

A Simple Neural Network Generating an Interactive Memory

JAMES A. ANDERSON
The Rockefeller University
New York, New York, 10021

Communicated by Donald H. Perkel

ABSTRACT

A model of a neural system where a group of neurons projects to another group of neurons is discussed. We assume that a trace is the simultaneous pattern of individual activities shown by a group of neurons. We assume synaptic interactions add linearly and that synaptic weights (quantitative measure of degree of coupling between two cells) can be coded in a simple but optimal way where changes in synaptic weight are proportional to the product of pre- and postsynaptic activity at a given time. Then it is shown that this simple system is capable of "memory" in the sense that it can (1) recognize a previously presented trace and (2) if two traces have been associated in the past (that is, if trace \bar{f} was impressed on the first group of neurons and trace \bar{g} was impressed on the second group of neurons and synaptic weights coupling the two groups changed according to the above rule) presentation of f to the first group of neurons gives rise to \bar{g} plus a calculable amount of noise at the second set of neurons. This kind of memory is called an "interactive memory" since distinct stored traces interact in storage. It is shown that this model can effectively perform many functions. Quantitative expressions are derived for the average signal to noise ratio for recognition and one type of association. The selectivity of the system is discussed. References to physiological data are made where appropriate. A sketch of a model of mammalian cerebral cortex which generates an interactive memory is presented and briefly discussed. We identify a trace with the activity of groups of cortical pyramidal cells. Then it is argued that certain plausible assumptions about the properties of the synapses coupling groups of pyramidal cells lead to the generation of an interactive memory.

All that we are is the result of what we have thought: it is founded on our thoughts, it is made up of our thoughts.—*Dhammapada*, I.I.

INTRODUCTION

Memory is a mental function which seems in its generality to be a central problem in neurophysiology and neuropsychology. Experiments on memory are difficult to perform and good data is scanty, although

considerable information about memory function in higher mammals has accumulated in recent years.

In two previous papers I discussed a very simple model for the organization of long term memory storage [3]. Basically, the model assumes, first, that memory "traces" (the items which are to be stored—the basic elementary unit of memory) are composed of a complex pattern of individual activities shown by a large spatial array of elements, and, second, that memory storage is formed by constructing a storage array which is the sum of many of the basic "traces."

More precisely, the trace was represented by an N element vector, \vec{f}_k , where N is very large. Elements of \vec{f}_k can take any values. (They are not binary elements.) Assume we have K traces ($\vec{f}_1, \vec{f}_2, \dots, \vec{f}_K$) to be stored. Then, we represented the memory array in this model by a storage vector, \vec{s} , defined as

$$\vec{s} = \sum_{k=1}^{k=K} \vec{f}_k.$$

\vec{s} represents our sole information about the system.

Storage of traces is simple in such a system and could be performed by a simple physiological mechanism; for example, strengthening of synaptic contacts by an amount proportional to activity in the presynaptic cell would be sufficient in some neuron models. It is not clear, however, that it is possible to extract much information from such a system. The retrieval problem was discussed in the previous papers and it was shown that it was indeed possible for such a simple system to perform many of the functions which would be expected of a biological memory. Some of the properties shown by the system, particularly the kinds of mistakes and distortions it makes, are reminiscent of those made by a human memory [3].

The following discussion is an extension of this simple idea to another model. It is an attempt to model a very common kind of anatomical configuration in the central nervous system where one large group of neurons projects to another large group of neurons, or projects to itself via recurrent collaterals having a significantly long conduction time. Examples of this type of highly parallel projection are legion: there are many projections of the thalamic nuclei to cortex, for example, and extensive intracortical projection systems.

If a group of neurons projects to another we shall show that strengthening or weakening the synaptic connection between the two groups according to a simple multiplicative function of activity in pre- and postsynaptic cells automatically generates an interactive memory akin to those discussed in the previous paper. The properties of this model will be discussed and a model of mammalian cortex embodying some of these ideas will be sketched in the last section of the paper.

THE MODEL

Let us consider two groups of neurons, α and β , where α sends projections to β (Fig. 1). We shall assume that (1) Exactly M neurons of α project to each neuron of β , (2) Exactly M neurons of β are projected to by each neuron of α , (3) Synaptic weights add linearly and the activity of a neuron in β is proportional to the sum of the synaptic weights times activity of neurons in α , and (4) α and β have the same number of neurons, N . Assumptions (1), (2) and (4) are approximations of the real situation for mathematical convenience.

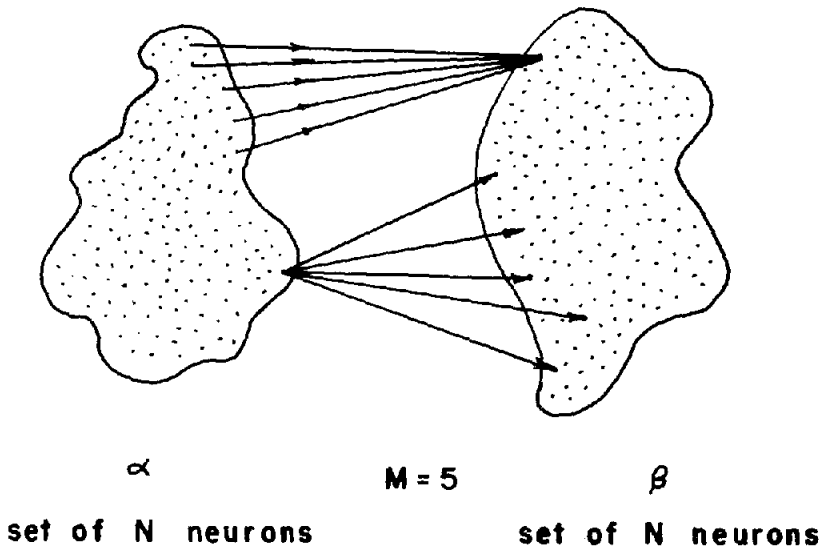


FIG. 1. A group of neurons, α , projects to another group of neurons, β . M neurons of α project to each neuron of β , and M neurons of β are projected to by each neuron of α . Both groups have N neurons.

Assumption (3) requires justification. There is an extensive literature on neural models which assumes that the nervous system is basically digital, that is, the presence or absence of a spike at a given moment is the matter of most importance to the nervous system. This approach to neural modeling is interesting mathematically, and has led to some elegant and important results in automata theory (Ref. 23, chap. 3). However, a great deal of experimental evidence now suggests that an analog model of the mammalian central nervous system is correct in many cases.

Evidence indicates that for most systems in mammals, particularly the "higher" systems, what is significant to the system is the behavior of the cell over a short period of time—the average firing frequency, for example—and not the presence or absence of a single spike. Perkel and Bullock

[26] have made an extensive list of biological systems and the apparent kinds of neural codes used, many of which depend on temporal spike patterns over a period long in relation to the duration of a single spike.

