

NOTICE WARNING CONCERNING COPYRIGHT RESTRICTIONS:
The copyright law of the United States (title 17, U.S. Code) governs the making of photocopies or other reproductions of copyrighted material. Any copying of this document without permission of its author may be prohibited by law.

CMU-CS-86-172

**A Distributed Connectionist
Production System**

**David S. Touretzky
Geoffrey E. Hinton**

December 1986

**Computer Science Department
Carnegie Mellon University
Schenley Park
Pittsburgh, PA 15213**

Copyright (C) 1986 David S. Touretzky and Geoffrey E. Hinton

This research was sponsored by the National Science Foundation under grants IST-8516330 and IST-8520359, and by the System Development Foundation.

Abstract

DCPS is a connectionist production system interpreter that uses distributed representations. As a connectionist model it consists of many simple, richly interconnected neuron-like computing units that cooperate to solve problems in parallel. One motivation for constructing DCPS was to demonstrate that connectionist models are capable of representing and using explicit rules. A second motivation was to show how "coarse coding" or "distributed representations" can be used to construct a working memory that requires far fewer units than the number of different facts that can potentially be stored. The simulation we present is intended as a detailed demonstration of the feasibility of certain ideas and should not be viewed as a full implementation of production systems. Our current model only has a few of the many interesting emergent properties that we eventually hope to demonstrate: it is damage resistant, it performs matching and variable binding by massively parallel constraint satisfaction, and the capacity of its working memory is dependent on the similarity of the items being stored.

Table of Contents

1. Introduction
2. The Structure of Working Memory
 - 2.1. Receptive Fields
 - 2.2. Properties of Coarse Coding
3. Selective Attention: Clause Spaces
4. The Rules
 - 4.1. Rule Format
 - 4.2. Representation of Rules
5. Variable Binding
 - 5.1. Constraints on Rules
 - 5.2. The Structure of Bind Space
6. The Match Process
 - 6.1. Hopfield Networks
 - 6.2. Matching as Parallel Constraint Satisfaction
 - 6.3. Boltzmann Machines
 - 6.4. Matching by Simulated Annealing
 - 6.5. Detecting Failed Matches
7. Rule Firing
 - 7.1. Variable-Free Actions
 - 7.2. Actions Requiring Instantiated Variable Values
 - 7.3. Functions on Variable Values
8. Experimental Results
 - 8.1. Measured Performance
 - 8.2. Difficult Match Cases
9. Discussion
 - 9.1. Alternative implementations of working memory
 - 9.2. Multiple interacting distributed representations
 - 9.3. Similarity and generalization
 - 9.4. Seriality and variable binding
- Appendix A. Model Parameters
- Appendix B. Generating Receptive Fields for Working Memory Units.

1. Introduction

DCPS is a connectionist production system interpreter that uses distributed representations. As a connectionist model (Feldman & Ballard, 1982), it consists of many simple, richly interconnected neuron-like computing units that cooperate to solve problems in parallel. One motivation for constructing DCPS was to demonstrate that distributed connectionist models are capable of representing and using explicit rules. Earlier connectionist models (Rumelhart & McClelland, 1986) have shown that many phenomena which appear to require explicit rules can be handled by using connection strengths that implicitly capture the regularities of the task domain without ever making these regularities explicit. However, we do not believe that this removes the need for a more explicit representation of rules in tasks that more closely resemble serial, deliberate reasoning.

The natural way to implement explicit rules is to apply a parallel best-fit search to the task of finding the rule whose left-hand side best matches the current contents of working memory. Connectionist networks are good at performing pattern-matching, especially when there is no perfect match and the aim is to find the best partial match. One difficulty with this approach is that the kind of matching required to implement a production system is more complex than simple template matching. The left-hand side of a production may contain several instances of the same variable, and matches are only valid if all instances of the variable receive the same binding. Ensuring consistent variable bindings in a parallel network is a difficult and important problem (Barnden, 1984) and one of the main aims of this paper is to demonstrate a feasible solution.

Ballard and Hayes have demonstrated that a rather elaborate connectionist network can decide whether two expressions can be unified (Ballard & Hayes, 1984; Ballard, 1986). DCPS uses a different solution which is based on earlier work (Hinton, 1981a) on viewpoint-invariant shape-recognition. In matching an object-model to a retinal image, it is essential to ensure that all the matches of a piece of the model to a piece of the image assume the same viewpoint. In matching the LHS of a rule to the contents of working memory, it is essential to ensure that all the matches of a clause in the LHS to a fact in working memory assume the same variable bindings.

A second motivation for DCPS is to show how "coarse-coding" or "distributed representations" can be used to construct a working memory that requires far fewer units than the number of different facts that can potentially be stored. The price of this economy is that only a small fraction of the potential facts can actually be present in working memory at any one time. Earlier analyses of coarse-coding have shown that it is efficient (Hinton, 1981b; Hinton *et al.*, 1986) but they have failed to demonstrate that it can be used effectively when many different groups of units must interact correctly. Coarse-coding "smears" the representation of a given item across many units, and when coarse-coded representations in several different groups of units interact during an iterative best-fit

search, there is a danger that the representation will become progressively more smeared with each iteration.

The simulation we present is intended as a detailed demonstration of the feasibility of certain ideas and should not be viewed as a full implementation of production systems. The production rules our model interprets are much simpler than those found in OPS5 or EMYCIN. Nevertheless, they do contain variables that get bound consistently by the connectionist network, and they are implemented using distributed representations throughout. This falsifies any strong claim that connectionist systems using distributed representations could not possibly implement symbol processing. However, it leaves us open to the alternative criticism that we have merely implemented a very simple production system in a peculiarly inefficient way.

