inputs via dynamically modifiable links to the input layer. Shown are the links for the leftmost node only. The links for the other nodes are the same, merely shifted. The strength of these links are determined by the control neurons, which make multiplicative couplings with the inputs.

where the term  $\Gamma_{ijk}$  determines how the control neurons are coupled with the inputs. The activities of both the inputs and the control neurons are assumed to be analog values between 0 and 1. It is sometimes helpful to think of Equation 2.1 in terms of two equations

$$I_i^{out} = \sum_j w_{ij} I_j^{in} \quad (2.2)$$

$$w_{ij} = \sum_k c_k \Gamma_{ijk} \quad (2.3)$$

where the term  $w_{ij}$  denotes the effective strength of the link from node  $j$  in the input to node  $i$  in the output. In this case, the term  $\Gamma_{ijk}$  can thus be thought of as the amount by which the  $k^{\text{th}}$  control neuron modulates  $w_{ij}$ . In general,  $\Gamma$  will be very sparse (i.e.,  $\Gamma_{ijk} = 0$  for the vast majority of combinations of  $i$ ,  $j$ , and  $k$ ).

In order to understand how the control neurons are to be coupled to the inputs for realizing translations and scalings, it is helpful to visualize the routing circuit in “connection space,” as shown in Figure 2.2a. Here, the horizontal axis represents the nodes constituting the input layer of the network and the vertical axis represents the nodes constituting the output layer. An  $\times$  at coordinate  $(j, i)$  in connection space

denotes that a physical connection exists from node  $j$  in the input to node  $i$  in the output; the lack of an  $\times$  at  $(j, i)$  implies that no connection pathway exists between those nodes. Note that for a 2D routing circuit the connection matrix would require four dimensions to display. We will use the 1D routing circuit for ease of illustration, but the concepts developed here are readily extendible to 2D.

If the window of attention is to be of a certain position and size, then the strength of each connection,  $w_{ij}$ , needs to be set appropriately. Figure 2.2*b* shows how this would look in connection space for an attentional window centered within the input array with a scale factor of one. The stippled area represents those connections that are enabled ( $w_{ij} > 0$ ); the remaining connections are effectively disabled by mechanisms discussed below. If the window of attention is to shift to the left or right, then the band of enabled connections must translate across the connection matrix. Changing the size of the window of attention corresponds to tilting the band of open connections, as shown in Figure 2.2*c*. Note that the band of open connections must also be widened as it is tilted (corresponding to blur); otherwise aliasing would occur, leading to spurious patterns in the output representation (Fig. 2.2*d*).

By viewing the routing circuit in this way, it can be seen that the problem of setting the position, size, and blur of the window of attention amounts to one of generating the proper patterns of active synapses in connection space. How this is to be accomplished by the control neurons depends on how they are connected to the feedforward synapses of the routing circuit. One possible scenario is for each control neuron to modulate the strength of a single physical connection  $(j, i)$ , as illustrated in Figure 2.3*a*. If a given control neuron were “on” ( $c_k > 0$ ) then its corresponding connection would be enabled, and if it were off ( $c_k \approx 0$ ) then the connection would be disabled. Nearly any remapping could then be accomplished by simply activating the control neurons corresponding to the connections we wish to enable. However, this scheme would require an enormous number of control neurons

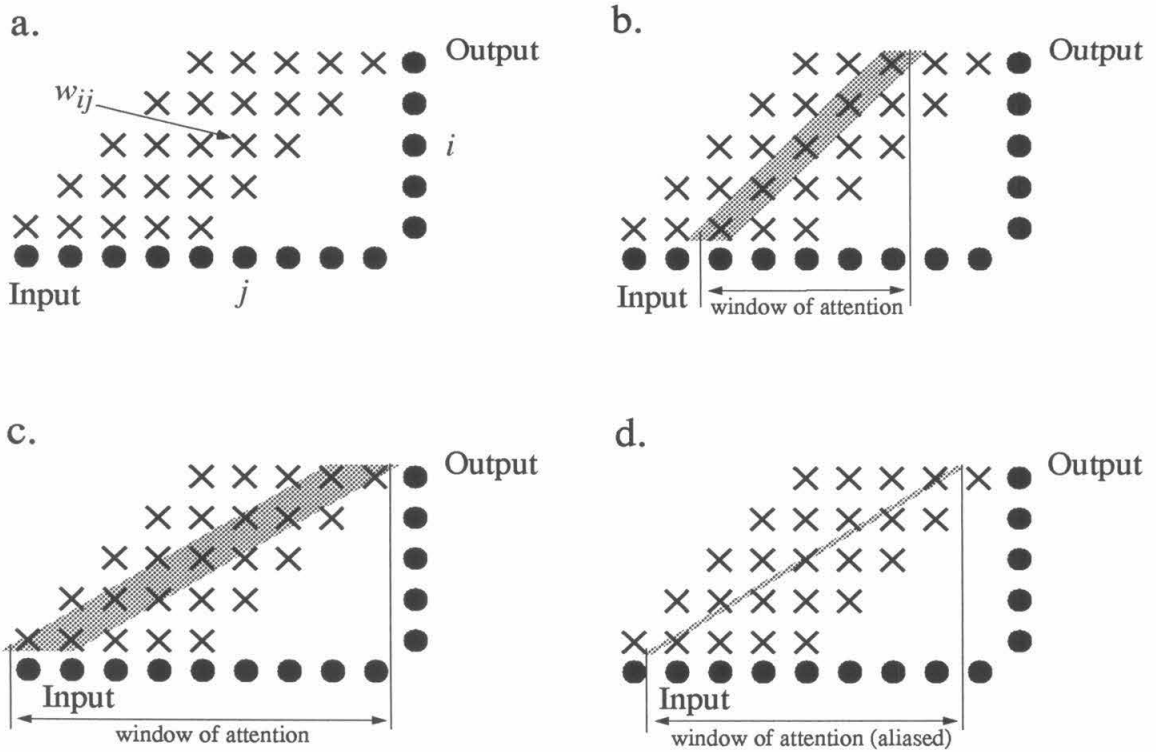


Figure 2.2: An illustration of “connection space.”

The input and output nodes in Figure 2.1 are represented on the horizontal and vertical axis, respectively. *a*, Each  $\times$  denotes a physical connection from an input node to an output node. The effective strength of the connection from node  $j$  in the input to node  $i$  in the output is denoted  $w_{ij}$ . *b, c*, The stippled region indicates those connections that need to be enabled ( $w_{ij} > 0$ ) in order to map the region within the window of attention onto the output nodes. *d*, If the width of the enabled region is too small, then aliasing will result; an exaggerated case is illustrated here (i.e., some output nodes will be lacking any input, leading to spurious patterns in the output).



for a scaled-up system. Also, since the set of remappings we wish to accomplish (translations and scalings) is but a minute fraction of all possible remappings, this scheme would arguably constitute a waste of computational resources.

Another possibility would be for the control neurons to gate connections globally so that each unit is responsible for effecting a single position and scale of the window of attention, as shown in Figure 2.3*b*. While this solution would be acceptable for a small circuit of this size, it has the potential disadvantage of requiring a large fan-out for each control neuron in a scaled-up system, which would render the circuit neurobiologically implausible.

A solution that attempts to simultaneously minimize both the number of control neurons and the fan-out required would have each control neuron modulate a local group of synapses—or a *control block* in connection space (Fig. 2.3*c*). The problem of forming the desired patterns in connection space then becomes an approximation problem, in which the control blocks form the basis functions and the activations of the corresponding control neurons form the coefficients. The connection strengths  $w_{ij}$  would then be determined according to

$$w_{ij} = \sum_k c_k \Psi_k(j, i) \quad (2.4)$$

where the function  $\Psi_k(j, i)$  specifies the shape of the  $k^{\text{th}}$  control block in connection space (note that this is merely an alternate form of expressing Equation 2.3). In order to facilitate their ability to approximate patterns in connection space, the control blocks should not have sharp boundaries; rather, they should have a Gaussian-like taper and overlap one another somewhat. For example, the control blocks shown in Figure 2.3*c* may be defined by

$$\Psi_k(j, i) = \exp\left[-\frac{(i - j - m)^2}{2\sigma^2} - \frac{(i - nW)^2}{W^2/2}\right] \quad (2.5)$$

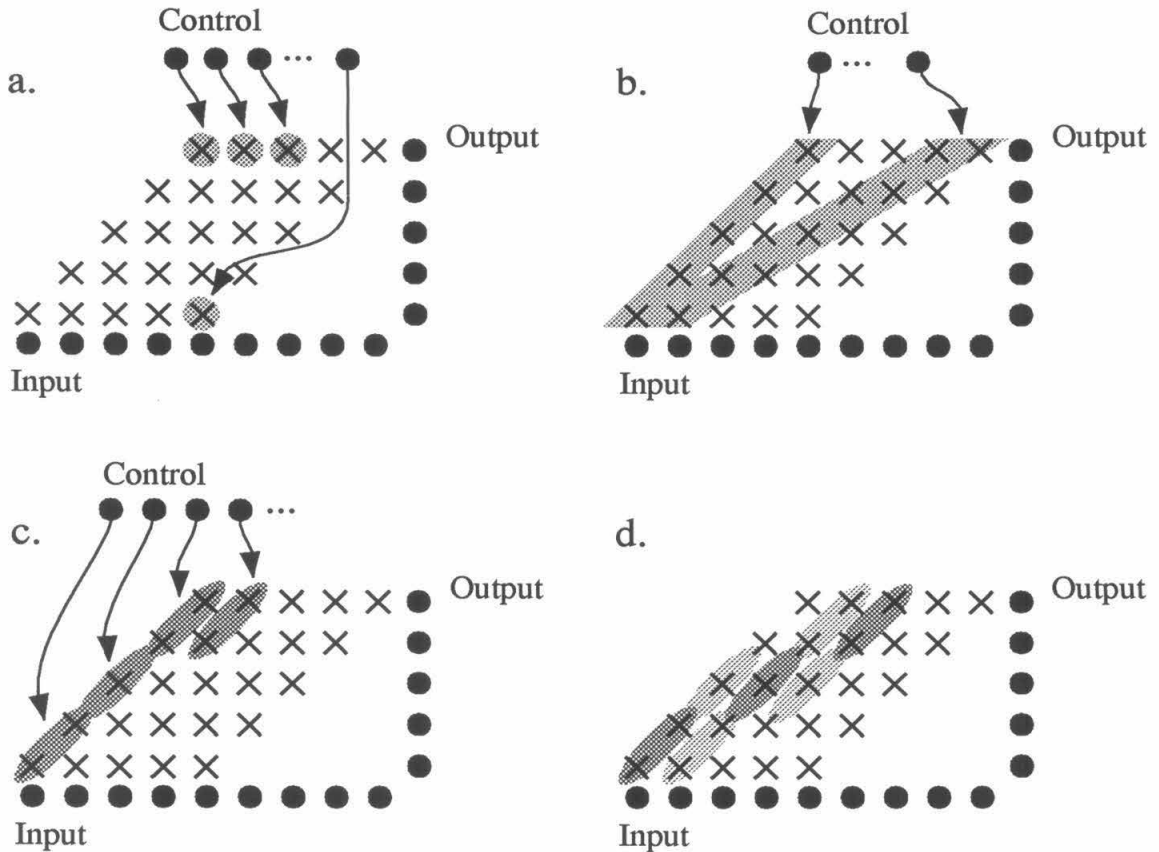


Figure 2.3: Some possible control scenarios.

*a*, Each control neuron modulates the strength of a single connection. *b*, Each control neuron modulates the strength of a large number of connections in order to effect a global position and scale of the window of attention. *c*, Each control neuron modulates a local group of connections, or a "control block." *d*, Approximating a desired position and scale of the window of attention using control blocks.

$$k = m N_B + n \quad (2.6)$$

where  $m$  denotes the index of translation (i.e., which diagonal in Fig. 2.3c),  $n$  is the index of the control block within a diagonal,  $N_B$  is the number of control blocks within a diagonal, and  $W$  is the spacing of the control blocks along the diagonal ( $W \approx \frac{\# \text{ of outputs}}{N_B}$ ). Shaping the control blocks as in Figure 2.3c would be most optimal for realizing translations, but could also be used to approximate scalings as well, as shown in Figure 2.3d. It may well be possible to optimize the shape of the control blocks using appropriate learning algorithms, but the strategy illustrated here will suffice for our immediate purposes.

## Accommodating larger convergence

If we wish to increase the size of the input array, then the fan-in on each output node must also increase if we wish the circuit to be capable of mapping any portion of the input onto the output. For example, increasing the number of input nodes to 33 would require a fan-in of 29 inputs on each output node. While this may be feasible for such a small model neural network, it will become neurobiologically implausible for input sizes on the order of 300,000 sample nodes, such as in the primate visual cortex.

In order to accommodate a larger input-output convergence without appreciably increasing the fan-in, it will be necessary to break the routing circuit into several stages. One possible design for a multi-stage circuit, initially conceived by Anderson (1993), is illustrated in in Figure 2.4. This circuit accommodates a total input-output convergence of 33:5 while maintaining a fan-in of 5 inputs on each node. An important feature of this circuit is that the spacing between inputs doubles at each stage of the circuit, which minimizes the number of stages required to meet

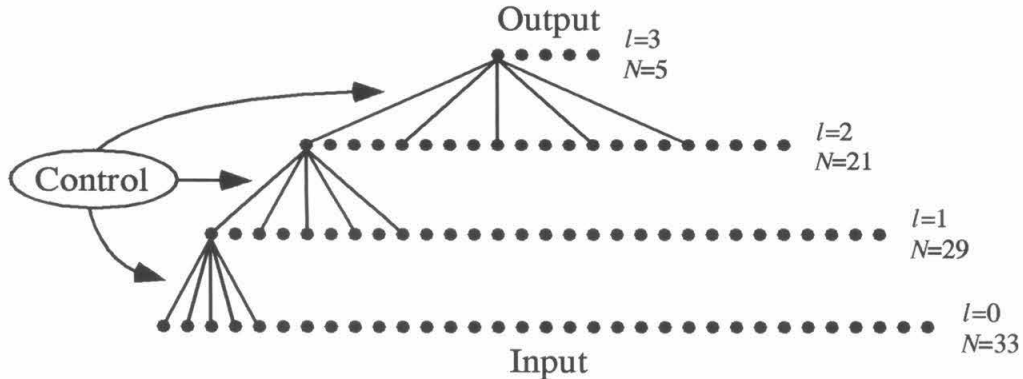


Figure 2.4: **A multistage dynamic routing circuit.**

As before, the connections are shown for the leftmost node in each layer.  $N$  denotes the number of nodes within each layer, and  $l$  denotes the layer number. The control neurons modulate connections at each interface to set the position and size of the window of attention. Because of the weight subsampling at higher stages, the lowest stages will be best suited for small, fine-scale adjustments to the position and size of the attentional window, while the upper layers will be able to handle large, coarse scale adjustments only (in chunks), as illustrated in the connection space diagram of Figure 2.5.

the desired convergence and fan-in constraint. A consequence of this architecture is that the lower stages of the circuit are best suited for performing small shifts and scale changes, whereas the higher stages are better suited for performing macro-shifts and scale changes only, because of the subsampling in connection space. Figure 2.5 illustrates the connection space for two possible settings of the position and size of the window of attention.

Since we are now working with a multistage circuit, some changes in notation and nomenclature are in order. We denote the nodes of layer  $l$  as  $I_i^l$ , and the weight from node  $j$  of level  $l$  to node  $i$  of level  $l+1$  as  $w_{ij}^l$ . Equations 2.2 and 2.3 are thus rewritten as

$$I_i^{l+1} = \sum_j w_{ij}^l I_j^l \quad (2.7)$$

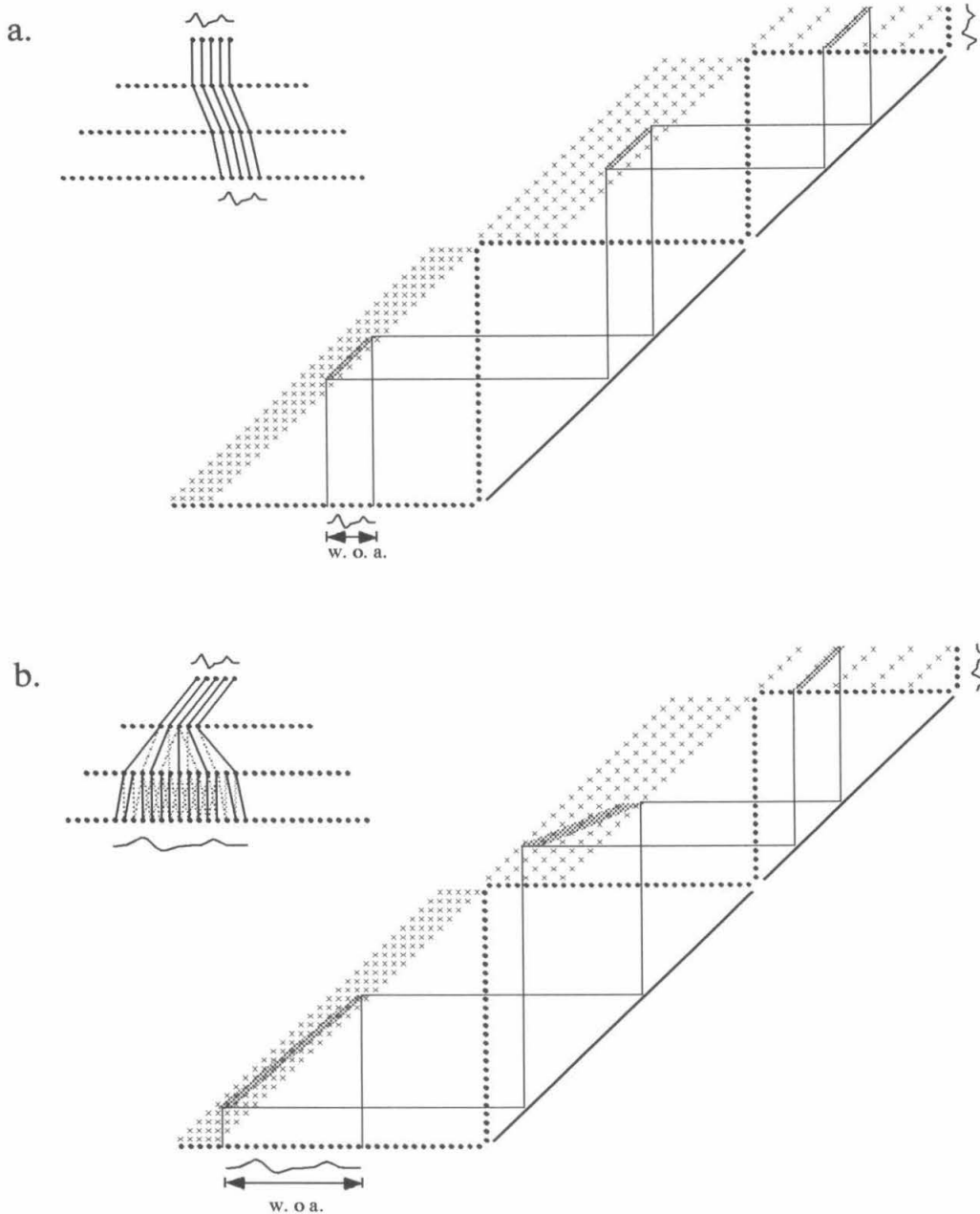


Figure 2.5: **Connection space for the multistage dynamic routing circuit.** *a*, For the smallest window size (5 nodes wide), connections are set 1:1 at each stage. Translating the window by a few steps could be accomplished at the first stage, while larger shifts will be accomplished at higher stages. *b*, For a larger window size, connections are set for a net convergence from input to output. An even larger window could be accomplished by tilting the band of connections at the top stage.

$$w_{ij}^l = \sum_k \Gamma_{ijk}^l c_k^l. \quad (2.8)$$

Two adjacent layers and the control neurons that modulate the connections between them shall be referred to as a *stage*. Thus, stage  $l$  refers to the interface between layer  $l$  and  $l + 1$ .

A consequence of the multistage architecture as specified in Equations 2.7 and 2.8 is that the nodes in the intermediate layers will be active only if the control neurons modulating the connections to them are enabled. Thus, those regions outside the window of attention in the intermediate layers ( $l > 0$ ) will not be active. For reasons that will become clear later, it will be desirable for these regions outside the window of attention in the intermediate layers to reflect the activity in the input in the “all-connections-open” state. This may be accomplished by one of two methods: 1) having the control neurons maintain a relatively low tonic activation ( $c_k \approx 0.1$ ), or 2) letting each connection,  $w_{ij}$ , have a default, resting value in the absence of any activity from the control neurons—i.e.,

$$w_{ij}^l = \sum_k \Gamma_{ijk}^l c_k^l + w_{rest}, \quad (2.9)$$

where  $w_{rest}$  is the “resting value” of the synapse with all control neurons off. In either of these cases, it would probably be desirable to renormalize the activities on the nodes in the intermediate layers so that the output is kept within a certain range, e.g.,  $[0,1]$ . This could conceivably be accomplished by local gain control networks, such as proposed by Grossberg (1976).

## Multiscale input representation

The circuit of Figure 2.4 requires that the control neurons as a whole be capable of dynamically blurring and scaling spatial information over a wide dynamic range in order to accommodate window sizes ranging from very small (5 nodes across) to very

large (33 nodes across). The complexity of control could be reduced somewhat if the representation at the input were to contain nodes that were preset, or hardwired, to integrate over different spatial scales. Much of the image smoothing could then be accomplished by switching between a set of fixed filters, rather than blurring dynamically: large objects would be attended by switching to the low resolution input array, while small objects would be attended by switching to the high resolution input array.

The routing circuit can be modified to accommodate a multiscale input representation, as shown in Figure 2.6. Here, the input is represented on three different sampling lattices separated in resolution by octaves. (This is essentially a Gaussian pyramid, as in Burt and Adelson (1983).) There are now three separate routing streams corresponding to each input lattice. Since the low resolution nodes are spaced more sparsely, the number of input nodes at these scales is fewer, and hence fewer intervening layers are required between the input and output in order to maintain the fan-in constraint. Each routing circuit for a particular scale performs translation over the entire range of its input, plus scaling within a factor of two. Scale changes greater than a factor of two are accomplished by switching between routing streams at the top stage of the circuit.

## Logarithmic spatial sampling

The final modification we must make to the routing circuit is to incorporate a logarithmic sampling lattice at the input. Koenderink and van Doorn (1978) have proposed a piecewise approximation to a logarithmic sampling lattice that meshes quite naturally with the circuit just described. In their *stack model*, illustrated in Figure 2.7, the image is represented by a stack of sampling lattices at different resolutions as before, except now each level of the stack comprises an equal number of sample nodes. The

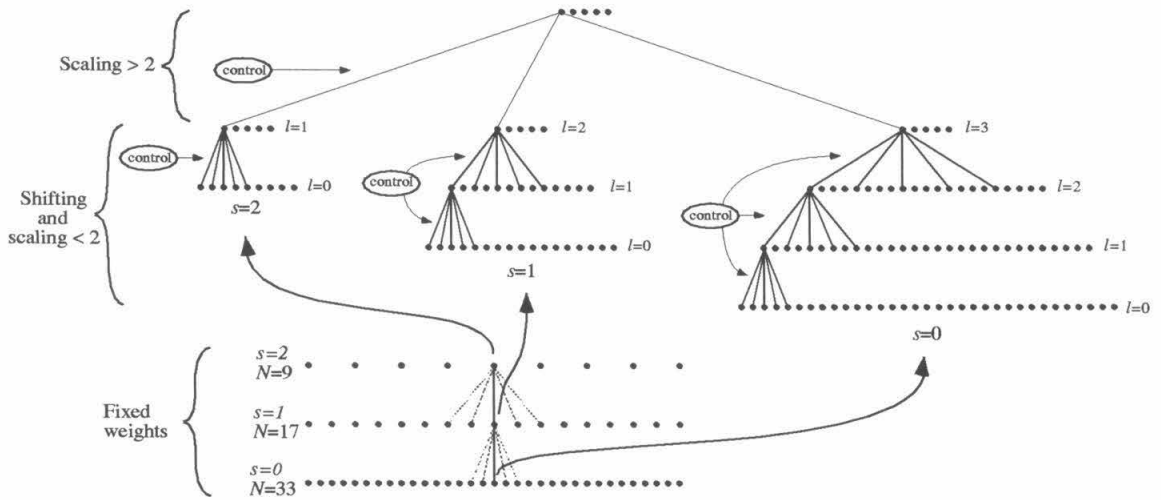


Figure 2.6: **Dynamic routing circuit with a multiscale input representation.** The input is decomposed into three different spatial scales via a set of fixed weights.  $s$  denotes the scale (resolution) of each sampling lattice, and  $N$  denotes the number nodes. Each lattice serves as input to a separate routing stream. Each routing circuit translates the window of attention within its input array, and rescales the window within a factor of two. Scale changes greater than a factor of two are accomplished by switching between routing streams at the top stage of the circuit.



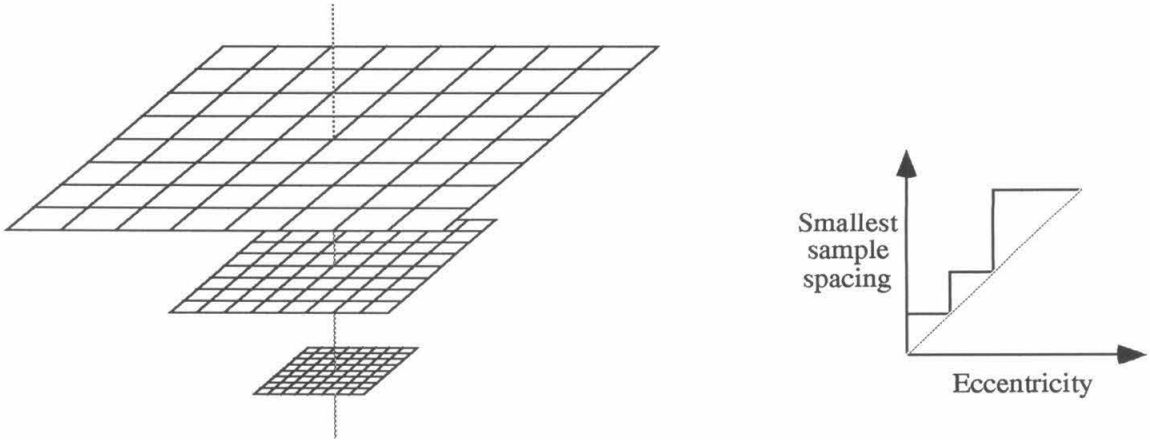


Figure 2.7: **The Koenderink “stack.”**

*a*, The input is represented by a stack of sampling lattices, each level covering a progressively greater extent, at lower resolution, than the level below it. *b*, Sample spacing as a function of eccentricity approximates the linear relationship found in primate vision.

consequence is that each level of the stack covers a progressively larger extent of visual space than the adjacent level below. When combined, the different levels of the stack provide a multi-resolution representation of the input image and also approximate the linear dependence of sample spacing on eccentricity found in the retina.

Figure 2.8 illustrates how the stack model can be incorporated with the routing circuit. In the particular implementation shown, the stack is composed of three levels, with the resolution changing by a factor of two between successive levels. The total number of nodes in each level of the stack is 29. As before, each level serves as input to a separate routing stream that translates the window of attention within the range its input and rescales within a factor of two. The final stage of the circuit switches between scales by selecting among the outputs of the different routing streams for each scale.

Each routing circuit for a particular scale is composed of three layers, or two stages, in order to accommodate the convergence from 29 input nodes to 5 output

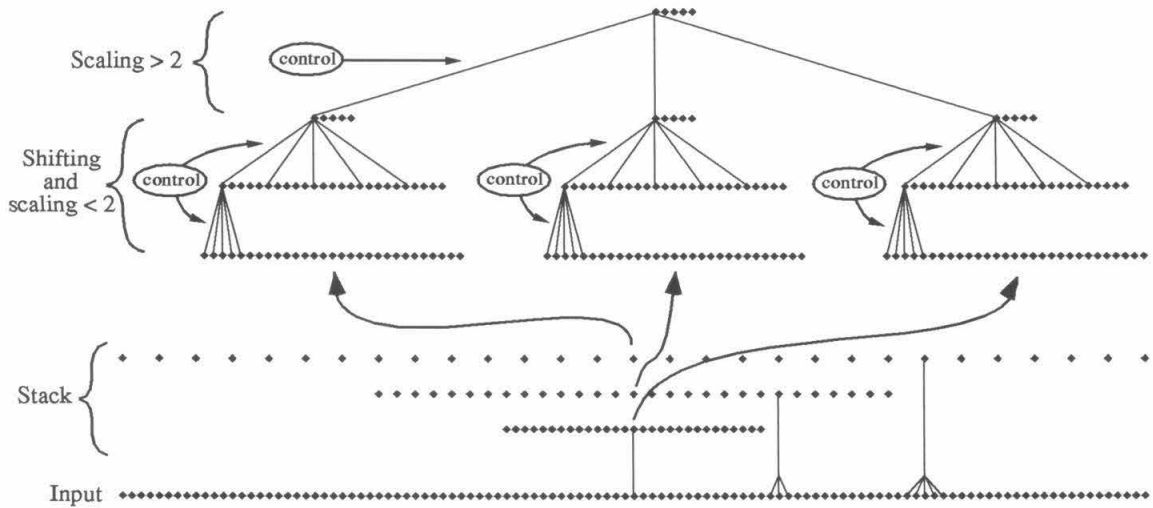


Figure 2.8: **Dynamic routing circuit with a “stack” input representation.** Each level of the stack serves as input to a separate routing stream that translates the window of attention within the range its input and rescales within a factor of two. The final stage of the circuit switches between scales by selecting among the outputs of the different routing streams for each scale.

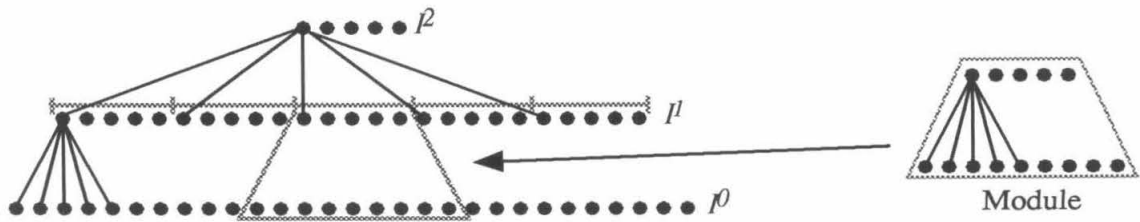


Figure 2.9: **Routing circuit for a single scale within the stack circuit.** It is helpful to think of the first stage as composed of modules of  $9 \rightarrow 5$  routing circuits, as in Figure 2.1, with the second stage switching between modules.

nodes. The first stage performs micro-shifting (in steps of a single node) in addition to scaling less than a factor two. The second stage performs macro-shifting in steps of 5 nodes. It may be helpful to think of the first stage as composed of modules of  $9 \rightarrow 5$  routing circuits (as in Fig. 2.1), with the second stage switching between modules, as in Figure 2.9. Note that this circuit has no redundant paths, in that there is one and only one path for routing a 5-9 node window from the input to the output. This is in contrast to the previous multi-stage routing circuit (Fig. 2.4) which had many redundant paths for the smallest window size. This, however, was a consequence of having to perform both shifting and scaling over a large range within the same circuit. Now that shifting has been separated from scaling for the most part, there is no need for this redundancy. It may in fact be desirable to include redundant paths for robustness to damage, which could be accomplished quite straightforwardly by overlapping the modules in the first stage somewhat. This would leave the first stage unaffected, but would increase the fan-in in the second stage, which may require an additional stage. In a redundant circuit, a cost function for the different paths should probably be included in order to ensure a uniquely optimum path for each window of attention.