One carefully studied sensory system in primates shows this kind of coding. Mountcastle has studied the tactile sensitivity of the monkey hand [24]. A stimulator probe indented the skin of the thenar eminence of the palm of a monkey. The number of impulses in 600 msec. was recorded from myelinated axons of the palmar branch of the median nerve and was found to be linearly related to the indentation of skin produced by the stimulator probe. The same linear relation to skin indentation was found in units in the ventrobasal thalamus to stimuli applied to the glabrous skin of the monkey hand if the spontaneous activity level of the cell was subtracted. The same linearity of response was found in a cortical neuron in the post-central gyrus of an unanesthetized monkey. Human observers gave a similar linear relation between subjective intensity and depth of mechanical stimulation by a probe tip applied to the pad of the middle finger.

A strictly linear relation between a stimulus (depth of indentation) and neuronal response is rare, however, what is often found and which encouraged Mountcastle to formulate it as a general rule for sensory systems, is that there is a linear relation between the output of the first order afferent fiber and sensory response of the nervous system. Most receptors are markedly nonlinear transducers; for example, response of single tactile receptors (average number of impulses to a stimulus) in the hairy skin areas of the hand gives a power function response to indentation with an exponent of about 0.5. Power functions as transducer outputs are common, and generally the exponent is preserved up to primary sensory cortex.

Maffei et al. have shown [20, 21] that the firing rate of lateral geniculate body (LGB) cells follows the sinusoidal intensity modulation of a light stimulus in quite linear fashion. Maffei has also shown that the LGB may use spatial averaging to improve signal to noise ratio. Their data indicates that spatial averaging preserves the linearity of the response of the cell, in consonance with our assumption (3).

It would be unrealistic to claim that neurons do not introduce substantial nonlinearities. (See Ref. 7 for examples of some of the significant nonlinearities encountered in LGB cells under conditions similar to those of Maffei). However, the overall picture is of much less nonlinearity than might at first be expected from a system incorporating such nonlinear threshold devices as neurons.

At this point it might be wise to make the distinction often made in circuit analysis between large and small signal characteristics. We would

expect a linear addition assumption to be quite accurate for cortical neurons if the stimulus increment was small, whereas significant nonlinearities would occur if the increment was large. The particular system under consideration is of paramount importance. Burns (Ref. 6, chap. 6) has shown that interactions between eyes can be quite nonlinear. Thus Burns finds that in primary visual cortex, where cells respond to stimuli applied to each eye, presentation of the same stimulus to both eyes in spatial register gives enhancement of firing rate far beyond what would be expected from the sum of the response to stimulation of either eye alone.

Before we go further, let us make some comments about the linear assumption. First, the model assumes, in general, small signal linearity in storage. Second, retrieval, in the model presented here, assumes as well large-signal linearity, but the model is not very sensitive to the details of the large-signal transfer function as long as monotonicity is preserved, that is, as long as increasing excitatory synaptic activity increases post-synaptic cell firing or increasing inhibitory activity decreases cell firing.

We will assume that what is of interest in our model is the simultaneous activity of the entire group of neurons, thus we define a trace to be this pattern of individual activities. We further assume that traces are "large" in that a single trace contains a good deal of information. We will make calculations with "traces" as our elementary units.

Since traces are assumed to be the simultaneous activities of large groups of neurons, we can formally represent a trace as a vector of N elements, where N is the number of neurons present. Thus, if \vec{f} is a trace, we define the "power" of \vec{f} , P , as the vector dot product,

$$P = \vec{f} \cdot \vec{f}.$$

We assume that traces add together in storage, they are not separated. Thus, our initial approach to calculations will be to place some statistical constraints on the set of allowable traces. We assume they have equal power, P . Since power, in some intuitive sense, stands for the "importance" of a trace, assumption of equal trace power seems unnatural. The calculations to follow could equally well have been carried out using P as the "average" power of a trace in the set of allowable traces. However, the results to be discussed are not basically changed and the additional complications in exposition seemed to make the equal power assumption a simplifying convenience.

We assume that different traces in a sum are uncorrelated. We assume that, on the average over sets of sums of allowable traces, that the statistics of every element will be the same.

By the Central Limit Theorem, we predict that the value of the sum of many uncorrelated traces approximates a normally distributed random variable.

Since we know nothing of the details of the traces involved in any particular sum, values we calculate are "averages," calculated over many sets of sums of allowable traces.

As discussed in Ref. 3, we assume that the mean value of elements in a trace is zero. This assumption can be shown to give nearly optimal properties (optimal in the sense that the signal power to noise power ratio is maximized) to the retrieval system to be discussed and greatly simplifies calculations. The zero mean assumption implies both positive and negative values for elements in a trace. However, neurons can only have positive average firing frequencies. We can meet the zero mean requirement by assuming neurons have a spontaneous activity level and then defining the activity comprising a trace as deflections, plus or minus, from this resting level. Other realizations of this requirement are possible, but this definition seems natural. Much cortical data shows transduction in both a positive and negative direction which often seems to be referred to the spontaneous activity level (see Ref. 24, p. 398), the well known data of Fox and O'Brien [13], as well as the data in other experimental papers. This implies that the neuron would be expected to act something like a limiter with less dynamic range in the negative than positive direction. Freeman [14, 15] in a well-developed and experimentally supported electrophysiological model of cat prepyriform cortex has developed a basically linear model with clipping to explain the results of his extensive experiments.

Calculations will be made as follows. We will be interested in the behavior of one trace in the sum. We will approximate the sum of the other traces by the values taken by a normally distributed random variable and with this additive "noise" can then make simple calculations.

Let us now return to the model shown in Fig. 1. Let us consider an input trace, \bar{f} , to α . We have assumed that the activity of a neuron in β will be given by the sum of the activities coupled to the neuron by the M neurons from α projecting to it [Assumption (3)]. If we denote by a_{ij} the value of the synapse connecting the i th element in α with the j th element in β , and if $g(j)$ is the value taken by the j th element in β , then we have as the fundamental relation for our system,

$$g(j) = \sum_{i=1}^{i=N} a_{ij}f(i).$$

We see that if \bar{f} is represented by a column vector and if a_{ij} are elements, representing synaptic weights in a "connectivity matrix", A , then a vector \bar{g} , representing activity in β is generated according to

$$\bar{g} = A\bar{f}.$$

We are interested in the following problem. We have trace \bar{f}_1 , representing activity in α , which we wish to have associated with another, trace \bar{g}_1 , representing activity in β , thus we wish to construct a connectivity matrix A so that

$$\bar{g}_1 = A\bar{f}_1,$$

as shown in Fig. 2.

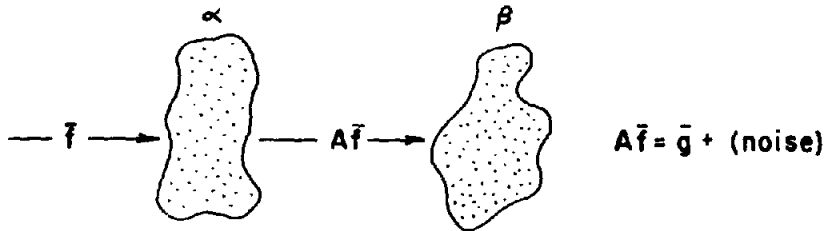


FIG. 2. The general scheme for associating two distinct traces, \bar{f} , and \bar{g} . Trace \bar{f} is impressed on group of neurons α . The connectivity matrix, A , coupling α and β gives rise to trace \bar{g} (which was associated with \bar{f} in the past) plus noise due to the presence of other associations in A .