One advantage of the implementation we present is that it is robust against the destruction of any small random subset of the units or connections, but the real advantage (which we have not demonstrated in this simulation) comes from the ability of a connectionist network to do a rapid best-fit match. This is potentially much more powerful than the standard implementations which find all exact matches and then do conflict resolution. In situations where no existing rule fits perfectly, it may be sensible to apply a plausible rule, particularly in a learning system that needs to explore the space of plausible actions in order to find a satisfactory one. The ability of a connectionist implementation to settle on plausible but imperfect matches could therefore be very helpful, but only if the matching apparatus is able to do more than simple, variable-free "template" matches. Our eventual aim is to exploit the best-fit ability of DCPS to allow it to do more of the computation in each match so that it can perform complex tasks with fewer rule-firings, and rules in one domain can be created by analogy with rules in other domains. But before we can do this we must demonstrate that it is possible to build a workable system that uses distributed representations and enforces consistent variable bindings during a match. So our current model only has a few of the interesting emergent properties that we eventually hope to demonstrate: it is damage resistant, and the capacity of its working memory is dependent on the similarity of the items being stored.

2. The Structure of Working Memory

The working memory elements of DCPS are triples of symbols, such as (F A B). We have chosen an alphabet size of 25 symbols, giving 25^3 or 15,625 possible triples. Only a few of these are present in working memory at any one time; typically there will be half a dozen. The sparseness of working memory is an important consideration in the design of the model.

The most straightforward representation for a set of triples, in a conventional architecture, would be a purely "localist" one, where every triple was represented by a dedicated unit. A unit in the active

state would then indicate that the corresponding triple was present. We have rejected this idea in favor of a distributed or "coarse-coded" representation (Hinton, 1981b; Hinton *et al.*, 1986). Localist representations require too many units and too many connections; they quickly succumb to combinatorial explosion as the alphabet size or the length of a sequence increases. This is because localist representations do not make efficient use of the units when the number of items that are simultaneously present in working memory is much less than the number of possible items. Distributed representations use the information-bearing capacity of the units more efficiently by making them active much more often.¹ In addition to the inefficiency of localist representations we think that a one-to-one mapping between individual neurons and symbolic structures is physiologically implausible; it is reminiscent of the grandmother cell idea. Recordings in the temporal lobe of the macaque cortex support the idea that neurons are tuned to very complex entities such as a face (Rolls, 1984) but they do not support the idea that a particular face is encoded by just one or just a few neurons. Each particular face is almost certainly encoded as a pattern of activity distributed over quite a large number of units, each of which responds to a subset of the possible faces. Using a distributed representation not only makes our model more efficient and neurally plausible, it also makes it tolerant of noise and occasional malfunctions.

2.1. Receptive Fields

The working memory space of DCPS, shown in figure 1, consists of 2000 binary state units. Each unit has a receptive field table such as the one in figure 2. A unit's receptive field is defined to be the crossproduct of the six symbols in each of the three columns, giving 6^3 or 216 triples per field. The unit described in figure 2 has the triples (C K R) and (F A B) in its receptive field, along with 214 others. Receptive field tables are generated randomly prior to beginning the simulation; they determine the connection pattern between units in the various spaces comprising DCPS. Once the connections have been built and the working memory units' states have been initialized, the tables are no longer needed; they are not consulted when running the model.

A triple may be stored in working memory by turning on all its receptors. With 2000 working memory units, triples will average $6^3/25^3 \times 2000$ or roughly 28 receptors. The number varies slightly from one triple to the next due to the random distribution of receptive fields. An external observer can test whether a particular triple is present in working memory by checking the percentage of active receptors for it. If this is close to 100%, the triple may be assumed to be present. For example, if the

¹If there are 15,625 possible items, but only 6 of these are present at any one time, the probability that a working memory unit is active in a localist scheme is only about 0.0004. The average information conveyed by the unit is therefore the entropy of the distribution {0.0004, 0.9996} which is about 0.005 bits. In DCPS, fewer units are used to encode the same information, and each unit is active much more often so it conveys much more information. The probability of an individual unit being active is about 0.08 and so the average information it conveys is about 0.4 bits. However, in DCPS the correlation between units cannot be ignored (as it can in the previous case) and so the average information conveyed per unit is actually only about 0.04 bits.

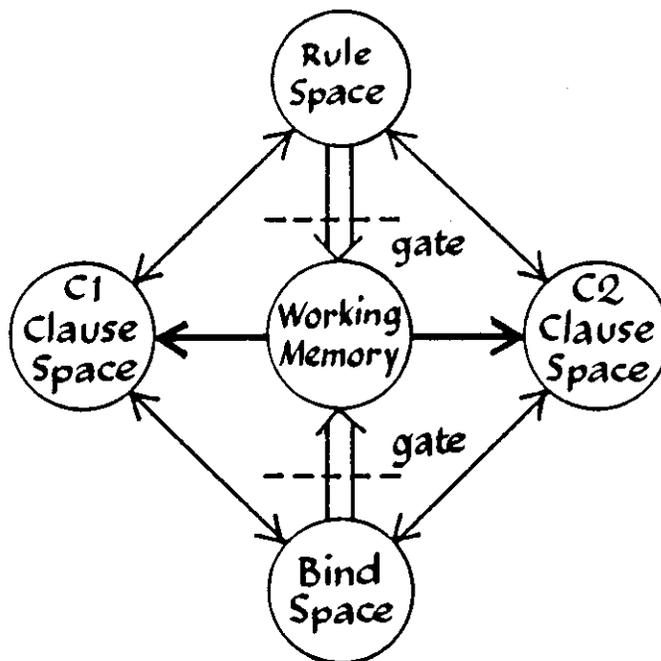


Figure 1: Block diagram of DCPS, a Distributed Connectionist Production System.

triple (F A A) were stored in working memory, the unit described in figure 2 would be active, along with about 27 other units. Although (C K R) also falls within the receptive field of this unit, the number of receptors two unrelated triples have in common is small; on average, it is less than one. Thus, while 100% of the (F A B) units become active when (F A B) is stored, only 1 out of roughly 28 (C K R) units would become active. To the external observer, (F A B) clearly is present in working memory and (C K R) clearly is not. But the network itself doesn't need to compute these percentages. It relies on the fact that triples that are present have strong effects and triples that are absent do not.

C	A	B
F	E	D
M	H	J
Q	K	M
S	T	P
W	Y	R

Figure 2: An example of a randomly generated receptive field table for a working memory unit. The receptive field of the unit is defined as the crossproduct of the symbols in the three columns.