An alternative means for arranging the routing circuit would be to perform only

shifts within the lower stages, leaving the top stage to perform rescalings less than a factor of two. However, this would necessitate routing up an image twice the size to the top stage. On the other hand, rotations and warps would best be performed at the top-stage, and since these operations will inevitably involve a moderate amount of rescaling, it would probably be desirable to route up an image slightly larger than the attentional window size to the top stage.

## **2.3 Autonomous control I: single stage**

Up to now we have described an essentially “open loop” routing circuit. That is, given a desired position and size for the window of attention, one could manually set the activity of the control neurons of the network so that the image within the window is remapped onto the output nodes of the network. We now describe how the network may be autonomously controlled when provided only with visual input and no external commands beyond the initial task specification. We begin in this section by describing the autonomous control for a single-stage routing circuit composed only of an input layer and an output layer. The next section will then extend the concepts developed here to a circuit composed of multiple stages.

### **System objective**

The purpose of attention in our model is to focus the neural resources for recognition on a specific region, or object, within a scene. Thus, it would make sense for the attentional window to be automatically guided to salient, or potentially informative areas of the visual input. Salient areas can often be defined on the basis of relatively low-level cues—such as pop-out due to motion, depth, texture, or color (e.g., Koch and Ullman, 1985; Anderson et al., 1985). Here, we utilize a very simple measure of

saliency based on luminance pop-out in which attention is attracted to “blobs” in a low-pass filtered version of a scene. (A blob may be defined simply as a contiguous cluster of activity within an image.) In reality, attention can also be directed via voluntary or cognitive influences, but these are not incorporated into the current model.

The following “algorithm” is proposed as a simple but useful strategy for an autonomous visual system (see Fig. 2.10):

1. Form a low-pass filtered version of the scene so that objects are blurred into blobs.
2. Select one of the blobs from the low-pass image—whichever is brightest or largest—and set the position and size of the window of attention to match the position and size of the blob.
3. Feed the high-resolution contents of the window of attention to an associative memory for recognition.
4. If a match with one of the memories is close enough (by some as yet unspecified criterion), then consider the object to have been recognized; note its identity, location, and size in the scene. If there is not a good match, then consider the object to be unknown; either learn it or disregard it.
5. Now inhibit this part of the scene from being attended and go to step 2 (find the next most salient blob).

The following three subsections describe the details for carrying out steps 2, 3, and 5. Step 1 is trivial, whereas step 4 is a high-level problem beyond the scope of this thesis (cf. Mumford, 1992; Carpenter and Grossberg, 1987a; Hinton, 1981b).

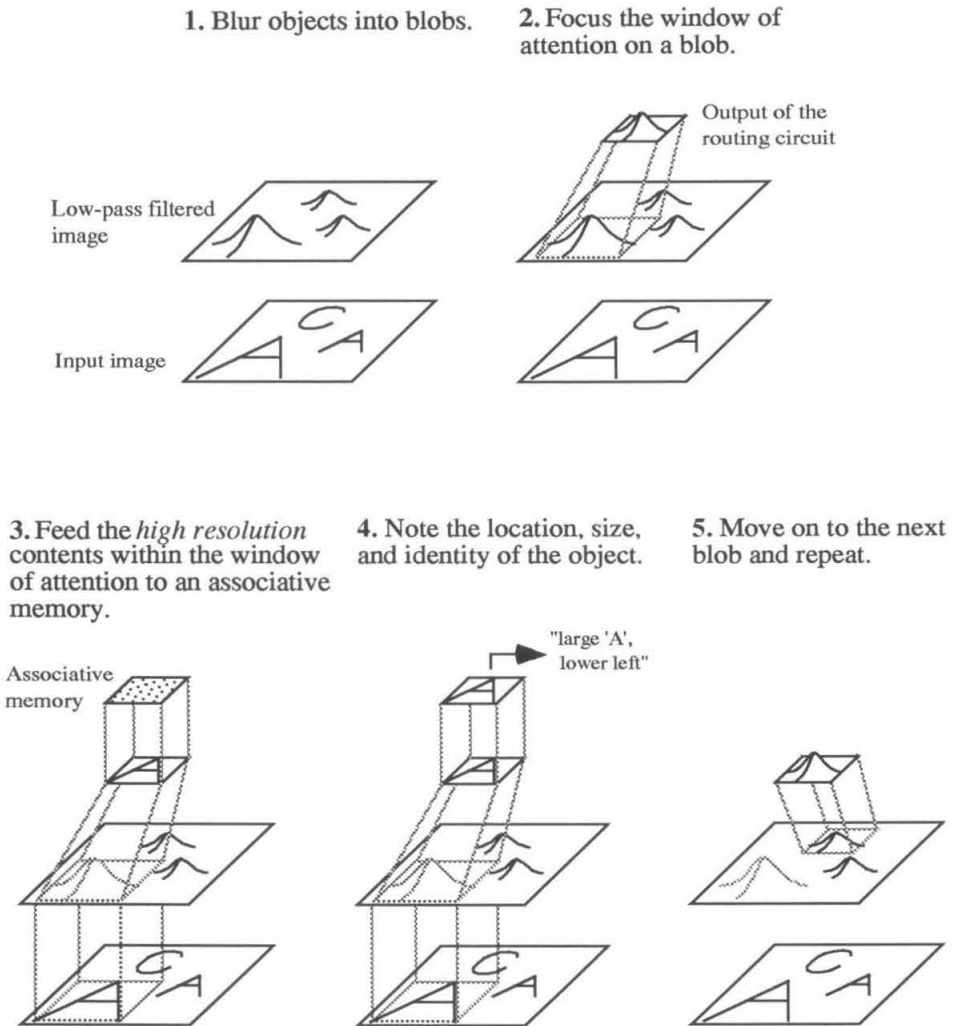


Figure 2.10: **A simple attentional strategy for an autonomous visual system.** Objects are preattentively segmented via lowpass filtering. Once an object has been localized, the contents of the window of attention are fed to an associative memory for recognition. This process is then repeated *ad infinitum*, or until all interesting locations have been attended.

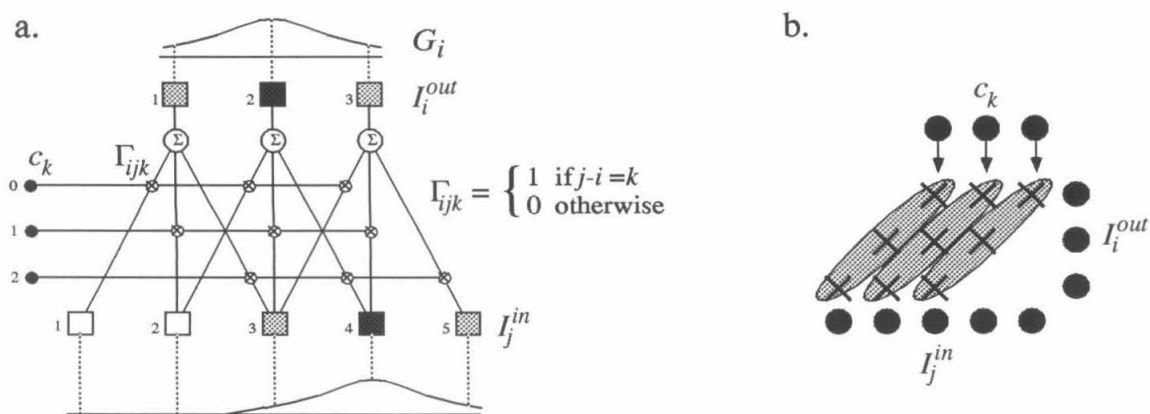


Figure 2.11: **A single-stage routing circuit with a Gaussian blob presented to the input units.**

a, Each control neuron corresponds to a different position of the window of attention: left ( $c_0$ ), center ( $c_1$ ), or right ( $c_2$ ). For example, in order to accomplish the remapping shown, the values on the control neurons should be  $c_2 = 1$  and  $c_0 = c_1 = 0$ . b, The connection-space diagram for the circuit.

## Focusing attention on a blob

We begin by formulating a solution for a simple single-stage routing circuit with one or more Gaussian blobs presented to the input units, as shown in Figure 2.11. In this circuit the  $\Gamma_{ijk}$  have been set so that each control neuron  $c_k$  corresponds to a global position of the window of attention, but in general this need not be the case.

In order to focus the window of attention on a blob in the input, the network's "goal" will be to fill the output units with a blob while maintaining a topographic correspondence between the input and output (Fig. 2.10, step 2). Since the dynamic variables in this network are the  $c_k$ , we need to formulate an equation governing the dynamics of  $c_k$  that accomplishes this objective. We can accomplish the first part of the objective by letting  $c_k$  follow the gradient of an objective function,  $E_{blob}$ , that provides a measure of how well a blob is focused on the output units. One possible

choice for  $E_{blob}$  is the correlation between the actual values on the output units,  $I_i^{out}$ , and the desired blob shape,  $\mathbf{G}$ . That is,

$$\begin{aligned} E_{blob} &= -\sum_i I_i^{out} G_i \\ G_i &= \exp(-(i-2)^2/4). \end{aligned} \tag{2.10}$$

The second part of the objective (maintaining topography) can be accomplished by letting  $c_k$  follow the gradient of a constraint function,  $E_{constraint}$ , that favors valid control states—i.e., those corresponding to translations or scalings of the input-output transformation. One possible choice for  $E_{constraint}$  is

$$E_{constraint} = -\frac{1}{2} \sum_{k,l} c_k U_{kl} c_l, \tag{2.11}$$

where the constraint matrix  $\mathbf{U}$  is chosen so as to appropriately couple the control neurons. For the simple circuit of Figure 2.11, each control neuron corresponds to a different position of the window of attention, so we could define  $\mathbf{U}$  as

$$U_{kl} = \begin{cases} -1 & k \neq l \\ 0 & k = l. \end{cases}$$

This has the effect of punishing any state in which two or more control neurons are active simultaneously, and thus forces a winner-take-all solution. (The more general case using control blocks is described below.)

A dynamical equation for  $c_k$  that performs gradient descent on both  $E_{blob}$  and  $E_{constraint}$  is given by

$$c_k = \sigma(u_k) \tag{2.12}$$



$$\frac{du_k}{dt} + \tau^{-1} u_k = \eta \sum_i \sum_j G_i \Gamma_{ijk} I_j^{in} + \eta \beta \sum_l U_{kl} c_l, \quad (2.13)$$

where the constants  $\tau$  and  $\eta$  determine the rate of convergence of the system, and the constant  $\beta$  determines the contribution of  $E_{constraint}$  relative to  $E_{blob}$ . A sigmoidal squashing function ( $\sigma$ ) is used to limit  $c_k$  to the interval  $[0,1]$ . (See Appendix A for derivation.)

A neural circuit for computing Equations 2.12 and 2.13 is shown in Figure 2.12. The first term on the right of Equation 2.13 is computed by correlating the Gaussian,  $G$ , with a shifted version of the input (the amount of shift depends on the index  $k$ ). The second term is computed by forming a weighted sum of the activities on the other control neurons. These two results are then summed together and passed through a leaky integrator and squashing function to form the output of the control unit,  $c_k$ . Thus, each control neuron essentially has a Gaussian receptive field in the input, and competition among the control neurons allows only the unit with the strongest input to prevail.

Figure 2.13 shows a computer simulation of a 2-D version of this circuit. The input is composed of a 9x9 array of sample nodes and the output is a 5x5 array of sample nodes. There are 25 control units, corresponding to the 5x5 possible positions of the window of attention. Figure 2.13a shows the window of attention centered on a blob in the input. When the blob moves, the window of attention subsequently moves (via a discontinuous jump) to track it (Fig. 2.13b,c). If two blobs are present in the input, the window of attention is attracted to the brightest blob, since it provides the greatest input to its corresponding control neuron (Fig. 2.13d).

The circuit of Figure 2.12 could be modified to allow for different sizes of the window of attention by adding another set of control neurons for each desired size of the window of attention. The control neurons corresponding to a large window

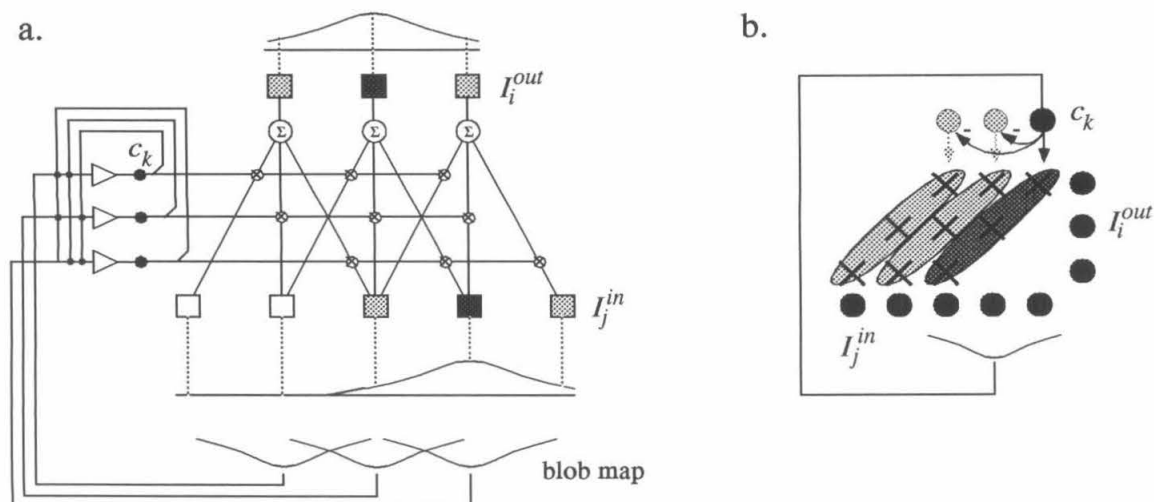


Figure 2.12: **Autonomous control.**

*a*, The circuit of Fig. 2.11 *a* with control circuitry added to autonomously focus the window of attention on a blob in the input. Each control neuron has a Gaussian receptive field in the input layer. The control neurons then compete among each other, via negatively weighted interconnections, so that only the control neuron corresponding to the strongest blob in the input prevails. The combined leaky integrator and squashing function (Equations 2.12 and 2.13) are denoted by the amplifier symbol. *b*, The circuit as depicted in connection space.

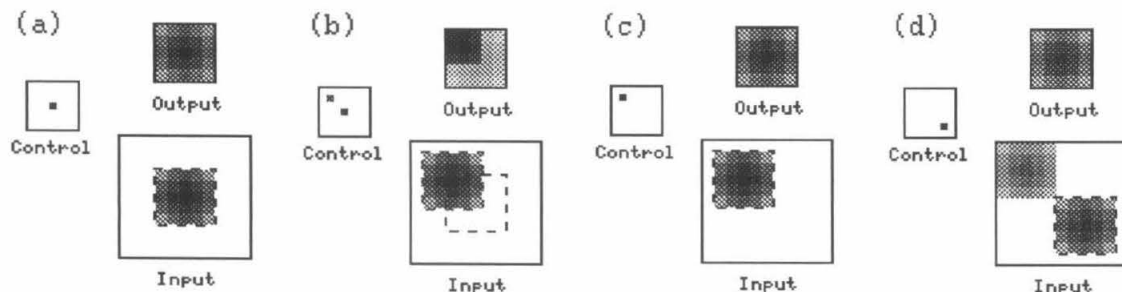


Figure 2.13: **Computer simulation of the autonomous routing circuit.** The dashed outline denotes the window of attention. The parameters chosen for the simulation are  $\eta = 0.04$ ,  $\alpha = 0.5$ , and  $\beta = 1.2$ . The window takes about half a time constant (or about 20 iterations) to shift position. (The time constant is defined as  $\frac{1}{\eta\alpha}$ .)

of attention would then have large Gaussian receptive fields, while control neurons corresponding to a small window of attention would have small receptive fields. (The strength of the receptive fields will also need to be normalized so that a large window control neuron out-competes a small window control neuron only when the entire extent of its receptive field is activated.) All of these units would then compete with one another so that the window of attention is constrained to a single position and scale.

In a more flexible scheme, the control neurons would be configured into control blocks (as in Fig. 2.3*c,d*). In this case, Equation 2.13 states that the input to each  $c_k$  would be computed by correlating the Gaussian values,  $G_i$ , and the input values,  $I_j^{in}$ , that are “connected” via that control unit (specified by  $\Gamma_{ijk}$ ). Note that since the  $G_i$  are fixed, the term  $\sum_i G_i \Gamma_{ijk}$  can essentially be considered a fixed weight. Thus, each control neuron will have a Gaussian-like receptive field that covers only those inputs that the control neuron couples with. Another consideration is that the constraint matrix,  $\mathbf{U}$ , would need to be modified in this case so that those control neurons corresponding to a common translation or scale reinforce each other ( $U_{kl} > 0$ ), while control neurons that are not part of the same transformation inhibit each other ( $U_{kl} < 0$ ). This is illustrated in Figure 2.14. This control scheme is demonstrated in the simulation of the recognition circuit below.

## Recognition

Once the window of attention has been focused on a blob, the underlying high-resolution information can also be fed through the routing circuit and into an associative memory for recognition. However, it is likely that the initial estimation of position and size made during blob search will be only approximately correct, and this may cause problems for matching the high-resolution information. Thus,

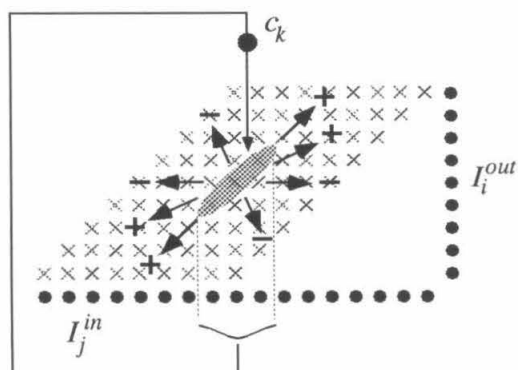


Figure 2.14: **Control neuron interactions when configured into control blocks.**

The control neuron corresponding to the block shown (stippled region) will have a receptive field covering only those inputs that the control neuron couples with. This control neuron should have excitatory connections ( $U_{kl} > 0$ ) to other control neurons whose blocks form a consistent position and size of the window of attention - i.e., those blocks lying along the '+' directions. Inhibitory connections ( $U_{kl} < 0$ ) should be formed with control units whose blocks are inconsistent with this one - i.e., those along the '-' directions. This scheme is somewhat analogous to the way constraints are imposed in the Marr/Poggio stereo algorithm (Marr and Poggio 1976).

it will be desirable to have the associative memory adjust the position and scale of the attentional window as it converges. How, then, shall the associative memory be incorporated into the control of the routing circuit?

If a Hopfield associative memory (Hopfield, 1984) is used for recognition, then we can replace  $E_{blob}$  with the associative memory's "energy" function,  $E_{mem}$ , which is defined as

$$E_{mem} = -\frac{1}{2} \sum_i \sum_j T_{ij} V_i V_j + \sum_i \frac{1}{R_i} \int_0^{V_i} g_i^{-1}(V) dV - \sum_i V_i I_i^{mem}. \quad (2.14)$$

In this equation the  $V_i$  denote the output voltages on the associative memory neurons,  $T_{ij}$  denotes the connection strength between neurons  $i$  and  $j$ ,  $I_i^{mem}$  denotes the inputs to the memory, and  $g_i$  is a squashing function such as  $\tanh(x)$ . Normally, the only dynamic variables are the  $V_i$ , which evolve by following a monotonically increasing function,  $g_i$ , of the gradient of the energy. That is,

$$\begin{aligned} V_i &= g_i(u_i^m) & (2.15) \\ C_i \frac{du_i^m}{dt} &= -\frac{\partial E_{mem}}{\partial V_i} \\ &= \sum_j T_{ij} V_j - \frac{u_i^m}{R_i} + I_i^{mem}, & (2.16) \end{aligned}$$

where  $C_i$  and  $R_i$  are constants that determine the integration time constant of each neuron. The dynamics of Equations 2.15 and 2.16 can be implemented in simple, neural-like circuitry. Note that the effect of minimizing  $E_{mem}$  is to simultaneously maximize (a) the similarity between the neuron voltages,  $V_i$ , and one of the stored patterns superimposed in the  $T_{ij}$  matrix (first term of  $E_{mem}$ ), and (b) the similarity between the  $V_i$  and the inputs  $I_i^{mem}$  (last term of  $E_{mem}$ ). (The second term of  $E_{mem}$  is the "leaky integrator term," which is unimportant for now. See Appendix A).

Since the inputs to the associative memory are obtained directly from the outputs

of the routing circuit ( $I_i^{mem} = I_i^{out}$ ), then the control neurons,  $c_k$ , become additional dynamic variables hidden in the last term of  $E_{mem}$ . By letting the  $c_k$  follow the gradient of  $E_{mem}$ , along with the  $V_i$ , the combined associative memory/routing circuit should relax to the closest stored pattern and to the correct position and size of the window of attention simultaneously.

A dynamical equation for  $c_k$  that performs gradient descent on both  $E_{mem}$  and  $E_{constraint}$  is given by

$$c_k = \sigma(u_k) \quad (2.17)$$

$$\frac{du_k}{dt} + \tau^{-1} u_k = \eta \sum_i \sum_j V_i \Gamma_{ijk} I_j^{in} + \eta \beta \sum_l U_{kl} c_l, \quad (2.18)$$

(See Appendix A for derivation.)

A neural circuit for computing Equations 2.17 and 2.18 is shown in Figure 2.15. The first term on the right of Equation 2.18 is computed by correlating the inputs,  $I_j^{in}$ , and outputs,  $V_i$ , whose connection pathways are influenced by control neuron  $c_k$  (specified by  $\Gamma_{ijk}$ ). The other terms are computed as before. Thus, the main qualitative difference between this circuit and the “blob finder” (Fig. 2.12) is that the control is guided by the interaction between top-down and bottom-up signals rather than purely bottom-up sources.

In order to avoid local minima, it will be advantageous to perform the combined process of pattern matching, shifting and scaling in a coarse-to-fine manner by utilizing information at multiple scales (e.g., Witkin and Terzopoulos, 1987; Buhmann et al., 1990). In this way, the low-pass information can be used to initially send the memory into the right part of its search space; the initial output of the associative memory can then be used to better refine the position and scale of the window of attention before allowing in higher-resolution information. A crude form of such a coarse-to-fine strategy has been utilized in the following computer simulation.

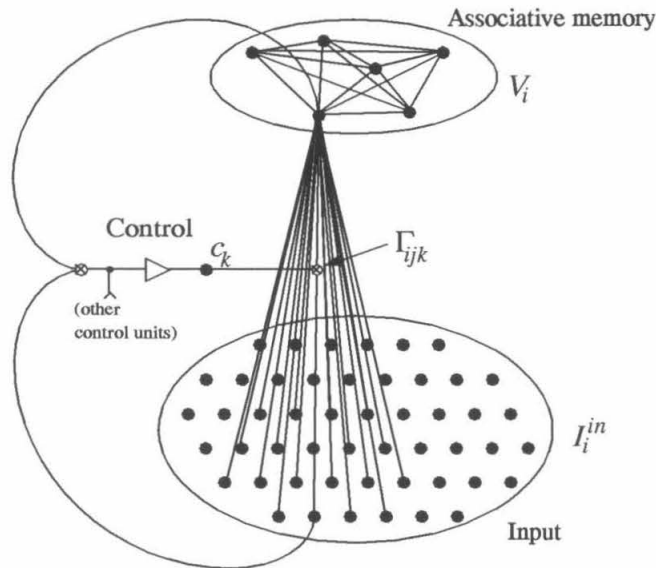


Figure 2.15: **An autonomous routing circuit for recognition.**

Each node of the associative memory receives its external input from an output node of the routing circuit. Hence, each node of the associative memory has dynamic connections to many input nodes. The outputs of the associative memory are then fed back and correlated with the inputs to drive the control neurons, as specified by Equation 2.18.

Figure 2.16 shows a computer simulation of an attentional system for recognizing objects. The network begins in blob search mode, attempting to fill the output of the routing circuit with something interesting. Rather than prefiltering shapes into blobs, the network attempts to find blobs directly from the original image. Thus, during blob search, an object ends up being low-pass filtered into the output of the routing circuit (Fig. 2.16*a*). This blurring function is facilitated by setting the constraint matrix,  $\mathbf{U}$ , so that control neurons corresponding to neighboring positions of the window of attention only weakly inhibit each other. After a fixed amount of time (one or two time constants), the network switches into recognition mode (Equation 2.18). Two patterns - 'A' and 'C' - have been stored in the associative memory using the outer product rule (Hopfield 1982). The blurred version of the object initially drives the inputs of the associative memory to begin the pattern search (Fig. 2.16*b*). If the position of the window of attention is slightly off, the low-pass version of the object will not be affected much and will still send the memory searching in the correct direction. As the associative memory converges, control neurons compute the correlation between memory outputs and retinal inputs and set their activation correspondingly. This tends to maximize the similarity between the outputs of the memory and the outputs of the routing circuit, which will also refine the position of the attentional window so that the high-resolution components can be properly matched (Fig. 2.16*c*). After allowing a fixed amount of time for the associative memory to converge (another time constant or two), the simulation states the position, size and presumed identity of the object.

It should be noted here that the particular form of associative memory used here has not been included as a model of the recognition process per se. Any number of schemes (e.g., population coding of shape) could conceivably be used. This example is provided only to give one an idea for how the control neurons may be driven top-down during the process of recognition.



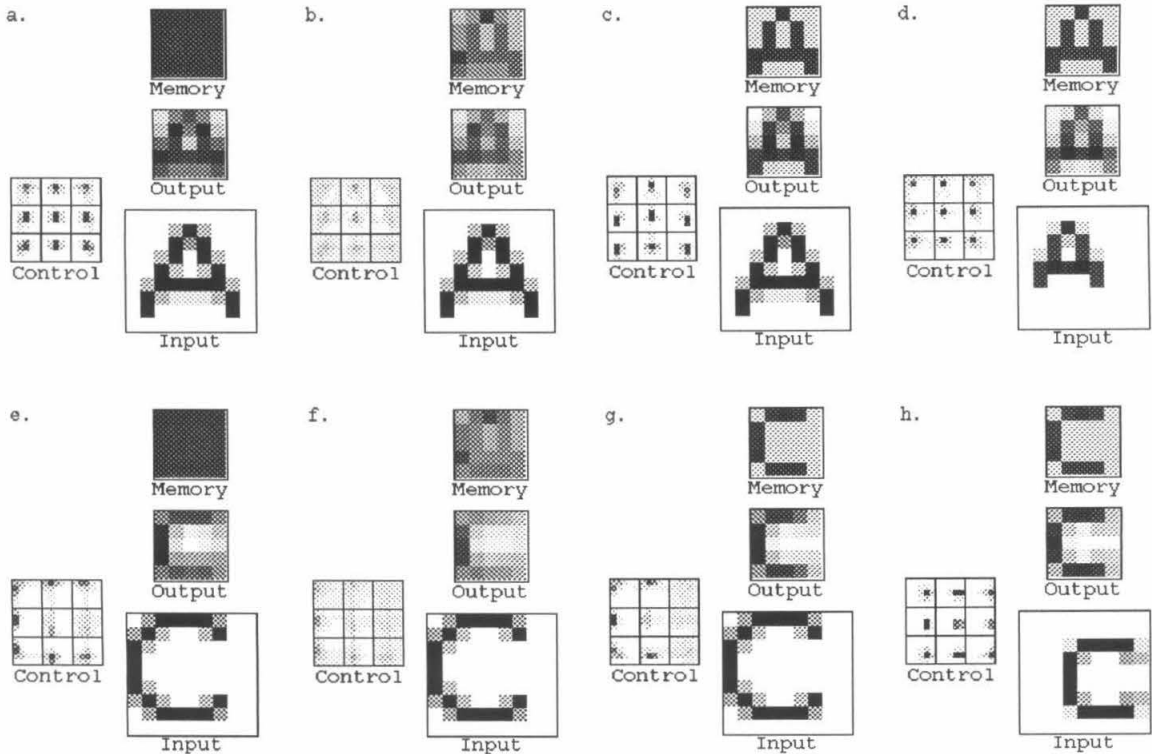


Figure 2.16: **Computer simulation of the recognition circuit.**

In this simulation, the control neurons have been configured into overlapping blocks composed of 3x3 synapses. control neurons compete within a block and cooperate and compete with other control neurons in neighboring blocks. Otherwise, the circuit is same as that shown previously. The Hopfield network ('Mem output') is composed of 25 units, fully interconnected and arranged into a 5x5 grid (i.e., one node for each output of the routing circuit). *a*, In blob search mode, the network rescales the large 'A' (7x7) into the window of attention. *b,c*, As the associative memory converges, the control neurons are driven top-down to refine the position and scale. *d*, If the size and position are changed, the control neurons update the connections to track it, acting as an "A finder." This same sequence is repeated for the letter 'C' in *e-f*.

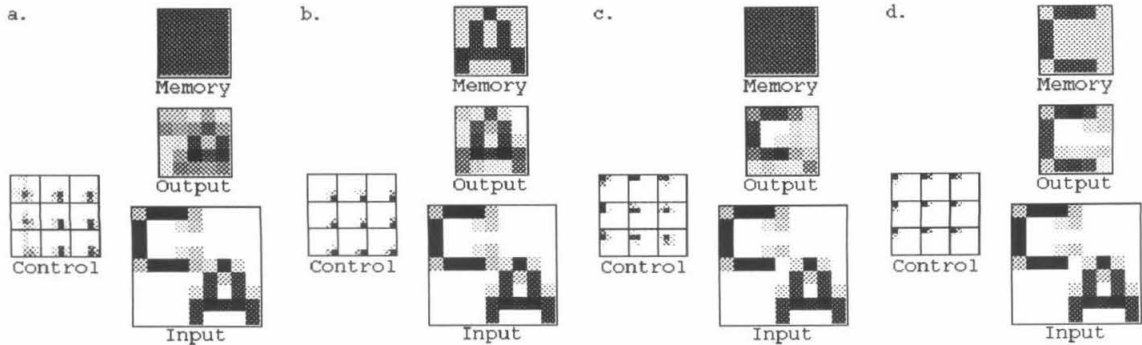


Figure 2.17: **Shifting attention.**

*a*, In blob search mode, the network has converged on the ‘A’ since it has the greatest overall brightness. *b*, After the associative memory converges, the current control state is inhibited, and the network is switched back into blob search mode. *c*, The ‘C’ is now at a competitive advantage in attracting the window of attention, and is subsequently recognized (*d*).

