

# The attraction of parallel distributed processing for modelling cognition

---

*Representation and processing of information in connectionist networks is distributed. Decisions are reached by consensus of a large number of simple components taking place in parallel as stimulus information interacts with stored information. In consequence, connectionist memories display many human characteristics. They are relatively immune to damaged components within the system or to loss of individual components; they allow retrieval by content; they are likely to retrieve typical information categories.*

The last decade has seen an explosive growth in the connectionist modelling of cognitive processes, with simulation of most of the classical experimental paradigms of cognitive psychology. One reason for this enthusiasm is that, independent of their specific modelling human performance at any particular cognitive task, all connectionist models exhibit some general characteristics which are shown by human memory systems and distinguish them from non-biological computational systems and computer programs: They still perform reasonably well after minor damage to components of the system; they still perform reasonably well if their input is noisy or incomplete; they allow memory retrieval by content.

In this chapter we will look at two aspects of connectionist systems which are particularly relevant for these characteristics. Like the principles of interneuronal communication described in chapter 1, these are based on general observations of brain function. First, knowledge representation is *distributed* across many processing units. Second, computations take place in *parallel* across these distributed units. The result is that conclusions are reached on the basis of a consensus of many units rather than depending on any particular one.

These principles put connectionist models in direct contrast to many traditional models of cognitive psychology or artificial intelligence where knowledge representation is serial and computation is serial. In general, such models are not immune to

damage or resistant to noisy input. So a traditional model of, say, reasoning, might give as good a fit to the experimental data as a connectionist model, but it would do so without exhibiting the full range of human characteristics performed the task.

## The representation of knowledge in connectionist networks is distributed

Traditional models of cognitive processing usually assume a local representation of knowledge. That is, knowledge about different things is stored in independent locations. In a traditional model of reading aloud, for example, information about how to pronounce the letter string DOG is stored in one location, information about how to pronounce the string CAT in another. What is more natural? The two pieces of information are independent and would be stored at different times. So storing them independently makes obvious sense for information storage systems we are familiar with in everyday life—directories, telephone directories, computer discs—use local representation. Each piece of information is stored separately. How else could it be done?

In connectionist models information storage is not local, it is *distributed*. There is no one place where a particular piece of knowledge can be located. Consider a segment of network at the bottom of figure 1.2 in chapter 1, part of a larger network which is learning to read aloud. Any input, such as the letter string DOG, excites units and connections all over the network. Learning takes place by adjusting the weights of the connections leading to all output units which have an effect on activity. The knowledge of how to pronounce the input DOG is distributed across many different connections in different parts of the system. It is the effect of all these connections which produces the pronunciation, not any one of them.

The concept of distributed storage may be difficult to grasp at first because it is counter to our everyday experience of information storage systems. The connections which contain the system's knowledge about how to pronounce DOG are distributed as those with the knowledge about how to pronounce any other letter string. The knowledge that the network contains is superimposed on the same connections. Intuitively this may seem entirely implausible. How can the system weights store independent and even contradictory pieces of information? I can see, it can be done, and some of the emergent properties of such systems are intriguingly similar to properties of human cognitive processes. But for this to work this will have to be taken on trust. There are no familiar information storage systems which use distributed coding, so analogy to a familiar system is not possible.

## uted representations are damage resistant and fault

considers the structure of the brain it is remarkable that it ever manages to ne correct conclusion about anything. By any conventional standards e e an entirely unsuitable medium for computation: they die throughout the , causing random loss of stored information; they have a finite probability , even when they are not engaged in signal processing; the response of a any particular input is probabilistic, not fixed.

back at the firing pattern of a single neuron the problem that probabilistic cause for the system will become clear. The upper part of figure 2.1 shows e response of a single neuron in the visual cortex to a stimulus presented to he stimulus is presented at time 0. Time after the presentation of the shown on the horizontal axis. The vertical axis shows the neuron's firing ; an average background firing rate of a few spikes per second. When the presented a signal is superimposed on this. About 50 ms after the stimulus e neuron fires strongly for about 30 ms. 100 ms later it fires again, rather ely, for about 50 ms. This is the signal that the neuron transmits to the ons to which it is connected, indicating what pattern of stimulation it has

ms fairly straightforward. But the histogram was obtained by summing the rms over a number of presentations of the stimulus. The lower part of the s 12 different occasions on which the same stimulus was presented. Each neuron fires there is a vertical spike. If we look at these individual trials the hich emerges is much less clear than that suggested by the overall average at On trial 5 the initial burst was missing. On trial 10 the second burst was On trial 11 the neuron did not respond at all. It is clear that the 'signal' histogram at the top shows is an idealised average. On any given trial the ll only approximate this, sometimes quite closely, sometimes not at all. the system produce a reliable response on every trial when the individual ts only produce their signal on average across trials?

rocessing components in a conventional digital computer produced random ous output, a different response to the same stimulus on different occasions ed from random component drop-out, the system would be totally un- e! Sometimes it would work correctly, but if a computation required access ents of a missing memory unit, or a burst of noise obliterated a signal, the uld be garbage. Although the components in the brain can fire and die at he computations performed by the brain are not unpredictable. With minor becomes a little slower and less accurate, but it still produces roughly the ver. It has to suffer serious damage before it produces nonsense.