In general, we will wish to couple many pairs of traces by way of our projection system. Let us assume we have a set of K pairs of traces

$$[(\bar{f}_1, \bar{g}_1); (\bar{f}_2, \bar{g}_2); \dots; (\bar{f}_K, \bar{g}_K)],$$

that we wish to couple, that is, if we present pattern of activity \bar{f}_k to the set of neurons α , we wish the set of neurons, β , to show a pattern of activity "close to" \bar{g}_k . Let us first try to construct a matrix A which is somehow optimal and, second, calculate just how accurately the projection system reconstructs \bar{g}_k at β when \bar{f}_k is presented at α .

The first problem is considered here in a slightly different manner than in Ref. 3. The second problem is considered later.

We will assume that we code the incoming trace, \bar{f} , according to a well defined set of rules, generating a vector $h(\bar{f})$. We assume $h(\bar{f})$ is a reasonable function—nonzero, continuous. Since $[\bar{f}_j]$ are uncorrelated, $[h(\bar{f}_j)]$ will be uncorrelated. Since we have assumed that synaptic increments or decrements due to different traces add together, we can write for the activity of a single element in β when \bar{f} is impressed on α ,

$$g(i) = \bar{f} \cdot \sum_j h(\bar{f}_j).$$

Let us assume that a trace corresponding to \bar{f} is present, that is, \bar{f} has appeared to the system and the coding of \bar{f} , $g(\bar{f})$ is represented in the synapses. Then

$$g(i) = \bar{f} \cdot h(\bar{f}) + \bar{f} \cdot \sum_j h(\bar{f}_j).$$

Remembering that $h(\vec{f}_j)$ are vectors, we see the second term on the average over allowable sets of traces is a function of the average length of $h(\vec{f}_j)$ since the \vec{f}_j are uncorrelated. We will assume an optimal coding, $h(\vec{f})$ is one which on the average maximizes the first term while keeping the second term as small as possible. The geometry of the dot product indicates that this maximum is obtained when \vec{f} and $h(\vec{f})$ point in the same direction, that is

$$h(\vec{f}) = c\vec{f},$$

where c is a constant.

Let us point out the implications of this simple result. Many models of long term memory assume that permanent changes in synaptic weight are made when the trace is laid down [12]. *We will show that if, when an association is learned by a set of neurons projecting to another set of neurons, synaptic weight is changed by an amount proportional to the product of activity in the presynaptic neuron and of activity in the postsynaptic neuron, a memory formed of interacting traces is generated which (1) is optimal in the sense just considered and (2) can be shown to effectively recognize and associate traces.* It is thus possible to generate a psychologically global memory system which is formed by physiologically local changes produced according to simple local rules, requiring the synapse to be affected only by the product of pre- and postsynaptic activities. Note that this does not correspond to a simple strengthening-by-use learning change, although synaptic change dependent only on the activity of the presynaptic cells is capable of generating a simple recognition memory model of an interactive type [3].

RECOGNITION

We will first consider the problem of recognition in this model where recognition is defined as the ability to state, with some calculatable degree of certainty, whether or not a trace presented to α has been presented to the system before. For this calculation we will assume that previous

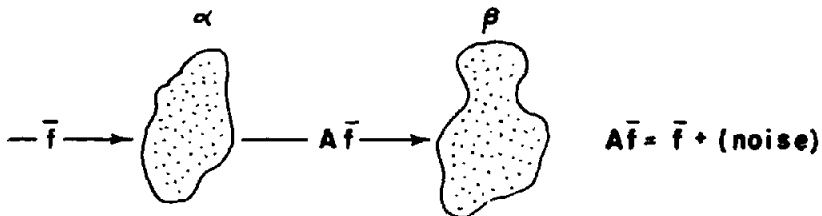


FIG. 3. The general scheme for recognition by self-association. Presentation of trace \vec{f} to set of neurons α gives rise to trace \vec{f} at set of neurons β , plus noise due to other stored associations in A .

storage implies that the trace is associated with itself, that is, presentation of the trace to α gives rise to the trace in β , thus

$$\bar{f} = A\bar{f},$$

(see Fig. 3). This may seem like an artificial definition if α and β are assumed to be two distinct sets of neurons, but if we were to consider a system, common in cortex, where a set of neurons projects to itself, we see that this scheme for recognition arises naturally.

We assume we have K pairs of associated traces stored in our system according to the rule

$$A_k \bar{f}_k = \bar{g}_k,$$

where A_k is the connectivity matrix generated when only (\bar{f}_k, \bar{g}_k) are associated by the system. Then we form the sum

$$A = \sum_{k=1}^{k=K} A_k.$$

We assume a trace \bar{f}_0 is presented whose connectivity matrix may be present in A .

We wish to calculate a statistic which will let us decide whether (\bar{f}_0, \bar{f}_0) is present or absent from the associations stored in A . A reasonable kind of statistic to use, since it is exceptionally easy to calculate with strictly local interactions, and which is also the optimal linear filter is the "matched filter" given in this case by

$$V = A\bar{f}_0 \cdot \bar{f}_0.$$

We will now proceed to calculate the part of V due to the presence of (\bar{f}_0, \bar{f}_0) and the amount due to the noise generated by the other stored associations. We will use as our parameter of interest the output signal to noise ratio $[(S/N)_0]$ defined as is usual by

$$(S/N)_0 = (\text{output due to signal})^2 / (\text{mean square output due to noise}).$$

We assume M neurons in α project to each neuron in β . We further assume diagonal elements of the matrix are zero. To assume otherwise for the recognition problem would lead to trouble since self-associations will always give rise to positive diagonal elements. For now we will assume that the correlation between elements in a trace and between traces is zero.

By our optimal coding scheme we know that nonzero elements of the elementary matrix coupling two traces must be proportional to a constant times the value of the input to α , \bar{f}_k , thus for the association (\bar{f}_k, \bar{g}_k) we obtain the following representation:

$$A_k = c \begin{bmatrix} g_k(1)\bar{f}_k^{(1)} \\ g_k(2)\bar{f}_k^{(2)} \\ \vdots \\ g_k(N)\bar{f}_k^{(N)} \end{bmatrix}.$$

Here we have assumed that M is large enough so that the normalizing constant for each element can be assumed to be the same. (This is equivalent to assuming that the trace power is approximately the same for every group of neurons in α that projects to a single neuron in β .) The $\bar{f}_k^{(i)}$ are row vectors defined so that $\bar{f}_k^{(i)}(j) = 0$ if neuron j in α does not connect with neuron i in β and

$$f_k^{(i)}(j) = f_k(j),$$

if neuron (j) in α is connected to neuron (i) in β . We can easily calculate the normalizing constant so that for K pairs of associations.