## Shifting attention

Once an object has been recognized, the window of attention should move on to another interesting part of the scene. One way this could be accomplished would be for the control neurons to be self-inhibited through a delay. Thus, when a group of control neurons are active for some time (long enough for recognition to take place) they should begin to shut off. This will then allow other blobs or interesting items to compete successfully for control of the window of attention. (See also Koch and Ullman, 1985.)

Figure 2.17 shows an example in which two shapes are presented in the input. The network initially settles on one of the shapes and tries to recognize it (Fig. 2.17*a,b*). Once this has been accomplished, the current control state is self-inhibited and the network switches back into blob search mode (Equation 2.13). This then puts the next object at a competitive advantage in attracting the window of attention so that it may be recognized (Fig. 2.17*c,d*).

## 2.4 Autonomous control II: multiple stages

We now formulate the autonomous control dynamics for a multistage routing circuit, which as we shall see introduces a whole host of new considerations. We first treat the case of two-stage routing circuit, and then turn to the case where this circuit is embedded in a multiscale system with competition among scales.

### Focusing on a blob

Consider the multistage routing circuit of Figure 2.18, with one or more Gaussian blobs presented to the input. We assume for starters here that the  $\Gamma_{ijk}^l$  have been set so that each control neuron in the first stage corresponds to a different position of the window of attention, and each control neuron in the second stage corresponds to a different module. More generally, each module in the first stage could perform scaling as well if the control neurons were broken up into control blocks, as demonstrated previously.

As before, we express the objective of finding a blob by the function  $E_{blob}$ , which measures the similarity between the output of the routing circuit,  $I^2$ , and a blob function  $G$ .

$$E_{blob} = - \sum_i I_i^2 G_i. \quad (2.19)$$

The input term for the control neurons of the top stage is derived by simply taking the derivative of this function (as derived in Appendix A)

$$- \frac{\partial E_{blob}}{\partial c_k^1} = \sum_i \sum_j G_i \Gamma_{ijk}^1 I_j^1. \quad (2.20)$$

This is essentially the same as the first term of Equation 2.13, and states that  $c_k^1$  has a Gaussian receptive field in layer 1 whose position corresponds to that of the  $k^{\text{th}}$

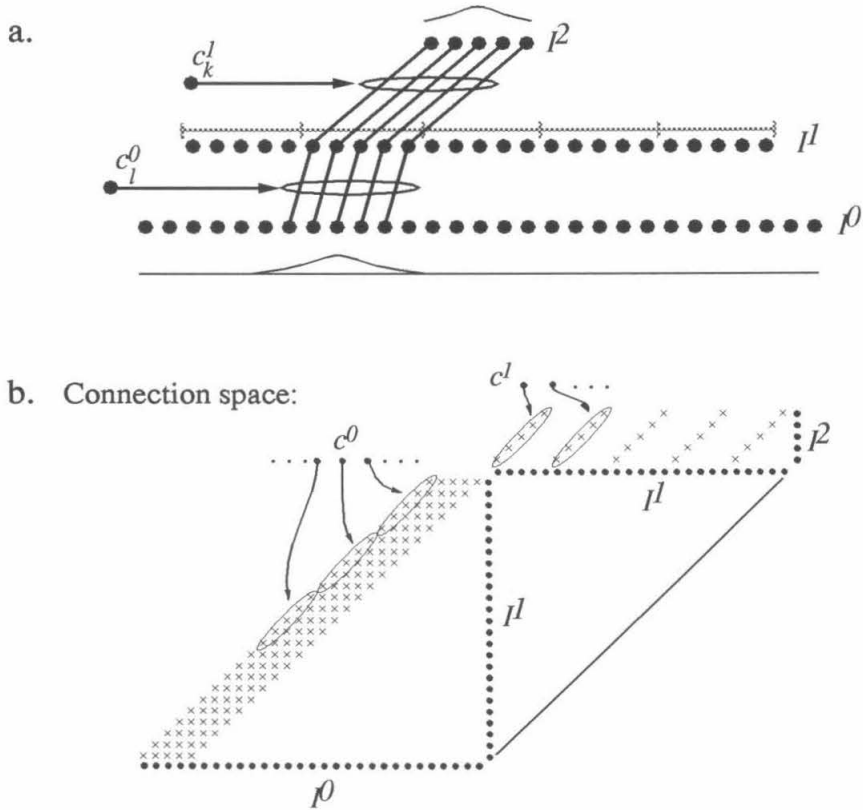


Figure 2.18: **A two-stage routing circuit and its control.**

*a*, The circuit of Figure 2.9 (composed of modules) with control made explicit. Each control neuron in the first stage corresponds to a different position of the window of attention, while each control neuron in the second stage corresponds to a different module in the first stage.

*b*, Connection-space diagram.

module in layer 1, as illustrated in Figure 2.19*a*.

The input term for the control neurons in the next lower stage is derived by using the chain rule to take the derivative one step further down

$$-\frac{\partial E_{blob}}{\partial c_k^0} = \sum_i \sum_j \sum_l G_i c_l^1 \Gamma_{ijl}^1 \frac{dI_j^1}{dc_k^0} \quad (2.21)$$

$$= \sum_i \sum_j \sum_l \sum_m G_i c_l^1 \Gamma_{ijl}^1 \Gamma_{jmk}^0 I_m^0. \quad (2.22)$$

This equation essentially states that  $c_k^0$  has a Gaussian receptive field in layer 0, gated by the control neuron in stage 1 that corresponds to the module to which  $c_k^0$  belongs. This has the effect of ensuring that paths through successive stages of the routing circuit are concatenated, as shown in Figure 2.19*b*.

The constraint term is the same as before, except now the control neurons need only be constrained locally within each module within each stage:

$$E_{constraint} = \sum_{l,m,n} c_m^l U_{mn}^l c_n^l \quad (2.23)$$

$$U_{mn}^l = \begin{cases} -1 & m \neq n \text{ and } m, n \text{ member of same module of stage } l \\ 0 & \text{otherwise.} \end{cases} \quad (2.24)$$

In the first stage, this has the effect of adding inhibitory connections between control neurons belonging to the same module. In the second stage, control neurons compete globally; but since the total number of control neurons at this stage is much fewer, the competition is still among a small number of neurons (5). Thus, a global winner-take-all is effected through local competition in a hierarchically organized control circuit, with the top-stage control neurons selecting which module (or chunk) of the image to attend to, and the bottom stage control neurons selecting a position within this module. In order to allow scaling ( $< 2$ ) in the lowest stage, the control neurons would be arranged into control blocks, and the constraint matrix for the bottom stage,  $\mathbf{U}^0$ ,

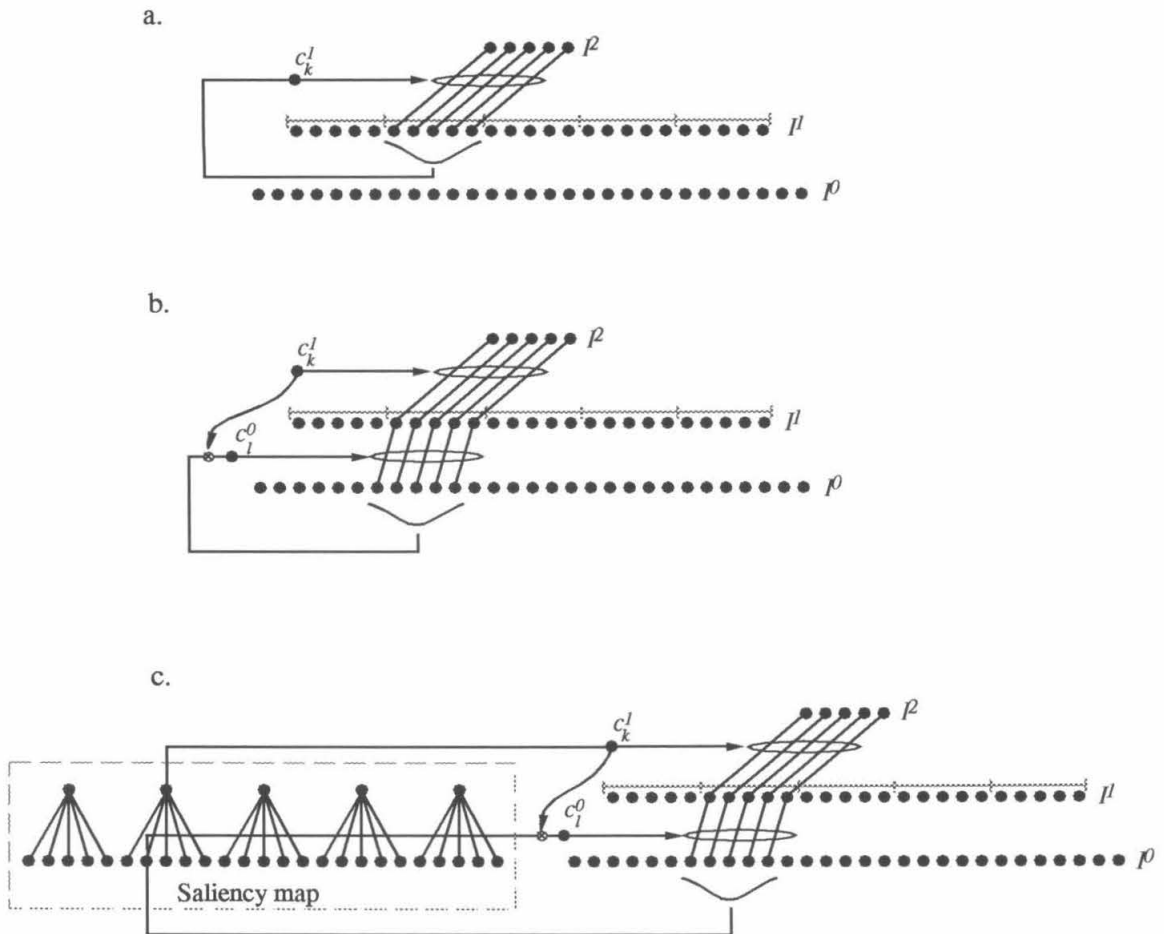


Figure 2.19: **Autonomous control of a multistage routing circuit.**

*a*, In the top stage, each control neuron,  $c_k^1$ , has a Gaussian receptive field in layer 1 whose position corresponds to that of the  $k^{\text{th}}$  module in layer 1. *b*, In the first stage, each control neuron,  $c_k^0$ , has a Gaussian receptive field in layer 0, gated by the control neuron in stage 1 that corresponds to the module to which  $c_k^0$  belongs. All five control neurons in the top stage compete among each other, whereas control neurons in the first stage compete in local groups of five within each module. *c*, A hierarchical saliency map for driving the control neurons. Each node in the first layer corresponds to a position of the window of attention, whereas nodes in the second layer correspond to modules, or “chunks,” of the input.

would be set as described previously (Fig. 2.14).

Now equations 2.20 and 2.22 seem to present something of a chicken-egg problem: The activity of the stage 1 control neurons depends on the activities of the layer 1 nodes, which in turn depend on the activity of the stage 0 control neurons. But the stage 0 control neurons are gated by the activity of the stage 1 control neurons. So which comes first? Since the stage 1 control neurons mediate macro-shifts, it makes sense for these to get set first, with the stage 0 control neurons initially off, or in the resting state. Thus, the activity of the layer 1 units will initially be determined by blurring  $I^0$  in the “all connections open” state—that is, with the  $c^0$  at their tonic resting state, or with  $w_{ij} = w_{rest}$ , as discussed previously. The control neurons of stage 1 will then settle on the most salient module in  $I^1$ . At this point, the winning  $c^1$  will enable the stage 0 control neurons belonging to the corresponding module in stage 0, and these control neurons will then compete and cooperate locally among each other to position and scale the window of attention within this module.

In the case where  $w_{ij} = w_{rest}$  when the control neurons are off, we can alternatively think of the control neurons as being driven by a hierarchical saliency map, as illustrated in Figure 2.19c. Each node in the first layer of the saliency map has a Gaussian receptive field in the input, while the second layer forms a summary by summing and subsampling this map.

Figure 2.20 illustrates a simulation of a 2D version of this circuit. Initially, the top stage control neurons are driven by the layer 1 saliency map to choose a module in the first layer (Fig. 2.20a). Once this module has been chosen, the control neurons within that module in the first stage compete with one another to choose a position within the module (Fig. 2.20b).

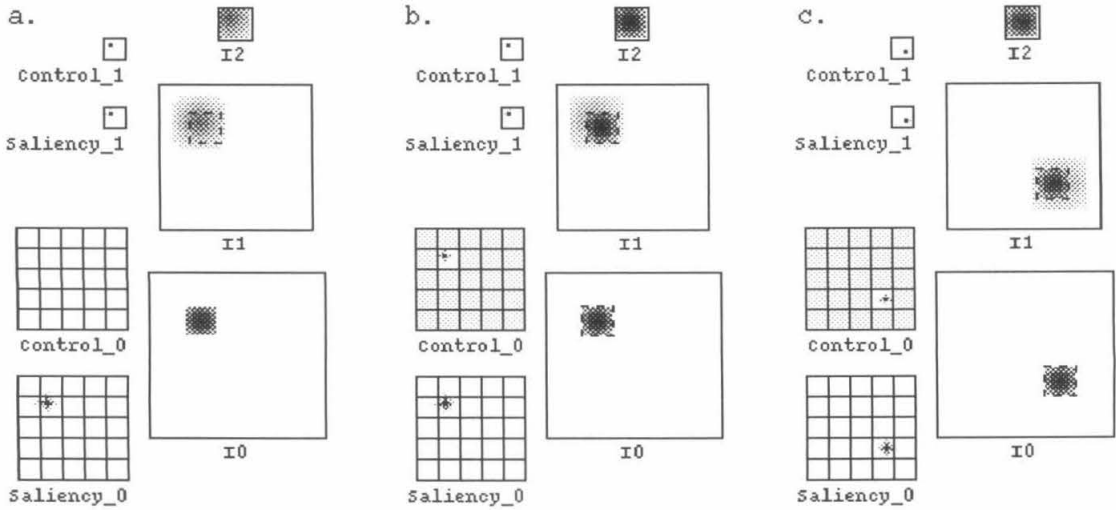


Figure 2.20: **Simulation of an autonomous, multistage routing circuit.** The input is composed of an array of  $29 \times 29$  nodes, the middle layer is composed of  $25 \times 25$  nodes, and the window of attention is  $5 \times 5$  nodes. The fan-in at each stage of the routing circuit is  $25:1$ . There are  $25 \times 25$  control neurons for the first stage ( $5 \times 5$  modules times  $5 \times 5$  control neurons per module), and  $5 \times 5$  control neurons for the second stage (one per module in layer 1). The value of  $w_{rest}$  was set to  $.04$  ( $1/25$ ), and the values on the output nodes,  $I_i^l$ ,  $l > 0$ , were renormalized locally within a layer. *a*, The top-stage control neurons converge initially, which then allows the first stage control neurons to converge (*b*). *c*, If the position of the blob changes, the control neurons change to track it.



## Recognition

In recognition mode, the output of the routing circuit is fed to an associative memory, and so the main objective function changes from  $E_{blob}$  to  $E_{mem}$ .

$$E_{mem} = - \sum_{ij} V_i T_{ij} V_j - \sum_i I_i^2 V_i. \quad (2.25)$$

The input term for the top stage control neurons is given by

$$- \frac{\partial E_{mem}}{\partial c_k^1} = \sum_i \sum_j V_i \Gamma_{ijk}^1 I_j^1, \quad (2.26)$$

which states simply that control neuron  $c_k^1$  is driven by the correlation between the layer 1 nodes,  $I_j^1$ , and memory outputs,  $V_i$ , that are connected via that control neuron (Fig. 2.21a). However, it is probably not necessary for the control neurons of this stage to be driven by the memory, because they set the position in such coarse chunks that fine-scale adjustments during recognition would have no impact. So we can keep them in the same state they settled on during blob search.

The input term for the control neurons in stage 0 is derived as before using the derivative chain rule

$$- \frac{\partial E_{mem}}{\partial c_k^0} = \sum_i \sum_j \sum_l \sum_m V_i c_l^1 \Gamma_{ijl}^1 \Gamma_{jmk}^0 I_m^0. \quad (2.27)$$

Note that this is just the same as Equation 2.18, except with  $G_i$  replaced with  $V_i$ . The major difference, though, is that the  $V_i$  are dynamic variables, and thus we cannot simply incorporate their multiplicative effect into a fixed weight as we did for the  $G_i$  previously. It is helpful to rewrite Equation 2.27 in a different form that leads to an interesting neural architecture for implementing it. Substituting  $w_{ij}^l$  for  $\sum_l c_l^1 \Gamma_{ijl}^1$ , we

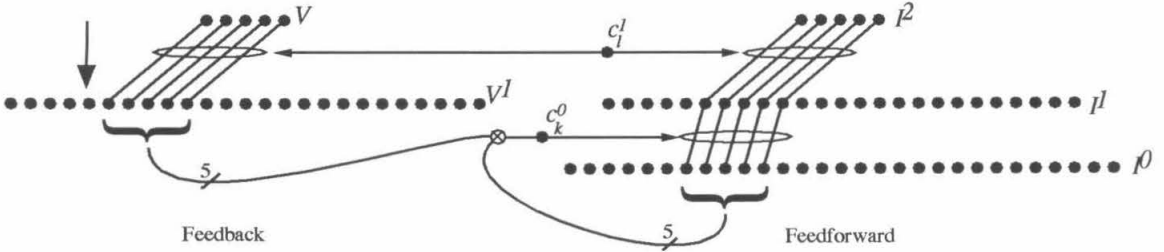


Figure 2.21: **Multistage recognition circuit.**

The output of the associative memory,  $V$ , is routed “backwards” into a separate population of neurons in layer 1 (*left*).  $c_k^0$  is then driven by the correlation inputs,  $I^0$ , and feedback nodes,  $V^1$ , as specified in Equation 2.29. The  $/5$  symbol denotes 5 parallel lines.

obtain

$$-\frac{\partial E_{mem}}{\partial c_k^0} = \sum_i \sum_j \sum_m V_i w_{ij}^1 \Gamma_{jmk}^0 I_m^0 \quad (2.28)$$

$$= \sum_j \sum_m V_j^1 \Gamma_{jmk}^0 I_m^0, \quad (2.29)$$

where

$$V_j^1 = \sum_i V_i w_{ij}^1. \quad (2.30)$$

$V^1$  is the result of routing the output of the associative memory,  $V$ , “backwards” into a separate population of neurons in layer 1. Thus,  $c_k^0$  is driven by correlating the inputs,  $I_m^0$ , and the fed back signals in the layer above,  $V_j^1$ , corresponding to those nodes in layer 1 connected with  $I_m^0$  via  $c_k^0$ . This is illustrated in Figure 2.21. A simulation of this circuit is demonstrated in the following subsection.

## Shifting attention

The method we used previously for shifting attention was to simply inhibit those control neurons that were “on” once recognition had taken place. This will not be

generally applicable for a multistage circuit, though, since inhibiting a control neuron in the higher stages will prevent all locations within the corresponding modules below from being attended. The method used for shifting attention in the multistage circuit will depend on which of the two methods is used to achieve the all-connections-open state in the intermediate stages. In the first case, where the control neurons have a low, tonically active resting state and blob search is done in  $I^1$ , we can simply inhibit the control neurons in the first stage only. This will then prevent any activity from showing up in  $I^1$  and subsequently being used to attract attention. In the other case, where  $w_{ij} = w_{rest}$  when all the control neurons are disabled, we cannot simply inhibit the first stage control neurons because the saliency in  $I^1$  is being computed independently of the control neurons and will still register these locations as interesting. Thus, the saliency nodes in the first layer must receive a delayed inhibition signal from the currently active control neuron. A third alternative is that the top-stage control neurons may be self-inhibited weakly, or with a fast time constant, and the bottom-stage control neurons self-inhibited strongly, or with a slow time constant. This way, attention would be more likely to be drawn to an object that is far away from (or a different size than) the currently attended object, but would go back to revisit neighboring objects after a short time.

Figure 2.22 shows an example of a simulation that utilizes the second method above for shifting attention. After an object has been recognized, the first-stage saliency nodes are inhibited in the location corresponding to the currently active control neuron of that stage. This then prevents that location from being attended to when the circuit is subsequently switched back into blob search mode.

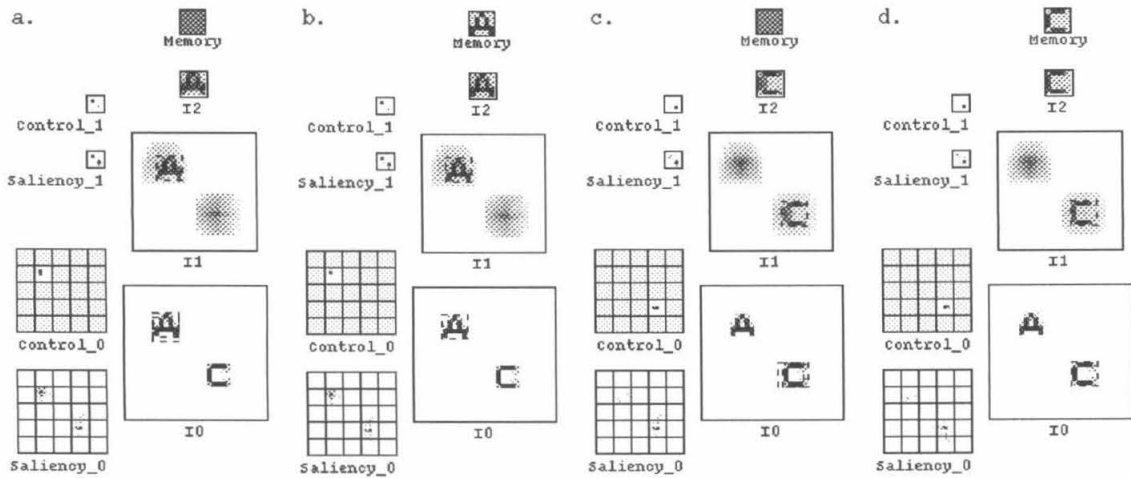


Figure 2.22: **Shifting attention in the multistage routing circuit.**

The circuit is the same as previously. In *a*, the circuit has converged on the 'A' since it has the greatest overall brightness. *b*, After allowing enough time for recognition, the first stage saliency map is inhibited in the vicinity corresponding to the currently active control neuron of that stage. *c*, The circuit is then switched back into blob search mode, allowing the next most salient object, 'C', to be attended and subsequently recognized (*d*).

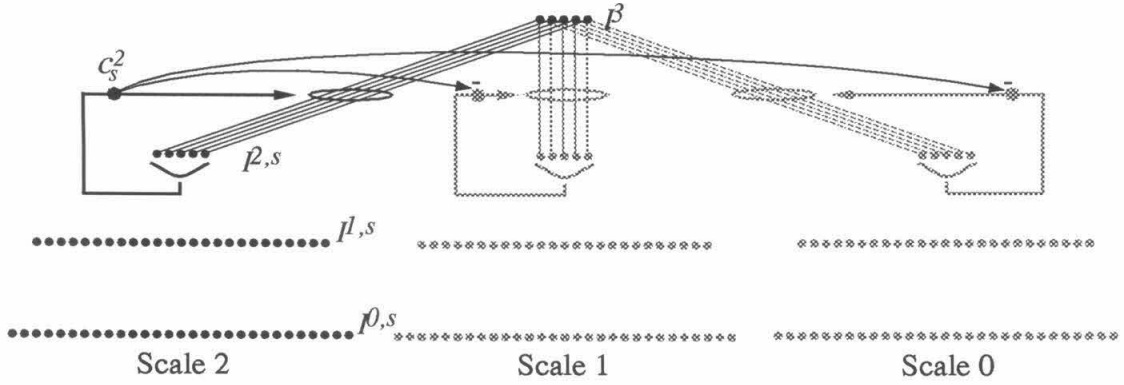


Figure 2.23: **Autonomous control for the multiscale stack routing circuit.** Each control neuron of stage 2 has a Gaussian receptive field in the top layer of the scale corresponding to the control neuron. These control neurons then compete among either other to decide on the scale to attend to.

## Competition among scales

We now revisit the multiscale stack circuit in which there are three different routing streams corresponding to different spatial scales (Fig. 2.8). Here, there are three control neurons at the top stage for gating the output of the routing circuit for each scale into a final output layer,  $I^3$ , as shown in Figure 2.23.

Following the same steps as before, the input term to  $c_s^2$  in blob search mode is given by

$$-\frac{\partial E_{blob}}{\partial c_s^2} = \sum_i \sum_j G_i \Gamma_{ijs}^2 I_j^{2,s}, \quad (2.31)$$

where  $I^{l,s}$  denotes the nodes of layer  $l$  of scale  $s$ . Equation 2.31 states that  $c_s^2$  will have a Gaussian receptive field in the top layer of scale  $s$  (Fig. 2.23). Or, in terms of the saliency map,  $c_s^2$  is driven by the sum of activity in the top-stage saliency map for scale  $s$ .

The constraint matrix is set so that control neurons of different scales compete

$$U_{st}^2 = \begin{cases} -1 & s \neq t \\ 0 & s = t. \end{cases} \quad (2.32)$$

Thus, extending the previous scheme, competition will begin at the top stage (stage 2); control neurons here will compete to select the scale with the most salience, and the winning control neuron will enable control neurons below to compete for the most salient module within that scale, and finally the most salient position within that module.

In order to make the comparison between saliencies at different scales meaningful, the shape of the saliency function,  $\mathbf{G}$ , will need to be changed so that it is selective for a particular scale. As it stands, the saliency nodes for the smallest scale will respond equally well or better to part of a large object as compared to a small object alone. The actual objective we seek during blob search is to just fill the window of attention with a blob that is confined within the bounds of the window. This objective can be expressed by adding an inhibitory surround to  $G$ , as in Figure 2.24. This way, a small object that stands alone in the high resolution array will be registered with higher salience than a high luminance region that is part of a larger object. It may also be desirable to build-in a precedence for global (low-resolution) over local (high-resolution) information by providing the low-resolution circuits with faster time constants. This would have the effect of “canceling out” the larger objects before attending to the small objects.

Figures 2.25-2.28 illustrate some examples from a computer simulation of an autonomous, 2D version of the model of Figure 2.8. The saliency function used here was +1 within a 7x7 window (with a gaussian taper) and -2 within a two-pixel wide perimeter. The method of computing the saliency for stages 1 and 2 was changed



slightly from that described above, because simply summing the total saliency in a module can lead to deceptively high saliencies in the higher stages. For example, there may be a single strong node within one salience module, but if another salience module has many weak nodes the total may actually be greater than that in the module with a single strong node. This effect was ameliorated somewhat by squaring the values of the saliency nodes, which has the effect of attenuating the low salience nodes. In Figures 2.25 and 2.26, the circuit is shown first attending to the ‘A’ at the lowest resolution level of the stack, and then to the ‘C’ at the highest resolution level several fixations later. Figures 2.27 and 2.28 demonstrate that the circuit is capable of discerning both local and global structure. In all figures, the “warper” window displays the state of the attended  $9 \rightarrow 5$  module in the first stage, where scaling (and warping) less than a scale factor of two is performed. (This is just the same circuit as demonstrated in Figure 2.16.) The circuit as a whole is capable of continuously scaling the window over a factor of eight (from  $5 \times 5$  to  $40 \times 40$  nodes).

## 2.5 Summary of the model

In this chapter we developed a model circuit that adheres to the important neurobiological constraints of fan-in, fan-out, multiscale representation, and logarithmic spatial sampling. In order to specify how the control neurons are driven autonomously, we assumed that a useful strategy would be to focus attention on interesting regions within a scene and then attempt to recognize whatever is there. From this basic assumption, we derived equations for governing the dynamics of the control neurons in both “pre-attentive” (blob search) and “attentive” (recognition) modes. Breaking the routing circuit into multiple stages results in a hierarchical control circuit that is capable of setting the position and size of the window of attention globally with only local interactions among the control neurons. Although these circuits have been



greatly scaled-down for the purpose of illustration and simulation, the basic principles can be extended to larger, scaled-up routing circuits. In the next chapter we show how such circuits may be implemented in the brain.

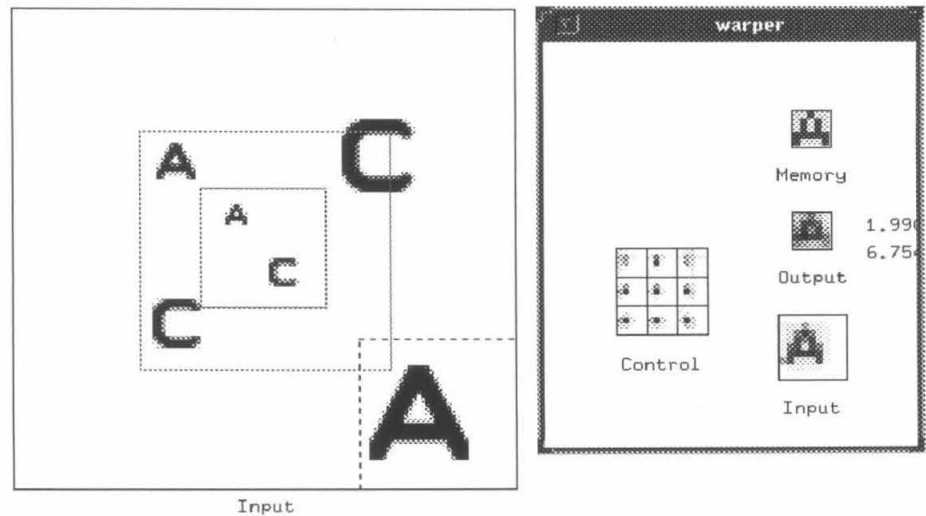
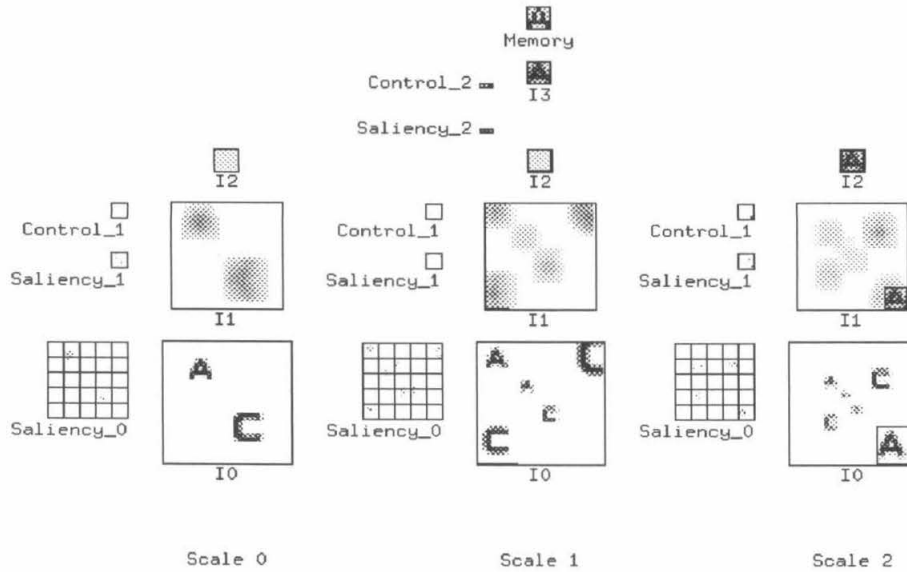


Figure 2.25: **Simulation of the stack circuit.**

The dotted square outlines in the input array denote the boundaries of the different levels of the stack. The dashed line indicates the window of attention. The circuit is shown attending to the 'A' at the lowest resolution level of the stack.