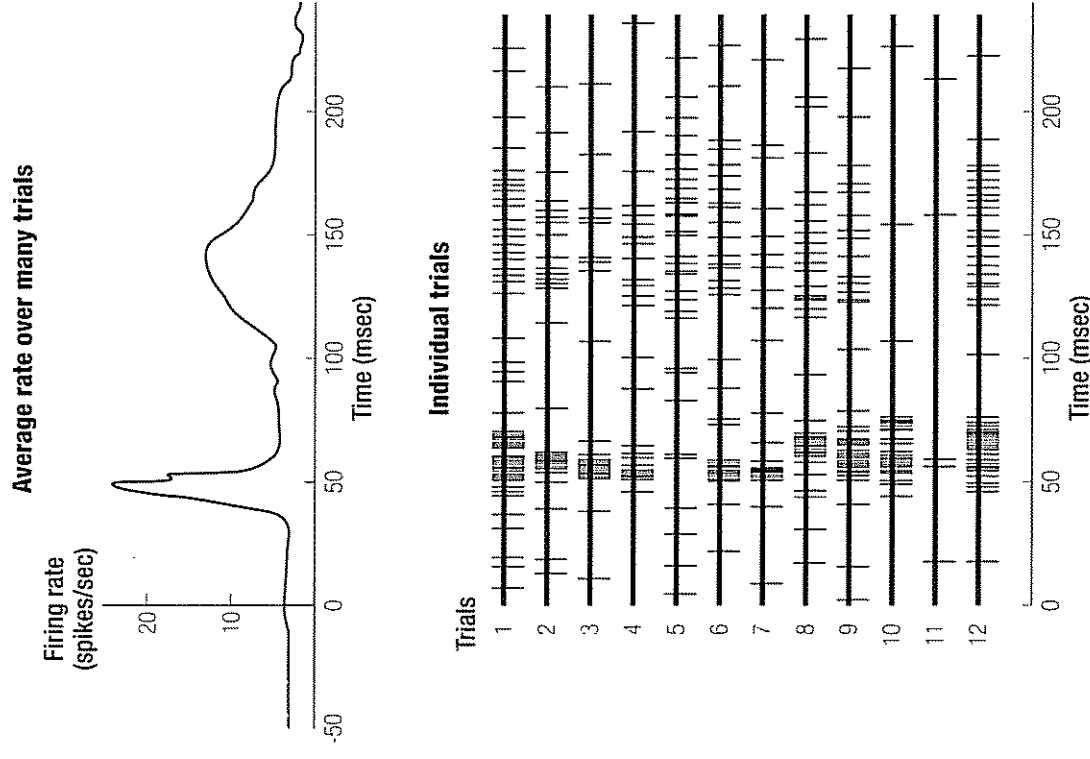


Figure 2.1 The response of a neuron in the visual cortex for 250 ms after stimulus presentation. The firing rate of the neuron summed over a number of presentations of the stimulus. Lower: on 12 of the trials which contributed to the average response shown at the top. Each vertical tick represents the time at which the neuron fired (Based on Morrell 1972 )

The brain escapes the consequences of the unpredictable behaviour of individual neurons because its computations are performed in parallel on representations are distributed over many neurons. No one neuron plays a crucial role in processing. The overall result is the outcome of many distributed sub-computations. The individual components of the calculation are not accurate, the ensemble

theless give an answer which is accurate enough. When the memory required for a calculation in a localist information storage system is the result is disaster. If the first page of a dictionary is missing, there is no checking whether *aardvark* is really spelt like that. But, in a connectionist there is no such thing as 'the memory location required for a calculation' and calculation are spread across the network. If one unit or connection work is damaged, others can make up for the missing part.

System will, of course, be slightly less accurate if a connection is lost. But the loss is quite different in localist and distributed systems. Damage to a system causes some information to be lost totally while other information is lost partially. In a distributed system any damage causes partial loss of a range of information. As damage increases, the performance of the system inevitably begins to drop. A small amount of damage may have no noticeable effect on the output of the system. The ability of brains and connectionist models to continue to produce a correct answer following damage, rather than suffering catastrophic failure, is an example of fault tolerance referred to as *degradation*.

## Connectionist networks allow memory access by content

Memory is content-addressable. That is, you can access a memory by using the content of the information contained in the memory (the *content*) as a retrieval key. Unlike retrieval from familiar forms of information storage, such as directories, telephone directories or computer discs. In these, the place where the information is stored has an *address*. The only way to access information for a specific address is the address. For example, in a dictionary the address is the spelling of the word; the information stored with this address is the definition of the word. In a connectionist system a typical address is the name of a file; the information which can be retrieved with this address is the contents of the file.