$$A = \sum_{k=1}^{k=K} A_k,$$

$$A = \frac{N}{MP} \sum_{k=1}^{k=K} \begin{bmatrix} g_k(1)\bar{f}_k^{(1)} \\ g_k(2)\bar{f}_k^{(2)} \\ \dots \\ g_k(N)\bar{f}_k^{(N)} \end{bmatrix}.$$

We will assume that K is sufficiently large so that we can make our random variable approximation. We can see that individual terms of this matrix can be approximated by the sum of K of the products of two uncorrelated random variables of mean zero and variance P/N . Thus, variance of an element of A_k averaged over sets of allowable traces is

$$E[f(i)^2 f(j)^2] = E[f(i)^2]E[f(j)^2] = P^2/N^2.$$

K of these elementary connectivity matrices go to form A with

$$\text{var}(a_{ij}) = KP^2/N^2.$$

We can now form the $(S/N)_0$ for recognition. We will assume we have $K + 1$ pairs of associated traces stored in A . We present a trace \bar{f}_k which is assumed to be identical to a stored self-association. Then

$$(S/N)_0 = \frac{(A_k \bar{f}_k \cdot \bar{f}_k)^2}{E\{[(A - A_k)\bar{f}_k \cdot \bar{f}_k]^2\}},$$

where the average is to be taken allowable sets of traces. Then,

$$(S/N)_0 = \frac{P^2}{E\{[(A - A_k)\bar{f}_k \cdot \bar{f}_k]^2\}}.$$

We now evaluate

$$E\{[(A - A_k)\bar{f}_k \cdot \bar{f}_k]^2\}.$$

We know that $(A - A_k)$ is a matrix composed of elements either zero or random variables approximated above. We have the useful formula, if X is a random variable

$$\text{var}(cX) = c^2 \text{var } X. \quad (1)$$

Then the vector approximating $(A - A_k)\bar{f}_k$ is given by a set of random variables of mean zero with

$$\text{var}[(A - A_k)f_k(j)] = (MP/N)(KP^2/N^2) = (KMP^3/N^3).$$

This vector is multiplied by a constant MP/N . Then

$$\text{var}[(A - A_k)\bar{f}_k \cdot \bar{f}_k] = \sum_{i=1}^{i=N} f(i)X = \text{var } X \sum_{i=1}^{i=N} f(i)X = (KMP^4)/N^3.$$

And for the $(S/N)_0$ we find, squaring the normalizing constant according to Eq. 1

$$(S/N)_0 = \frac{P^2}{(N^2/M^2P^2)(KMP^4/N^3)} = MN/K.$$

This result should now be discussed. Let us first note that this result is consistent with those obtained in previous papers, in particular, it shares with them the important property that an interactive memory works better as it gets larger, and, in this case, more highly interconnected. This finding is suggestive since it is known that in mammals the amount of cortex associated with a function appears to be determined by the relative importance of the function in the animal's behavior. Two examples are given by Thompson (Ref. 32, p. 317). First, he notes that anatomical studies indicate that in the dog, a relatively small amount of cortical tissue is devoted to representation of forepaw. In the racoon, which makes extensive use of its forepaws, there is a far larger amount of cortex devoted to forepaw representation. Second, and perhaps the best known example of this, is the grotesque little man representing the relative size of parts of motor cortex (determined by noting response to electrical stimulation) in humans (Ref. 32, p. 318). He has huge fingers and face and relatively tiny body and feet. Also, the simple increase in size of the human brain during evolution suggests the presence of a significant mass effect where the cortex appears to work better (generating a selective advantage presumably due to more complex or appropriate behavior) as it gets larger (Ref. 22, p. 634).

Second, let us point out that the recognition model postulated here, where a self-association is assumed to be stored, is also capable of recognizing temporal sequences (Fig. 4). If we assume that there is a significant delay between the arrival of a pattern of activity \bar{f} at α and the generation of pattern of activity \bar{g} at β , we can see that if the input pattern changes, say, from \bar{f} to \bar{g} during this time delay, the formal scheme for recognition of the sequence and for recognizing a self-association are identical, giving rise to identical $(S/N)_0$. (Systems with long conduction times (20-200 msec) are characteristic of cortex.)

Third, let us briefly try to justify the choice of a matched filter recognition scheme. Although it has the maximum $(S/N)_0$ of any simple filter and requires only local operations to form, this, of course, is no guarantee that it is used by the nervous system. The presence of a matched filter would have some implications.

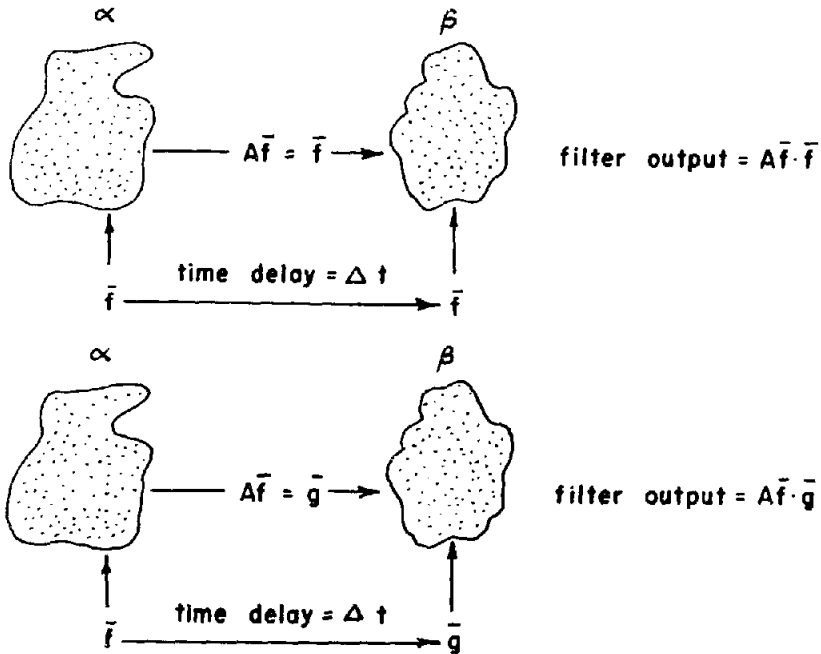


FIG. 4. A system recognizing traces by self-association also allows recognition of temporal sequences. When trace \bar{f} is impressed on set of neurons α it may give rise to output \bar{f} (plus noise) if the self-association (\bar{f}, \bar{f}) has been stored in A . The filter output $A\bar{f} \cdot \bar{f}$ indicates whether or not the self-association (\bar{f}, \bar{f}) was stored. If we assume that the connection between α and β requires a time delay, Δt , then if the input to the system, \bar{f} , changes to \bar{g} during the same period that the association $A\bar{f}$ is being generated at β , the filter output is identical to that for recognition.