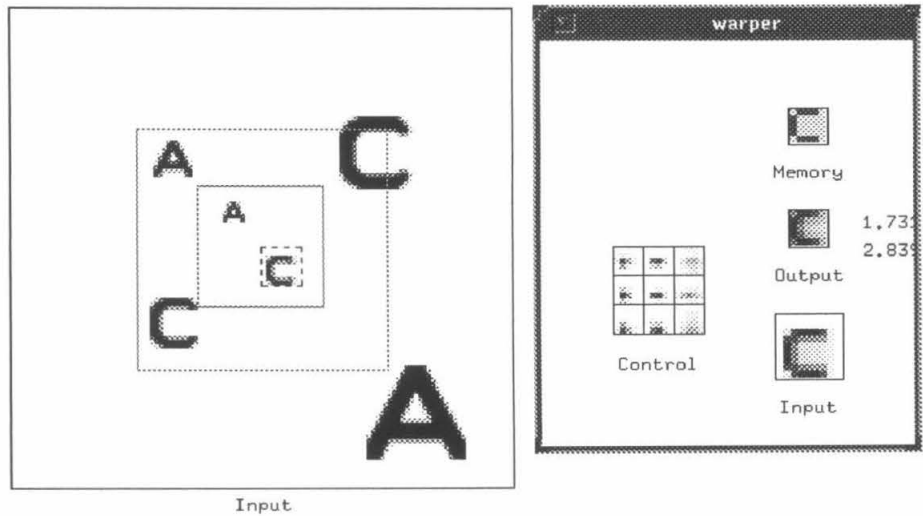
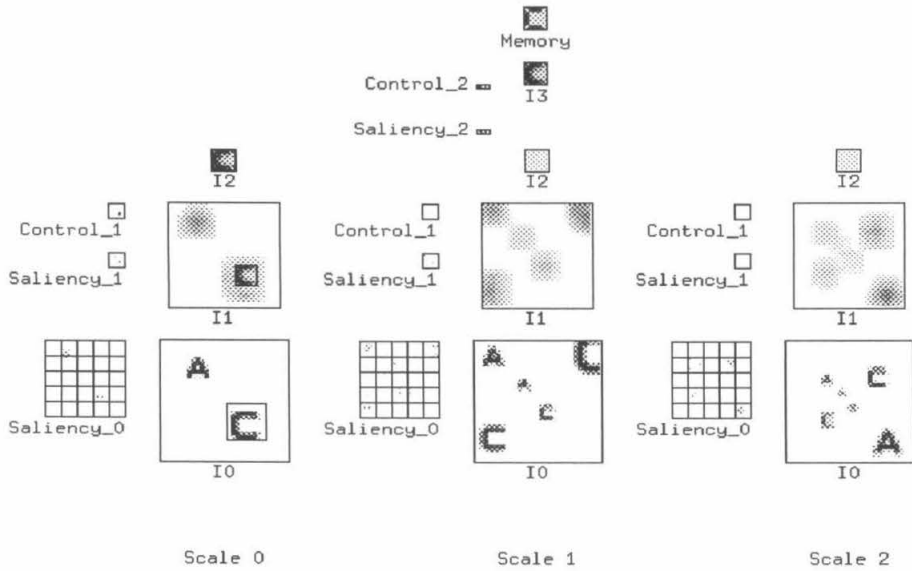


Figure 2.26: **Simulation of the stack circuit.**  
 The circuit is shown attending to the 'C' at the smallest scale several attentional fixations later.

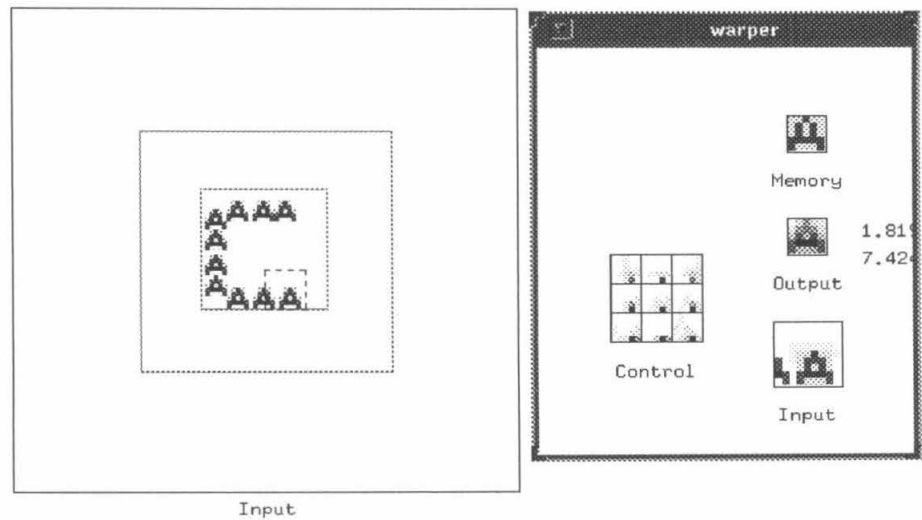
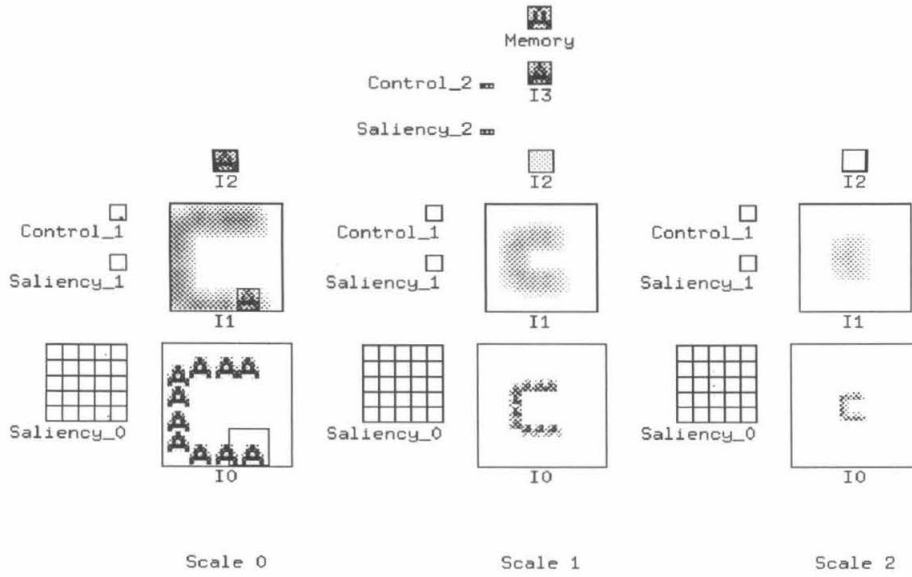


Figure 2.27: **Simulation of the stack circuit (local vs. global).**  
 The circuit is shown attending to one of the small A's that make up a global 'C' shape.

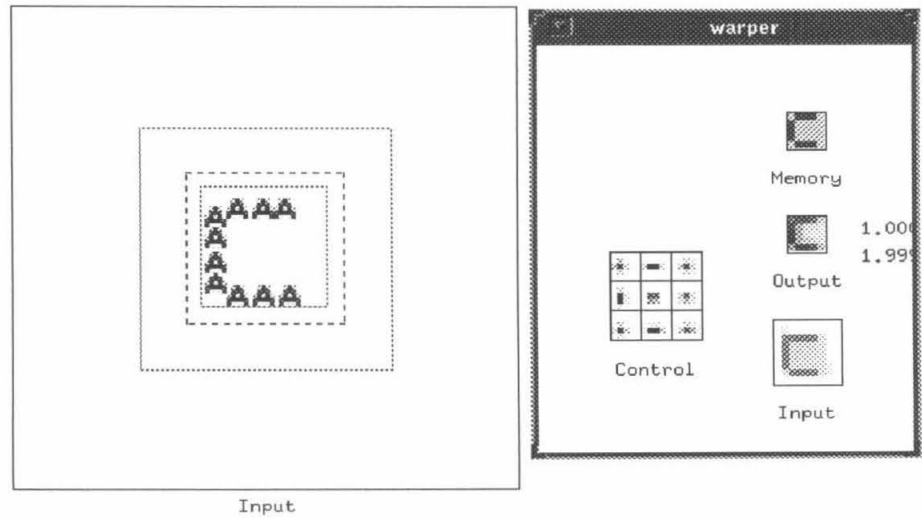
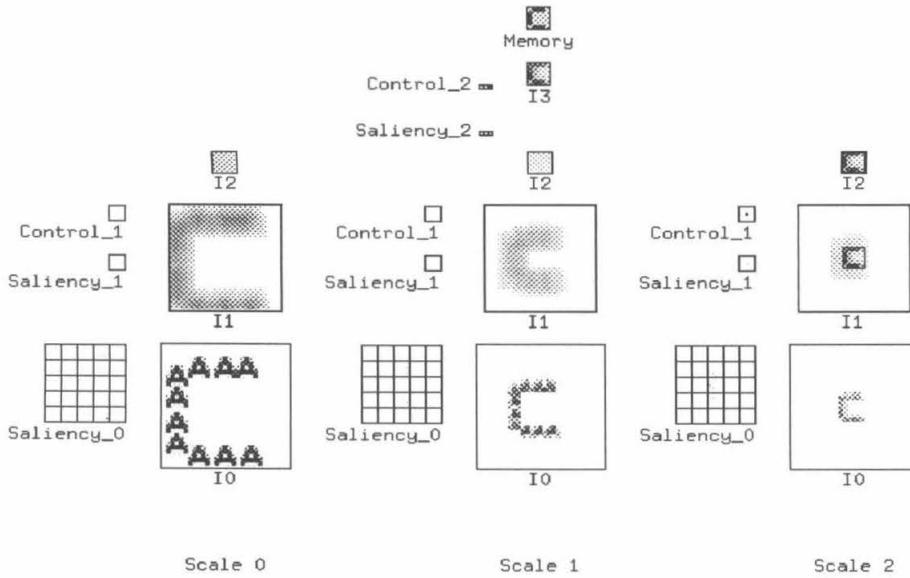


Figure 2.28: **Simulation of the stack circuit (local vs. global).**  
The circuit is shown attending to the global 'C' shape.

## Chapter 3

# Neurobiological substrates and mechanisms

We now turn to the issue of how the routing circuit we built up in the previous chapter may be implemented in the brain of the macaque monkey. The major areas that will be of interest to us, along with their anatomical relationships, are shown in Figure 3.1. Information from the retino-geniculo-striate pathway enters the visual cortex through area V1 in the occipital lobe and proceeds through a hierarchy of cortical visual areas that can be subdivided into two major functional streams (Ungerleider and Mishkin, 1982). The so-called “form” pathway leads ventrally through V4 and inferotemporal cortex (IT) and is mainly concerned with object identification, regardless of position or size. The so-called “where” pathway leads dorsally into the posterior parietal complex (PP), and seems to be concerned with the locations and spatial relationships among objects, regardless of their identity. The pulvinar, a sub-cortical nucleus of the thalamus, makes reciprocal connections with all of these cortical areas (cf. Robinson and Petersen, 1992) and also receives a projection from the superior colliculus, a midbrain structure involved in directing overt visual attention (i.e., eye movements).

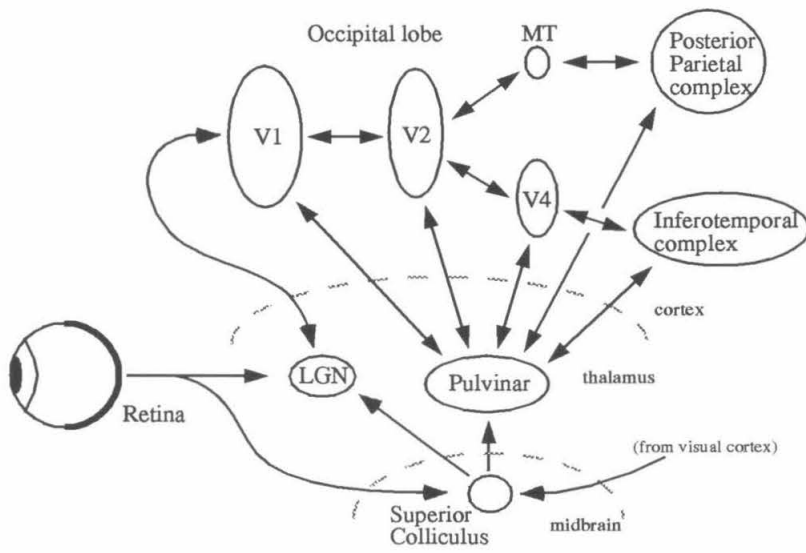


Figure 3.1: Neurobiological substrates.

The following sections describe how the dynamic routing circuit may be mapped onto this collection of neural hardware. We begin by describing how the multiscale stack and the intermediate stages of the routing circuit may be represented in the cortex. We then turn to the proposed neural substrates for control, and possible gating mechanisms for modulating connection strengths.

### 3.1 Routing circuit substrates

Since we are primarily interested in how invariant representations are formed for object recognition, we will focus on the “form” pathway from V1 to IT.<sup>1</sup> Figure 3.2 shows a schematic of the main cortical areas in this pathway. Each area is drawn to scale according to its relative size in one dimension (square root of the area). The fan-out shown at each stage is derived from the fact that receptive field size and visual field overlap between hemispheres approximately doubles at each stage in this pathway (Gattass et al., 1985), in addition to the anatomical observation that areal inter-connections increase their divergence and become more patchy in higher stages of the visual cortical hierarchy (Rockland, 1992; Van Essen et al., 1986; Van Essen et al., 1990; Van Essen and DeYoe, 1993; DeYoe and Sisola, 1991). This is only a rough depiction, however, and more data are needed to construct a firmer quantitative picture.

In previous work (Olshausen, Anderson, and Van Essen, 1993), it was described how the multistage circuit of Figure 2.4 could be mapped onto these stages of cortical visual processing. Here, we shall propose that area V1 forms a multiscale “stack” representation of the retinal image, and that the major intermediate visual areas in

---

<sup>1</sup>More generally, we can conceive of routing taking place in the other visual processing streams as well—for example, in the motion pathway for making fine discriminations of motion (Nowlan and Sejnowski, 1993; Van Essen and Anderson, 1990).



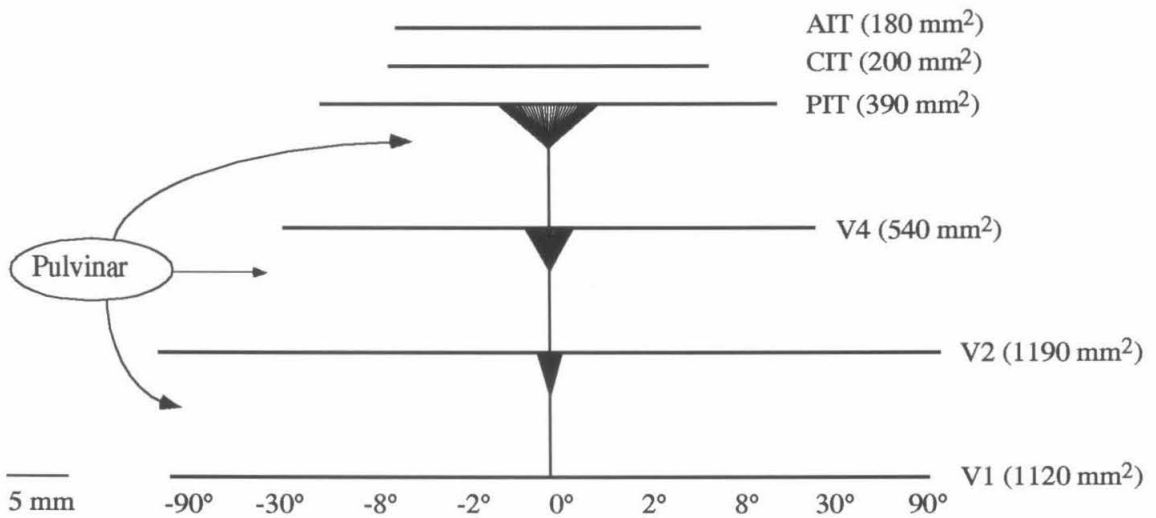


Figure 3.2: **Schematic of the main cortical areas in the “form” pathway.** The size of each area is drawn to scale ( $\sqrt{\text{area}}$ ). The fan-in at each stage is drawn to comprise 30 outputs ( $\sim 1000$  in 2D), and the divergence increases by a factor of two at each stage. The pulvinar makes connections with all these visual areas, and is proposed to be a source of the control signals for modulating connection strengths (see Section 3.2 below).

the form pathway (i.e., V2, V4) serve as intermediate stages for routing visual information from the stack in V1 into a progressively more position and scale invariant representation at higher stages. This process culminates with information being represented within an object-centered reference frame at the initial stages of IT (i.e., PIT). The size of each area would be expected to correspond roughly to the number of “sample nodes” in each area. Each node would correspond to a vector of features extracted in that area for a particular position and size of visual space (e.g., a full set of orientations in V1). Cells at higher stages of IT (i.e., CIT, AIT, or in the STS) would then perform their analyses on the contents of the window of attention, with variations in position and scale removed (e.g., face cells). Although we do not yet fully understand the nature of form processing occurring in this pathway—especially in the intermediate stages—we will assume we can leave this as an unknown and deal with issues of routing independently.

### **The multiscale “stack”**

In order to propose a quantitative model for a multiscale stack representation in V1, we need to specify 1) the highest resolution available as a function of eccentricity, and 2) the resolution ratio between adjacent levels of the stack. For the primate visual system, the highest resolution available at eccentricity  $E$  is given by

$$\delta(E) = .01(E + 1.3) \text{ deg}, \quad (3.1)$$

where  $\delta$  gives the one-dimensional spacing between sample nodes in the retina (Van Essen and Anderson, 1990). In two dimensions, each sample node would cover an area of approximately  $\delta^2$ . The resolution ratio between adjacent levels can be inferred from the spatial-frequency bandwidths of V1 cells, since an efficient coverage of the spatial frequency domain would require that the spacing in spatial frequency be related to

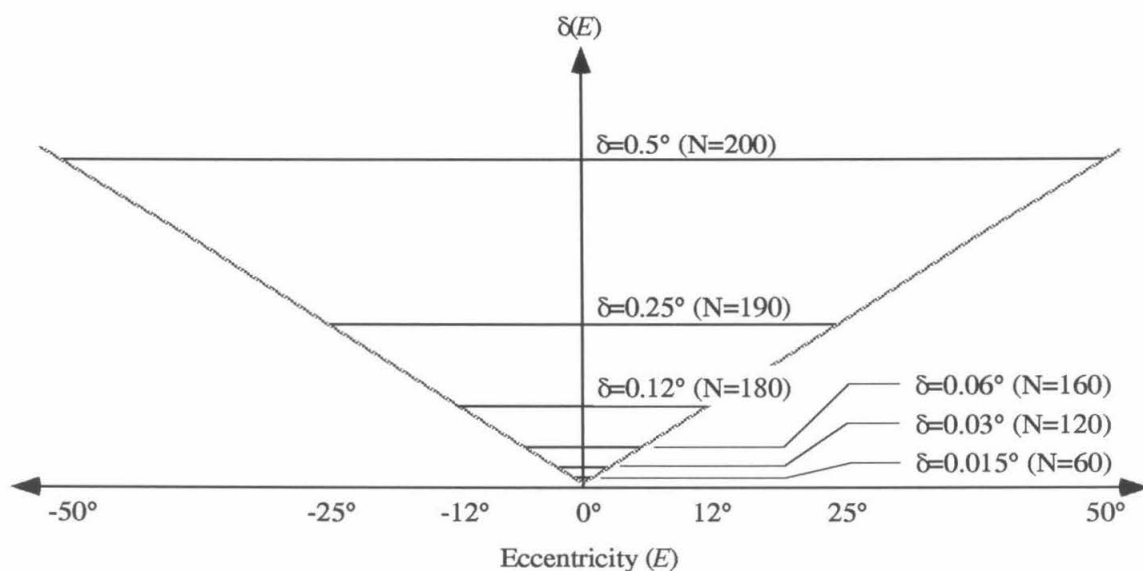


Figure 3.3: A six-level “stack” model for V1.

$\delta$  denotes the sample spacing and  $N$  denotes the number of nodes in each level of the stack.

bandwidth. Since the bandwidths of V1 cells tend to hover in the range of 1 to 1.5 octaves (De Valois et al., 1982), we will assume that resolution approximately doubles for each successive level of the stack.

Given these constraints, a stack comprising approximately 6 levels would suffice to cover the visual field up to  $\pm 50^\circ$  eccentricity (a total of  $100^\circ$ ), as shown in Figure 3.3. Beyond this eccentricity, retinal ganglion cell sample spacing is no longer scale invariant (i.e., no longer adheres to the linear relationship of Equation 3.1), and it also begins to be somewhat unreasonable to suppose that objects beyond this size<sup>2</sup> would be recognizable as a whole. The number of sample nodes in 1D for each level

<sup>2</sup>This would be the extent of a fully stretched hand when held 3 inches from the eye.

of the stack is given by

$$N = \frac{2E}{\delta(E)} = \frac{2E}{.01(E + 1.3)} \quad (3.2)$$

which will equal approximately 200 for  $E \gg 1.3^\circ$ . At eccentricities near or below  $1.3^\circ$  the number of nodes within a level will be fewer. The total number of sample nodes for the entire stack will thus be on the order of  $6 \times 200^2 = 240,000$ , which is about equal to the total number of sample nodes delivered by the optic nerve for the central  $100^\circ$  (80% of the total) when one takes into account the fact that information is divided into on- and off-channels, magno and parvo streams, and different spectral bands (Van Essen and Anderson, 1990). The highest resolution level of the stack will have a mean sample spacing of about  $.015^\circ$ , which implies a peak spatial-frequency tuning for each sample node in the range of 15 cy/deg, assuming that the optimum spatial-frequency for a given level is somewhat less than half the sampling frequency. The lowest frequency nodes will have a spacing of about  $0.5^\circ$  with an expected peak spatial-frequency tuning of about 1 cy/deg or lower.

Figure 3.4 illustrates the stack in cortical dimensions, where space has been logarithmically compressed according to

$$X_c(E) = 10 \log\left(\frac{E + 0.8}{0.8}\right) \quad (3.3)$$

where  $X_c$  gives the cortical distance in millimeters from the origin (i.e., fovea) of V1. (Equation 3.3 was obtained by integrating the formula for cortical magnification factor,  $10(E + .8)^{-1.1}$  mm/deg, as given by Van Essen et al. (1984), and rounding the exponent down to  $-1.0$  to make the integration simple.) In foveal V1, the highest resolution nodes would be spaced by about  $200\mu$ , and the lowest frequency nodes would be spaced by about 6 mm. With increasing eccentricity, the spacing between low resolution nodes will decrease, and the total number of levels will decrease as well until only the lowest resolution level is left. The spacing between the lowest resolution

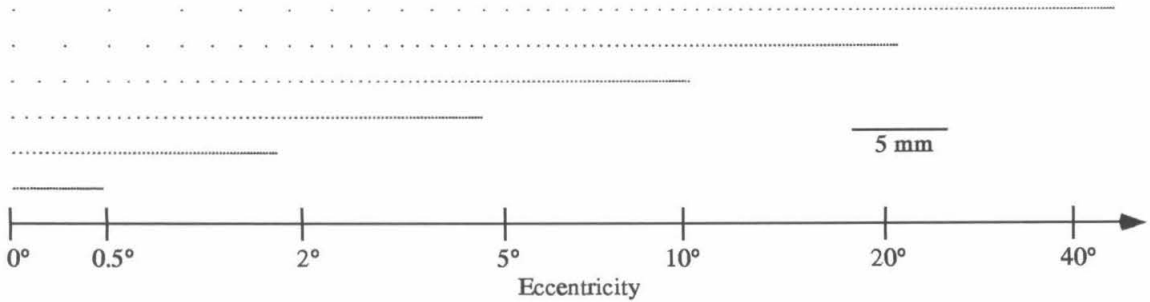


Figure 3.4: **The stack in cortical dimensions.**

Visual space has been logarithmically compressed according to Equation 3.3. Sample nodes for the lower resolution levels are distributed sparsely in the foveal cortex, but densely in the periphery.

nodes at the largest eccentricity ( $50^\circ$ ) will be the same as the spacing between the highest resolution nodes in the fovea ( $\approx 200\mu$ ). Note that since the density of sampling nodes decreases by a factor of four (in 2D) for each octave decrease in resolution, the total density does not vary appreciably with eccentricity, even though there are many more levels of the stack represented in the fovea than in the periphery.

Are these characteristics consistent with the physiological and anatomical data on V1? Determining the number and range of spatial-frequency tuned cells in V1 is difficult because of the disparate and conflicting data, so a more in depth discussion of this issue is provided in Appendix B. To summarize, though, it appears that within foveal V1, the necessary range of spatial frequencies exists: Most cells have very small receptive fields with central excitatory zones on the order of 2-4 minutes in diameter, which puts their peak spatial-frequency in the range of 8-15 cy/deg (Parker and Hawken, 1988). Cells with larger receptive fields become fewer in number with increasing diameter, bottoming out with a very small number on the order of 1 cy/deg (DeValios et al., 1982; Tootell et al., 1988). There is also evidence that the range of peak spatial-frequencies is largest in the fovea, and that this range

progressively decreases with eccentricity, as would be expected of a multiscale stack representation. It is not possible to determine whether the relative spacing of sample nodes at different resolutions is consistent with a stack model, but one can infer from the relative numbers of cells at each spatial-frequency that lower spatial frequency tuned cells would be spaced farther apart, assuming that cells are spread uniformly.

### **Intermediate stages and wiring constraints**

The number and size of intermediate stages of routing required depends on the total input-output convergence and the maximum allowable fan-in (number of inputs per neuron). The input to the routing circuit for each scale will be a 2D array comprising approximately 200x200 nodes (as described above), and the output of the routing circuit is hypothesized to be a relatively small array, comprising on the order of 30x30 sample nodes. (This estimate is based largely on spatial acuity and recognition studies that provide hints about the resolution of the window of attention—see Discussion, Section 4.1). Thus, the total convergence for the routing circuit for each scale will be about 40,000:1. Since the maximum allowable fan-in is on the order of 1000 inputs per neuron (Cherniak, 1990; Douglas and Martin, 1990a), the routing circuit for each scale must be broken into several stages.

A nominal configuration would be for each routing circuit to be broken into two stages, as shown in Figure 3.5*a*. This circuit is simply a scaled-up version of the circuit described in the previous chapter (Fig. 2.9), where the middle layer is now composed of modules of size 30. On the right, the circuit is pictured in terms of its fan-out, which is more neurobiologically relevant. Although this architecture does not appear to pose any anatomical problems as drawn, it must be kept in mind that the routing circuits for each each scale will be superimposed in register in the cortex. Thus, the low resolution nodes will need to have a very great divergence in terms of cortical

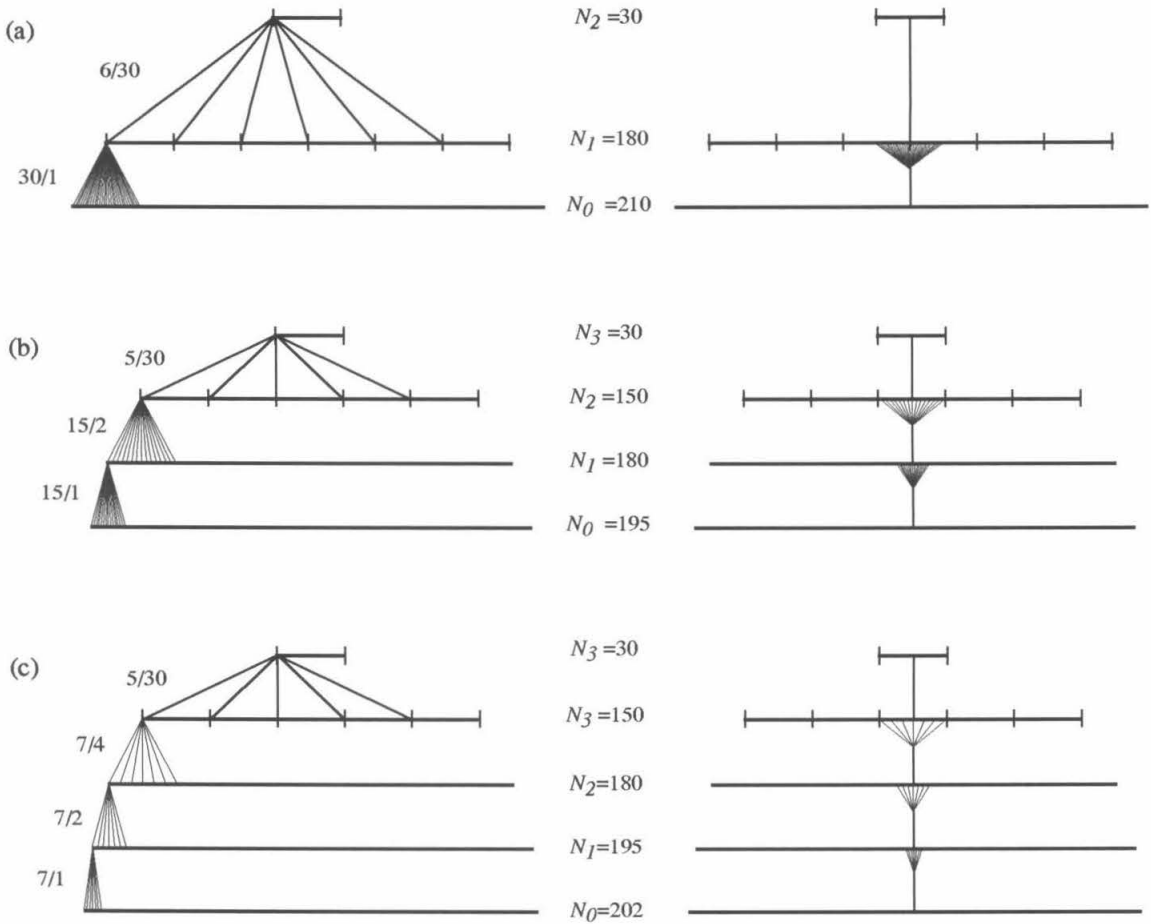


Figure 3.5: **Some possible multistage routing circuits.**

The number to the left of each stage ( $x/y$ ) denote the fan-in ( $x$ ) and the spacing between inputs ( $y$ ) for the stage.  $N_l$  denotes the number of sample nodes in 1D for each layer. The vertical tick marks in the next to last layer denote the modules, of size 30 (equal to the size of the window of attention). *a*, a nominal configuration, corresponding to the circuit of Figure 2.9. The fan-in/fan-out in the first stage can be reduced by adding more intervening stages, as in *b, c*.

distance in order to span 30 nodes in the next cortical area (cf. Fig. 3.4), which is rather implausible. We can reduce the fan-in/fan-out at each stage by breaking the routing circuit into more stages, as shown in Figure 3.5*b,c*.