To understand the difference between accessing a memory by address and by content, contrast retrieving information from a dictionary with obtaining the same information from a person. Imagine that you want to know the name of a man-made dam across a valley, to contain water to build up a head of pressure to generate electricity. With a dictionary, there is no way to access the location where this information is stored and extract the missing piece, the word 'dam'. If you start with a Dam Address Memory (DAM) you can access all the information stored at the address. Without access to the DAM you can access only part of the information would be able to retrieve the rest. Unlike a dictionary, human memory allows access to any part of the information that forms the memory. One of the reasons connectionist models of human memory are attractive is that content address-

ability follows as a natural consequence of their distributed structure. Content-addressability can be built into localist storage but only by adding a complex referencing system.

Any information processing system which works in the brain must be fault-tolerant because the signals it has to work with are seldom perfect. There is a random element to neuronal firing; speech is usually heard against a background of other objects rarely present the same image on different occasions. An attractive content-addressable memory is that it is inherently fault-tolerant. Imagine you asked you to guess who they were thinking of: 'This man was a British Conservative politician. He became Prime Minister in 1978 and was Prime Minister during the Falklands War. He was ousted from office by his own party, being held responsible for the fiasco of the Poll Tax. He was eventually replaced as Prime Minister by a Major.' You could probably suggest 'Margaret Thatcher' as an answer despite the fact that some of the information is incorrect. Mrs Thatcher did not become Prime Minister until 1979, of course. With content-addressable memory the wrong evidence pointing to one answer can overcome other evidence that is inconsistent. The best fit solution can be chosen even if it is not perfect. This is unlike memory in which access by address is the only possibility. Any error in the address will lead to failure. A search of *Who's Who* using the address 'Margaret Thatcher' would find nothing, despite the fact that most of the search term fits an existing entry.

## Retrieving information from a distributed database

To see how a distributed system with parallel processing works in practice, look at retrieval from a simple connectionist memory described by McClelland (1981).<sup>1</sup> This memory demonstrates content addressability and fault tolerance. It also shows typicality effects in retrieval—if asked to retrieve a random member of a category it will produce a typical member. Considering the simplicity of the simulation it demonstrates a remarkable range of human characteristics in memory retrieval.

Imagine that you live in a neighbourhood where many of your male acquaintances belong to one of the two rival local gangs, the Jets or the Sharks. Your knowledge about these characters will come from a succession of independent episodes. One night Fred emerges from behind a bush and offers you some white powder. Another night Dave and his wife trading insults as she drives off with a car full of stolen goods. Everyone in the bar is laughing because Don has been admitted to college on the basis of forged examination results. Nick hangs out with Karl whom you know as a Shark. All these pieces of information about your neighbours are gathered together in your memory.

<sup>1</sup>The Jets and Sharks memory system is not implemented in tearn but its properties and the simulation described in the text can be explored using the iac program in McClelland and Rumelhart (1981).

table 2.1, address + contents, is a logical way of storing the information name *Alan* is the key which binds one set of information together, *Clayde* another set, and so on.

This form of storage is efficient for retrieving information in response to a question like 'Is Fred a pusher?'. The question contains the address, and the address directly to the place where the information which provides the answer is stored. It is not so good for answering other enquiries. If you are asked 'Do you know the name of a pusher?', the only way is to search through the list of addresses until you find one where the information *Pusher* is stored under **Occupation**. Although the information *Pusher* is stored at many locations, it does not form part of the answer. So an answer to this question cannot be extracted directly from the memory.

Admittedly, with this particular database it would not take long to find a particular address at random. But a more realistic representation of knowledge of these people would include many unique pieces of information, such as that Fred's grandparents came from Ballylickey. The only way to store a system like that shown in table 2.1 is as the fact 'grandparents came from Ballylickey' at the storage location with the address *Fred*. The question 'Whose grandparents came from Ballylickey?' could only be answered by random search of the addresses until the one containing that information was found. This might take a long time although you would get there in the end. But a human memory would not be like that. If you could remember the information at all you would usually be able to answer reasonably quickly. This is because human memory can be accessed by *content*—any part of the knowledge base can be used to access any other part. 'Ballylickey' can be used as a cue, and will lead to *Fred*. The information shown in table 2.1 is perfectly logical. Indeed, it is probably the sort of system you would use if you were asked to store the information about the Jets and Sharks. But human memory cannot be organised like this. A memory organised like this does not allow content-addressable retrieval; human memory does.

Another way of seeing why human memory cannot be organised so that it is only possible by address is to consider how you would answer the question 'Who is the Jets like?'. Table 2.1 allows easy access to information about individual Jets, but it offers no simple way to answer questions requiring generalisations about the number of entries. Human memory does allow generalisations across memory. A person who knew these two gangs could probably tell you that the Jets were younger than the Sharks without having to think very hard. Although the method of information storage shown in table 2.1 seems natural, it cannot be a way that human memory is organised.