A simple matched filter (a variant of a template matching scheme) runs into the problem of the so-called perceptual invariants. The problem can be avoided by assuming that storage takes place after invariant transformations have occurred. However, even a simple template scheme is more realistic than it might seem at first. In the visual system, it is clear that matched filter detection of a trace that still retains a good deal of topographic organization is not rotation invariant. Although some literature on visual pattern recognition makes the assumption that human

pattern recognition is largely rotation invariant, this is not the case, as sensory psychologists have known since the last century [18]. (A simple experiment will demonstrate the effect: it is very difficult to recognize even familiar faces if the observer's head is tilted 45° .) Relatively little quantitative work has been done on this effect. Dearborn (1899, described in Ref. 18) required his subjects to detect repeated presentations of forms displayed with various degrees of rotation. Dearborn reported that cards repeated in their original orientation were recognized 70% of the time while forms rotated by 9° were recognized only 43% of the time.

The matched filter is translation invariant if even elementary centering systems are permitted, as head or eye movements in vision. It is pertinent to point out that animals and humans have an extensive repertoire of orienting and gaze directing behaviors. Thus a snake, with an eye with a slit pupil, can be seen to keep its pupil opening vertical, preserving a constant orientation of the retinal image, no matter what the angle of the snake's body. Size invariance is a more difficult problem. Simple schemes have been proposed (Ref. 6, p. 110) which preserve size invariance over a limited range in ways which are compatible with a matched filter detection system. It is known [31] that brain damage which leaves a large central scotoma ("blind spot") does not interfere with recognition of figures, such as a large triangle, whose contours ring the scotoma, a result which might be expected from a matched filter. In any case, there often appears to be an optimal size of retinal image for many complex items. Thus the size of type will unconsciously determine the distance at which printed matter will be held, both very large and very small type being "hard to read". In vision, a complex system is present to give an output which is directly related to size invariance.

A MORE COMPLEX SYSTEM

More detailed calculations are possible on the model. The most interesting involve associations between unrelated traces. In this section an example of such a system will be considered.

Assuming that the trace, \bar{g} , associated with \bar{f} , is a trace that can be recognized by itself allows us to use a recognition and an association system together to improve the $(S/N)_0$. Figure 5 shows such a system. We assume that K_A traces are stored in the synapses connecting α and β and that K_B traces are stored in the recognition system attached to β . We assume that α and β each have N neurons and that all stored traces have power P . A is the connectivity matrix associated with connections between α and β ; B is the connectivity matrix representing the recognition system associated with β .

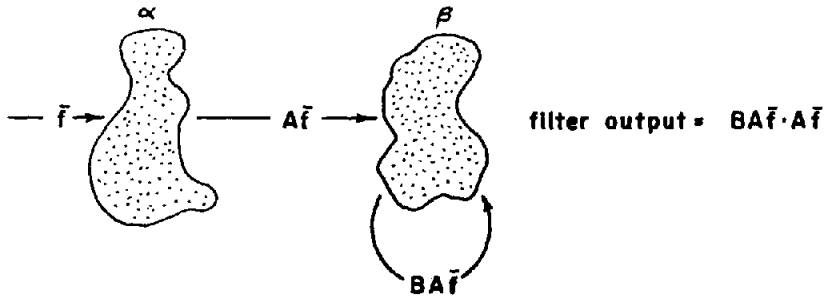


FIG. 5. A more complex system allows the reconstruction of output traces uncorrelated with the input traces. It combines a projection system A , coupled with a projection system which couples β with itself.

We will first calculate the equivalent of the recognition $(S/N)_0$ for this system and then show how this result can be generalized to partially reconstruct the associated trace.

We wish to know whether a trace \tilde{f} presented to α generates a trace, \tilde{g} , which is one of the traces recognized by the recognition system.

For these calculations we will introduce the notation \tilde{n}_A to indicate the noise added to the associated trace \tilde{g} due to the presence of the other stored associations, and \tilde{n}_B to represent the added noise produced by the recognition system, thus,

$$A\tilde{f} = \tilde{g} + \tilde{n}_A, \quad B\tilde{g} = \tilde{g} + \tilde{n}_B.$$

Since we only have noisy information about \tilde{g} , we form our statistic (equivalent to the recognition statistic discussed previously when $A = I$, the matrix with all diagonal elements equal to one) as

$$V = A\tilde{f} \cdot B A\tilde{f}.$$

Then

$$\begin{aligned} A\tilde{f} \cdot B A\tilde{f} &= (\tilde{g} + \tilde{n}_A) \cdot (B\tilde{g} + B\tilde{n}_A) \\ &= (\tilde{g} + \tilde{n}_A) \cdot (\tilde{g} + \tilde{n}_B + B\tilde{n}_A) \\ &= \tilde{g} \cdot \tilde{g} + \tilde{g} \cdot \tilde{n}_B + \tilde{g} \cdot B\tilde{n}_A + \tilde{g} \cdot \tilde{n}_A + \tilde{n}_A \cdot \tilde{n}_B + \tilde{n}_A \cdot B\tilde{n}_A. \end{aligned}$$

We assume as before that traces are uncorrelated, trace elements are uncorrelated and the mean of elements in a trace is zero. We will sketch the calculation of these quantities, averaged over sets of allowable traces.

We wish to calculate the variance of $\tilde{g} \cdot \tilde{n}_B$. In the calculation in the previous section we calculated $(S/N)_0$ for recognition, thus

$$\tilde{g} \cdot B\tilde{g} = \tilde{g} \cdot \tilde{g} + \tilde{g} \cdot \tilde{n}_B = P + \tilde{g} \cdot \tilde{n}_B.$$

We found

$$(S/N)_0 = (\tilde{g} \cdot \tilde{g})^2 / \text{var}(\tilde{g} \cdot \tilde{n}_B) = M_B N / K_B.$$

Solving

$$\text{var}(\tilde{g} \cdot \tilde{n}_B) = (K_B P^2) / (M_B N).$$

Similarly,

$$\text{var}(\bar{g} \cdot \bar{n}_A) = (K_A P^2)/(M_A N).$$

The term $\text{var}(\bar{n}_A \cdot \bar{n}_B)$ is also easy to calculate. We can show that the power

$$P_{n_A} = \bar{n}_A \cdot \bar{n}_A = (K_A P)/M_A,$$

$$P_{n_B} = \bar{n}_B \cdot \bar{n}_B = (K_B P)/M_B.$$

Then we see that since the variance of an element of $\bar{n}_A = K_A P/M_A N$ and of $\bar{n}_B = K_B P/M_B N$ and there are N elements,

$$\text{var}(\bar{n}_A \cdot \bar{n}_B) = (K_A K_B P^2)/(M_A M_B N).$$

We must go to the details of matrix B for calculation of the variance of the last two terms. We can approximate elements of B by random variables with variance $K_B P^2/N^2$; \bar{n}_A is approximated by elements with variance $K_A P/M_A N$. Variance of the product vector is given by a vector with elements with variance

$$(N/M_B P)^2 M_B (K_B P^2/N^2) (K_A P/M_A N) = (K_A K_B P/M_A M_B N).$$

The first term in the above expression arises from the normalizing constant which multiplies B and which appears as the square in the variance.