Note that since the divergence is not as great a problem for the smallest scales, fewer stages of routing would be required for these levels of the stack, which is consistent with the direct projections observed from V1 to V4 for the central visual field only (Yukie and Iwai, 1985). For example, a high resolution level of the stack may pass merely through V1-V4-IT, while a low resolution level may be routed through V1-V2-V4-IT. Also, since the spacing between inputs doubles at each stage, the circuit is consistent with the observed doubling in receptive field size and progressively greater divergence in connectivity patterns at each stage, as mentioned earlier. (There is a big jump at the top, though, but this is not a necessity, just a minimal configuration.)

The particular scheme adopted here of dedicating a distinct routing stream for each scale and switching between scales at the top stage is not a strict requirement. It is quite possible that switching between scales could be done at earlier stages of the circuit, thus allowing the resources at the top stages to be shared between scales. The scheme proposed here has the advantage that the routing circuits for each scale can work independently, which allows the possibility that the window of attention could be set in advance for the upcoming attentional fixation at that scale while attending somewhere else at a different scale. More importantly, it allows for a wide variety of features to be computed in parallel at different scales in the intermediate stages (V2,V4), which would provide a richer set of “primitives” for object matching within the window of attention.



## 3.2 Control substrates

Postulating the substrates for control requires that we specify 1) the anatomical location of the control neurons, 2) bottom-up sources for driving the control neurons during blob search—i.e., the “saliency map,” and 3) top-down sources for driving the control neurons during recognition. We consider here each of these issues in turn.

### Control neurons

The pulvinar nucleus of the thalamus is proposed as a major source of the control signals for routing information through the cortex. Since the pulvinar is reciprocally connected to all areas in the form pathway, it is a good candidate for modulating information flow from V1 to IT. The pulvinar also receives a massive projection from the superior colliculus, which is known to encode the direction of saccade targets and may also be involved in setting up attentional targets (Gattass and Desimone, 1991, 1992; Posner and Petersen, 1990). In addition, neurophysiological studies (Petersen et al., 1985, 1987), lesion studies (Desimone et al., 1990; Bender, 1988; Rafal and Posner, 1987) and PET studies (LaBerge and Buchsbaum, 1990; Corbetta et al., 1991) of the pulvinar suggest that it plays a role in engaging visual attention, or filtering out unattended stimuli. A detailed analysis of these studies and their relevance to the model is provided in Appendix C.

A subcortical nucleus such as the pulvinar also has the important property of being spatially localized while at the same time being able to communicate with vast areas of the visual cortex. The relative proximity of pulvinar neurons to each other would facilitate the competitive and cooperative interactions among the control neurons which are necessary for choosing a single position and size of the attentional window and for maintaining spatial relationships within the window. Although it is not known whether such interactions exist among pulvinar neurons, Ogren and Hendrick-

son (1979) have reported the existence of interneurons with elaborate dendritic trees approaching  $600\mu$  in diameter, which could mediate communication among pulvinar neurons. In addition, neuropharmacological experiments by Petersen et al. (1987) have shown that enhancing or depressing inhibition within the pulvinar can respectively slow down or speed up attentional shifts, which is suggestive of lateral inhibitory connections within the pulvinar. An analogous function might also be served by the reticular nucleus of the thalamus, which is an inhibitory structure through which pulvinar neurons project on their way to the cortex. One study in *Galago* (Conley and Diamond, 1990) has shown that the pulvinar projects quite diffusely into the reticular nucleus, which would be desirable for a winner-take-all type circuit.

To first order, it would make sense for each stage of the routing circuit to have its own set of control neurons. The anatomical subdivisions of the pulvinar correspond roughly with this scheme, insofar as the inferior pulvinar projects mainly to lower areas (V1, V2) and the lateral and medial pulvinar to higher areas (V4, IT). (A small fraction of inferior pulvinar neurons (10%) have been shown to project to both V1 and V2 (Kennedy and Bullier, 1985).) The control neurons for the lower stages would need to compete only locally, since these stages would be more concerned with making local adjustments in the position and scale of the window of attention. Control neurons at the highest stage would need to compete globally, since these stages are setting the position and scale of the window of attention for the entire scene.

The number of control neurons that would be required for the routing circuit depends on how many cortical synapses are modified by each control neuron. Theoretically, the minimal number of control neurons is given by

$$\# \text{ of control neurons} = \frac{(\# \text{ of output nodes}) \times (\text{fan-in per node})}{(\# \text{ of synapses per control block})}.$$

Assuming that the control blocks comprise a maximum of 1000 synapses each, then

the minimum number of control neurons required for each stage for each scale of the routing circuit would be on the order of the number of output nodes of each stage (since the maximum fan-in per node is about 1000). The total number of control neurons required for each stage across all scales would be obtained by multiplying this number by 6 (the number of levels in the stack). Thus, for the circuit of Figure 3.5*b,c*, approximately  $6 \times 40,000 = 240,000$  control neurons would be required for the first stage, approximately  $6 \times 30,000 = 180,000$  for the second stage, etc., which is well within the estimated number of neurons in the pulvinar.<sup>3</sup> However, each output node in the circuit actually corresponds to a multitude of neurons representing various features, such as local orientation, texture, etc. Thus, each pulvinar control neuron would require an additional fan-out for controlling the inputs to all the neurons corresponding to an output node. Since there may be hundreds of neurons for each node, the pulvinar neurons would need to amplify their fan-out via other neurons (a fan-out of 100,000 for pulvinar neurons is probably too large to be plausible). This could possibly be subserved by neurons residing in the deeper layers (5 and 6) of the cortex, as proposed previously by Van Essen and Anderson (1990). Control might then be implemented in a hierarchical fashion, with each pulvinar neuron specifying how information is routed between nodes, and cortical control neurons specifying how information is routed between the neurons belonging to each node.

### **Bottom-up control sources (saliency map)**

As discussed in the previous chapter, the control neurons may be driven by bottom-up, or low-level signals—such as “pop-out” in motion, color, texture, etc.—in order to direct the window to salient regions of the input. Since each of the saliency measures

---

<sup>3</sup>The pulvinar has somewhat lower neuronal density than the LGN, but also is several times larger. Since the LGN contains  $\sim 10^6$  projection neurons, this would constitute a reasonable lower bound for the number of neurons in the pulvinar.

may be computed in a separate cortical area, it would be advantageous to fuse them together into a single representation of saliency—such as the “saliency map” proposed by Koch and Ullman (1985). Two possible anatomical substrates for such a saliency map are the posterior parietal complex and the superior colliculus.

The posterior parietal complex (PP) is known to play an important role in attentional processes. Some studies have reported that neurons in this area show an enhanced response to attended targets within their receptive fields, even when no eye movements are made (Bushnell et al., 1981). Others have reported a 3-fold enhancement for *unattended* targets when the animal is in an attentive state (Mountcastle et al., 1981), or even a relative suppression for attended targets as opposed to unattended targets (Robinson et al., 1991; Steinmetz et al., 1992). In addition, lesion studies show that damage to the parietal lobe in humans hinders the ability of other objects in the field of view to attract the attentional window away from the currently attended location (Posner et al., 1984). Taken together, these results suggest that PP may be representing the locations of potential attentional targets, as opposed to targets already being attended. This is exactly the property we would expect of a saliency map. If this is the case, then these neurons would drive the control neurons in the pulvinar which compete to select the locus of the window of attention.

This proposal contains at least two potential weaknesses, however. One possible drawback is that PP neurons typically have relatively long latencies— $\sim 100$  ms (Robinson et al., 1978; Duhamel et al., 1992)—which is hard to reconcile with psychophysical data that imply that attention takes  $\sim 50$  ms to move to a new location in the visual field (Saarinen and Julesz, 1991; Nakayama and Mackeben, 1989). A possible solution to this dilemma is that the superior colliculus may supplement PP by acting as a crude saliency map, but with a quicker response time due to its direct retinal input (the latency of neurons in the superficial layers of the superior colliculus is in the range of 40-50 ms; Goldberg and Wurtz, 1972). The other drawback of using

PP as a saliency map is that the currently available anatomical data seem to offer relatively few direct pathways by which PP could influence those pulvinar neurons that would be able to modulate connection strengths in the “form” pathway, since PP and V4/IT connect with rather segregated portions of the pulvinar (Baleydier and Morel, 1992). However, there do exist indirect pathways, such as through the superior colliculus, that may provide viable alternatives. Another possibility of course is that the salience measures made along different feature dimensions in different cortical areas could drive the control neurons in the pulvinar directly via cortico-fugal pathways.

### **Recognition-guided control sources**

During recognition, top-down influences will need to take over to refine the position and size of the attentional window for object matching, as depicted in Figures 2.15 and 2.21. The pulvinar would thus need to alternate between bottom-up and top-down sources of input as attention moves from one object to the next. Top-down guidance during recognition would presumably be propagated to lower cortical areas via cortico-cortical feedback pathways from IT, or IT may influence control neurons in the pulvinar directly via its diffuse projections to many nuclei within the pulvinar. Alternatively, IT could supply top-down guidance primarily to *cortical* control neurons via the feedback pathways. Under this scenario then, the pulvinar’s role would be analogous to that of a general in an army—coarsely specifying a plan of action, which the cortical control neurons refine into a concise remapping under top-down, or object-based guidance from IT.

### 3.3 Gating mechanisms

In order for the control neurons to modulate connection strengths we must postulate some possible neuronal gating mechanisms. This is not at all a new idea, as neural gating mechanisms are believed to play an important role in many aspects of nervous system function. For example, the extent to which a noxious stimulus is perceived as painful varies greatly as a function of one's emotional state and other external factors. This is subserved at least in part by gating mechanisms in the spinal cord, where descending fibers from the raphe nuclei form part of a control system that modulates pain transmission via presynaptic inhibition in the dorsal horn (Fields and Basbaum, 1978). Gating mechanisms are also thought to play an important role in sensori-motor coordination; for example, there are many instances in which spinal cord central pattern generators gate sensory inputs according to the phase of the movement cycle in which the input occurs (Sillar, 1991). A somewhat different form of gating seems to take place in the LGN, where thalamic relay cells exhibit two distinct response modes: a *relay* mode, in which cells tend to more or less faithfully replicate retinal input, and a non-relay *burst* mode, in which cells burst in a rhythmic pattern that bears little resemblance to the retinal input (Sherman and Koch, 1986). In this instance, the reticular nucleus of the thalamus is thought to be the source of the signal that switches the LGN into the non-relay burst mode.

Although there is as yet no explicit evidence for gating mechanisms in the visual cortex, there are several possible biophysical mechanisms that would allow control neurons to gate synapses along the V1-IT pathway. Pre-synaptic inhibition, as in the spinal cord, would probably provide the most localized gating effect. However, to date there exists no morphological evidence for this type of synapse in the visual cortex (Berman et al., 1992). Postsynaptically, a control neuron could decrease or possibly nullify the efficacy of a cortico-cortical synapse via shunting inhibition. Evidence for

this type of mechanism playing a role in orientation or direction tuning is mixed, with some for (Volgushev et al., 1992; Pei et al., 1992) and some against (Douglas et al., 1988). Another possible post-synaptic gating mechanism could be realized via the combined voltage- and ligand-gated NMDA receptor channel, which has been shown to play an important role in normal visual function (Nelson and Sur, 1992; Miller et al., 1989). In this case, a control neuron could effectively boost the gain of a cortico-cortical synapse by locally depolarizing the membrane in the vicinity of the synapse. Also, there exist voltage-gated  $\text{Ca}^{++}$  channels in dendrites (Llinas, 1988) that could provide non-linear coupling between inputs. Evidence for non-linear interactions of this type have been reported for synaptic inputs into layer 1 of neocortex (Cauller and Connors, 1992). All of these mechanisms, and possibly others, offer a multiplicative-type effect that is suitable for gating information flow through the cortex (see also Koch and Poggio, 1992).

Under an inhibitory gating scheme, such as shunting or pre-synaptic inhibition, the control neurons would need to become active only when attention is actively engaged on an object. The finer the resolution desired within the window of attention, the more the control neurons would need to be engaged. The absence of any activity on the control neurons would correspond to the all-connections-open, or inattentive state, in which neurons in IT would exhibit the very large receptive fields observed in anesthetized or inattentive animals (Gross et al., 1972; Desimone et al., 1984).

Under an excitatory gating scheme, such as via NMDA receptors, one would need to hypothesize the existence of a gain control mechanism working in concert with the control neurons. When no control signals are provided, cortical input would be rather weak, and the firing threshold of pyramidal cells should be lowered to let all information through. When control signals are present to boost the gain of individual synapses, however, the threshold should be raised. This way, the unboosted synapses will be essentially suppressed to a relatively low strength. Threshold adjustment could

perhaps be subserved by chandelier cells, which make strong inhibitory connections exclusively onto the axon initial segment of pyramidal cells (Douglas and Martin, 1990b). Evidence that gain control mechanisms indeed exist in visual cortex has been established in previous physiological studies (Ohzawa et al., 1982; Pettet and Gilbert, 1992).

From a computational viewpoint, gating of inputs within individual dendrites provides a much higher degree of flexibility than would merely gating the outputs of pyramidal cells. Since the output of a pyramidal cell may branch to several cortical areas and make synaptic connections to a multitude of neurons, any modulation of the cell's output will simply be duplicated at all these subsequent input points. Gating inputs within the dendrites, on the other hand, allows the non-linear computation of many intermediate results ( $\sum_k c_k \Gamma_{ijk} I_j^{in}$ ) within the post-synaptic membrane, which can then be summed together within a single cell. This results in a computational structure that is orders of magnitude richer (Mel, 1992), and provides a higher degree of flexibility in sculpting patterns in connection space (cf. Fig. 2.3). The demonstrable computational advantage of dendritic gating mechanisms for visual processing motivates the need to specifically look for such mechanisms experimentally. (See also Desimone, 1992, for a discussion of output vs. input gating mechanisms.)



## Chapter 4

# Discussion

We shall see in this chapter that the proposed neurobiological correlates of the model lead to a number of neurophysiological, neuroanatomical, and psychophysical predictions that can be tested experimentally. We shall also describe what distinguishes this model from other network models that have been proposed for visual attention and invariant pattern recognition, and how the model can be generalized to be understood in a broader context. Finally, we shall briefly examine some of the unresolved issues that remain as topics for future research.

### 4.1 Predictions

#### Neurophysiology

The most obvious prediction of the dynamic routing circuit model is that the receptive fields of cortical neurons should change their position or size as attention is shifted or rescaled. This effect should be especially pronounced in higher cortical areas. Some support for this prediction comes from the neurophysiological findings of Moran and Desimone (1985) in areas V4 and IT of primate visual cortex. As schematized in

Figure 4.1, they found that if two bar-shaped stimuli were placed within the classical receptive field (CRF) of a V4 cell, and the animal was trained to attend to only one of them, then the cell's response to the unattended stimulus was substantially attenuated. This is what one would expect from our routing circuit, since the pathways between the cell and the unattended stimulus would be effectively disabled in this case (Fig. 4.1*c*). They also found that the V4 cell responded to an unattended stimulus anywhere within its CRF when the animal attended a stimulus outside the CRF. This effect is also predicted by the model, because once a V4 cell lies outside the region of interest in V4 it no longer needs to restrict its inputs (Fig. 4.1*d*). Indeed, other targets of V4, such as those in the posterior parietal cortex, would presumably be interested in the information from regions lying outside of the attentional beam.

While Moran and Desimone's findings offer some support for attentional modulation effects predicted by the model, they did not attempt to map receptive fields under different attentional conditions with any precision; thus, their results do not address the more specific effects predicted by the model. For a cell in one of the intermediate stages of a particular scale, one would expect a cortical receptive field to shift as the attentional window is translated, or to expand or shrink somewhat as the attentional window is made slightly larger or smaller respectively. In the highest cortical stages, beyond scale selection, we would expect to see dramatic size changes, as well as a shift in the spatial frequency tuning, for large changes in the size of the attentional window. These predictions can be tested by giving an animal a task that forces it to attend to a region of a specific size and location, and then probing the receptive field with a neutral (behaviorally irrelevant) stimulus to measure its extent. Preliminary results using such a paradigm suggest that the receptive fields of V4 cells do indeed translate toward attentional foci in or near the classical receptive field (Connor et al., 1993). However, the extent of the observed shift is modest (at most about  $2^\circ$  for an approximate  $5^\circ$  shift in window position). One would not expect a

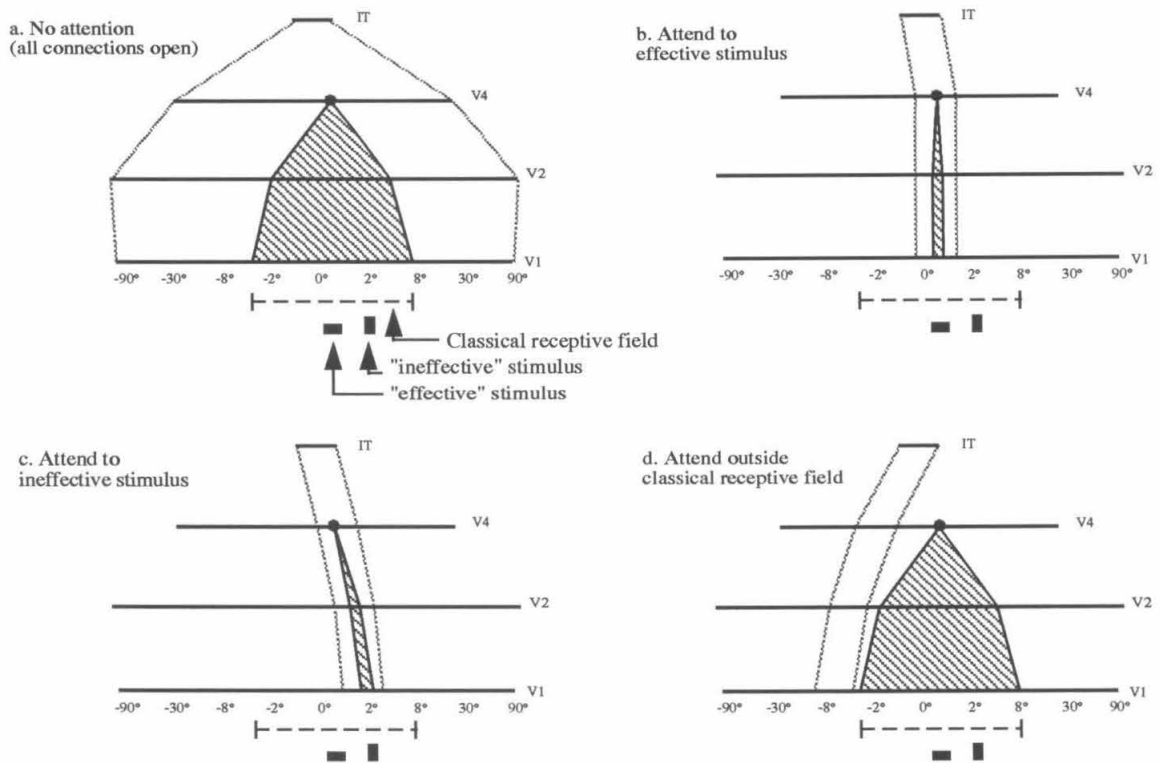


Figure 4.1: **The dynamic routing circuit interpretation of the Moran and Desimone (1985) experiment.**

The node in layer V4 indicates the cell under scrutiny. The hashed region indicates those connections to the cell that are enabled; the others are disabled. The bounds of the window of attention in each area are shown by the stippled lines. (a) In the non-attentive state, all connections will be open and the effective stimulus can excite the cell anywhere within its classical receptive field. (b) When attending to the effective stimulus, the cell's response should be unaltered since the neural pathways to the stimulus are still open. (c) When attending to the ineffective stimulus, the cell's response should decrease substantially since the neural pathways to the effective stimulus are gated out. (d) When attending outside the cell's classical receptive field, there is no need to gate the cell's inputs since it is no longer taking part in the process of routing information within the window of attention.

V4 cell to cover the entire shift, since it is situated in an intermediate stage of routing (see Figure 4.1), but still for a cell in the third layer of a three-stage routing circuit the average shift should come out to about half the total window shift. The lack of shift could possibly be due to cell relaxing back to its all-connections-open state, since due to various experimental considerations the receptive field is probed 200 ms after the presumed attentional fixation and the cell's response is integrated over a 300-500 ms period. Also, there seems to be little appreciable size change when the attended object is quadrupled in size, but since this is an intermediate area, it is unclear what if any change would be expected for such a large size change outside the scope of a routing circuit for a particular scale.

Another physiological prediction of the model is that lesions to the pulvinar, the hypothesized control center, should dramatically degrade attention and pattern recognition abilities. While there is substantial evidence linking pulvinar lesions to attentional defects (Desimone et al., 1990; Bender, 1988; Rafal and Posner, 1987), some pattern recognition abilities appear to be relatively unimpaired by pulvinar lesions (Bender and Butter, 1987; Nagel-Leiby et al., 1984; Chalupa et al., 1976; Mishkin, 1972). One possible reason for the apparent sparing of pattern recognition is that the tasks used in these studies generally were very simple, such as distinguishing a large 'N' from a 'Z' (Chalupa et al., 1976). It is conceivable that such a task could be carried out even when the fidelity of the remapping process has been compromised. A more rigorous test using stimuli that demand the full spatial resolution capacity of the window of attention would be better suited to test the effect of pulvinar lesions on recognition abilities. Pulvinar lesions would also be expected to diminish the result found by Moran and Desimone (1985) and Connor et al. (1993), and it would be interesting to repeat these experiments while reversibly deactivating the pulvinar.

The physiological responses to be expected from pulvinar neurons depend on how they are configured to gate information flow in the cortex. In an inhibitory gating

scheme, one would expect enhanced responses from pulvinar neurons projecting to areas of the cortex within and immediately surrounding the attentional beam, and little or no response from pulvinar neurons projecting to those areas of the cortex substantially outside the attentional beam. In an excitatory gating scheme, one would expect to find enhanced responses from pulvinar neurons projecting to areas of the cortex within the attentional beam only. Petersen et al. (1985) have reported such an enhancement effect for neurons in the dorsomedial portion of the pulvinar (which is connected with PP), but not in the inferior or lateral portion (which is connected to V1-IT). The lack of enhancement in these latter areas may be due to the fact that the task used in this experiment was very simple (detecting the dimming of a spot of light). Again, a more appropriate task would be one that fully taxes the capacity of the attentional window, as this would require the greatest participation from the control neurons in gating out irrelevant information.

## Neuroanatomy

The convergence and wiring constraints discussed in the previous chapter (Section 3.1) suggest that the anatomical divergence of intra-cortical connections should increase by roughly a factor of two in the intermediate stages, and that there should be evidence of progressively larger “modules” at higher cortical areas (Fig. 3.5). While there is considerable evidence in support of progressively greater divergence, patchiness, and modularity in ascending stages of the form pathway (Rockland 1992; Van Essen et al., 1986; Van Essen et al., 1990; Van Essen and DeYoe, 1993; DeYoe and Sisola, 1991; Felleman et al., 1992; Felleman and McClendon, 1991), more quantitative and higher resolution data are needed in order to confirm or contradict the proposed routing architecture. Most interesting of all would be to determine the topography and size (in terms of “sample nodes” of visual field) of the the observed modules in V2, V4,

and PIT.

Another anatomical prediction of the model is that the terminations of pulvinar-cortical projections should be suitably positioned for effective modulation of inter-cortical synaptic strengths. The pulvinar is known to project to the output layers (2,3) of V1 and to both the input and output layers (3,4) of extrastriate areas V2, V4, and IT (Ogren and Hendrickson, 1977; Rezak and Benevento, 1979; Benevento and Rezak, 1976). These synapses are suspected to be excitatory since they are of the asymmetric type (in layers 1 and 2, Rezak and Benevento, 1979). However, it is not known whether the pulvinar afferents make synapses with inhibitory interneurons or directly onto the dendrites of pyramidal cells.

Finally, the model predicts that there should exist lateral inhibitory and excitatory connections within the pulvinar in order to enforce the constraint of preserving spatial relationships within the window of attention. This prediction is partially supported by the existence of interneurons within the pulvinar (Ogren and Hendrickson, 1979), but it remains to be seen if the axons of projection neurons have collaterals that spread horizontally within the pulvinar, or to what extent the reticular nucleus of the thalamus might subserve this role.

## **Psychophysics**

The fixed window size at the top layer of the routing circuit implies that the spatial resolution of the window of attention is limited. Thus, a large window of attention should have rather poor spatial resolution, whereas a small window of attention should have rather high spatial resolution. It has been proposed by C. Anderson that the size of this window may be on the order of approximately 30x30 sample nodes, an estimate based on psychophysical studies of spatial acuity and pattern recognition

(Van Essen et al., 1991; Campbell, 1985).<sup>1</sup> However, one problem with this analysis is that the critical data were derived from experiments in which visual attention was not explicitly controlled. In particular, most of the experiments had display times long enough to permit multiple shifts of attention (although we doubt that this would have been a major contaminating factor in most cases). On the other hand, those experiments that have been directed at studying the amount of “resources” allocated during visual attention have largely ignored the issue of spatial resolution. For example, various studies have reported evidence for a “zoom lens” model of attention in which the density of processing resources decreases as the size of the attentional window increases (Eriksen and St. James, 1986; Shulman and Wilson, 1987). However, these experiments were not designed to measure spatial resolution explicitly. Also, Verghese and Pelli (1992) have attempted to measure the information capacity of the window of attention, which they conclude to have an upper bound of about 50 bits. However, they studied only two tasks—detecting a non-moving target among moving distractors, or detecting a non-flashing square among flashing squares—neither of which is well suited for measuring spatial resolution. In a more recent study, Farrell and Pelli (1993) have reported that localization, but not identification, suffers as the window size increases. The fact that identification is not affected may be attributable to the fact that simple or overlearned forms were used in this task (e.g., black or white checks, or digits among letters), thereby reducing identification to an essentially pre-attentive task. It would be interesting to do this with more complicated, non-overlearned shapes. The fact that localization becomes worse as the window size increases is consistent with the model, and it would be interesting to see quantitatively how the size of the positional errors correlates with attentional window size.

---

<sup>1</sup>This prediction shares a basic similarity to Nakayama’s (1991) “iconic bottleneck” theory, although his estimate ( $\sim 100$  pixels total) is somewhat lower.

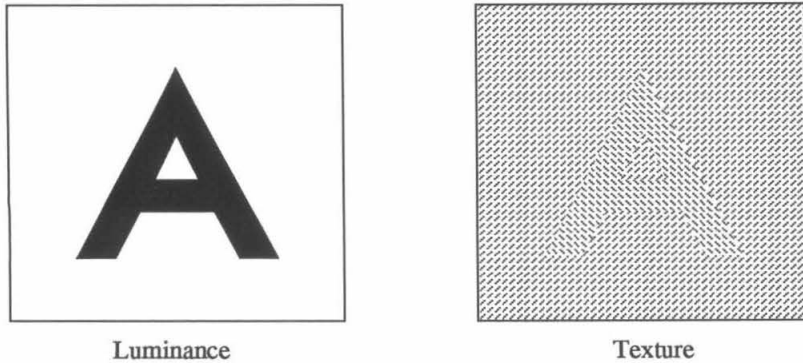


Figure 4.2: **The meaning of “cycles per object.”**

An object composed of high frequency texture elements is still recognizable even if the frequency of the elements is beyond 15 cy/window. What is important here is the structure of the ‘A’ and the amount of spatial detail required to specify it. When this exceeds 15 cy/object, discrimination and/or recognition should suffer.

The kind of experiment that would be most useful in determining the resolution of the attentional window would be one that tested pattern discrimination ability as a function of the position, size, and resolution of an object. In this case, the model predicts that performance would drop off sharply once the spatial frequency content of the stimulus exceeded approximately  $15 \times 15$  cycles per object (as dictated by the Nyquist theorem). However, it is important not to interpret “spatial frequency” literally here: A shape composed of high frequency elements would be perfectly discriminable due to non-linear or texture processing. What is important is the resolution of the overall structure, as illustrated in Figure 4.2.

The model also makes some interesting predictions with regard to the dynamics of visual attention. For example, once a location has been attended to in the visual field it should be difficult to stay there or immediately revisit the site, because the the control neurons and/or saliency map neurons corresponding to that part of the visual field would be transiently inhibited from firing. There is some evidence



for such a mechanism, in that involuntary attentional fixations tend to be transient (Nakayama and Mackeben, 1989) and appear to be inhibited from return (Posner and Cohen, 1984). The amount of time that it takes the attentional window to shift from one location to another would be expected to be roughly independent of the distance between locations. Unlike eye saccades, there is no obvious reason why the control neurons should sequence through all intervening positions of the attentional window. Rather, moving the locus of attention would require merely inhibiting the current control state and activating a new one. This prediction is most consistent with Remington and Pierce's (1984) study showing time-invariant shifts of visual attention, although other studies (e.g., Tsal, 1983) are in disagreement (see also Eriksen and Murphy (1987) and Cave (1991) for a critical commentary on these and other studies). On the other hand, if attention were to actually track a stimulus, then one would indeed expect a smooth transition of activity across the control neurons. It is interesting to note that Cavanagh (1992) has discovered some forms of visual stimuli that produce a motion percept only when tracked with attention. We speculate that the progression of activity across the control neurons is what underlies one's perception of motion in such cases.

## **4.2 Comparison with other models**

### **Control vs. synchronicity**

A number of other models of visual attention and pattern recognition have been proposed that rely on the synchronous firing of neurons in order to effectively change connection strengths (e.g., von der Malsburg and Bienenstock, 1986; Crick, 1984; Crick and Koch, 1990). We contend that a key disadvantage of such approaches is that information about the effective connection-state at any one point in time is

not explicitly encoded anywhere in the system. In the routing circuit model, this information is encoded explicitly in the activities of the control neurons, which then allows it to be utilized advantageously in a number of ways.

One way that information about connectivity can be utilized is in constraining the active connections between retinal- and object-based reference frames to be in accordance with a global shift and scale transformation. This constraint is incorporated in our model via the competitive and cooperative interactions among the control neurons (Equation 2.11). During object recognition, this constraint drastically reduces the number of degrees of freedom in matching points between the retinal and object-centered reference frames, because once a few point-to-point correspondences have been established, the number of potential matches between other pairs of points is greatly reduced. In machine vision, this is known as *viewpoint consistency constraint*, and it has proved to be a powerful computational strategy for object recognition systems (Lowe, 1987; Hinton, 1981b).