(1) *Setting up a distributed database for the Jets and Sharks base.* McClelland explored the consequences of storing the information about the Jets and Sharks in a distributed system. The architecture of his system is shown in the upper

Contents				
Gang	Age	Education	Marital Status	Occupation
Jets	30s	JH	Married	Burglar
Jets	40s	JH	Single	Pusher
Jets	40s	JH	Single	Bookie
Sharks	30s	HS	Divorced	Pusher
Sharks	30s	Col	Married	Burglar
Jets	30s	HS	Single	Bookie
Sharks	40s	HS	Married	Burglar
Jets	20s	HS	Single	Pusher
Jets	20s	Col	Single	Pusher
Jets	20s	JH	Divorced	Burglar
Jets	20s	HS	Married	Pusher
Sharks	30s	JH	Single	Bookie
Jets	20s	JH	Divorced	Burglar
Jets	20s	JH	Married	Burglar
Sharks	40s	HS	Married	Bookie
Sharks	20s	HS	Single	Burglar
Jets	20s	JH	Married	Burglar
Jets	30s	JH	Single	Bookie
Sharks	30s	HS	Single	Bookie
Sharks	30s	Col	Married	Bookie
Sharks	30s	HS	Single	Pusher
Sharks	30s	Col	Married	Pusher
Jets	20s	HS	Single	Bookie
Sharks	30s	Col	Married	Pusher
Jets	30s	JH	Single	Pusher
Sharks	30s	HS	Divorced	Burglar
Jets	20s	Col	Married	Bookie

McClelland 1981.)

over the years. Fred is a pusher; Dave is divorced; Don went to college; Shark.

2.1 shows how this information might be stored in a conventional distributed storage system. The system has a set of storage locations, corresponding to cards in an index file or a set of files on disc in a computer, for example, each card being headed by an address. The address is the name of the person. Each new piece of information relating to him is stored at the location headed by this address. In this way we can imagine that we have information about the Gang each person belongs to (*Jets* or *Sharks*), a rough idea of his Age (*20s*, *30s* or *40s*), the extent of his education (*Junior High*, *High School* or *College*), his Marital Status (*Married*, *Single* or *Divorced*), and his Occupation (*Bookie*, *Burglar* or *Pusher*). The format used in



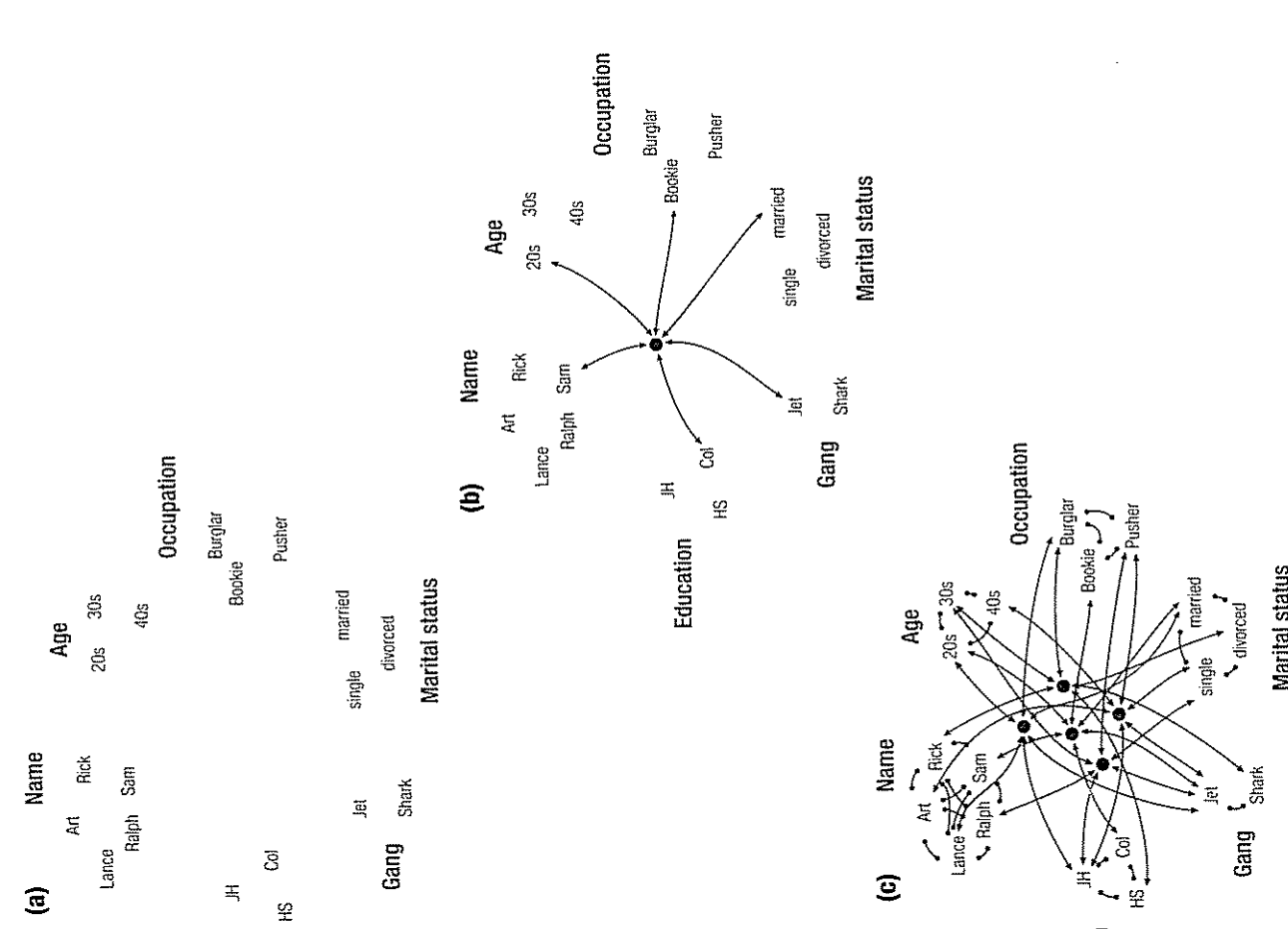
**Figure 2.2** The architecture of McClelland's system for storing the information about the Sharks shown in table 2.1 (a) Each cloud represents an area of knowledge about the members of a gang, with the nodes within a cloud representing possible instances (b) The information is represented by setting up excitatory connections between the facts that are known about a node, done by setting up a Person node (the black node in the central circle) and linking this to all nodes which represent his properties If any one of these nodes becomes active these links will activate the nodes representing his other characteristics will be activated (c) The excitatory connections to represent all the information about five members of the gang have been entered connections (links with filled circles on their ends) have been set up between competing instances within each area of knowledge When the model runs, any node which has a positive activation inhibit any other node to which it is connected by one of these links. (Based on McClelland 1981)