Figure 3 shows the state of working memory when the two triples (F A B) and (F C D) have been stored. The 2000 working memory units are arranged in a 40x50 array, with the 55 that are active indicated by black squares. The positions of these 55 units in the array are not significant, since units' receptive fields are generated randomly. However, if we were to examine the receptive fields of each of the active units we would see that every one contains either (F A B) or (F C D), or both.

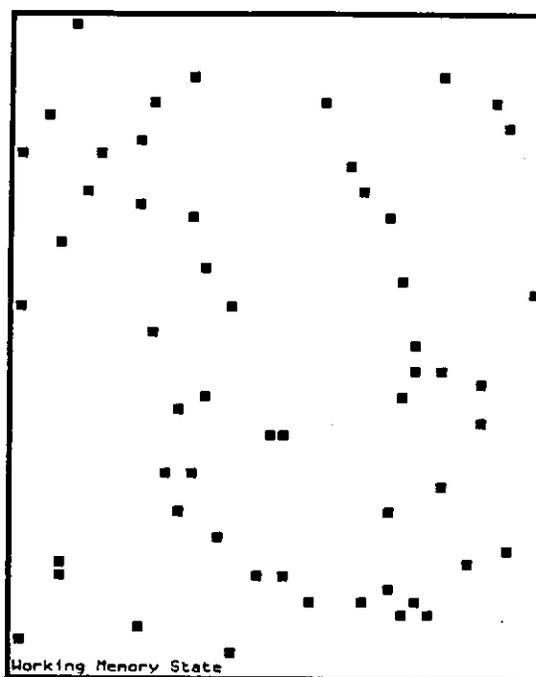


Figure 3: The state of working memory after the triples (F A B) and (F C D) have been stored. Active working memory units are indicated by black squares. 55 of the 2000 units are active.

Table 1 shows the first dozen triples with the strongest representations when working memory is in the state shown in figure 3. (F A B) and (F C D) each have 100% of their receptors active, while the next best represented triple, (F N B), has only 42% active. The average activity level over all 15,625 triples is much lower: only 2.7%. If we adopt the criterion that 75% of a triple's receptors must be active for it to be deemed present in memory, the division between present and absent triples in Table 1 is quite clear.

Figure 4 shows the levels of support for all 15,625 possible triples when working memory contains (F A B) and (F C D). In the figure, (A A A) is located in the upper left corner and (Y Y Y) in the lower right. The blobs in this figure are associated with triples, not units; the size of each blob indicates how many receptors are active for that triple. A simple thresholding operation yields figure 5, in which the (F A B) and (F C D) blobs stand out clearly and there is only a small amount of noise remaining.

Triple	Percent Active	Active Receptors	Total Receptors
(F A B)	100%	28	/ 28
(F C D)	100%	28	/ 28
(F A D)	40%	11	/ 27
(F B D)	38%	10	/ 26
(F A X)	37%	11	/ 29
(S A B)	37%	10	/ 27
(F Q D)	37%	10	/ 27
(F C N)	37%	10	/ 27
(F C B)	37%	10	/ 27
(F C M)	35%	10	/ 28
(F T D)	35%	10	/ 28
(N C D)	34%	10	/ 29

Table 1: The first dozen triples with the strongest representations when working memory is in the state shown in figure 3.

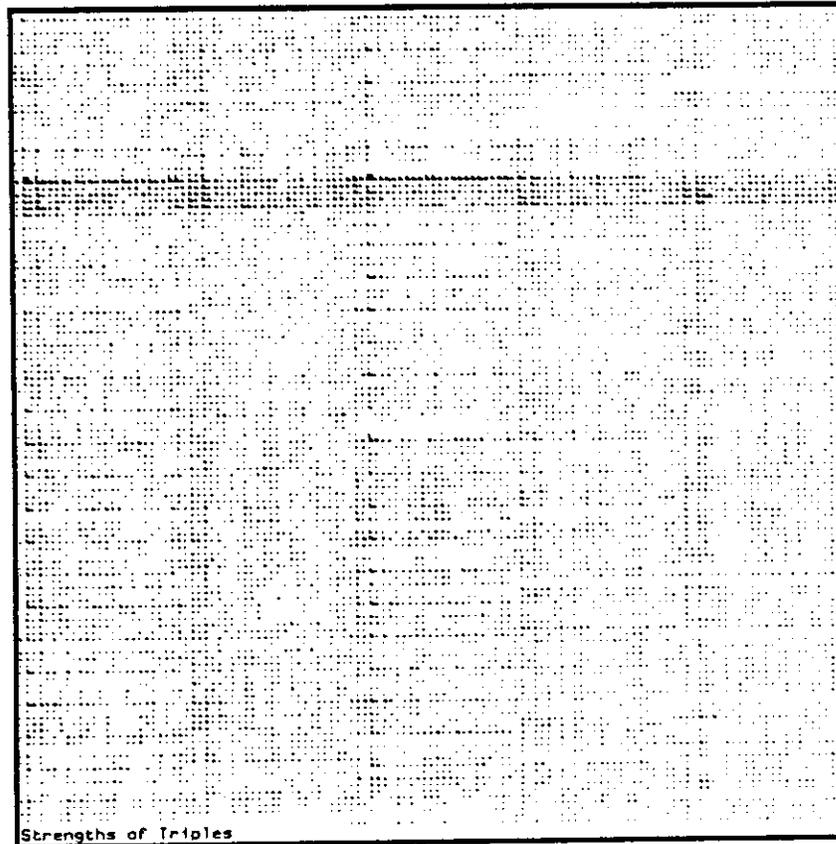


Figure 4: The levels of support for all 15,625 possible triples when working memory contains (F A B) and (F C D), represented by the 55 active receptors in figure 3. The size of each blob indicates the number of active receptors for that triple.