Another advantage of having knowledge of the active connection state readily available is that the ensemble of control neurons together form a neural code for the current position and size of the window of attention. Therefore, information about the position and size of an object can be obtained by simply reading out the state of the control neurons. In addition, it would also be possible for the control neurons to warp the reference frame transformation in order to form object representations that are invariant to distortion (e.g., hand written digits), in which case information about the particular shape of the object (e.g., its slant or style) could also be preserved. Note that such information is typically lost in networks that utilize feature hierarchies of complex cells (Fukushima, 1980, 1987; LeCun et al., 1990) or Fourier transforms (e.g., Pollen et al., 1971; Cavanagh, 1978, 1985) for forming position-, scale-, and/or distortion-invariant representations.

The routing circuit model can also explain how attention may be directed “at

will,” or by other modalities, to the extent that those areas of the brain having access to the control neurons (such as parietal cortex) can directly influence where attention is directed. This also provides a convenient format for mediating the access to control among various competing demands. While such forms of top-down control are not impossible to incorporate in models based on synchronicity-gated connections, its implementation would seem to be less straightforward.

### **Control-based network models**

A number of other network models of attention and recognition have also utilized the concept of control neurons for directing information flow. Niebur et al. (1993), Desimone (1992), Laberge (1990; 1992), Ahmad (1992), and Posner et al. (1988), among others, have proposed models that involve the pulvinar as a control site for routing information from a select portion of the visual scene. In addition, Tsotsos (1991) and Mozer (1992) have proposed somewhat more abstract connectionist models that utilize gating units to control attention. However, none of these models explicitly preserve spatial relationships within the window of attention, which is presumed here to be a critical component of the routing process.

Hinton and Lang (1985) and Sandon (1988; 1990) have proposed control-based models that do preserve spatial relationships within the window of attention and share the same basic principle as the model presented here—i.e., remapping object representations from retinal into object-centered reference frames via a third set of units (equivalent to control neurons). Although these models attempt to explain various psychophysical data, they do not contain the necessary level of neurobiological detail to give them strongly predictive value in biology.

Postma et al. (1992) have proposed a neural model based upon the original shifter circuit proposal (Anderson and Van Essen, 1987) to account for translational

invariance in visual object priming (Biederman and Cooper, 1992). This model shares many similarities to the routing circuit model, including top-down, or template-driven control, but it differs in the specifics of the gating and control structure. Postma’s circuit uses distinct control neurons for each synapse that are locally connected to each other in a “gating lattice” that settles upon a global shift for the circuit. The model also utilizes a series of stages of local, winner-take-all circuits to control the shift, which shares a basic similarity to the hierarchical control scheme proposed here.

### 4.3 Generalizations of the model

#### Bayesian interpretation

As discussed in the introduction, the routing circuit model can be viewed as a means for generating separate, independent representations of *what* and *where*, which results in an efficient usage of computational resources. As we shall see here, the energy functional that we use to express the “goal” of the network has an interesting Bayesian interpretation that illustrates another advantage of separating *what* and *where*.

The total energy functional of the network, in its most general form, is

$$E_{total} = -\beta_1 \mathbf{I}^{out} \mathbf{V} - \beta_2 \mathbf{V} \mathbf{T} \mathbf{V} - \beta_3 \mathbf{c} \mathbf{U} \mathbf{c}. \quad (4.1)$$

Note that we have switched into vector notation here to eliminate the bulky summation signs. The first term of Equation 4.1 measures the similarity between the output of the routing circuit,  $\mathbf{I}^{out}$ , and the “desired” output,  $\mathbf{V}$ . In blob search mode,  $\mathbf{V}$  is set equal to the blob function,  $\mathbf{G}$ , and in recognition mode  $\mathbf{V}$  is the output of the associative memory. The second term is the associative memory energy term, which only takes effect in recognition mode; in blob search mode,  $\beta_2 = 0$ . The third term is the constraint term for the control neurons,  $\mathbf{c}$ , that enforces a single scale and position

of the window of attention. We can write  $\mathbf{I}^{out}$  in terms of the input image,  $\mathbf{I}^{in}$ , and the control neurons as

$$\mathbf{I}^{out} = \mathbf{c}\mathbf{\Gamma}\mathbf{I}^{in} \quad (4.2)$$

where  $\mathbf{c}$  and  $\mathbf{\Gamma}$  are a concatenation of the control neurons and coupling coefficients,  $\Gamma_{ijk}^l$ , in all stages of the routing circuit. Thus, we can write  $E_{total}$  as

$$E_{total} = -\beta_1 \mathbf{V}\mathbf{c}\mathbf{\Gamma}\mathbf{I}^{in} - \beta_2 \mathbf{V}\mathbf{T}\mathbf{V} - \beta_3 \mathbf{c}\mathbf{U}\mathbf{c}. \quad (4.3)$$

Now, in plain English we can say that the goal of our model network is to infer the position, size, and identity of objects in an image, given the image data. Or, in probabilistic terms, we can say that that we wish to maximize

$$P(\text{WHAT,WHERE}|\text{IMAGE}) \quad (4.4)$$

where WHAT denotes the identity of the object, WHERE denotes the position and size of the object, and IMAGE denotes the image data. Expanding Equation 4.4 according to Bayes rule gives us

$$\begin{aligned} P(\text{WHAT,WHERE}|\text{IMAGE}) &\propto P(\text{IMAGE}|\text{WHAT,WHERE})P(\text{WHAT,WHERE}) \\ &= P(\text{IMAGE}|\text{WHAT,WHERE})P(\text{WHAT})P(\text{WHERE}). \end{aligned} \quad (4.5)$$

The last step assumes that WHAT and WHERE are statistically independent (i.e., that any given object is equally likely to appear at any location and size), which is going to be true for the most part.

In our network, the IMAGE is  $\mathbf{I}^{in}$ , WHAT is expressed in the ensemble of activity in  $\mathbf{V}$ , and WHERE is expressed in the ensemble of activity in  $\mathbf{c}$ . Thus, if we make the

following equivalences via the Gibb's distribution

$$P(\text{WHAT,WHERE}|\text{IMAGE}) \propto e^{-\beta E_{total}} \quad (4.6)$$

$$P(\text{IMAGE}|\text{WHAT,WHERE}) \propto e^{\beta_1 \mathbf{V} \mathbf{c} \Gamma \mathbf{I}^{in}} \quad (4.7)$$

$$P(\text{WHAT}) \propto e^{\beta_2 VTV} \quad (4.8)$$

$$P(\text{WHERE}) \propto e^{\beta_3 cUc} \quad (4.9)$$

then we can see that Equation 4.3 is just the logarithm of Equation 4.5, and so minimizing the energy function,  $E_{total}$ , will also tend to maximize  $P(\text{WHAT,WHERE}|\text{IMAGE})$ .

The terms  $P(\text{WHAT})$  and  $P(\text{WHERE})$  are the “priors,” which provide the prior probability of WHAT and WHERE before we know anything about the image. The term  $e^{\beta_2 VTV}$  gives a high probability to any state of the  $\mathbf{V}$  corresponding to a known object, and the term  $e^{\beta_3 cUc}$  gives high probability to any state of the  $\mathbf{c}$  corresponding to a unique position or size of the window of attention. The term  $P(\text{IMAGE}|\text{WHAT,WHERE})$  is the “likelihood,” which expresses how likely the IMAGE could have arisen from a particular WHAT and WHERE. The term  $e^{\beta_1 \mathbf{V} \mathbf{c} \Gamma \mathbf{I}^{in}}$  gives a high probability to that state of the  $\mathbf{c}$  and  $\mathbf{V}$  that can “explain”  $\mathbf{I}^{in}$  (measured by taking the inner product  $\mathbf{V} \mathbf{c} \Gamma \cdot \mathbf{I}^{in}$ ).

Finding the WHAT and WHERE that maximizes the posterior,  $P(\text{WHAT,WHERE}|\text{IMAGE})$ , requires optimizing over a huge search space. The strategy we adopt in the routing circuit for dealing with this dilemma is to first set  $\beta_2 = 0$  (no recognition) and let  $\mathbf{c}$  evolve while holding  $\mathbf{V} = \mathbf{G}$ . This essentially moves  $\mathbf{c}$  into a good initial state. From this point, we turn on  $\beta_2$  and let  $\mathbf{V}$  evolve. In terms of the probabilistic framework, the network hill climbs on  $P(\text{WHAT,WHERE}|\text{IMAGE})$  initially along the WHERE axes, and then along the WHAT axes. The reason we can do this is that the statistical independence of WHAT and WHERE allows the use of rather primitive, pre-attentive measures to guess the WHERE without knowing WHAT. More

generally though, the pre-attentive measures can also help guess WHAT. For example, measures such as color, texture, or convexity may be able to narrow down the class of objects to search among. In this framework, then, “attention” can be understood as a heuristic that exploits the statistical independence of WHAT and WHERE in order to make an extremely computationally intensive problem tractable with limited resources.

## The big picture

How are the various “snapshots” obtained by the window of attention eventually incorporated in order to form an overall percept of a scene? One possibility, initially suggested by Hinton (1981b), is that a compact representation of each object may be maintained in the form of the activities on a set of neurons within a “scene buffer.” My own rendition of this scheme is illustrated in Figure 4.3. One can essentially think of the scene buffer as a spatially indexed RAM (random access memory). Each attentional fixation writes its contents into a different part of the buffer, depending on the position and size of the attentional window as well as the orientation of the eyes, head, and body with respect to the environment. (See also Baron, 1987, for another variation on this theme.)

An important property of the model advanced here is that all the features belonging to an object, such as the eyes, nose, and mouth of a face, are bound together as an entity, and that a symbolic code for the object is entered into the spatial buffer. An alternative scheme that may be favored by a “feature *Gemischist*” would be to collect each feature with an attentional snapshot independently, and then enter these into the scene buffer with the proper spatial relationships. This would avoid the earlier mentioned problem of confusing the spatial relationships of features (Fig. 1.3), but it would not be very efficient because it would require accurate storage of the pointers

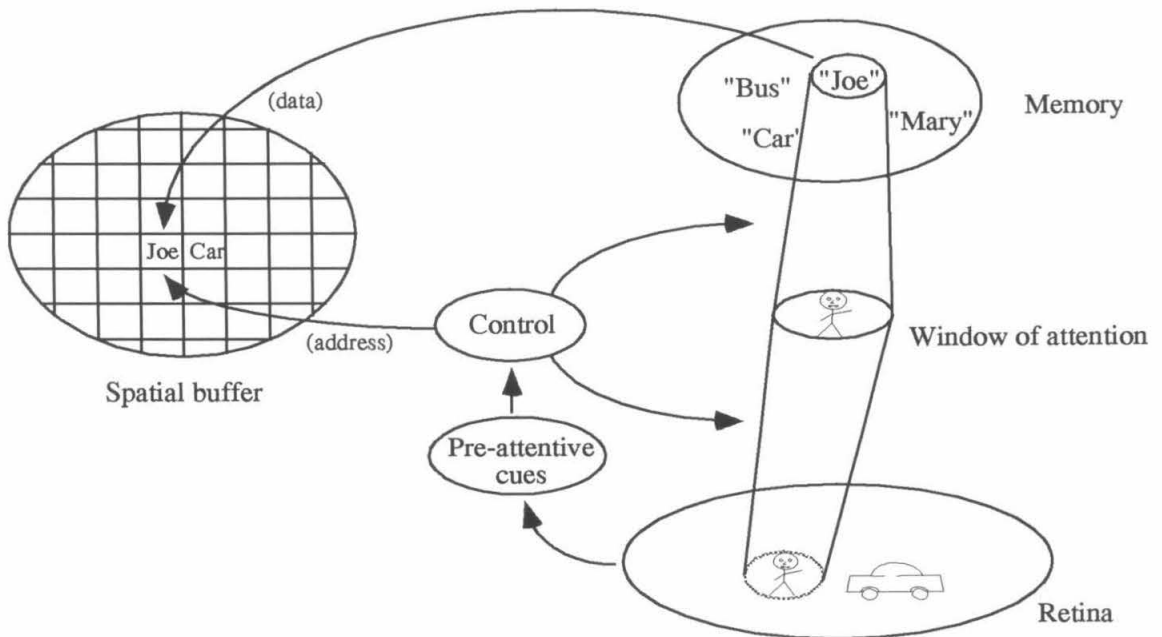


Figure 4.3: **The big picture.**

Pre-attentive cues drive the control neurons, which direct attention to a restricted region of visual space and “object space.” Once the object has been recognized, its label is written into the spatial buffer at the address indexed by the control neurons.



for each feature of an object. The advantage of the scheme proposed here is that the spatial relationships that recur repeatedly in the visual world are encoded explicitly (presumably by a population of cells in IT). The spatial buffer is then reserved only for essentially non-predictable structure. For example, Joe may only rarely appear with his car—indeed, he appears at home, in the lab, riding a bicycle, etc. What is the same in all these situations are the spatial relationships among the features of Joe’s face. But if Joe were to recur repeatedly with some other object, say, a particular shirt or style of clothing, then one might eventually adopt an explicit representation, or a new code, for this combination. Subsequent presentations of this combination would then result in the code being entered into the scene buffer, instead of the individual parts. The strategy advocated here, then, is that the natural redundancies in the input should be captured explicitly, leaving the dynamic representation of the spatial relationships (in the spatial buffer) for those combinations that rarely or spontaneously occur in our everyday visual experience.

## **Control as a general strategy for neural computation**

Another of the general themes advanced in this thesis is that the utilization of explicit control neurons is a useful computational principle for visual processing. This principle may be employed by the brain in other domains as well. A different perspective of dynamic control is illustrated in Figure 4.4. In most “static” neural network models, the output of a neuron is computed by forming the inner product of a weight vector,  $\bar{w}$ , with the inputs to the neuron, and then passing the result through a non-linearity. The weight vector may change on a slow time scale in order to optimize the network for performing a certain task, but typically  $\bar{w}$  remains fixed over the relatively short time in which the task is actually performed (e.g.,  $< 1$  sec). By having control neurons available to modify  $\bar{w}$  on a short time scale, the computation being carried out by

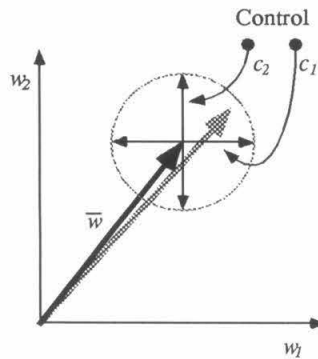


Figure 4.4: **A more general way of viewing control.**

A weight vector with two components,  $w_1$  and  $w_2$ , is shown. Control neurons  $c_1$  and  $c_2$  modulate each of these components, respectively, to dynamically change the weight vector. Thus, the weight vector may be able to occupy any region within the circular outline in order to optimize the network for the particular input and task at hand.

the network can be dynamically reconfigured and optimized for the particular task at hand. This added degree of flexibility reduces the neural resources required for solving a complicated task, since it is no longer necessary to have dedicated, specialized networks with fixed connections to deal with each variation of a task.

A particularly interesting example is in the motor system. As Lashley (1942) pointed out long ago, somehow we are able to route prototypical motor commands to any set of limbs (he termed this “motor equivalence”). For example, one may write with either hand, or even with the mouth, and the style of writing is more or less preserved. This would seem to imply that the trajectories for writing the name may be stored in a canonical reference frame without reference to a particular limb, and that these trajectories are dynamically routed out to any desired set of actuators. This problem could essentially be considered the inverse of the sensory routing problem (see Van Essen et al., 1994).

## 4.4 Unresolved issues

The dynamic routing circuit as described in this paper is intended as a “zero-th order” model, and as such many details have been neglected or oversimplified. Here we outline some of the more important unresolved issues that remain as topics for future research.

### Features instead of pixels

As already noted, one key neurobiological characteristic neglected in the present model is the known preponderance of feature selective cells in the visual cortex. V1, for example, is known to contain cells tuned for various orientations, and V2 and V4 contain cells that seem to be tuned for more complex stimuli (von der Heydt and Peterhans, 1989; Gallant et al., 1993). It has been assumed here that this will not significantly affect the general nature of the routing process, but in order to make our predictions more precise, it would be helpful to have a better idea of what features are computed. For example, what is role of position invariance in complex cells? Also, dynamic routing need not necessarily be restricted to the space domain, but could work over feature domains as well.

### Feedback pathways

Although the model relies upon feedback pathways for top-down control during recognition, the information processing role of these pathways has more or less been ignored. Mumford (1992) has sketched a theory proposing that the role of these feedback pathways is to relay the interpretations of higher cortical areas to lower cortical areas in order to verify the high-level interpretation of a scene. Such a mechanism would obviously be of use for Step 4 of our proposed strategy for an autonomous visual

system. Under this scenario, it would be necessary to route information flow within the feedback pathways as well in order to ensure that the high-level interpretation is matched against the appropriate region within the cortical area below (i.e., within the window of attention). Another possible attentional role for information flow in the feedback pathways may be to refine the tuning characteristics of lower-level cortical cells based upon the interpretations made in higher cortical areas (see, for example, Tsotsos, 1991).

## **Pop-out in multiple dimensions**

In the simple autonomous visual system we have proposed, “blobs” were the only salient features used to attract the window of attention. Presumably, other salience measures—such as pop-out due to motion or texture gradients—would provide a much richer and more robust system for pre-attentively guessing the size and position of potential objects. It will also be of interest to see how these measures could be used to position the window in object space, as illustrated in Figure 4.3.

## **Rotation and warp**

Our model accounts for how reference frames can be shifted and rescaled, but it does not address rotation and other distortions (e.g. hand-written characters). The ability to rotate or warp reference frames could probably be included in the model without much difficulty, since this would just involve another form of routing. Moreover, for foveated objects the log-polar representation in V1 would convert rotations into approximate linear shifts on the cortex (Schwartz, 1977), which may facilitate the routing.

### **3D objects**

How are 3D objects represented neurally, and how is information in the retinal reference frame transformed to match this representation? One possibility, as advanced by Poggio and Edelman (1990), is that 3D objects are actually represented by a few characteristic two-dimensional views, and that a match to the retinal representation is achieved by interpolating among these views. In this case, the routing circuit would be required to properly reposition and rescale the object so that the interpolation could take place.

### **Learning**

Although the model we have presented here is neurobiologically plausible in terms of the number of neurons, connectivity, and computational mechanisms required, it remains to be seen whether such a system can self-organize or fine tune itself with experience, beginning with only roughly appropriate connections. A hint as to how this may be accomplished has been described by Foldiak (1991), who has demonstrated how a complex cell can learn translation invariance using the objective function of “perceptual stability.” In our model, perceptual stability would be desired in IT, and the control neurons would need to learn how to configure themselves to maintain a stable percept as an attended object moves or changes size on the retina. More generally, there is a clear need to devise learning rules for networks with control-like structures, or three-way interactions, rather than simple perceptron-type networks with two-way interactions only. Recent work in this direction by Lee and Olshausen (1994) has shown that networks in which inputs interact in a local, non-linear fashion are capable of learning higher-order regularities—such as disparity—using a local Hebb rule. It is conceivable that such learning rules may be extended to the control networks proposed here.

## Chapter 5

# Conclusions

In order for us to make sense of the visual world, the brain must be capable of forming object representations that are invariant with respect to the dramatic fluctuations occurring on the retina. We have demonstrated here how this feat may be accomplished by simplified, model neural circuits that are largely consistent with our current knowledge of neurophysiology and neuroanatomy.

While the model has ignored the exact nature of the representation of visual information—for example, the preponderance of feature selective cells in the cortex—it still retains its essential predictive value. Feature processing will certainly affect the picture; but the need for routing still exists, because there is no known theory by which feature processing *per se* can take care of the invariance problem. Thus, the dynamic effects predicted by this model will be expected to be evidenced in some form if routing is indeed employed by the visual cortex.

Besides generating some basic predictions, the model also provides us with a more concrete understanding of what attention is (or what it may be). The action of attention in our model neural circuit can be understood as exploiting the statistical independence of *what* and *where* to make an otherwise enormous, computationally

intensive problem tractable with limited resources. More generally, such mechanisms could be used to take advantage of statistical independence along other dimensions as well, such as motion and color. Within this framework, attention can be seen not merely as a phenomenon, but rather as a computational strategy for dealing with complex visual tasks.

It is interesting (and sometimes amusing) to look back at the early theories, such as those of Lashley, and Pitts and McCulloch, because in certain respects their ideas seem ludicrous by today's standards. It is certainly possible that the theory proposed here will be viewed in retrospect with equal disbelief. However, to date there are no alternative, comparably detailed models that suggest a means for forming invariant representations in a manner that is more parsimonious with our current understanding of the brain. Only after the proper experimental evidence is collected can we alter our conviction in the model. As these experiments are carried out, the results will either help to increase our confidence in the model, or will suggest where it is wrong and how it might be revised. It is this combined process of computational modeling and experimentation that eventually will lead us to understand how visual attention and recognition are actually implemented in the brain.

# Appendix A

## Derivation of autonomous control dynamics

### A.1 Blob search

The total energy functional we wish to minimize is

$$E_{total} = E_{blob} + \beta E_{constraint}, \quad (\text{A.1})$$

where  $E_{blob}$  and  $E_{constraint}$  are defined in Equations 2.10 and 2.11, and  $\beta$  is a constant determining the relative contribution of the constraint term. Letting  $c_k$  follow the gradient of this functional, we obtain

$$\begin{aligned} \frac{dc_k}{dt} &= -\eta \frac{\partial E_{total}}{\partial c_k} \\ &= -\eta \frac{\partial E_{blob}}{\partial c_k} - \eta \beta \frac{\partial E_{constraint}}{\partial c_k}, \end{aligned} \quad (\text{A.2})$$

where  $\eta$  is a constant determining the rate of gradient descent.

As it stands,  $c_k$  is unbounded; hence  $E_{blob}$  and  $E_{constraint}$  will also be unbounded



and the network will not be guaranteed to converge. We can ameliorate this problem by letting  $c_k$  be a monotonically increasing function of another analog variable,  $u_k$ , that actually follows the gradient. That is,

$$c_k = \sigma(u_k) \quad (\text{A.3})$$

$$\frac{du_k}{dt} = -\eta \frac{\partial E_{total}}{\partial c_k} \quad (\text{A.4})$$

$$\sigma(x) = [1 + \exp(-\lambda x)]^{-1}. \quad (\text{A.5})$$

This has the effect of limiting  $c_k$  to the interval  $[0, 1]$ , but since we know *a priori* that the desired minimum of  $E_{blob}$  and  $E_{constraint}$  lies in this range, the limitation does not present a problem.

Taking the derivative of  $E_{blob}$  and  $E_{constraint}$  with respect to  $c_k$  yields

$$\frac{\partial E_{blob}}{\partial c_k} = -\sum_i \sum_j G_i \Gamma_{ijk} I_j^{in} \quad (\text{A.6})$$

$$\frac{\partial E_{constraint}}{\partial c_k} = -\sum_l U_{kl} c_l \quad (\text{A.7})$$

and so the dynamical equation for  $u_k$  is thus

$$\frac{du_k}{dt} = \eta \sum_i \sum_j G_i \Gamma_{ijk} I_j^{in} + \eta \beta \sum_l U_{kl} c_l. \quad (\text{A.8})$$

One remaining problem is that  $u_k$  must be computed via pure integration, which may cause implementation difficulties. We can convert the integrator to a more biologically plausible leaky integrator by adding to  $E_{total}$  the term

$$E_{leak} = \sum_k \int_{0.5}^{c_k} \sigma^{-1}(c) dc. \quad (\text{A.9})$$

The total energy functional is now defined as

$$E_{total} = E_{blob} + \beta E_{constraint} + \alpha E_{leak}, \quad (\text{A.10})$$

where the constant  $\alpha$  determines the relative contribution of  $E_{leak}$ . (The effect of adding this term is discussed in Hopfield's 1984 paper. It essentially pushes  $c_k$  slightly away from 0 and 1.0, depending on the value of  $\alpha$  and  $\lambda$ .)

Taking the derivative of  $E_{leak}$  with respect to  $c_k$  yields

$$\frac{\partial E_{leak}}{\partial c_k} = u_k \quad (\text{A.11})$$

and so the final dynamical equation for  $c_k$  is now

$$c_k = \sigma(u_k) \\ \frac{du_k}{dt} + \tau^{-1} u_k = \eta \sum_i \sum_j G_i \Gamma_{ijk} I_j^{in} + \eta \beta \sum_l U_{kl} c_l, \quad (\text{A.12})$$

where the time constant,  $\tau$ , is defined as  $\frac{1}{\eta\alpha}$ .

## A.2 Recognition

Now the total energy functional is

$$E_{total} = E_{mem} + \beta E_{constraint} + \alpha E_{leak}, \quad (\text{A.13})$$

where  $E_{mem}$  is defined as in Equation 2.14. Note that Equation A.13 is just the same as Equation A.10, except with  $E_{blob}$  replaced by  $E_{mem}$ .

Taking the derivative of  $E_{mem}$  with respect to  $c_k$  yields

$$\frac{\partial E_{mem}}{\partial c_k} = - \sum_i \sum_j V_i \Gamma_{ijk} I_j^{in} \quad (\text{A.14})$$

and so the new dynamical equation for  $c_k$  is thus

$$c_k = \sigma(u_k) \quad (\text{A.15})$$

$$\frac{du_k}{dt} + \tau^{-1} u_k = \eta \sum_i \sum_j V_i \Gamma_{ijk} I_j^{in} + \eta \beta \sum_l U_{kl} c_l. \quad (\text{A.16})$$

Note that this result is just the same as Equation A.12, with the exception that  $G_i$  is replaced with  $V_i$ .

## Appendix B

# Open questions about spatial frequency tuning

The multiscale stack model proposed in Chapter 3 raises a number of interesting issues with regard to the nature of the multiscale representation of visual space within V1. What is the range of peak spatial-frequency tuning observed among cells in foveal V1? Is the coverage—i.e., the number of cells at each spatial frequency and their relative spacing—sufficient to preserve spatial information at a level of detail that is consistent with our perceptual capabilities? This appendix examines the available evidence and points out some lurking mysteries that will need to be resolved in future experiments.

### B.1 Conflicting data

De Valois et al. (1982) and Tootell et al. (1988) report that the peaks of the spatial-frequency tuning curves of foveal V1 cells are in the range of 1-10 cy/deg, with the greatest number of cells in the range of 4-8 cy/deg. Few if any cells in their study have a peak spatial-frequency at 16 cy/deg or above.

On the other hand, Parker and Hawken (1988) report that the majority of foveal V1 cells can be fit by a difference-of-difference-of-Gaussians function with a central, excitatory zone of about 2-4 minutes in diameter. The number of cells with larger central diameters drops off rapidly, with few if any larger than 20 minutes. The cells with the smallest central diameters (2-4 minutes) would presumably have their peak spatial-frequency in the range of 8-15 cy/deg, while the cells with the largest diameters (20 minute) would be centered around 1.5 cy/deg (Hawken and Parker, 1987; this assumes the spacing between the peaks of the inhibitory flanks is about twice the central diameter).

The results of Parker and Hawken seem to be shifted upwards by about an octave from the results of the Tootell and De Valois groups. The apparent discrepancy between the data is not addressed by Parker and Hawken, although both experiments are in anesthetized monkeys and utilize similar methods. Some potential reasons for this difference are given below.

## **B.2 How do we see so clearly?**

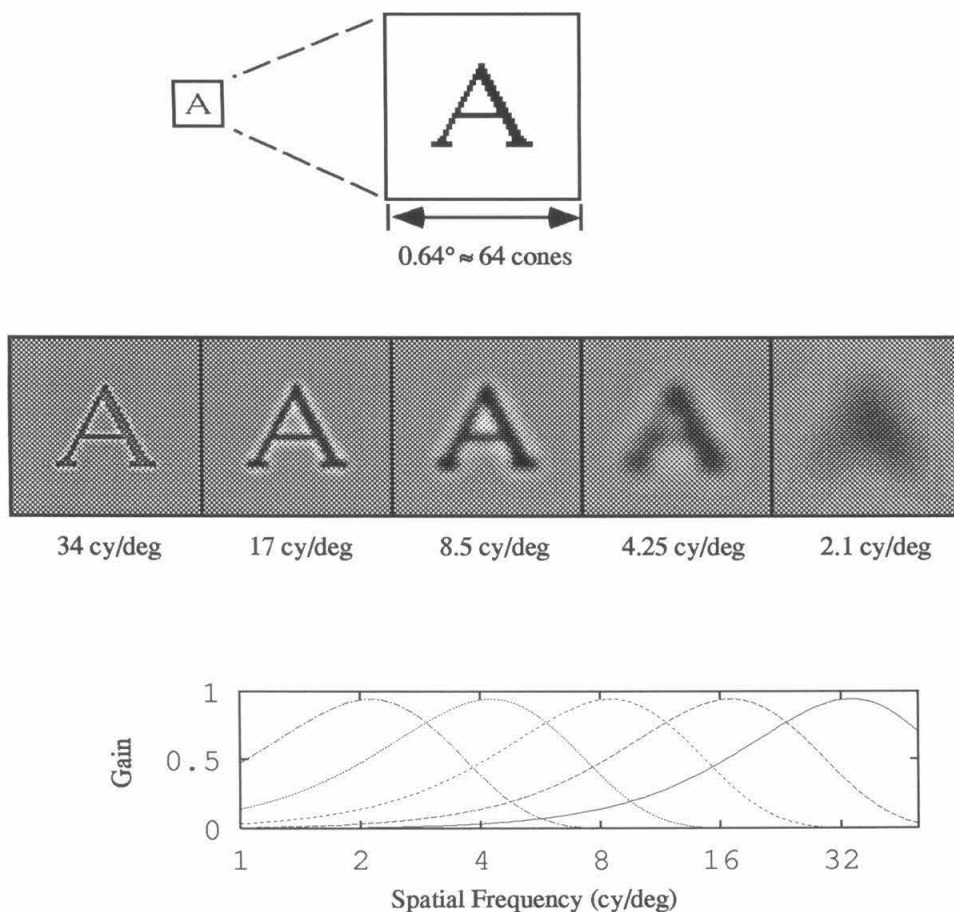
We can essentially consider each cortical cell as a sample node conveying certain properties (e.g., spatial-frequency, orientation) within a local region of visual space. In order to represent information within a certain spatial-frequency band without loss, cortical cells should be spaced by an amount proportional to their peak spatial-frequency, as dictated by the Nyquist sampling theorem. Low spatial-frequency cells on the order of 1 cy/deg would be expected to be spaced by about 0.5 degrees or less, while high spatial-frequency cells on the order of 15 cy/deg would be expected to be spaced by about 0.03 degrees or less. In a 1D array, there should be twice as many cells for each octave increase in spatial-frequency; in 2D, there should be a fourfold increase for each octave.