inputs are exactly balanced.<sup>4</sup> The Name node which is most strongly activated in the steady state is reached is the system's answer to the question.

Does such a system behave like human memory? Apart from being able to store information it had been given directly, by answering such questions as 'Fred does Fred do?', anyone who knew these people would find it easy to answer questions like 'Do you know the name of a pusher?' or 'What are the Joneses?' These are the sort of questions which it is difficult to answer with a localist store organised like table 2.1. Will this distributed, connectionist memory system be any easier?

With the aid of the bottom part of figure 2.2 it is possible to get some idea of what happens when the system runs. To see what answer the system will retrieve if asked: 'Can you remember the name of a pusher?' the *Pusher* node is activated, and activity passes along all the connections from the *Pusher* node. So the *Ralph* *Person* nodes become excited because there is a positive connection to them from *Pusher*. (In the real model all the *Person* nodes of pushers would be excited, for simplicity we will just follow two of them.) The *Bookie* and *Burglar* nodes are inhibited because there are negative connections to them from *Pusher*. In the processing cycle the *Ralph* *Name*, *Jet*, *30s*, *JH*, *Single* and *Pusher* nodes are excited by the *Ralph* *Person* node, and the *Art* *Name*, *Jet*, *40s*, *JH*, *Single* and *Pusher* nodes become excited by the *Art* *Person* node. At each succeeding cycle every node that is active influences every node to which it is connected by an amount which depends on its activity level and in a direction which depends on whether the connection to it is excitatory or inhibitory. So, for example, the fact that the *30s* node is excited by the *Ralph* *Person* node will in turn cause excitation of all the *Person* nodes connected to *30s*, and inhibition of the *20s* and *40s* *Age* nodes. At the same time, excitation of the *40s* *Age* node by the *Art* *Person* node may be sufficient to excite

<sup>41</sup>In McClelland's model the activity of each unit also decays on each cycle by an amount proportional to its activity level. This affects the dynamic behaviour of the net but to understand why the net reaches a steady state it can be considered as another negative input contributing to the balance between positive and negative inputs to each unit.



it excites all the nodes to which it has a positive connection and inhibits all which it has a negative connection. Eventually the system reaches a steady state in which the activity level of each node is constant, either because it has reached a maximum or minimum permitted value, or because its negative and positive

use excitation of all Person nodes connected to 40s and the inhibition of d 20s Age nodes. After several processing cycles the activity level of every e system is being influenced by a mixture of positive and negative inputs. it soon becomes impossible to keep track of the patterns of excitation and and predict whether the system will reach a stable state, and if so, what ited and what depressed. The only way to find out what the system will do e to stimulation of any of its nodes is to run a computer simulation of the em.

d now be clear that a connectionist memory system is totally unlike a nal memory such as a computer filing system. In a conventional system ent pieces of information are stored separately. When a specific piece of on is accessed, it and it alone is retrieved. But in a distributed connectionist n attempt to extract any information from the system leads to a flow of and inhibition throughout the system to everything which has any relation oration. This results in many different nodes becoming active. What is is the information which corresponds to the most active node(s) once this stabilised. If one is accustomed to information retrieval from a con- , non-connectionist system such as a telephone directory, this might seem

If you looked up Tom Brown's number in a connectionist telephone the number retrieved would be influenced by the entries of everyone with a me or a similar number to Tom's. This would not be useful. You do not number retrieved to be influenced by the fact that there happens to be called Tim Brown who has a telephone number quite unlike Tom's. You information stored at the location with the address 'Tom Brown' and else. But the interference which a connectionist system allows during between related items of stored information turns out to have some g and useful properties.