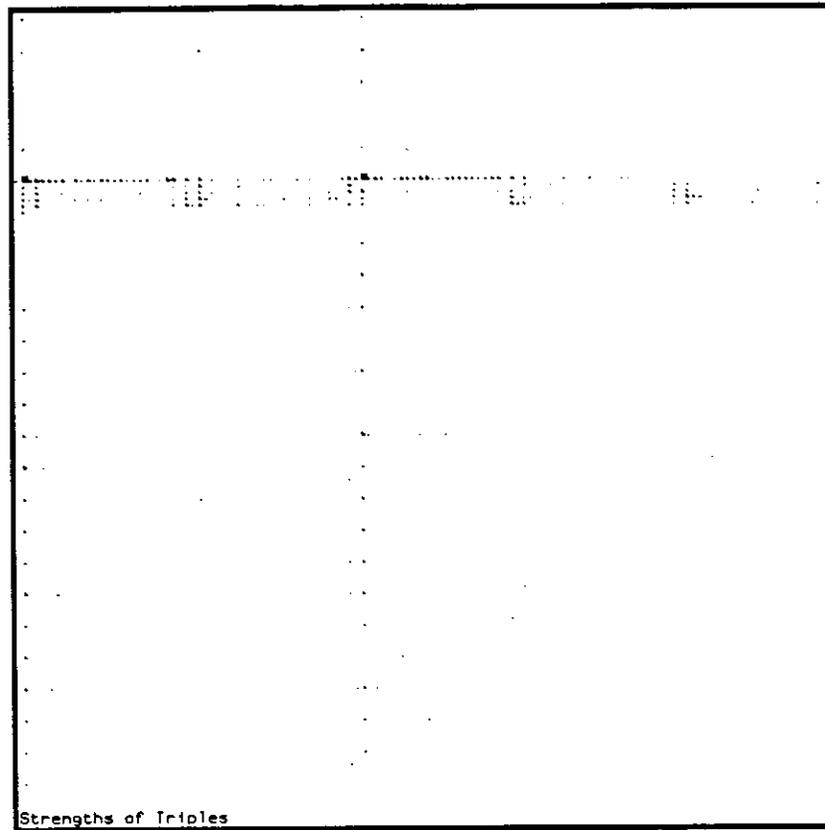


Figure 5: A moderately thresholded version of figure 4. The (F A B) and (F C D) blobs stand out clearly here.

2.2. Properties of Coarse Coding

Coarse coded representations have a number of interesting and useful properties. One of these is tolerance of noise. If after storing some triples in working memory a few units are flipped on or off at random, the perceived contents of working memory will not be affected at all.² Tolerance of noise is especially important when items will be deleted from the memory as well as added to it. A slight overlap in the receptor set of related triples causes deletion of a triple to affect any related ones previously stored. That is, if (F A B) and (F C D) were stored in memory and (F C D) were then deleted by turning off all its receptors, it is likely that only 27 of the 28 (F A B) receptors would remain active, the 28th having been shared between the two.

The contents of working memory remain reasonably persistent because the overlap between any two triples is small. A visual effect resulting from this overlap can be seen in figure 4. The dot pattern may appear completely random at first, but closer examination will reveal a regular series of thin

²Assuming, of course, that we do not require strictly 100% of a triple's receptors to be active for it to be considered present.

horizontal and vertical bands. These bands are formed by triples that have 2 out of 3 components in common with the stored triples (F A B) and (F C D); on average such triples have 7 of their receptors active, while triples with no components in common, such as (G K Q), average about 0.4 receptors active. Another effect that can be seen in the figure is the horizontal F band that is thicker and also somewhat darker than the other bands. Since both of the stored triples begin with F, all other triples beginning with F have a slightly higher number of active receptors.

Another interesting property of the coarse coded representation is that the memory has no fixed capacity; instead its ability to distinguish stored items from other items decreases gradually as the number of stored items increases. Each triple added to working memory raises the number of active units, thereby increasing the support for other triples that have not been stored. As working memory fills up, the fraction of active receptors for certain triples that are "close" to those that have been stored approaches 100%, and the dividing line between present and absent triples blurs. If many closely related triples are stored, such as (F A A), (F A B), (F A C), (F A D), etc., then the system may exhibit local blurring, where it can't tell whether (F A P) is present or not but it is certain that (G S Q) is absent. Figure 6 illustrates the local blurring that occurs when four closely related triples are stored.

Finally, triples stored early on in a coarse coded memory eventually fade away if production rules delete a large number of other triples. This gradual decay phenomenon is again an effect of the overlap of receptive fields. One way to counteract the decay effect is to recall a triple before it has completely faded away, and then store it again. Whenever a triple is stored all its receptors become active, so its representation in working memory is refreshed.

3. Selective Attention: Clause Spaces

Clause spaces, labeled C1 and C2 in figure 1, are a device for focusing the network's attention on particular triples from the set stored in working memory. Michael Mozer of UCSD independently invented a device similar to clause spaces, which he calls "pullout networks," that allow a perceptual system to attend to specific objects in a scene (Mozer, 1984). The matching problem in DCPS consists of selecting two triples in working memory (which may contain half a dozen or more) that together satisfy the left hand side of some production rule. Each clause space is responsible for pulling out one of these triples.

There is a one-one excitatory mapping between working memory units and units in C1 and C2 spaces, so that each working memory unit that is active tries to turn on its corresponding C1 and C2 units. What prevents the C1 and C2 spaces from exactly copying the activity pattern in working memory is the fact that clause units are mutually inhibitory within their space, i.e., each of the 2000 C1