We see that

$$\text{var}(\bar{g} \cdot B \bar{n}_A) = (K_A K_B P^2)/(M_A M_B N)$$

$$\text{var}(\bar{n}_A \cdot B \bar{n}_A) = (K_A^2 K_B P^2)/(M_A^2 M_B N).$$

Now we can find $(S/N)_0$:

$$\begin{aligned} (S/N)_0 &= \frac{P^2}{\frac{K_B P^2}{M_B N} + \frac{K_A K_B P^2}{M_A M_B N} + \frac{K_A P^2}{M_A N} + \frac{K_A K_B P^2}{M_A M_B N} + \left(\frac{K_A}{M_A}\right)^2 \frac{K_B P^2}{M_B N}} \\ &= \frac{N}{(K_B/M_B)(K_A/M_A + 1)^2 + K_A/M_A}. \end{aligned}$$

If we assume $K_A = K_B = K$; $M_A = M_B = M$ as a simplifying approximation, then

$$(S/N)_0 = \frac{MN}{K[(K/M + 1)^2 + 1]}.$$

If many traces are stored (K/M is large)

$$(S/N)_0 = M^3 N/K^3.$$

This result would signify that $(S/N)_0$ would drop off rapidly as K became very large, but indicates that there is still a linear dependence of $(S/N)_0$ on N . However, the $(S/N)_0$ is very dependent on M , indicating that a highly interconnected system should work far better than a weakly interconnected system.

This $(S/N)_0$ is related to the question as to whether a trace presented to the system has an associated trace present in the system but says nothing

about the details of the associated trace, a result of more interest. In general we would wish to extract as much information as possible from the association, in particular, we wish to have an idea of the structure of the associated trace—enough detail to be able to use it in further processing in the nervous system.

If we have N elements, a good strategy to recover information about the associated trace is to process the output of groups of elements, increasing or decreasing the number of elements in a group as needed to effect a compromise between amount of recovered detail and noise added by the memory.

As an example of this approach, we can use the formula just derived in a grouping scheme. We will assume that we are most interested in recovering the magnitude of the sum of squares of the values taken by elements. If we take V ,

$$V = A\bar{f} \cdot BA\bar{f},$$

and note that the dot product can be taken over groups of R elements in β , we see that for a group of R elements the $(S/N)_0$ is given by

$$(S/N)_0 = \frac{R}{(K_B/M_B)(K_A/M_A + 1)^2 + K_A/M_A}.$$

(This result holds if $\bar{g} \cdot \bar{g}$ over the R elements has the average value,

$$E(\bar{g} \cdot \bar{g}) = RP/N.$$

The $(S/N)_0$ for a given region depends on the power of $\bar{g} \cdot \bar{g}$ in the region.)

We see from the above that V will be an estimate of the sum of squares of the trace summed over the R chosen elements. By increasing R we obtain the intuitive conclusion that we can increase the $(S/N)_0$ while simultaneously losing detail on the level of single elements.

It should be noted that a system like this apparently exists in the retina where the size of the spatio-temporal light intergrating area can be increased or decreased depending on the average light intensity, sacrificing acuity for detectability in low light and attaining very high acuity in high light intensities (Ref. 8, chaps. IV, V, VI; Ref. 4).

By grouping elements, the $(S/N)_0$ of an average output can be varied over a range

$$\frac{1}{(K_B/M_B)(K_A/M_A + 1) + K_A/M_A} \leq (S/N)_0 \leq \frac{N}{(K_B/M_B)(K_A/M_A + 1) + K_A/M_A}.$$

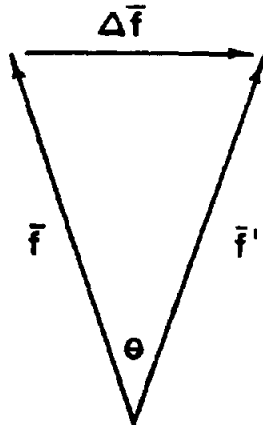
If we assume that groups of neurons will be chosen spatially close to each other, we see that by assuming even weak topographic organization of the

neurons involved, this will decrease the high spatial frequency response of the recovered associated trace while the low frequency outline will be more stable. Detail will be lost but major structure will be preserved. This is in accord with common sense.

SELECTIVITY

An important problem for a model where traces can easily be confused is the question of selectivity, that is, how "close" must a trace be to the stored trace to be recognized.

We can do some simple geometrical calculations to get an insight into this question. We have assumed all traces have power P , that is, they are vectors with tips lying on an N -dimensional hypersphere of radius $P^{\frac{1}{2}}$ (Fig. 6).



$$\bar{f} \cdot \bar{f} = \bar{f}' \cdot \bar{f}' = P$$

FIG. 6. A trace \bar{f} is perturbed by a vector $\Delta\bar{f}$, generating a new trace \bar{f}' . Both \bar{f} and \bar{f}' have the same power.

Let us consider \bar{f} perturbed by a vector $\Delta\bar{f}$ so that

$$\bar{f} + \Delta\bar{f} = \bar{f}'; \quad \bar{f} \cdot \bar{f} = \bar{f}' \cdot \bar{f}' = P.$$

Considering a simple recognition system with only one stored trace, \bar{f} , if we present our perturbed vector \bar{f}' , then

$$\begin{aligned} A\bar{f}' \cdot \bar{f}' &= (A\bar{f} + A\Delta\bar{f}) \cdot \bar{f}' \\ &= (A\bar{f} + A\Delta\bar{f}) \cdot (\bar{f} + \Delta\bar{f}) \\ &= P + \bar{f} \cdot \Delta\bar{f} + \bar{f}' \cdot A\Delta\bar{f} + A\Delta\bar{f} \cdot \Delta\bar{f}. \end{aligned}$$

We see $A\Delta\vec{f}$ can be calculated if we assume M is large enough so that for all (j, k)

$$\vec{f}^{(j)} \cdot \Delta\vec{f} = \vec{f}^{(k)} \cdot \Delta\vec{f} = (M/N)(\vec{f} \cdot \Delta\vec{f}),$$

is an adequate approximation. Then

$$A\Delta\vec{f} = (N/MP)(M/N)(\vec{f} \cdot \Delta\vec{f})\vec{f},$$

and

$$A\Delta\vec{f} \cdot \vec{f} = \vec{f} \cdot \Delta\vec{f}.$$

Similarly,

$$A\Delta\vec{f} \cdot \Delta\vec{f} = (\vec{f} \cdot \Delta\vec{f})^2/P.$$

By geometry (Fig. 6) we see

$$\vec{f} \cdot \Delta\vec{f} = P(\cos \theta - 1).$$

If we form a ratio of filter output as a function of θ over filter output when $\theta = 0$, where θ is the angle between \vec{f} and \vec{f}' , then,

$$(A\vec{f}' \cdot \vec{f}')/(A\vec{f} \cdot \vec{f}) = \cos^2 \theta.$$

We wish to find the angle θ at which the output of the filter is reduced by half; we see $\theta = 45^\circ$. This corresponds to the half-power point of engineering filter theory.