ntent-addressable memory in *Jets and Sharks*. To see whether this memory lows access by content, we can ask it the question 'Do you know the name er?'. To do this we activate the *Pusher* node, leave it on, and see which des becomes activated. Figure 2.3 shows the activity level of three of the des as a function of the number of processing cycles for which the system allowed to run. All the pushers names initially become activated. Most of e *Oliver*, quickly return to their resting level. But *Fred* and *Nick* both ncreasingly activated. After about 50 cycles *Fred* starts to dominate and system enters a stable state with *Fred* activated and all the other names back resting level. The system answers the question with the reply: 'Err... Fred.' So, e storage system of table 2.1, this system *does* allow information to be when it has been accessed by content rather than address.

lative activity of the name *Fred* compared to the name *Nick* over the last 50

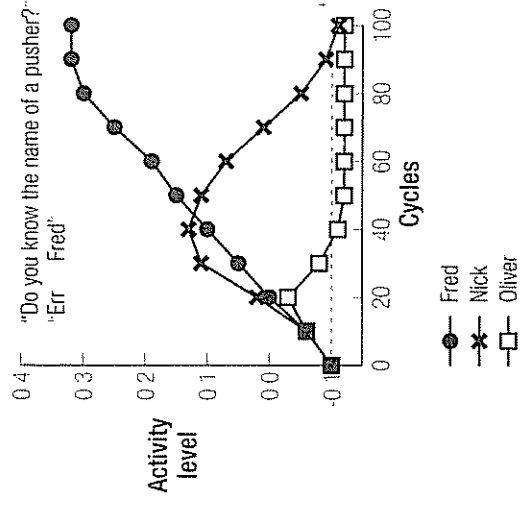


Figure 2.3 To see how the system responds to the question, the *Pusher* node in the Occasion area is held On and activity flows from through the system. The initial activity level of the nodes is set to -0.1. The activity level of the nodes in the Name area is plotted as a function of the number of cycles of activity passing around the system.

cycles demonstrates an important characteristic of models with mutual inhibition between competing responses. When alternatives are equally activated then each other equally and everything is balanced. But once one gets ahead it inhibits others more than they inhibit it. This reduces their activity and thus the system which they inhibit the one which is ahead. So it becomes more active and inhibits others yet more. This rapidly results in the one that is a little more active consolidating its position in the lead and completely inhibiting the alternatives. The effect is sometimes referred to as 'the rich get richer' or 'winner takes all'.

Building positive feedback into the system in this way makes it likely that the system will quickly come to a definite conclusion, even if the difference between evidence favouring one alternative rather than the other is small. But it is a decision process vulnerable to noise. A random disturbance may be magnified and treated as a signal. This would generally be considered a drawback in a making system but it has one useful consequence. Figure 2.3 suggests that the system would always answer the question 'Do you know a pusher?' with the reply 'Fred', so, it would be an indifferent model of human memory. People would give an answer to this question on different occasions, or if asked for an alternative could provide the names of other pushers. It is straightforward to add a response variability with the model. If random noise is added to the starting levels the system will produce a different answer. Positive feedback ensures that a node that gets ahead is likely to consolidate its advantage. So a small change in starting conditions, or during processing, can make a radical difference to the outcome. Figure 2.4 shows the result of setting the initial activity level of the *Pusher* node to -0.07 rather than -0.1 before activating the *Pusher* node. The system answers the question 'Do you know a pusher?' with the reply 'Err... Fred'.







al final state would be a set of activities for the individual units where all constraints were satisfied. The network would then be stable because no unit is trying to change the state of any of the units to which it was connected. A solution is unlikely to exist because most units are connected to some units trying to increase its activity and others which are trying to reduce it. A way of satisfying both. But if the system can find a state in which any the activity levels of the units reduces the overall number of satisfied constraints, it will stop changing activities. That is, it will have found a stable state. In networks there will be many possible stable states with a different pattern of activity, each one of which will be reached from a different pattern of input. The realisation that these stable states could be viewed as the network's response—that is, the set of possible states that it could reach in response to inputs—was an important step in the history of connectionism which will be discussed further in chapter 15.

In a conventional connectionist network where knowledge is distributed across many units, it is difficult to follow constraint satisfaction at work because it is not so obvious what role any particular unit is playing. In Jets and Sharks it is easy to see how concept coding is localist rather than distributed. Each node stands for a specific concept. The constraints on the system are the various facts in the network. The fact that a Jet in his 40s means that if either of the nodes representing these concepts is activated it will try to activate the other, and inhibit the *Shark*, *30s* and *20s* nodes. That there are also both Jets and Sharks in their 30s and 20s means that a contradiction exists. In other words, mutually contradictory constraints are influencing the way that the network changes the activity level of the units in response to any particular pattern of input. The key point is that the changes in activity on each cycle will increase the number of constraints which are satisfied.