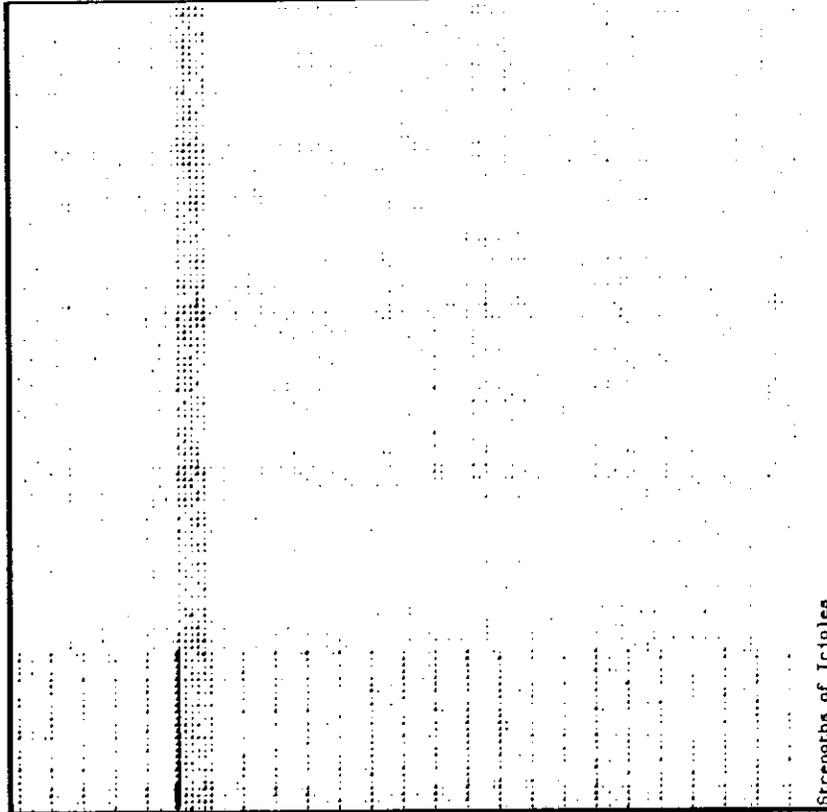


Figure 6: An illustration of the local blurring that occurs when several related triples are stored. Here, (F A A), (F A B), (F A C), and (F A D) have been stored. As a result, similar triples receive a high degree of support, as shown by the dark (F A x) line at the beginning of the F band and the weaker (x A y) lines in other bands. Moderate thresholding was applied.

units inhibits the other 1999 units, and similarly for C2; working memory units do not inhibit each other. See figure 7. The inhibition level in clause space is carefully adjusted so that only about 28 units per space can remain active simultaneously, i.e., just enough to represent a single coarse coded triple. Exactly which triple is selected depends on various outside influences imposed on the clause space by units in the Rule and Bind spaces. Briefly, a clause unit will be able to remain active despite inhibition from its siblings only if it receives support from rule and bind units that are also active.

The apparent requirement that a clause space have $(N^2-N)/2$ bidirectional inhibitory connections might seem a flaw in the design, since as the number of units grows the number of connections quickly becomes unreasonable. With 2000 clause units there would have to be 1,999,000 connections. But these connections need not actually be built. The inhibition function can be accomplished more economically by $2N$ unidirectional connections: N excitatory connections from clause units to a special regulatory unit with a graded or integer-valued rather than binary response,³

³These regulatory units resemble inhibitory inter-neurons which probably play a similar role in cortex.

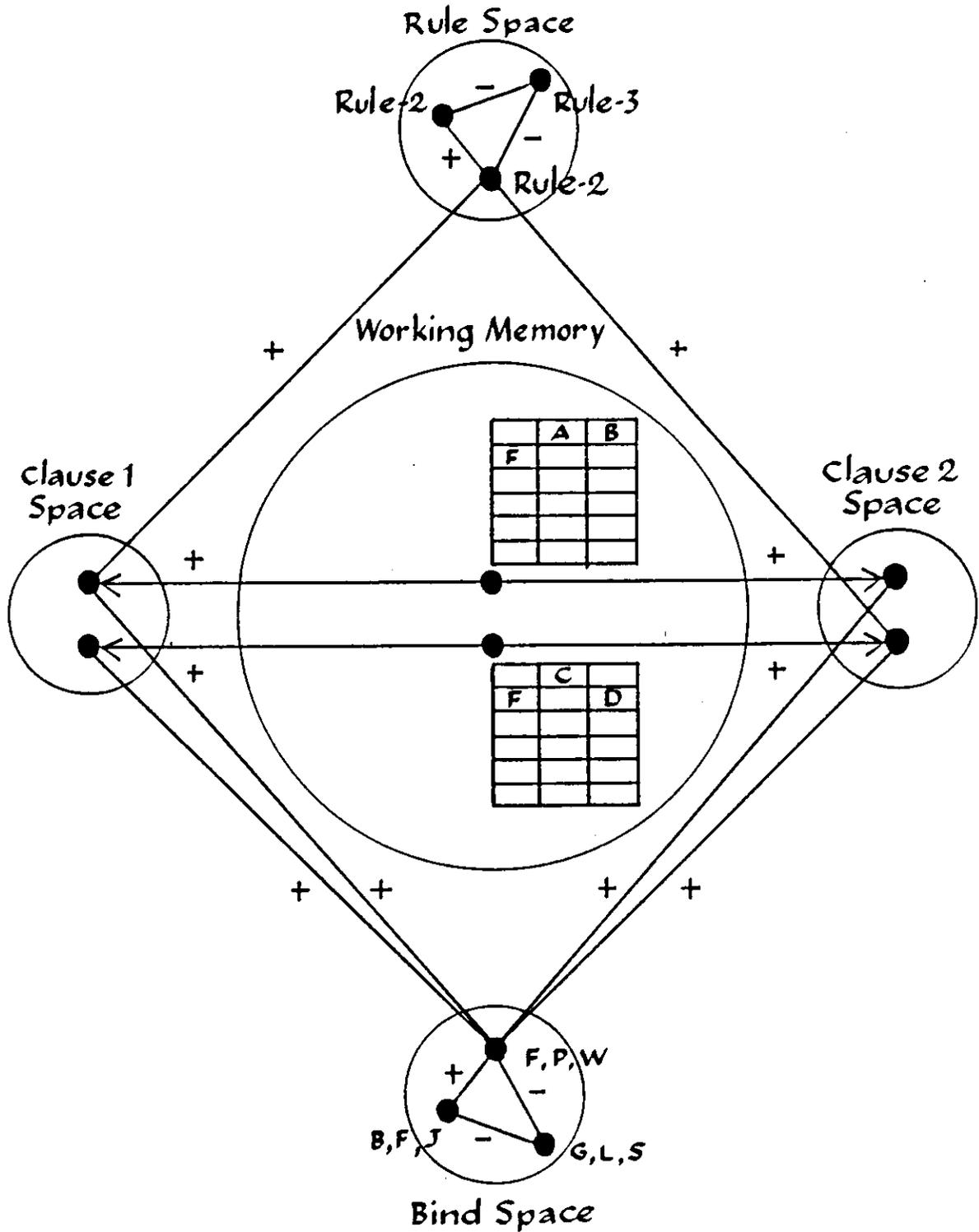


Figure 7: Connection pattern between clause units and working memory, rule, and bind units.

plus N inhibitory connections in the opposite direction. To exactly mimic the effect of $N(N-1)$ pairwise connections we would also need one excitatory connection from each unit to itself to cancel out the inhibitory effect it has on itself via the regulatory unit, giving $3N$ total connections. However, in practice these recurrent connections may be omitted with negligible effect.