We now ask what are the chances of a random trace falling into the region where the filter output is greater than half of its maximum value. Restated, this is the ratio of the surface content of the N -dimensional hypersphere within 45° of a given trace to the surface content of the N -dimensional hypersphere [30]. We can calculate this quantity by observing we can establish a recursion formula, when S_N is the surface-content of an N -dimensional hypersphere and V_N is the volume-content of an N -dimensional hypersphere and when R is the radius,

$$S_N = 4 \int_0^{\pi/2} V_{N-2} R d\theta.$$

By a similar argument, if SC_N is the surface-content of the hypersphere contained in a cone of base angle θ and VC_N is the volume content of such a cone (assuming the cone extends in positive and negative directions, like an hourglass)

$$SC_N = 4 \int_0^\theta VC_{N-2} R d\theta.$$

Thus the desired ratio is given, for N elements in the vector by

$$SC_N/S_N = \frac{4RVC_{N-2} \int_0^\theta d\theta}{4RV_{N-2} \int_0^{\pi/2} d\theta} = (2\theta/\pi)(VC_{N-2}/V_{N-2}).$$

If N is even

$$(SC_N/S_N) = (2\theta/\pi)^{N/2},$$

and if N is odd,

$$(SC_N/S_N) = (2\theta/\pi)^{(N+1)/2}$$

When $\theta = 45^\circ$, $(2\theta/\pi) = \frac{1}{2}$. This quantity becomes extremely small when N is large. Thus we see that selectivity as well as $(S/N)_0$ increases as N increases.

CORTICAL MODEL

The preceding discussion suggests a simple cortical model which would be capable of generating an interactive memory. The model will be sketched here.

Cortex is similar in histological structure over much of its extent. The most frequently occurring kind of cell—up to 80% of neurons in rabbit cortex [16]—are variants of what are called “pyramidal cells” with characteristic pyramid-shaped cell bodies and dendrites perpendicular to the cortical surface, extending almost to it, and then branching extensively parallel and close to the surface. These cells show rich synaptic contacts on all parts of the dendrites and have an array of “spines” which apparently correspond to synaptic contacts on the dendrites, although not all synapses are associated with spines. These cells also show extensive collateral branches of their axons, often projecting both intra- and extracortically. It is hard to estimate the number of different axons synapsing on pyramidal cells. Measurements quoted by Brodal (Ref. 5, p. 658) indicate that in monkey as many as 60,000 synapses are found on a motor area neuron, in visual cortex, about 7,000. In any case, several thousand cells may synapse on a single pyramidal cell. Origins of these synapses differ: some are sensory afferents, others intracortical fibers of various types—callosal, short projections, long projections. Much is known about the details of interconnection, for this brief discussion we will merely assert that the cortex is very highly interconnected.

Most electrophysiological recordings from cortex have been made from pyramidal cells, simply because they are the largest and most numerous cells present. This data suggests that cells behave in many respects as unique individuals, each cell taking its particular sample of the surrounding afferent inputs. Since afferents are organized topographically in some cases (vision, for example) this gives common features to neighboring cells, but when cells are investigated in detail, each cell appears different from its neighbors in its particular blend of properties. Creutzfeldt and Ito [10] and Creutzfeldt [9] found that primary visual cortical cells appeared to receive large inputs from only a few lateral geniculate fibers (2 to 4). They

suggested that the variety of forms of receptive field displayed by primary visual cortex cells could be explained by the possible permutations of the types of lateral geniculate cells. Goldstein et al. [17] found that auditory cortex is weakly organized tonotopically with nearby cells showing quite different tuning curves. Single units in primary auditory cortex responded in an "individualistic and variegated" manner. Hubel and Wiesel [19] suggested that in monkey visual cortex, many overlapping mosaics of sensory parameters are present, the sample received by a single cell depending on the afferents in its particular area.

The viewpoint that cortical cells respond as individuals finds confirmation in a paper by Noda and Adey [25]. They recorded single units chronically from cat association cortex (parietal cortex). In over 70 cases they recorded two units with the same electrode. They separated and analyzed these units. They found that when the animal was alert and awake, the cells, spatially close together, were uncorrelated in their discharge. This was true as well when the animal was in REM sleep, presumably a time of intense subjective experience. When the animal was awake but drowsy, there was a weak correlation between cells and when the animals were in deep sleep, cell discharges were highly correlated.

Let us now assume that we can identify the "traces" discussed earlier with patterns of increased or decreased pyramidal cell discharge in a given cortical area.

We assumed in the mathematical model that elements of a single trace were uncorrelated. The Noda and Adey findings [25] lend direct support to this assumption. Observations on EEG suggest a similar conclusion more indirectly. The generators apparently giving rise to the EEG tend to desynchronize (i.e. the resulting EEG amplitude distribution approaches a normal distribution more closely) during REM sleep and the awake state than in deep sleep [2]. This evidence indicates that during times when memory may be presumed to be functioning (alert, awake state and dream state) cortical elements tend to show uncorrelated activity.

We must consider the interconnections of the cortex. I would like to suggest that there is some evidence indicating the presence of two major classes of synapses significant for memory in the pyramidal cells. More probably, there may be a spectrum of types of synapses with the classes to be described forming two ends of a spectrum, but the analysis is not affected by assuming two discrete classes instead of a spectrum.

First, there is the familiar type of synapse, with relatively large PSPs with relatively short rise times (1.5–10 msec). This type of synapse would be characteristic of incoming sensory afferents or strong intracortical projection systems. The work of Rall [27] suggests that PSPs can be classified on the basis of shape as to their electrical distance from the cell

body, since distant synapses will give slowly rising long PSPs due to the cable properties of the dendrites. Closer synapses will give more rapidly rising and falling PSPs, other properties of the PSPs being identical. Thus the first system would consist primarily of synapses close to the cell body.

However, we have noted that many thousand synapses occur on pyramidal cells. Since there are generally only a few clearly recognizable PSPs [11], one might wonder what the other synapses are doing.

Intracellular recordings from cortical neurons generally show spontaneous fluctuations of the membrane potential [1]. Although spontaneous PSPs are often present, this activity usually appears superimposed on some other fluctuating activity. Elul has determined the frequency spectrum of neuronal intracellular activity and finds it to have a frequency spectrum similar to that of the EEG with most energy in the frequencies below 10 cps. Although the distribution of amplitudes of the intracellular slow wave activity is not Gaussian (it is somewhat asymmetrical skewed in the positive direction) it is sufficiently close to the histogram expected of a Gaussian process to suggest some interesting modeling possibilities to see if the known asymmetries of the cortical neurons coupled with various input probability distributions might lead to the observed distribution.

In any case, I would like to suggest, as a second system, that this slow activity is generated by the activity of other synapses on the pyramidal cells. These synapses would be electrically remote from the cell body (on spines?) and would be severely low-pass filtered. Their influence would appear as a biasing of cell activity where the activity of any given input would be unnoticeable but the action of many thousand weak, long time course inputs would be highly significant. Since we know that most cells in cortex are spontaneously active and are interconnected, assuming weak, extensive interconnections seems the most reasonable way of explaining these ubiquitous low frequency membrane fluctuations.

Burns (Ref. 6, p. 64) has presented evidence which suggests that isolated cat forebrain acts very much like a highly interconnected random network.