The system which works by constraint satisfaction has a number of desirable characteristics for modelling human cognition. The main one is that it allows a solution to be reached by a consensus of evidence, a reasonable fit between input and output, rather than requiring an exact match. We have already seen this as a virtue of the model of human cognition because the nature of the nervous system requires a degree of fault tolerance in the information processing system (remember figure 2.1). It is desirable because of the nature of the input which the cognitive system has with it in the real world. Consider what happens when you listen to one speaker in a crowded room. The signals arriving at your ear contain the sounds made by the person you are listening to, but superimposed on these are a host of sounds from different speakers, the whole thing obliterated from time to time by bursts of laughter and other noises. And yet, most of the time, what you hear is words. The signal you receive bears some relationship to a prototypical

representation of the word that you perceive but will be far from an exact match. The fact that you *perceive* words shows that the word recognition system must be able to find a *best fit* to the word patterns it has stored rather than for an exact match. This effect has been studied in the laboratory with an experimental paradigm known as 'phoneme restoration'. In the original study by Warren (1970) the sound of a sentence containing the word 'legislature' and replaced with a cough. People then heard. People reported the sentence correctly, adding that there was a cough or after the word 'legislature'. In other words, the perceptual system necessarily gives a veridical account of the stimulus, it gives a plausible interpretation of the input, given its knowledge of English words.

The same effect can be seen in the Jets and Sharks system. Figure 2.5 shows what happens when the system is probed with a variety of retrieval cues. The fill shows what happens when the net is asked: 'Do you know a Shark, in his went to High School, who is single and a burglar?' (i.e. the *Shark*, *20s*, *HS*, and *Burglar* nodes are all switched On). The circles joined by solid lines the activity of Ken's Name node and the circles joined by dashed lines the activity of the next most activated Name node. Not surprisingly, the net answers 'Ken'.

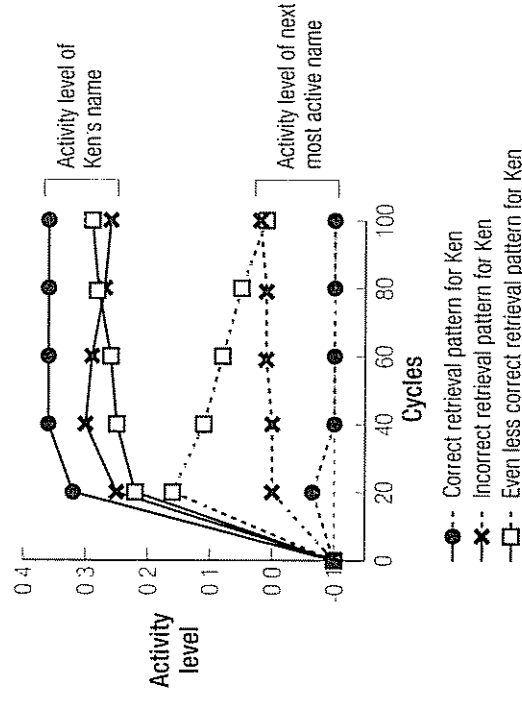


Figure 2.5 Constraint satisfaction in operation. Given a correct description of Ken as a retrieval cue, the system retrieves Ken's name. Given progressively less accurate retrieval cues, which are, none the less, closer to a description of him than to anyone else in the database, it still produces his name. In levels:

- (●) *Shark* = 1; *20s* = 1; *HS* = 1; *Single* = 1; *Burglar* = 1.
- (×) *Shark* = 1; *20s* = 1; *JH* = 1; *Single* = 1; *Burglar* = 1.
- (□) *Shark* = 1; *20s* = 1; *JH* = 1; *Burglar* = 0.5; *Pusher* = 0.25; *Bookie* = 0.25.

the equivalent of presenting a listener with a clear and unambiguous example word 'legislature' and asking her what word she heard. The crosses show what happens when *High School* is changed to *Junior High* in the input activity pattern. But there is no longer an accurate description of anyone in the database (the result of presenting a listener with '*legi<cough>laturre*'). A system which tried to find an exact match to the input would fail. There is no person who matches that pattern in the database. But the network has no problem. It takes slightly more time to respond (i.e. Ken's name takes more time to become activated, and does not reach a high level) but *Ken* still comes out as the clearly preferred item retrieved from memory. The squares show what happens with an even more ambiguous input. In the *Bookie* and *Pusher* nodes have all been activated as well as *Burglar*. As the input pattern is closer to a description of Ken than to any alternative, the network makes a clear decision in favour of the response 'Ken'.<sup>55</sup>

The fact that connectionist systems work by constraint satisfaction is the reason they exhibit fault tolerance. No part of the input uniquely determines the response. The network's response is the best fit it can make between the current input information it has acquired in the past. It would be possible to devise a system, in which an address + contents information storage like table 2.1, which, if given an input that failed to match any stored information, could compute a best fit. However, this would be a time consuming process once possible inputs achieved any degree of complexity. The parallel distributed computation in the Jets and Sharks networks automatically computes a best fit between input and stored information. It would continue to do so whatever the degree of complexity of the patterns describing individual entries without taking any more time.