For analysis purposes we will treat DCPS as an instance of a Hopfield network, and later, a Boltzmann machine. In order to meet those definitions we will ignore the regulatory unit solution and adopt the pretense, for the remainder of this article, that $(N^2-N)/2$ bidirectional inhibitory connections are actually built where required.

Note that although clause spaces are constrained to have roughly 28 units active at a time, not all patterns of 28 active units correspond to a valid triple. Clause spaces can sometimes be in an intermediate state where there are, say, 15 receptors for $(F A B)$ active, 10 for $(G K Q)$, and 5 for something else. In other words, the clause-space units can divide their attention among several partially represented triples simultaneously. At higher temperatures (more relaxed constraints), more than 28 units can be active, which increases the chance that multiple triples will be partially represented. There is nothing analogous to this in conventional computers, where symbol structures remain discrete and must be considered one at a time (Derthick & Plaut, 1986).

4. The Rules

4.1. Rule Format

Production rules in DCPS consist of two left hand side clauses that specify triples and any number of right hand side actions that modify working memory by adding or deleting triples. We first consider rules without variables. A typical rule would be:

Rule-1: $(F A B) (F C D) \rightarrow +(G A B) +(P D Q) -(F C D)$

This rule can fire if $(F A B)$ and $(F C D)$ are both present in working memory. If it does fire, the triples $(G A B)$ and $(P D Q)$ will be added to memory and $(F C D)$ will be deleted.

4.2. Representation of Rules

Each rule is represented by a population of 40 Rule units; the pattern of connections between these units and the clause units is determined by the left hand side of the rule. For example, Rule units that represent Rule-1 above will have bidirectional excitatory connections to C1 units whose receptive field includes $(F A B)$ and C2 units whose receptive field includes $(F C D)$, as shown in figure 7. If a sufficiently large number of these C1 and C2 units become active, indicating that the triples $(F A B)$ and $(F C D)$ are present in working memory, the rule unit will also become active.

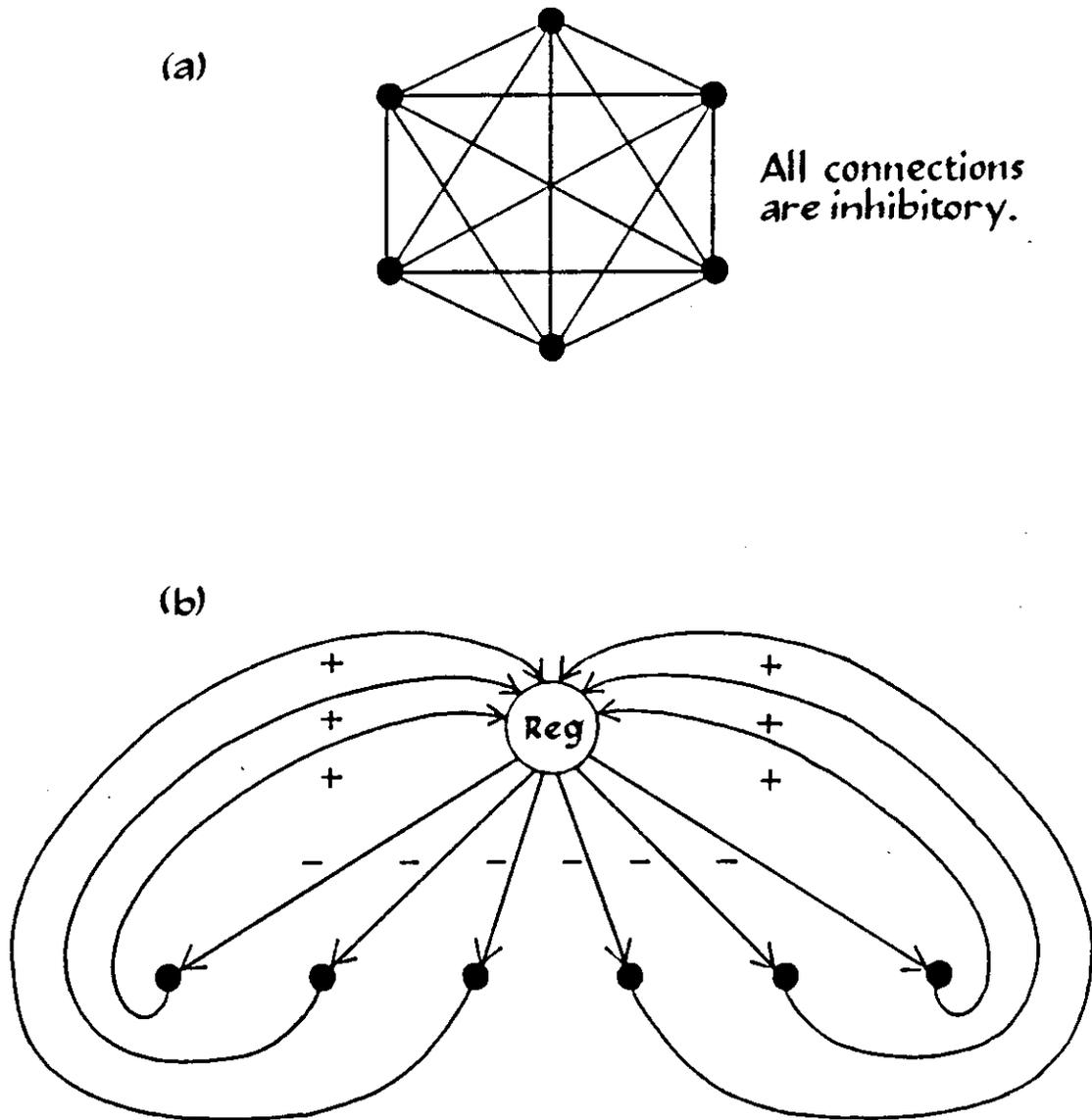


Figure 8: (a) Space of binary state units whose activity is limited by mutual inhibition using $(N^2 - N)/2$ bidirectional connections. (b) Introduction of a regulatory unit with graded response accomplishes the same effect with only $2N$ one-way connections.