The fact that the De Valois and Tootell groups find the most number of cells in the range of 4-8 cy/deg is disturbing, because it implies that the representation of spatial structure at 8 cy/deg or above is incomplete or non-existent. Grating detection tasks or simple hyperacuity tasks could still be carried out at high frequencies by utilizing information in the tails of the low frequency tuning curves, but pattern recognition or other tasks requiring a veridical representation of spatial structure at 8 cy/deg or above will suffer. To make this more concrete, Figure B.1 illustrates how a 12-point Times-font 'A' would appear on the retinal sampling lattice in the fovea when viewed at a distance of two feet from the eye, and how it would appear when subsequently filtered by hypothetical cortical cells centered at different spatial frequencies. Needless to say, our perception of this shape corresponds more closely with those filters centered above 8 cy/deg, and not with those reported by the De Valois and Tootell groups. Not only can we recognize the letter, but we are also perfectly capable of discerning the details pertaining to its font, style, etc. The apparent discrepancy becomes even more obvious when one considers that faces can be recognized when they are reduced to a size spanning as few as 19x19 cones (Campbell, 1985). Such faces would appear as little more than amorphous blobs when filtered at 8 cy/deg or below. Thus, there is clearly something awry with the range of peak spatial-frequencies reported by the De Valois and Tootell groups. The data of Parker and Hawken is certainly more consistent with our perceptual capabilities, but one would still expect to find four times as many cells with widths of 2 minutes than of 4 minutes, and the ratio is only 1 : 1 at best.

Thus, we are left with a discrepancy not only between the existing data, but also between the data and our perceptual capabilities. <sup>1</sup> One possible explanation for

---

<sup>1</sup>While the neurophysiological data have been collected on monkeys, it can be assumed that the resolution of their visual perception is nearly equivalent to ours, since their foveal acuity is about 3/4 that of humans (Merigan and Katz, 1990) and their contrast sensitivity functions are practically identical to those of humans (De Valois and De Valois, 1988).



**Figure B.1: Relation of spatial-frequency tuning to perception.** Viewing the small 'A' in the box at upper left from a distance of two feet will result in it being projected onto an array of about 64x64 cones. Shown below this are the results of filtering this image with a difference-of-Gaussians filter centered at various peak spatial frequencies, as shown in the plot. Each filter has approximately a 1 octave bandwidth, as with the cells in V1.

the discrepancy is that the receptive fields of V1 cells are dynamic, as demonstrated by Pettet and Gilbert (1992); each cell may have a very large “potential” receptive field, and depending on the nature of the visual stimuli and the task at hand the cell’s receptive field may either expand or contract. It is conceivable that putting the animal into an anesthetized state and exposing it only to gratings in a darkened room allows the receptive fields to expand somewhat, and thus lowers their peak spatial-frequency tuning. It would be desirable then to repeat these assays in awake animals. Another consideration is that the stimuli used to locate and isolate a cell may bias the experimenter to record from particular types of cells. For example, it is typical in these experiments to use fairly large bars to locate a cell, and these stimuli will preferentially excite lower frequency tuned cells. Interestingly, Parker and Hawken also used a variety of high-frequency stimuli to locate their cells, such as lines and spots, which may be why they found more of the small, high frequency cells. It is also important to keep in mind that the low spatial-frequency cells will fire most often and most prolonged in response to a bar passing over their receptive fields, whereas the high frequency cells with small receptive fields will register only blips of activity by comparison, and will thereby be harder to find. In an analogous situation in the hippocampus, where cells are believed to form a sparse population code, new recording techniques using “stereotrodes” have been successfully used to reliably record from many of the cells that fire only rarely (Wilson and McNaughton, 1993). It would be highly desirable then to utilize these same techniques to determine the number and range of spatial-frequency tuned cells in the visual cortex.

### **B.3 Very low-frequency cells**

The observation of very low-frequency cells ( $\leq 1$  cy/deg) in foveal V1 is likely to be correct, since it is hard to conceive of ways this measurement could have been made



unless a cell actually does possess the capability to integrate information in this way. How these cells achieve such low frequency tuning is an interesting question, since they would need to integrate information over a diameter of about 12 mm on the cortical surface, and the largest observed diameter of cortical axonal arbors is on the order of 6-8 mm (Lund, 1988; Gilbert and Wiesel, 1989). One possibility is that these cells integrate over a large region by collecting the responses of cells two or three synapses removed via intermediary cells acting as relays. Another possibility is that they may expand the range of their connectivity via reciprocal pathways to V2. In any case, it would be interesting to compare the extent of the input field for low vs. high spatial-frequency cells, perhaps by using viruses that are transported retrogradely and transynaptically.

## Appendix C

# Details of pulvinar anatomy and physiology

### C.1 An overview of the pulvinar

The pulvinar occupies the posterior 2/5 of the thalamus and is the largest nucleus of the thalamus in man. It has evolved along with the primate brain, growing in size with increasingly complex subdivisions as the extrastriate visual cortex becomes more complex. Figure C.1 shows a first-pass attempt at constructing a comprehensive connectivity diagram of the primate pulvinar and its relationship with the cortex and other brain structures. Roughly, the pulvinar can be subdivided into four separate subnuclei on cytoarchitectural grounds: inferior (PI), lateral (PL), medial (PM), and oral (PO). The connectivities shown have been accumulated from many different anatomical studies. Recently, a graphical anatomical database of pulvinar connectivity has been constructed by Press and Olshausen (1993) that allows these connections to be analyzed and compared to each other in a detailed fashion.

There are at least two distinct maps of visual space in the pulvinar, which are

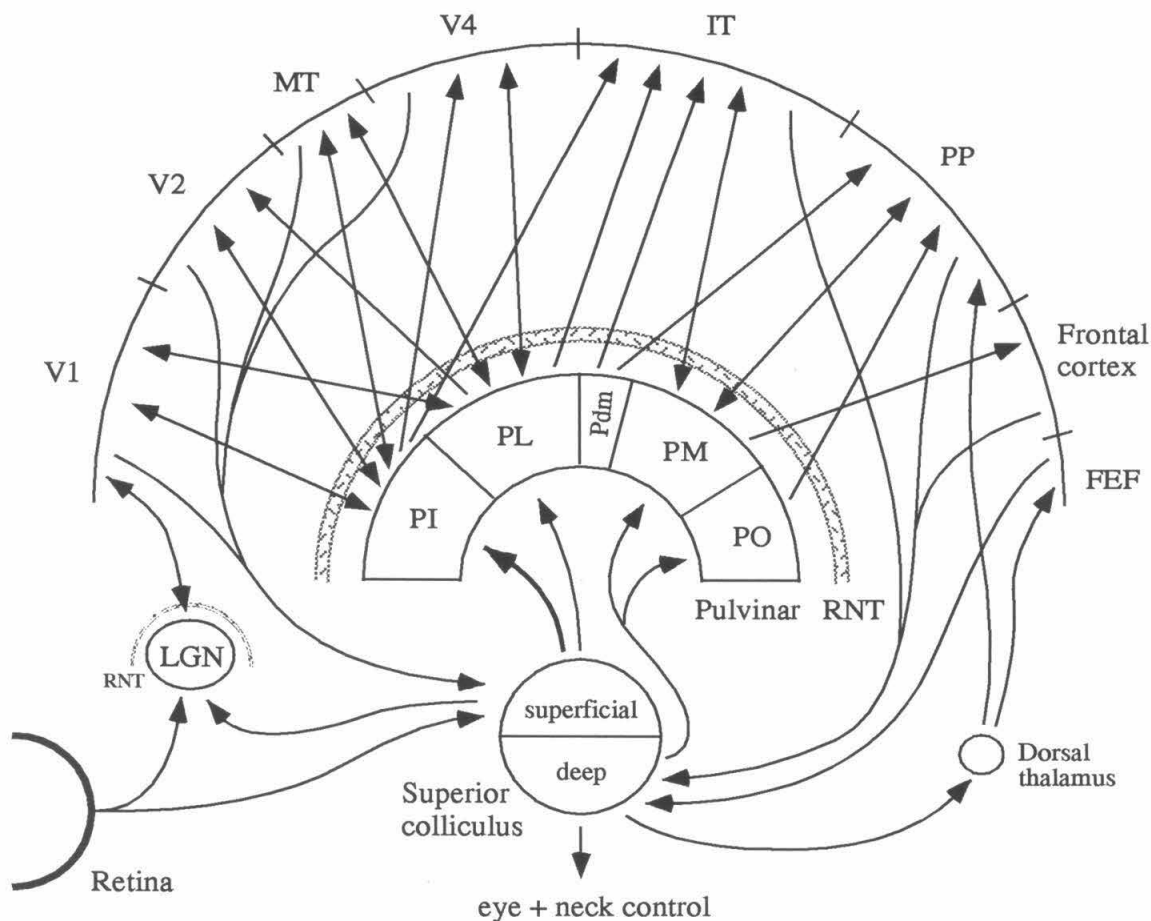


Figure C.1: **A schematic of the known connections of the pulvinar.**

These connections were compiled from the following references: Baleyrier and Morel (1992), Hardy and Lynch (1992), Yeterian and Pandya (1991), Schmahmann and Pandya (1990), Robinson and McClurkin (1989), Benevento and Davis (1977), Benevento and Rezak (1976), Trojanowski and Jacobsen (1976), Ogren and Hendrickson (1977). Anatomical distinctions within subnuclei of the pulvinar are not shown in this diagram (e.g., the projections from PM to IT and PP arise from distinctly different zones within PM, and they intermingle only at the PM/PL border). The pulvinar anatomy database of Press and Olshausen (1993) allows for a more detailed viewing of the topography of connections revealed by these and other studies.

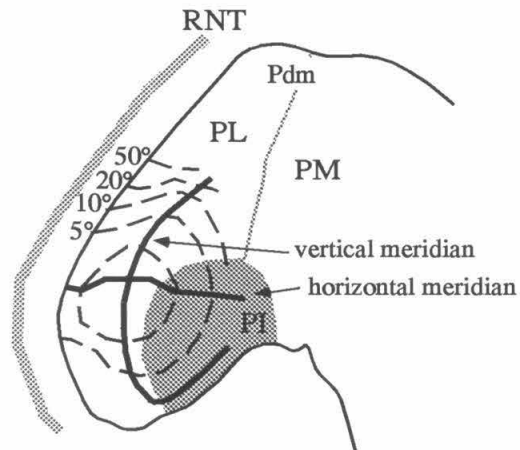


Figure C.2: **The two maps of visual space within the pulvinar.** Shown is a coronal cross-section of the pulvinar. The PI/PL map fills all of PI and extends a bit into PL. The PL map lies entirely within PL. (Adapted from Robinson and McClurkin (1989) and Bender (1981).)

illustrated in Figure C.2. One map occupies all of PI and some of PL and is referred to as the PI/PL map. The other map lies entirely within PL and is referred to as the PL map. (Notice that the PL map wraps around the PL/PI map the same way that V2 wraps around V1! - i.e., by splitting the map along horizontal meridian.) A separate region within PL lying dorsal and medial to the PL map is referred to as Pdm, which does not seem to exhibit much topography - or at least no map has been made of it. The RF sizes in the PI/PL and PL maps as a function of eccentricity are comparable to those in V2. RF sizes in Pdm are a mix of very large and normal (V2-like) sizes.

## C.2 Doubts about pulvinar

The routing circuit model proposes that the pulvinar is a major source of control signals for gating intra-cortical connection strengths along the occipito-temporal pathway. During visual attention, this allows information from a windowed region of retinal inputs to be selectively routed to higher cortical areas concerned with pattern analysis and recognition (namely, IT). Thus, one may ask, why is the pulvinar not just interconnected to areas in the occipito-temporal pathway, but to nearly every area within the visual cortex and to many other parts of the cerebral cortex as well? Indeed, this massive interconnectivity with such diverse areas may seem more suggestive of a general integrative role for the pulvinar, rather than the very specific role we propose for it in pattern recognition. In fact, its connectivity is reminiscent to that of the claustrum, which some people have proposed as an area for integrating information across multiple modalities. Perhaps, then, the pulvinar acts something like a blackboard - a la Mumford (1992) - where various areas of the cortex “write” their interpretation of the world, and this information gets fed back to cortex to further enhance or nullify the interpretation.

Some of the physiological studies of the pulvinar are no less discouraging. The apparent lack of any distinct pattern recognition or visual search deficits after pulvinar lesions, such as in the experiments of Bender (1988), would seem to cast doubt on the specific role we propose for the pulvinar in pattern recognition. In fact, it is the lesions to SC that seem to have the effect we would expect from lesions to pulvinar! Also, Gattass has reported that stimulating pulvinar does not result in any behavioral modulation in an attentional task (personal communication), whereas stimulating SC does. Even in the physiological experiments of Petersen et al. (1985), attentional enhancement was obtained only in Pdm, which is interconnected with PP, and not in PL or PI, which interconnect to the occipito-temporal pathway. Perhaps then the

enhancement observed here is merely a reflection of activity in PP, which has been previously reported to show such attentional enhancements.

In addition, the activity of pulvinar neurons seems to be commonly enhanced or inhibited as a function of eye position or eye movements. This raises the possibility that the pulvinar may be more involved in bringing information about eye movements to bear on cortical visual processing, rather than directing covert visual attention.

A first glance at this evidence might easily lead one to the conclusion that the pulvinar has little or no role in directing attention. But the brain is a complex beast, and so deducing its function requires a deeper analysis of all the factors involved in these and other experiments (see Section C.5).

### C.3 The case for pulvinar

If we make the assumption that the process of visual attention involves control of information flow in the visual cortex (which every model of attention does), then we might suspect that some area (or areas) of the brain would serve as the controller of attention. The properties we would expect of such a controller are that

- A. it should be able to integrate, or bring together, information over a large region of visual space in order to be informed of what's "out there" to attend to;
- B. it should be able to arbitrate, or mediate, the various demands for visual attention and choose *one* locus of the visual field for focusing attention (assuming for now that we attend to only one locus at a time);
- C. it should be able to affect information flow in the cortex in such a way that information within the chosen attentional window is processed preferentially.

The pulvinar seems to have the capacity to satisfy each of these important properties. Cells in Pdm and PM tend to have very large receptive fields, which would imply that

these subnuclei could bring together information over a large region of visual space (thus satisfying **A**). On the other hand, neurons within PI and PL tend to have fairly localized receptive fields, which probably implies more local processing in these areas (although cells with extremely large, diffuse RF's are scattered throughout PI and PL). Such a mixture of large and small RF's is actually what one would expect of a hierarchical routing circuit, since control neurons in the initial stages would receive their input from V1 and V2, whereas control neurons in the top stages would receive their input from V4 and IT, and hence have larger RF's). The fact that PI and PL contain nearly equal RF sizes may or may not be consistent with this scheme, but it is hard to say at this stage of formulation of the model.

Property **B** could be satisfied by the pulvinar since a relatively large portion of the visual field is brought within the extent of horizontal connections within the pulvinar. For example, in the inferior pulvinar,  $50^\circ$  of visual space is traversed in about 3 mm. Given that the arbors of interneuron dendritic fields are on the order of 600 microns in diameter, this would provide a good degree of crosstalk among large parts of visual space, thus facilitating mediation. This crosstalk would become even more global in Pdm and PM, which project to IT. Moreover, a study in *Galago* (Conley and Diamond, 1990) has shown that pulvinar neurons (presumably in the homologue of the PL/PM area) project very diffusely into the RNT, which is in stark contrast to the very localized projections of LGN afferents into the RNT. This combined with the strongly inhibitory nature of the RNT would provide a suitable mechanism for a global winner-take-all function, thus enabling the pulvinar to choose one locus for attention. Note that such diffuse connectivity with the RNT is inconsistent with the notion of a blackboard, since it would tend to screw-up the topology of any pattern of activation fed into it.

Finally, the fact that the pulvinar projects to all areas in the occipito-temporal pathway means that it could perform the requisite modulation of information flow

(property **C**). The fact that the pulvinar projects to other visual areas outside the occipito-temporal pathway, such as MT, may be to focus attention for other aspects of visual processing—such as motion. Indeed, our current focus on pattern vision is probably too restrictive. The broader role proposed for routing by Van Essen and Anderson (1990), in which control operates in the “M stream” in addition to the form pathway, is more parsimonious with the observed connections between the pulvinar and other visual areas. The connections with areas outside the visual cortex, such as frontal cortex, may play the role of informing other brain regions of the current focus of attention, or allowing them to control where attention is focused.

It is also worth noting that nearly every neurobiological theory of visual attention proposes nearly identical roles for the pulvinar. For example, Posner et al. (1988), on the basis of lesion studies, assign the pulvinar with the operation of “engaging” attention, which would presumably mean implementing the modulation of cortical activity. Laberge (1990), on the basis of his PET study, assigns the pulvinar with the task of filtering out distracting information in a cluttered field. And Niebur et al. (1993) propose the pulvinar as the site of origination of oscillatory activity in the cortex for gating connections. It may well be that the commonality in all these theories is an artifact of the sociology of science, but nevertheless it is interesting that no one else has yet proposed a better alternative for control of visual attention.

In any case, the increasing size of the pulvinar throughout evolution would seem to imply that it plays some sort of an important role in visual function. Otherwise, the cost of wiring and neuronal resources would surely dictate its demise.

## C.4 Other alternatives

What are the alternatives to the pulvinar as a control site? Desimone, on the basis of his lesion study, has suggested a distributed control system. However, what exactly



he means by this is unclear. I have trouble with the concept of having more than one controller, because if you did, you would then have to posit the existence of another controller to decide which of the control systems takes precedence over the other. At this point you are back to having a single control site. Posner also has proposed a sort of distributed control system, with PP disengaging attention, SC moving attention, and pulvinar engaging attention. However, even this scheme centralizes the process of engaging - which is the quintessential operation of attention - in the pulvinar.

One form of distributed control that does make sense is a hierarchical control scheme, such as in the hierarchical routing circuit of Figure 2.19. In this scheme, the top-stage control neurons would serve as the global “master controller,” and the bottom-stage control neurons would control information flow over more localized regions of visual space. Thus, if control neurons near the top stage were disabled, as presumably was done with the PL lesions of Desimone et al. (1990), one would observe an asymmetry for global vs. local attentional effects, just as Desimone did in his study (i.e., there was no impairment for distractors within a hemifield, but distractors in the opposite hemifield did produce a decrement in performance). On the other hand, if control nodes near the bottom were disabled, one would be able to focus attention only roughly (without much precision) for a limited region of visual space.

Alternatively, the control neurons may be distributed between the pulvinar and the cortex. In this scheme, the pulvinar control neurons may set up the more global context for attention (rough position and size) at each stage of the routing circuit, and then control neurons in the cortex (presumably in layer 6, as described by Van Essen and Anderson (1990)) would provide refinements. (Such a distribution of control circuitry would indeed make sense, because the function of control is so important that it would be risky to put all the eggs in one basket.) It is conceivable under such a scheme that lesioning the pulvinar control neurons would have little effect on

recognition when a single object is presented alone in the visual field, since the cortical control neurons themselves may be able to handle the routing for such a simple scene. The pulvinar control neurons may only be necessary to filter out other information when multiple objects are present.

It seems unlikely that the parietal cortex could play a role in controlling information flow because it is so lacking in connectivity to the occipito-temporal pathway. Also, it does not appear to have the necessary structure for choosing a single locus for attention, such as the RNT, although this remains an unknown. The superior colliculus also does not appear as a good candidate for control for the same reasons.

## **C.5 Which data are important?**

The routing circuit model was formulated to solve a specific problem - namely, the recognition of objects in a complex environment independent of translation and scale on the retina. We have hypothesized that the process of visual attention provides a solution to this problem by reducing the amount of incoming information to a manageable size and by bringing this information to the pattern recognition system in the correct format (i.e., normalized for position and scale). In the psychophysics community, on the other hand, attention is largely viewed as a phenomenon - something that gives you faster reaction times or better performance when you are pre-cued to a particular location (or feature dimension). This view naturally leads to differences between the way attention is commonly investigated and the way one would investigate attention if testing our model. However, we must make do with the data that are available, and so we need some means of judging which data to weigh seriously and which to weigh less seriously.

Since the model is directed at the task of pattern recognition, then those experiments involving the processing of patterns, such as the experiments of Bergen and

Julesz (1983) and Triesman (1988), are probably ones that we can most directly relate to our model. On the other hand, those experiments involving such tasks as detecting the onset or dimming of a spot of light, or color or orientation discrimination, are more difficult to relate to our model. Although these tasks may qualify as “attention-related” since certain stimuli can be processed quicker or more efficiently when computational resources are focused on them, they can be solved with much simpler models that do not require a routing circuit for shifting and scaling information. It may well be that the capabilities of our model could transfer to such tasks, but the outcome of these experiments, either way, cannot be weighed as serious evidence for or against the model.

Lesion studies, in particular, need to be interpreted with extreme caution. Previous lesion studies of other parts of the brain have revealed that it is a highly distributed, redundant system. This makes it very difficult to discount a particular area taking part in some function just because one doesn’t obtain a substantial behavioral effect. For example, lesions of SC or FEF alone seem to have only a mild effect on eye movements, but simultaneous lesions to both areas produce a dramatic effect (eye movements are virtually abolished). Also, as Ungerleider and Mishkin have pointed out, initial lesion studies of the prestriate cortex seemed to indicate that it was not the principal pathway between V1 and IT, but later studies made it clear that sparing even small amounts of the prestriate cortex left viable pathways between V1 and IT that enabled the animals to still perform object recognition tasks. Even prosopagnosics seem to retain more or less normal visual function in most respects, except that they may occasionally confuse one object for another.

In this context, I would evaluate each of the physiological and lesion studies of the pulvinar to date as follows:

- Petersen et al.’s (1985) physiology experiment utilizes the task of detecting the

dimming of a spot of light. It is assumed that the spot of light was attended if its dimming was correctly detected. This task is unrelated to pattern recognition and can be solved by many other means besides the dynamic routing circuit, so we would be hard pressed to use the outcome as direct evidence. They find an enhancement in 50% of the cells recorded in Pdm when the attended spot of light falls within its RF. Attentional enhancements are absent in PL and PI, but one does find cells here, and in Pdm, that are modulated by saccades. Cells in Pdm are connected with both IT and PP. Perhaps the Pdm cells are exerting an attentional effect on IT, or on PP? Or perhaps the Pdm cells are merely reflecting the enhancement of PP cells - i.e., they are an effect, and not the cause, of attention.

- Petersen et al.'s (1987) neuropharmacological experiment utilizes the task of detecting the onset of a spot of light in pre-cued and un-cued conditions. Reaction times are faster in the pre-cued condition than the un-cued condition. Again, this task is unrelated to pattern recognition, so it provides relatively weak evidence. Injecting GABA agonists and antagonists into Pdm seems to respectively slow down or speed up shifts of attention, insofar as reaction times get slower and faster. The results of this experiment would seem to support a causal role for Pdm in visual attention, although whether the effect is exerted on PP or IT cannot be resolved. Presumably, either of these areas could be involved in performing this task.
- Rafal and Posner's (1987) lesion experiment also utilizes the task of detecting the onset of a spot of light (actually, an asterix). Subjects were precued to one of two locations on opposite sides of the visual field and asked to respond as quickly as possible when an asterix was presented in one of the locations. The subjects with thalamic lesions (i.e., pulvinar + presumably other parts of

the thalamus) consistently had slower reaction times on the contralesional side, although the improvement obtained with increasing SOA times was the same for both sides. This is in contrast to what one obtains with SC lesions, in which targets presented on the contralesional side result in a much slower improvement with increasing SOA times, suggesting an impairment in the “move” operation of attention. Thus, it is concluded that patients with pulvinar lesions are moving their attention OK, but not “engaging” it as thoroughly as without the lesion. Although it is noted that all the thalamic lesioned patients had “no clinical evidence of visual impairment,” it is still hard to rule out the possibility that the thalamic lesions merely screw up the visual system in some general way unrelated to attention.

- Desimone et al.’s (1990) lesion experiment utilizes the task of color discrimination. The animal is pre-cued to a certain location in the visual field and must make a judgement of the color of a bar subsequently presented at that location (maintaining fixation) with and without a distractor. Although this task does not involve pattern recognition, it does at least involve a distractor stimulus - albeit one - that would seem to force the animal to process one location preferentially over the other. They find that complete deactivation of PL (the posterior portion, connected to V4 and IT) of one hemisphere results in an increase in errors ( $\sim 30\%$ ) when a distractor is present in the opposite hemifield. One does not get this effect, however, without a distractor, or when the distractor is placed in the same hemifield as the target. On the other hand, local deactivation of the superior colliculus does produce an increase in errors when the target and distractor are in the same hemifield. Thus, it would seem that both pulvinar and SC play some role in directing information processing to a particular part of the visual scene; but note the only way that SC could

actually affect cortical information flow is through the pulvinar, or perhaps via its direct connection to LGN.

- Chalupa et al.'s (1976) tachistoscopic pattern discrimination experiment is getting closer to an appropriate experiment for testing the model, in that he utilizes spatial patterns with brief presentation times. They find that lesions to PI result in a dramatic impairment in pattern discrimination ability, whereas lesions to PL and PM do not. However, his lesions to these latter two areas are incomplete. Also, the patterns being discriminated are an "N" vs. a "Z", which are probably not substantially complex to tax the monkey's pattern recognition system. Bender attempted to repeat Chalupa's experiment, however since he used cynomolgus monkey's, which are slower learners of visual discriminations, he had to scale up the patterns so they were 20°x20° in size! So it does not come as a great surprise that his lesions produced no effect. Both Chalupa and Bender utilized RF lesions, which damage fibers of passage. Thus, it is difficult to say whether Chalupa's effect is due to destroying the cortico-tectal pathway or the pulvinar itself.
- Bender and Butter's (1987) study also included a visual search task in which the monkey scanned a display (with eye movements) to look for a particular target item - for example, a circle - placed in a field of distracting items - squares, triangles, crosses, etc. The patterns were 2°-4° in size. Although this task relies on overt instead of covert attention, it does require the animal to discriminate among patterns. Bender's lesions produced massive bilateral damage to PI and PL - and in one monkey much of PM too. However, only a very mild increase in errors was observed. This result is not inconsistent with our model, since it is conceivable that the animal could simply be homing in on a particular feature in the target, and hence there would be no need for a routing circuit. However,

it is still a bit disturbing that a stronger effect was not observed.

- Nagel-Leibey et al.'s (1984) lesion study utilizes the task of discriminating a plus vs. square pattern. Both RF and kainic acid lesions are used (the latter leaves the cortico-tectal fibers of passage intact). Both lesions result in mild deficits in the retention of the ability to learn these tasks. The lesions involved large portions of PI and PL. The effect obtained here seems to be suggestive of some role for the pulvinar in pattern discrimination, however the patterns used here are still quite simple. Nor is the task particularly demanding of attention (long presentation times, no distractors, etc.).
- Laberge and Buchsbaum's (1990) PET study utilizes the task of discriminating the shape of a target item when it is surrounded by a field of distractors vs. when it is displayed alone. This is probably the best designed task of any experiment to date for testing the model, in that it involves both pattern discrimination and the filtering out of irrelevant information. They report that the pulvinar shows a higher differential activation when the target item is surrounded by distractors than when it is displayed alone. However, his data are less than convincing. Only the left pulvinar shows enhanced activity when the distractors are present. In addition, V1 shows enhanced activity too when the distractors are present. Thus, it could well be the case that the presence of the distractors causes more activation overall in V1, which then excites the pulvinar differentially. This experiment would have been more compelling had he done a control in which the distractor stimuli were viewed passively, as opposed to attentively.
- Finally, it should be noted that Robinson et al. (1986;1990) have carried out several studies investigating the responses of pulvinar neurons in relation to eye movements. Most cells are enhanced before an eye movement, irrespective of the direction of the eye movement; others are selective for direction. Some cells

are suppressed during eye movements, others not. Also, some cells seem to show activity signaling the end of a saccade. It has been determined that most of these effects are due to an extraretinal signal - i.e. proprioceptive information. In general, it makes sense that the control for covert attention would somehow have to be coordinated with eye movements, but since the model does not yet address this issue it is hard to say what one would expect of pulvinar neurons with respect to eye movements. At this stage, my best guess is that the general enhancement effect observed before and after saccades could possibly be a way of resetting the pulvinar before and after an eye movement (this effect is also found in V1 and other visual cortical areas). The suppression during eye movements would also be desirable in order to filter out irrelevant information as the eye sweeps over a scene (such an effect has also been observed in the LGN).

## C.6 Design of experiments

Consider the problem a monkey faces in recognizing another monkey in the jungle, among the many other objects—coconuts, flowers, foliage, etc.—that clutter the natural visual environment. This task demands enormous amounts of sensory information processing, yet it probably comes quite naturally and effortlessly to the monkey. Now compare this with the situation a monkey typically faces in any one of the above experiments—for example, detecting the dimming of a spot of light in the periphery, or judging whether the orientation or color of two bars is the same or different. These latter tasks demand trivial amounts of visual computation by comparison, yet I would guess that they are conceptually much more difficult for the animal to perform.

The model we have proposed is meant to address a problem of the former type—i.e., computationally intensive, yet effortless. Thus, experiments designed to test the



routing circuit hypothesis should be designed along these lines. The task used in the experiment should be *computationally* challenging, but not necessarily conceptually difficult for the animal to perform. Indeed, designing tasks that naturally fit an animal's capabilities would be an advantage in training. An example might be to train an animal to discriminate among two or more complex shapes - for example, different faces - and then test the animal's performance with the same shapes presented at different sizes and locations on the retina, or among a field of distracting items. If an animal could still perform such a task with extensive pulvinar lesions, or if no modulation of activity were observed in pulvinar neurons during such a task, then we would have to reconsider what role, if any, the pulvinar plays in visual attention.

Another important consideration in the design of experiments is to determine where data are most needed to test a certain theory. Here, it would be helpful to have an alternative theory against which to judge the evidence - i.e., does the data fit model A or model B? Finding data that don't fit model A will tend to lower the likelihood of model A somewhat, but if it could be shown that the data better fit an alternative model B that solves the same problem as model A, then this makes the case against model A even stronger. Unfortunately, few models exist that are designed to solve the same problem as ours. One popular theory that is often touted is that recognition is achieved through a succession of feature selective and complex-like cells that allow for some variability in the position of the feature. The features become increasingly more complex and less position sensitive as one proceeds through the hierarchy of visual areas until - voilà! - one achieves face-selective and hand-selective cells in IT. However, no one has yet spelled out a coherent neurobiological theory that goes beyond this rather intuitive, two-sentence description. In order to qualify as a valid neurobiological theory, one must be able to show (a) that the method actually works for recognizing shapes over wide ranges of position and scale, (b) that it does so in a manner consistent with known neuroanatomy, neurophysiology and

psychophysics, and (c) that it makes solid, experimentally testable predictions about what one would expect to find in various cortical areas.

On the other hand, there are plenty of other models of attention in general that could be tested against ours. The model of Niebur et al. (1993), for example, would predict 40Hz oscillatory activity in the pulvinar and within attended regions of V1, and non-oscillatory activity during states of inattention. Their model also predicts shifts and enlargements of V4 receptive fields, but our model would predict that an RF would cover only a fraction of the attentional window, whereas in their model the RF could easily cover the entire attentional window, since spatial relationships within the attentional window are not preserved. Another consequence of preserving spatial relationships is that the patterns of activation on the cortical surface should translate, or dilate and contract, as attention is shifted and scaled over the scene. This prediction could probably be tested most directly with optical recording techniques, but this method still seems to be in its infancy.

## C.7 Conclusion

The available data on the pulvinar leave us with a mixed bag of evidence. All of the physiological and lesion study evidence is rather weak, with most of it for, and some (namely, Bender's search task) against the pulvinar playing an important role in attention and pattern recognition. Thus, the strongest argument for the pulvinar as a control site comes from the way that it is anatomically situated with respect to the rest of the cortex. Until additional data are obtained to the contrary - from experiments designed in the appropriate manner - it would be premature to discount the pulvinar's role in controlling attention with what little negative data exist.