## There is no distinction between 'memory' and 'processing' in connectionist models

A general point to note about distributed representations is that they blur the distinction between memory and processing. Traditional models of cognitive systems often distinguish between 'memory', a store of learnt information, and 'processing', operations which enable the system to interpret incoming information. Processing operations may use information from memory, but the conceptual distinction is clear. Indeed, in many models this is made explicit with separate parts of the model labelled 'memory' and 'processor'. Such models exploit the analogy to

the generation of information in the Jets and Sharks system produces some strange and unpredictable results which can only be appreciated by playing with the tac model. For example, figure 2.5 shows that if the network is long enough the *less* accurate description of Ken produces a stronger preference for his name than the more accurate description.

conventional digital computers where there are independent systems for memory information and for processing it.

There is no such distinction in connectionist models. All the information in the network has—its memory—is stored in the weights of the connections between nodes. All the processing that the net can do is determined by the same set of weights.

## Problems for distributed representations

We have emphasised the advantage of distributed representations over localist representations in allowing the system some degree of resistance to damage and tolerance of noisy inputs. Given the unreliable nature of the matter which uses for computation, it seems inevitable that it would use distributed representations. However, there are two properties of human memory which would be easy to account for with connectionist models but which would be difficult to account for with localist information storage: First, the addition of new information does not necessarily cause the loss of old. Second, learning can be immediate.

In a distributed system, any new information has to be added to the contents which already carry the system's current store of knowledge. To add new information the strength of connections must be changed. If this is done incrementally, addition of new information is likely to lead to some loss of old information. In a localist system, in contrast, the addition of new information is no problem. Simply adding to new storage locations and does not affect old information. At this point, however, we will just suggest that at an intuitive level there are differences between the two kinds of human learning.

At one extreme it seems clear that some sorts of knowledge can be immediately recalled, without interfering with other information. All young chess players show the smothered mate sequence with a queen sacrificed to a rook on g1 by mate of the king on h1 by a knight moving from h3 to f2. If you understand this, you only have to see this sequence once to remember it for ever, despite the fact that you may never have an opportunity to use it in a game. Similarly, if you understand anything about the British ex-Prime Minister Mrs Thatcher, and were told her nickname at school was 'Bossy Roberts', you would be unlikely to forget it. A piece of information which you find interesting or amusing, in a domain with which you already have sufficient knowledge to understand its significance, is likely to be remembered after a single presentation. And it can be retrieved as a specific item of information in future, independent of any other facts in the database. It is difficult to believe that such acquisition is accompanied by the loss of any other information. Quick, cost-free, addition of new information to existing databases characterises certain sorts of human knowledge acquisition. It is natural with localist representations of knowledge—you just add another entry to the database. But it is difficult

### 3 Pattern association

---

it can happen with a distributed system. (That connectionist models *can* do one trial learning will be shown in chapter 13 which demonstrates a model of the hippocampus in episodic memory formation.)

In other extreme there are many areas of knowledge acquisition, such as going to play tennis or learning to talk, where acquisition of new knowledge is gradual, and accompanied by the modification or loss of previous patterns. As your performance improves or you learn to pronounce the language correctly you *want* to forget some aspects of your old response patterns because they were inaccurate. Later it is difficult to recall when a specific piece of information was added to the database. In a pattern, where new information is inextricably interwoven with old, occurs a difficulty with a distributed system, but not with a localist one. Distinctions between different sorts of knowledge representation occur in many models of the cognitive system.

Most of the connectionist learning algorithms we will look at are more appropriate for modelling the latter sort of acquisition and representation than the

*This chapter will introduce the architecture and properties of one specific network, a pattern associator, and the operation of one particular learning rule, the Hebb rule. During training a pattern associator is presented with pairs of patterns. Learning is successful then the network will subsequently recall one of the patterns in the output when the other is presented at input. After training, a pattern associator will respond to novel inputs, generalising from its experience with similar patterns. Pattern associators are tolerant of noisy input and resistant to internal damage. They are capable of extracting the central tendency or prototype from a set of similar examples.*

The first two chapters introduced general principles which are shared by all connectionist networks. In this and the next four chapters we will look in detail at the structure and properties of a variety of specific network architectures: pattern associators, autoassociators, competitive nets and recurrent nets. Networks can be trained in a variety of ways. In this chapter we will look at one particular learning rule, the Hebb rule. In chapter 4 we will look at the Delta rule and in chapter 5 at backpropagation.

The examples and the networks in these early chapters have deliberately been very simple so that the principles involved in their operation can be seen at once. It may seem that the problems they solve are so trivial that they have little to do with human cognition. Don't worry. In chapters 8–14 we will look at networks that have been scaled up to the point where they can simulate realistic aspects of human behaviour. For example, in chapter 8 we will look at a model with roughly 1000 processing units and 200 000 connections which learns to pronounce monosyllabic words in the English language. In chapter 13 we will look at a model of episodic memory formation in the hippocampus with about 4000 processing units and about half a million connections. Exactly the same principles will be at work in these networks as in the examples which follow.

### The architecture and operation of a pattern associator

A fundamental task for the nervous system is to discover the structure of the world by finding what is correlated with what. That is, to learn to associate one