Conversely, since the connections are bidirectional, when a Rule-1 unit becomes active it provides support for units in C1 and C2 space that support that rule.

The 40 units representing one production rule are connected so as to form a clique. Each active unit provides a slight excitatory stimulus to the other units in its clique and a slight inhibitory stimulus to units in all the other cliques. Thus, Rule space is organized as a "winner take all" network (Feldman and Ballard, 1982); when the network settles, all the units in one clique will be active and all the remaining units will be inactive. This is how the system decides which rule to fire.

There are several reasons for implementing rules as collections of units rather than as individual units. First, it is damage resistant. Second, it allows binary units to give a graded response.⁴ If, during the settling phase, there is a weak match between one rule and working memory, this will be indicated by only some of the corresponding rule units being active. If another rule matches more strongly, more of the units in its clique will be active, and they will eventually overpower the units in the other cliques. The implementation of rules in DCPS is "semi-distributed:" rules are represented by the collective activity of a set of units, but each unit codes for only one rule.

A further reason for implementing rules with multiple units is that it frees any one unit from having to represent the entire pattern associated with a rule's left hand side. Each rule unit is connected to a random subset of all the clause units associated with the rule's left hand side; only the clique as a whole has a complete representation for the rule. This is a more plausible organization than one in which rules are represented by single units, since it allows us to limit the connectivity of rule units without limiting the complexity of rules.

As in the case of clause spaces, the problem of building $O(N^2)$ connections among rule units can be solved by the use of regulatory cells with graded outputs and a combination of one-way and bidirectional connections, as shown in figure 9. Each rule unit excites its clique's "pro" regulatory unit which in turn excites all its siblings in the clique; the unit also receives inhibition from its clique's "con" regulatory unit. The regulatory units of the various cliques are in turn connected to a master regulatory unit that controls the entire rule space. Each clique's pro unit also has an inhibitory connection to the corresponding con unit, to counterbalance the tendency for a clique to inhibit itself via the master regulatory unit. As in figure 8, the recurrent connections from rule units to themselves, which are needed for absolute equivalence to the original network, have been omitted.

⁴One could implement rules as individual units with continuous rather than binary outputs, but the resulting network would not be a Hopfield net or Boltzmann machine. The fact that our hypothesized regulatory units have graded (either continuous or integer-valued) activation levels can be ignored because those units are merely used to simulate an equivalent Hopfield net composed solely of binary state units, with $O(N^2)$ rather than $O(N)$ connections.

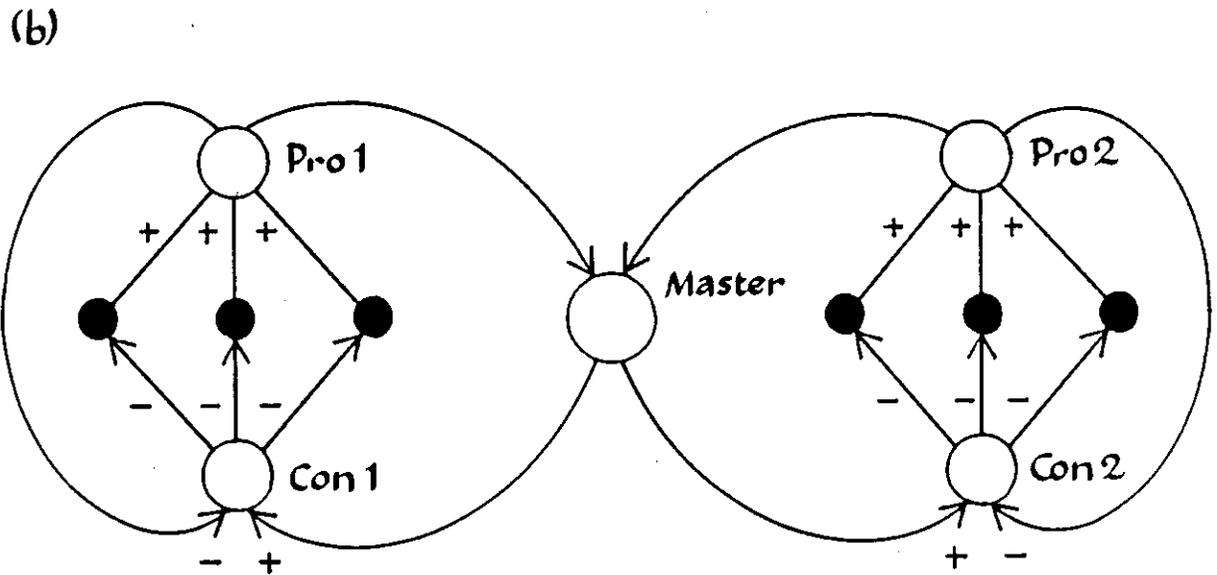
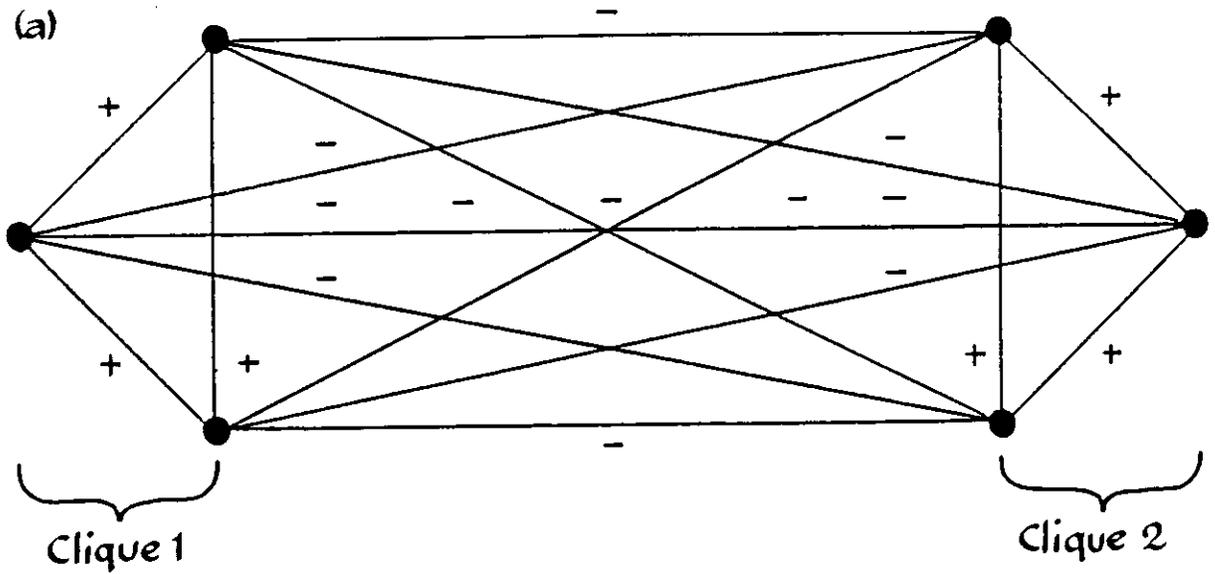