Another line of evidence indicates the presence of a system electrically remote from the cell body. Smith and Smith [29] in an extensive study of the statistics of spontaneous cell activity in cortex detected the presence of two distinct systems giving rise to the spontaneous activity of cortical cells in their preparation, the unanesthetized, isolated cat forebrain. They found that the statistics of spontaneous activity could be explained as the result of two Poisson processes, one a "shower" of spikes following Poisson statistics and the other a process which "gated" the shower on and off at random intervals. They found when they passed a weak polarizing

current through the tip of their electrode that the "gating" process was greatly affected but the average frequency of the Poisson shower was not changed drastically. Their interpretation of this is that the gating process is located near the cell body and the system giving rise to the shower is located farther from the soma. One might like to identify the "shower" process with the extensive weak interactions and the "gating" process with the first system but this is premature. Creutzfeldt, et al. [11] suggest in interpreting their stimulation data that "nonspecific" thalamic afferents are electrically more remote than "specific" thalamic afferents and that most inhibitory synapses appear to be located on the soma. Scheibel and Scheibel [28] suggest on anatomical grounds that nonspecific afferents exert a temporally diffuse biasing control on the pyramidal cell because of the pattern and location of the terminal arborization of these fibers.

Let us consider how an idealized model which assumes a relatively strong, fast system, and a weak, slow, highly interconnected system might work to generate memory. We have assumed that in learning, as we assumed in the mathematical model, synaptic strength coupling two neurons in the weak system is increased or decreased proportional to the product of pre- and postsynaptic activity during a given period—when the two groups of cells are "associating" the desired traces—then we can see that the weak system will give rise to an interactive memory system of the type discussed in previous sections.

The system would work roughly like this. A sensory input is impressed on the first group of neurons by the fast, strong system. A pattern of activity, corresponding to the input trace, is established in the first group of neurons. This pattern then "filters through" the weak system generating associations in a second group of neurons by the mechanisms discussed in the earlier sections of this paper. These associations can interact with the input trace to produce the recognition statistics. The input trace might be made available to the second group of neurons by a connection between first and second groups by the fast strong system or by another means, such as a direct sensory projection. Or the trace produced by the weak system could then be processed further by the second group of neurons and its connections.

A slow, weak system seems like a good candidate for memory since it would have the properties of extensive interconnection (high " M ") and long time course (allowing complex temporal interactions at times consistent with the time courses of psychological events).

A common pattern of response of cortical neurons is to show a brief burst of activity when a new stimulus is presented and then to show a long inhibition. It might be possible that this inhibition serves to quiet the cells for the memory readout so that the diffuse and noisy memory trace trickling

through the weak system would not be submerged by a strong sensory input.

I would like to thank Professor W. Ross Adey and the Space Biology Laboratory, Brain Research Institute, UCLA, for providing financial support for this work. I would also like to thank Professor Alan Grinnell, Department of Zoology, UCLA for providing greatly appreciated assistance.

REFERENCES

- 1 W. R. Adey, *Neurosciences Research Program Bulletin* 7(2), 75 (1969).
- 2 W. R. Adey, *The Neurosciences*, Vol. 2 (F. O. Schmitt, Ed.) Rockefeller University, New York (1970), p. 224.
- 3 J. A. Anderson, *Kybernetik* 5, 113 (1969); *Mathematical Biosciences* 8, 137 (1970).
- 4 M. A. Bouman, in *Sensory Communication* (W. A. Rosenblith, Ed.) MIT Press, Cambridge, Mass. (1961), p. 377.
- 5 A. Brodal, *Neurological Anatomy*, 2nd ed., Oxford University Press, New York (1969).
- 6 B. D. Burns, *The Uncertain Nervous System*, Arnold, London (1968).
- 7 B. Cleland and C. Enroth-Cugell, *Acta Physiol. Scand.* 68, 365 (1966).
- 8 T. N. Cornsweet, *Visual Perception*, Academic, New York (1970).
- 9 O. D. Creutzfeldt, in *The Neurosciences*, Vol. 2 (F. O. Schmitt, Ed.) Rockefeller University, New York (1970).
- 10 O. D. Creutzfeldt and M. Ito, *Experimental Brain Research* 6, 324 (1968).
- 11 O. D. Creutzfeldt, K. Maekawa and L. Hösl, in *Progress in Brain Research*, Vol. 31 (K. Akert and P. G. Waser, Eds.) Elsevier, Amsterdam (1969).
- 12 J. C. Eccles, in *Brain and Conscious Experience* (J. C. Eccles, Ed.) Springer, Berlin (1966), p. 314.
- 13 S. S. Fox and J. H. O'Brien, *Science* 147, 888 (1965).
- 14 W. J. Freeman, *J. Neurophysiol.* 31, 337 (1968).
- 15 W. J. Freeman, *Mathematical Biosciences* 2, 181 (1968).
- 16 A. Globus and A. B. Scheibel, *J. Comp. Neurol.* 131, 155 (1966).
- 17 M. H. Goldstein, Jr., J. L. Hall, II, and B. O. Butterfield, *J. Acoust. Soc. Am.* 42, 444 (1968).
- 18 H. W. Hake, quoted in *Pattern Recognition* (L. Uhr, Ed.) Wiley, New York (1966), p. 151.
- 19 D. H. Hubel and T. N. Wiesel, *J. Physiol.* 195, 215 (1968).
- 20 L. Maffei, *J. Neurophysiol.* 31, 283 (1968).
- 21 L. Maffei and G. Rizzolatti, *J. Neurophysiol.* 30, 333 (1967).
- 22 E. Mayr, *Animal Species and Evolution*, Harvard, Cambridge, Mass. (1963).
- 23 M. L. Minsky, *Computation: Finite and Infinite Machines*, Prentice-Hall, Englewood Cliffs, New Jersey (1967).
- 24 V. B. Mountcastle, in *The Neurosciences* (G. C. Quarton, T. Melnechuk and F. O. Schmitt, Eds.) Rockefeller University, New York (1967).
- 25 H. Noda and W. R. Adey, *J. Neurophysiol.* 33, 572 (1970).
- 26 D. H. Perkel and T. H. Bullock, *Neurosciences Research Program Bulletin* 6, 221 (1968).
- 27 W. Rall, in *Neural Theory and Modeling* (R. F. Reiss, Ed.) Stanford University, Stanford, Calif. (1964), p. 73.

- 28 M. E. Scheibel and A. B. Scheibel, in *The Neurosciences*, Vol. 2 (F. O. Schmitt, Ed.), Rockefeller University, New York (1970), p. 443.
- 29 D. R. Smith and G. K. Smith, *Biophys. J.* **5**, 47 (1965).
- 30 D. M. Y. Sommerville, *An Introduction to the Geometry of N-Dimensions*, Dutton, New York (1929), chap. 8.
- 31 H. L. Teuber, in *Brain and Conscious Experience* (J. C. Eccles, Ed.) Springer, Berlin (1966), p. 182.
- 32 R. F. Thompson, *Foundations of Physiological Psychology*, Harper and Row, New York (1967).