The brain has been wonderfully constructed to perform computational feats that are beyond our most powerful supercomputers. In order to deduce how the brain

carries out these operations, it is important to observe the system during periods of intense computation. Simple probes have traditionally been used in order to isolate various aspects of the stimulus, and this approach has been useful at giving us a first cut estimate of what may be going on in different brain areas. However, revealing the more complex operations of the brain - such as visual attention and complex pattern recognition - will require experimental designs that better challenge an animal's natural computational capabilities.

## References

- Ahmad S (1992) VISIT: A neural model of covert visual attention. In: *Advances in Neural Information Processing Systems 4* (Moody JE, Hanson SJ, Lippman RP, eds), San Mateo, CA: Kaufmann, pp 420-427.
- Anderson CH, Burt PJ, van der Wall GS (1985) Change detection and tracking. *SPIE Vol. 579—Intelligent Robots and Computer Vision*, 72-78.
- Anderson CH, Van Essen DC (1987) Shifter circuits: A computational strategy for dynamic aspects of visual processing. *Proceedings of the National Academy of Sciences, USA*, 84:6297-6301.
- Anderson CH, Van Essen DC (1993) Dynamic neural routing circuits. In: *Visual Search 2* (Brogan D, Gale A, Carr K, eds), London: Taylor & Francis.
- Anderson JA, Rosenfeld E (1988) *Neurocomputing: Foundations of research*. Cambridge, MA: MIT Press. p.29.
- Anstis SM (1974) A chart demonstrating variations in acuity with retinal position. *Vision Research*, 14:589-592.
- Atick JJ, Redlich AN (1990) Towards a theory of early visual processing. *Neural Computation*, 2: 308-320.

- Attneave, F (1954) Some informational aspects of visual perception. *Psychological Review*, 61: 183-193.
- Attneave, F (1965) Triangles as ambiguous figures. *American Journal of Psychology*, 81: 447-453.
- Baleydier C, Morel A (1992) Segregated thalamocortical pathways to inferior parietal and inferotemporal cortex in macaque monkey. *Visual Neuroscience*, 8: 391-405.
- Baron RJ (1987) *The Cerebral Computer*. Erlbaum.
- Bender DB (1981) Retinotopic organization of macaque pulvinar. *J Neurophysiol*, 46: 672-693.
- Bender DB (1988) Electrophysiological and behavioral experiments on the primate pulvinar. In: *Progress in Brain Research*, 75 (Hicks TP, Benedek G, eds), New York: Elsevier, pp 55-65.
- Bender DB, Butter CM (1987) Comparison of the effects of superior colliculus and pulvinar lesions on visual search and tachistoscopic pattern discrimination in monkeys. *Exp Brain Res*, 69: 140-154.
- Benevento LA, Rezak M (1976) The cortical projections of the inferior pulvinar and adjacent lateral pulvinar in the Rhesus monkey: An autoradiographic study. *Brain Research*, 108:1-24.
- Benevento LA, Davis B (1977) Topographical projections of the prestriate cortex to the pulvinar nuclei in the macaque monkey: An autoradiographic study, *Brain Research*, 30: 405-424.
- Bergen JR, Julesz B (1983) Parallel versus serial processing in rapid pattern discrimination. *Nature*, 303:696-698.

- Berman NJ, Douglas RJ, Martin KAC (1992) GABA-mediated inhibition in the neural networks of visual cortex. In: *Progress in Brain Research*, 90 (Mize RR, Marc RE, Silito AM, eds), New York: Elsevier, pp 443-476.
- Biederman I, Cooper EE (1992) Evidence for complete translational and reflectional invariance in visual object priming. *Perception*, 20:585-593.
- Buhmann J, Lades M, von der Malsburg C (1990) Size and distortion invariant object recognition by hierarchical graph matching. In: *Proceedings of the International Joint Conference on Neural Networks*, San Diego, June, pp. 411-416.
- Burt PJ, Adelson EH (1983) The Laplacian pyramid as a compact image code. *IEEE Transactions on Communications*, 31: 532-540.
- Bushnell C, Goldberg ME, Robinson DL (1981) Behavioral enhancement of visual responses in monkey cerebral cortex. I. Modulation in posterior parietal cortex related to selective visual attention. *Journal of Neurophysiology*, 46(4): 755-772.
- Campbell FW (1985) How much of the information falling on the retina reaches the visual cortex and how much is stored in the visual memory? In: *Pattern Recognition Mechanisms* (Chagas C, Gattass R, Gross C, eds), Berlin: Springer, pp 83-95.
- Carpenter G, Grossberg S (1987a) A massively parallel architecture for a self-organizing neural pattern recognition machine. *Computer Vision Graphics and Image Processing*, 37:54-115.
- Carpenter GA, Grossberg S (1987b) Invariant pattern recognition and recall by an attentive self-organizing ART architecture in a nonstationary world. In: *Proceedings of the IEEE First International Conference on Neural Networks*.

- Cauler LJ, Connors BW (1992) Functions of very distal dendrites: Experimental and computational studies of layer I synapses on neocortical pyramidal cells. In: Single Neuron Computation (McKenna T, Davis JL, Zornetzer SF, eds), Academic Press, Cambridge, MA, pp 199-229.
- Cavanagh P (1978) Size and position invariance in the visual system. *Perception*, 7:167-177.
- Cavanagh P (1985) Local log polar frequency analysis in the striate cortex as a basis for size and orientation invariance. In: Models of the Visual Cortex (Rose D, Dobson VG, eds), New York: Wiley, pp 85-95.
- Cavanagh P (1992) Attention-based motion perception. *Science*, 257:1563-1565.
- Cave KR (1991) What makes a spotlight a spotlight? (unpublished manuscript—Dept. of Psychology, Vanderbilt University)
- Chalupa LM, Coyle D, Lindsley DB (1976) Effect of pulvinar lesions on visual pattern discrimination in monkeys. *Journal of Neurophysiology*, 39(2):354-369.
- Cherniak C (1990) The bounded brain: Toward quantitative neuroanatomy. *Journal of Cognitive Neuroscience*, 2(1):58-68.
- Conley M, Diamond IT (1990) Organization of the visual sector of the thalamic reticular nucleus in *Galago*. *European Journal of Neuroscience*, 2(3):211-226.
- Connor CE, Gallant JL, Van Essen DC (1993) Effects of focal attention on receptive field profiles in area V4. *Soc. Neurosci. Abstr.*, 19.
- Corbetta M, Miezin FM, Dobmeyer S, Shulman GL, Petersen SE (1991) Selective and divided attention during visual discriminations of shape, color, and speed:

- Functional anatomy by positron emission tomography. *The Journal of Neuroscience*, 11(8): 2383-2402.
- Crick F (1984) Function of the thalamic reticular complex: The searchlight hypothesis. *Proceedings of the National Academy of Sciences*, 81:4586-4590.
- Crick F, Koch C (1990) Towards a neurobiological theory of consciousness. *Seminars in the Neurosciences*, 2:263-275.
- Desimone R, Albright TD, Gross CG, Bruce C (1984) Stimulus-selective properties of inferior temporal neurons in the macaque. *The Journal of Neuroscience*, 4(8):2051-2062.
- Desimone R, Schein SJ (1987) Visual properties of neurons in area V4 of the macaque: Sensitivity to stimulus form. *Journal of Neurophysiology*, 57(3): 835-868.
- Desimone R, Wessinger M, Thomas L, Schneider W (1990) Attentional control of visual perception: Cortical and subcortical mechanisms. In: *Cold Spring Harbor Symp Quant Biol*, 55:963-971.
- Desimone R (1992) Neural circuits for visual attention in the primate brain. In: *Neural Networks for Vision and Image Processing* (Carpenter GA, Grossberg S, eds), Cambridge, Mass.: MIT Press, pp 343-364.
- De Valois RL, De Valois KK (1988) *Spatial Vision*. New York: Oxford.
- De Valois RL, Albrecht DG, Thorell LG (1982) Spatial frequency selectivity of cells in macaque visual cortex. *Vision Res*, 22: 545-559.
- DeYoe EA, Sisola LC (1991) Distinct pathways link anatomical subdivisions of V4 with V2 and temporal cortex in the macaque monkey. *Society for Neuroscience Abstracts*, 17: 1282.



- Douglas RJ, Martin KAC, Whitteridge D (1988) Selective responses of visual cortical cells do not depend on shunting inhibition. *Nature*, 332:642-644.
- Douglas RJ, Martin KAC (1990a) Neocortex. In: *Synaptic Organization of the Brain* (Shepard GM, ed), New York: Oxford UP, pp 389-438.
- Douglas RJ, Martin KAC (1990b) Control of neuronal output by inhibition at the axon initial segment. *Neural Computation*, 2:283-292.
- Duhamel J, Colby L, Goldberg ME (1992) The updating of the representation of visual space in parietal cortex by intended eye movements. *Science*, 255:90-92.
- Eriksen CW, Murphy TD (1987) Movement of attentional focus across the visual field: A critical look at the evidence. *Perception and Psychophysics*, 42(3):299-305.
- Eriksen CW, St. James JD (1986) Visual attention within and around the field of focal attention: A zoom lens model. *Perception and Psychophysics*, 40 (4):225-240.
- Farell B, Pelli DG (1993) Can we attend to large and small at the same time? *Vision Research*, 33: 2757-2772.
- Felleman DJ, McClendon E (1991) Modular connections between area V4 and temporal lobe area PITv in macaque monkeys. *Society for Neuroscience Abstracts*, 17: 1282.
- Felleman DJ, Van Essen DC (1991) Distributed hierarchical processing in the primate cerebral cortex. *Cerebral Cortex*, 1(1):1-47.
- Felleman DJ, McClendon E, Lin K (1992) Modular segregation of visual pathways in occipital and temporal lobe visual areas in the macaque monkey. *Society for*

- Neuroscience Abstracts, 18: 390.
- Field, DJ (1987) Relations between the statistics of natural images and the response properties of cortical cells. *J Opt Soc Am, A*, 4: 2379-2394.
- Fields HL, Basbaum AI (1978) Brainstem control of spinal pain-transmission neurons. *Annu Rev Physiol*, 40:217-248.
- Fischer B (1973) Overlap of receptive field centers and representation of the visual field in the cat's optic tract. *Vision Research*, 13:2113-2120.
- Foldiak P (1991) Learning invariance from transformation sequences. *Neural Computation*, 3:194-200.
- Fukushima K (1980) Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, 36:193-202.
- Fukushima K (1987) Neural network model for selective attention in visual pattern recognition and associative recall. *Applied Optics*, 26(23):4985-4992.
- Gallant JL, Braun J, Van Essen DC (1993) Selectivity for polar, hyperbolic, and Cartesian gratings in macaque visual cortex. *Science*, 259:100-103.
- Gattass R, Sousa APB, Covey E (1985) Cortical visual areas of the macaque: Possible substrates for pattern recognition mechanisms. In: *Pattern Recognition Mechanisms* (Chagas C, Gattass R, Gross C, eds), Berlin: Springer, pp 1-20.
- Gattass R, Desimone R (1991) Attention-related responses in the superior colliculus of the macaque. *Soc Neurosci Abstr*, 17: 545.
- Gattass R, Desimone R (1992) Stimulation of the superior colliculus (SC) shifts the focus of attention in the macaque. *Soc Neurosci Abstr*, 18: 703.

- Gilbert CD, Wiesel TN (1989) Columnar specificity of intrinsic horizontal and corticocortical connections in cat visual cortex. *The Journal of Neuroscience*, 9: 2432-2442.
- Goldberg ME, Wurtz RH (1972) Activity of superior colliculus in behaving monkey. I. Visual receptive fields of single neurons. *Journal of Neurophysiology*, 35:542-559.
- Gross CG, Rocha-Miranda CE, Bender DB (1972) Visual properties of neurons in inferotemporal cortex of the macaque. *Journal of Neurophysiology*, 35:96-111.
- Grossberg S (1976) Adaptive pattern classification and universal recoding: I. Parallel development and coding of neural feature detectors. *Biological Cybernetics*, 23: 121-134.
- Hardy SGP, Lynch JC (1992) The spatial distribution of pulvinar neurons that project to two subregions of the inferior parietal lobule in the macaque, *Cerebral Cortex*, 2: 217-230.
- Hawken MJ, Parker AJ (1987) Spatial properties of neurons in the monkey striate cortex. *Proc R Soc Lond B*, 231: 251-288.
- Hinton GE (1979) Some demonstrations of the effects of structural descriptions in mental imagery. *Cognitive Science*, 3: 231-250.
- Hinton GE (1981a) A parallel computation that assigns canonical object-based frames of reference. In: *Proceedings of the Seventh International Joint Conference on Artificial Intelligence 2*, Vancouver B.C., Canada.
- Hinton GE (1981b) Shape representation in parallel systems. In: *Proceedings of the Seventh International Joint Conference on Artificial Intelligence 2*, Vancouver B.C., Canada.

- Hinton GE, Lang KJ (1985) Shape recognition and illusory conjunctions. In: Proceedings of the Ninth International Joint Conference on Artificial Intelligence, Los Angeles.
- Hopfield JJ (1982) Neural networks and physical systems with emergent collective computational abilities. *Proc Natl Acad Sci USA*, 79:2554-2558.
- Hopfield JJ (1984) Neurons with graded response have collective computational properties like those of two-state neurons. *Proc Natl Acad Sci USA*, 81:3088-3092.
- Kanerva P (1988) *Sparse Distributed Memory*. Cambridge, Mass.: MIT Press.
- Kennedy H, Bullier J (1985) A double-labeling investigation of the afferent connectivity to cortical areas V1 and V2 of the macaque monkey. *The Journal of Neuroscience*, 5: 2815-2830.
- Koch, C. and Poggio, T. (1992) Multiplying with synapses and neurons. In: *Single Neuron Computation* (McKenna T, Davis JL, Zornetzer SF, eds), Cambridge, MA: Academic, pp 315-345.
- Koch C, Ullman S (1985) Shifts in selective visual attention: towards the underlying neural circuitry. *Human Neurobiology*, 4:219-227.
- Koenderink JJ, van Doorn AJ (1978) Visual detection of spatial contrast; Influence of location in the visual field, target extent and illuminance level. *Biological Cybernetics*, 30: 157-167.
- Kosslyn SM, Schwartz SP (1978) Visual images as spatial representations in active memory. In: *Computer Vision Systems* (Hanson AR, Riseman EM, eds), New York: Academic Press.

- LaBerge D (1990) Thalamic and cortical mechanisms of attention suggested by recent positron emission tomographic experiments. *Journal of Cognitive Neuroscience*, 2(4): 358-372.
- LaBerge D, Buchsbaum MS (1990) Positron emission tomographic measurements of pulvinar activity during an attention task. *The Journal of Neuroscience*, 10(2):613-619.
- LaBerge D, Carter M, and Brown V (1992) A network simulation of thalamic circuit operations in selective attention. *Neural Computation*, 4: 318-331.
- Larsen A, Bundesen C (1978) Size scaling in visual pattern recognition. *Journal of Experimental Psychology: Human Perception and Performance*, 4: 1-20.
- Lashley, KS (1942) The problem of cerebral organization in vision, *Biol Symp*, 7: 301-322.
- LeCun Y, Boser B, Denker JS, Henderson D, Howard RE, Hubbard W, Jackel LD (1990) Backpropagation applied to handwritten Zip code recognition. *Neural Computation*, 1:541-551.
- Lee CW, Olshausen BA (1993) A nonlinear hebbian network that learns to detect disparity in random-dot stereograms. (submitted)
- Li Z, Atick JJ (1994) Towards a theory of the striate cortex. *Neural Computation*, 6.
- Llinas RR (1988) The intrinsic electrophysiological properties of mammalian neurons: Insights into central nervous system function. *Science*, 242:1654-1663.
- Lowe DG (1987) The viewpoint consistency constraint. *International Journal of Computer Vision*, 1:57-72.

- Lund JS (1988) Anatomical organization of macaque monkey striate visual cortex. *Annual Review of Neuroscience*, 11: 253-288.
- Marr D, Poggio T (1976) Cooperative computation of stereo disparity. *Science*, 194:283-287.
- Marr D (1982) *Vision*. New York: W.H. Freeman.
- Mel, B.W. (1992) NMDA-based pattern discrimination in a modeled cortical neuron. *Neural Computation*, 4: 502-517.
- Merigan WH, Katz LM (1990) Spatial resolution across the macaque retina. *Vision Research*, 30: 985-991.
- Miller KD, Chapman B, Stryker MP (1989) Visual responses in adult cat visual cortex depend on N-methyl-D-aspartate receptors. *Proceedings of the National Academy of Sciences*, 86:5183-5187.
- Mishkin M (1972) Cortical visual areas and their interactions. In: *Brain and Human Behavior* (A.G. Karczmar, J.C. Eccles, eds), New York: Springer, pp 187-208.
- Moran J, Desimone R (1985) Selective attention gates visual processing in the extrastriate cortex. *Science*, 229:782-784.
- Mountcastle VB, Andersen RA, Motter BC (1981) The influence of attentive fixation upon the excitability of the light-sensitive neurons of the posterior parietal cortex. *The Journal of Neuroscience*, 1(11):1218-1235.
- Mozer MC, Behrmann M (1992) Reading with attentional impairments: A brain-damaged model of neglect and attentional dyslexias. In: *Connectionist Approaches to Natural Language Processing* (Reilly RG, Sharkey NE, eds), Hillsdale, NJ: Erlbaum, pp 409-460.

- Mumford D (1992) On the computational architecture of the neocortex. II The role of cortico-cortical loops. *Biological Cybernetics*, 66:241-251.
- Nagel-Leiby S, Bender DB, Butter CM (1984) Effects of kainic acid and radiofrequency lesions of the pulvinar on visual discrimination in the monkey. *Brain Res*, 300: 295-303.
- Nakayama K, Mackeben M (1989) Sustained and transient components of focal visual attention. *Vision Res*, 29 (11):1631-1647.
- Nakayama K (1991) The iconic bottleneck and the tenuous link between early visual processing and perception. In: *Vision: Coding and Efficiency* (Blakemore C, ed), Cambridge: Cambridge UP, pp 411-422.
- Nelson SB, Sur M (1992) NMDA receptors in sensory information processing. *Current Opinion in Neurobiology*, 2:484-488.
- Niebur E, Koch C, Rosin C (1993) An oscillation-based model for the neural basis of attention. *Vision Research*, 33: 2789-2802.
- Nowlan SJ, Sejnowski TJ (1993) Filter selection model for generating visual motion signals. *Advances in Neural Information Processing Systems*, 5 (Hanson SJ, Cowan JD, Giles CL, eds), San Mateo, CA: Morgan-Kaufmann, pp. 369-376.
- Ogren MP, Hendrickson AE (1977) The distribution of pulvinar terminals in visual areas 17 and 18 of the monkey. *Brain Research*, 137:343-350.
- Ogren MP, Hendrickson AE (1979) The structural organization of the inferior and lateral subdivisions of the *Macaca* monkey pulvinar. *J Comp Neur*, 188:147-178.
- Ohzawa I, Sclar G, Freeman RD (1982) Contrast gain control in the cat visual cortex. *Nature*, 298:266-268.

- O'Kusky J, Colonnier M (1982) A laminar analysis of the number of neurons, glia, and synapses in the visual cortex (area 17) of adult macaque monkeys. *Journal of Comparative Neurology*, 210: 178-290.
- Olshausen BA, Anderson CH, Van Essen DC (1993) A neurobiological model of visual attention and invariant pattern recognition based on dynamic routing of information. *The Journal of Neuroscience*, 13: 4700-4719.
- Palmer SE (1983) The psychology of perceptual organization: A transformational approach. In: *Human and Machine Vision* (Beck J, Hope B, Rosenfeld A, eds), Orlando: Academic, pp 269-339.
- Parker AJ, Hawken MJ (1988) Two-dimensional spatial structure of receptive fields in monkey striate cortex. *Journal of the Optical Society of America A*, 5: 598-605.
- Pei X, Volgushev M, Creutzfeldt O (1992) A comparison of directional sensitivity with the excitatory and inhibitory field structure in cat striate cortical simple cells. *Perception*, 21, supplement 2, p. 26.
- Petersen SE, Robinson DL, Keys W (1985) Pulvinar nuclei of the behaving rhesus monkey: Visual responses and their modulation. *Journal of Neurophysiology*, 54(4):867-886.
- Petersen SE, Robinson DL, Morris JD (1987) Contributions of the pulvinar to visual spatial attention. *Neuropsychologia*, 25(1A):97-105.
- Pettet MW, Gilbert CD (1992) Dynamic changes in receptive-field size in cat primary visual cortex. *Proc Natl Acad Sci USA*, 89:8366-8370.
- Pitts W, McCulloch WS (1947) How we know universals: The perception of auditory and visual forms. *Bulletin of Mathematical Biophysics*, 9:127-147.



- Pollen DA, Lee JR, Taylor JH (1971) How does the striate cortex begin the reconstruction of the visual world? *Science*, 173:74-77.
- Poggio T, Edelman S (1990) A network that learns to recognize three-dimensional objects. *Science*, 343: 263-266.
- Posner MI, Cohen Y (1984) Components of visual orienting. In: *Attention and Performance X: Control of Language Processes*. (Bouma H, Bouwhuis DG, eds), Hillsdale, NJ: Erlbaum, pp 531-556.
- Posner MI, Walker JA, Friedrich FJ, Rafal RD (1984) Effects of parietal injury on covert orienting of attention. *The Journal of Neuroscience*, 4(7):1863-1874.
- Posner MI, Petersen SE, Fox PT, Raichle ME (1988) Localization of cognitive operations in the human brain. *Science*, 240: 1627-1631.
- Posner MI, Petersen SE (1990) The attention system of the human brain. *Annu Rev Neurosci*, 13:25-42.
- Postma EO, van den Herik HJ, Hudson PTW (1992) The gating lattice: A neural substrate for dynamic gating. In: *CNS\*92 Proceedings*, July 26-29, San Francisco, California. Kluwer Academic Publishers.
- Press WA, Olshausen BA, Van Essen DC (1993) Analyzing connections between the pulvinar and visual cortex: An interactive graphical database. *Soc Neurosci Abstr*, 19: 331.
- Rafal RD, Posner MI (1987) Deficits in human visual spatial attention following thalamic lesions. *Proc Natl Acad Sci USA*, 84: 7349-7353.
- Remington R, Pierce L (1984) Moving attention: Evidence for time-invariant shifts of visual selective attention. *Perception and Psychophysics*, 35(4):393-399.

- Rezak M, Benevento LA (1979) A comparison of the organization of the projections of the dorsal lateral geniculate nucleus, the inferior pulvinar and adjacent lateral pulvinar to primary visual cortex (area 17) in the macaque monkey. *Brain Research*, 167:19-40.
- Robinson DL, McClurkin JW (1989) The visual superior colliculus and pulvinar. In: *The Neurobiology of Saccadic Eye Movements* (Wurtz and Goldberg, eds), Elsevier Science Publishers BV, pp. 337-359.
- Robinson DL, Petersen SE (1992) The pulvinar and visual salience, *Trends in Neuroscience*, 15(4):127-132.
- Robinson DL, Goldberg ME, Stanton GB (1978) Parietal association cortex in the primate: Sensory mechanisms and behavioral modulations. *Journal of Neurophysiology*, 41(4):910-932.
- Robinson DL, Petersen SE, Keys W (1986) Saccade-related and visual activities in the pulvinar nuclei of the behaving rhesus monkey. *Exp Brain Res*, 62: 625-634.
- Robinson DL, McClurkin JW, Kertzman C (1990) Orbital position and eye movement influences on visual responses in the pulvinar nuclei of the behaving macaque. *Exp Brain Res*, 82: 235-246.
- Robinson DL, Bowman EM, Kertzman C (1991) Covert orienting of attention in macaque: II. A signal in parietal cortex to disengage attention. *Society for Neuroscience Abstracts*, 17: 442.
- Rock I (1973) *Orientation and form*. New York: Academic Press.
- Rock I (1988) On Thompson's inverted-face phenomenon (Research Note). *Perception*, 17:815-817.

- Rockland KS (1992) Configuration, in serial reconstruction, of individual axons projecting from area V2 to V4 in the macaque monkey. *Cerebral Cortex*, 2:353-374.
- Rolls ET, Baylis GC (1986) Size and contrast have only small effects on the responses to faces of neurons in the cortex in the superior temporal sulcus of the monkey. *Exp Brain Res*, 65: 38-48.
- Saarinen J, Julesz B (1991) The speed of attentional shifts in the visual field. *Proceedings of the National Academy of Sciences*, 88:1812-1814.
- Sandon PA, Uhr LM (1988) An adaptive model for viewpoint-invariant object recognition. *Proceedings of the 10th Annual Conference of the Cognitive Science Society*, Montreal, Canada, August, pp. 209-215.
- Sandon PA (1990) Simulating visual attention. *Journal of Cognitive Neuroscience*, 2(3):213-231.
- Schmahmann JD, Pandya DN (1990) Anatomical investigation of projections from thalamus to posterior parietal cortex in the rhesus monkey: A WGA-HRP and fluorescent tracer study, *The Journal of Comparative Neurology*, 295: 299-326.
- Schwartz EL (1977) Spatial mapping in the primate sensory projection: Analytic structure and relevance to perception. *Biological Cybernetics*, 25: 181-194.
- Shepard RN, Metzler J (1971) Mental rotation of three-dimensional objects. *Science*, 171: 701-703.
- Sherman SM, Koch C (1986) The control of retinogeniculate transmission in the mammalian lateral geniculate nucleus. *Experimental Brain Research*, 63:1-20.
- Shulman GL and Wilson J (1987) Spatial frequency and selective attention to local and global information. *Perception*, 16: 89-101.

- Sillar KT (1991) Spinal pattern generation and sensory gating mechanisms. *Current Opinion in Neurobiology*, 1:583-589.
- Sperling G, Landy MS, Cohen Y, Pavel M (1985) Intelligible encoding of ASL image sequences at extremely low information rates. *Computer Vision Graphics and Image Processing*, 31: 335-391.
- Steinmetz MA, Connor CE, MacLeod KM (1992) Focal spatial attention suppresses responses of visual neurons in monkey posterior parietal cortex. *Society for Neuroscience Abstracts*, 18: 148.
- Toet A, van Eekhout P, Simons HLJJ, Koenderink JJ (1987) Scale invariant features of differential spatial displacement discrimination. *Vision Res*, 27: 441-451.
- Tootell BH, Silverman MS, Hamilton SL, Switkes E, De Valois RL (1988) Functional anatomy of macaque striate cortex. V. Spatial Frequency. *The Journal of Neuroscience*, 8: 1610-1624.
- Treisman AM (1988) Features and objects: The fourteenth Bartlett memorial lecture. *The Quarterly Journal of Experimental Psychology*, 40A(2): 201-237.
- Treisman AM, Schmidt H (1982) Illusory conjunctions in the perception of objects. *Cognitive Psychology*, 14: 107-141.
- Trojanowski JQ, Jacobson S (1976) Areal and laminar distribution of some pulvinar cortical efferents in rhesus monkey, *J Comp Neur*, 169: 371-392.
- Tsal Y (1983) Movements of attention across the visual field. *Journal of Experimental Psychology: Human Perception and Performance*, 9:523-530.
- Tsotsos JK (1991) Localizing stimuli in a sensory field using an inhibitory attention beam. Technical Report, RBCV-TR-91-37, Dept. of Computer Science,

University of Toronto.

- Ungerleider LG, Mishkin M (1982) Two cortical visual systems. In: *Analysis of Visual Behavior* (Ingle DJ, Goodale MA, Mansfield RJW, eds), Cambridge, Mass.: MIT Press, pp 549-586.
- Van Essen DC, Newsome WT, Maunsell JHR (1984) The visual field representation in striate cortex of the macaque monkey: Asymmetries, anisotropies, and individual variability. *Vision Res*, 24: 429-448.
- Van Essen DC, Newsome WT, Maunsell JHR, Bixby JL (1986) The projections from striate cortex (V1) to areas V2 and V3 in the macaque monkey: Asymmetries, areal boundaries, and patchy connections. *The Journal of Comparative Neurology*, 244:451-480.
- Van Essen DC, Felleman DJ, DeYoe EA, Olavarria J, Knierim J (1990) Modular and hierarchical organization of extrastriate visual cortex in the macaque monkey. In: *Cold Spring Harbor Symp Quant Biol*, 55:679-696.
- Van Essen DC, Anderson CH (1990) Information processing strategies and pathways in the primate retina and visual cortex. In: *An Introduction to Neural and Electronic Networks* (Zornetzer SF, Davis JL, Lau C, ed), New York: Academic, pp 43-72. (2nd Edition in press)
- Van Essen DC, Olshausen B, Anderson CH, Gallant JL (1991) Pattern recognition, attention, and information bottlenecks in the primate visual system. In: *Proc SPIE Conf on Visual Information Processing: From Neurons to Chips*, Vol 1473 (Mathur BP, Koch C, eds), Bellingham, WA: SPIE, pp 17-28.
- Van Essen DC, Anderson CH, Olshausen BA (1994) Dynamic routing strategies in sensory, motor, and cognitive processing. In: *Large Scale Neuronal Theories of*

- the Brain (Koch C, Davis J, eds) MIT Press.
- Van Essen DC, DeYoe EA (1993) Concurrent processing in the primate visual cortex. In: *The Cognitive Neurosciences* (Gazzaniga MS, ed), Cambridge, MA: MIT Press.
- Vergheze P, Pelli DG (1992) The information capacity of visual attention. *Vision Research*, 32(5):983-995.
- Volgushev M, Pei X, Vidyasagar TR, Creutzfeldt OD (1992) Orientation-selective inhibition in cat visual cortex: an analysis of postsynaptic potentials. *Perception*, 21, supplement 2, p. 26.
- von der Heydt R, Peterhans E (1989) Mechanisms of contour perception in monkey visual cortex. *The Journal of Neuroscience*, 9(5):1731-1763.
- von der Malsburg C, Bienenstock E (1986) Statistical coding and short-term synaptic plasticity: A scheme for knowledge representation in the brain. In: *Disordered Systems and Biological Organization* (NATO ASI Series, Vol. F20) (Bienenstock E, Fogelman Soulie F, Weisbuch G, eds), Berlin: Springer, pp 247-272.
- Wechsler H, Zimmerman GL (1988) 2-D invariant object recognition using distributed associative memory. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 10: 811-821.
- Wilson MA, McNaughton BL (1993) Dynamics of the hippocampal ensemble code for space. *Science*, 261: 1055-1058.
- Witkin AP, Terzopoulos D, Kass M (1987) Signal matching through scale space. *International Journal of Computer Vision*, 1(2):133-144.

- Yeterian EH, Pandya DN (1991) Corticothalamic connections of the superior temporal sulcus in rhesus monkeys, *Exp Brain Res*, 83: 268-284.
- Yukie M, Iwai E (1985) Laminar origin of direct projection from cortex area V1 to V4 in the rhesus monkey. *Brain Research*, 346: 383-386.