Figure 9: (a) A winner take all network composed of two cliques with three units each requires $(N^2 - N)/2$ connections. (b) Use of regulatory units with graded response produces the same effect with only $2N + 3C$ connections, where C is the number of cliques.

5. Variable Binding

5.1. Constraints on Rules

The first version of DCPS, called DCPS1, did not allow rules to contain variables. In developing DCPS2, which allows a limited form of variable binding, there were three distinct binding problems to consider:

1. Left hand sides in which variables impose intra-clause constraints, e.g., the clause $(=x R =x)$ can only match triples such as $(F R F)$ or $(G R G)$.
2. Left hand sides in which variables impose inter-clause constraints. The pair of clauses $(=x A B)$ and $(=x C D)$ can match pairs of triples such as $(F A B)$ and $(F C D)$ or $(G A B)$ and $(G C D)$, but not $(F A B)$ and $(G C D)$.
3. Right hand side actions in which variables appear. Variable binding requires a memory so that the variable's value can be instantiated into right hand side actions when the rule fires.

Each of these problems requires a different type of wiring pattern. Intra-clause constraints are the least interesting, and so they were not included. DCPS2 does allow a limited form of inter-clause constraint: each rule must have a variable in the first position of *both* left hand side clauses.⁵ DCPS2 also permits unrestricted use of variables on the right hand side. A typical DCPS2 rule is:

Rule-2: $(=x A B) (=x C D) \rightarrow +(G =x P) -(=x R =x)$

If this rule fires by matching $(F A B)$ and $(F C D)$, so that $=x$ is bound to F , its right hand side will add $(G F P)$ to working memory and delete $(F R F)$.

5.2. The Structure of Bind Space

Variable binding, which refers both to the imposition of constraints on rule matching and the instantiation of bound variables, is handled by the fifth space of units in figure 1, the Bind space.⁶ The units in this space form a winner take all network with 25 cliques, one for each of the 25 symbols of the alphabet. The space is coarse coded, so that each unit belongs to three cliques (votes for three distinct symbols) rather than one. Since Bind space contains a total of 333 units, each symbol falls in the receptive field of $(3/25) \times 333$ or 40 bind units, except for Y which has only 39.

Each bind unit has a set of bidirectional excitatory connections to units in C1 and C2 space whose receptive field table contains one or more of the letters the bind unit votes for. An $F/J/W$ bind unit,

⁵This choice was arbitrary; we could have chosen to require that the variable appear, say, in the first position of clause 1 and the third position of clause 2. The important constraint is that the variable be in the same position in all rules.

⁶These bind units are similar to the mapping units used for object recognition by Hinton (1981a).

for example, connects to a randomly chosen set of 240 C1 units: 80 that are receptors for triples beginning with F, 80 for J, and 80 for W. The same bind unit would also connect to a similar but independently chosen set of 240 C2 units. If a C1 unit that is a receptor for (F A B) and is connected to this bind unit becomes active, it will excite the bind unit, which in turn will excite other C1 and C2 units that code for triples beginning with F, J, or W. With many units in the F bind clique active, C2 space is more likely to adopt an activity pattern representing a triple beginning with F. The global effect of bind space is that it forces the C1 and C2 spaces to select triples beginning with the same symbol; that is how the "variable binding constraint" is imposed.

The inhibitory connections between cliques in bind space prevent the number of active bind units from growing much above 40, which is just enough to activate all the units that vote for a particular symbol as the value of the bound variable. The stable states of this network (considered in isolation) each consist of one active clique of 40 units, with the remaining units inactive. But because each unit is a member of three cliques, in a stable state the winning symbol receives 40 votes whereas the 24 remaining symbols receive 3 to 4 votes each.⁷ Even when Bind space has settled on a value for the variable, it is still giving some slight consideration to other values. This consequence of the coarse coded representation may help the network avoid getting trapped in local minima when searching for a globally optimal rule match, though this issue needs further research.

6. The Match Process

So far we have described a network consisting of five spaces of units: working memory, C1, C2, rule, and bind. Working memory units are essentially latches; they do not perform computation, but their activity pattern drives the rule match process. C1, C2, rule, and bind units are wired up in complex but principled ways. Ignoring the possible use of regulatory units, all units have binary states, and all connections between units are bidirectionally symmetric. The important questions to ask at this point are:

1. What are the stable states of such networks?
2. Under what conditions will a network eventually settle into one of its stable states?
3. Do stable states bear any relation to valid rule matches?

The first two questions have already been answered by Hopfield (1982); we will try to present a convincing argument for the third.

⁷Each symbol is voted for by 40 units, and each unit votes for 3 symbols, so in a stable state there are 120 votes to be had. Since 40 go to the winner, the losers average $(120-40)/(25-1) = 3.333$ votes apiece.

6.1. Hopfield Networks

A Hopfield network is a neural network composed of binary threshold units, all of whose connections are symmetric. Hopfield proved that if units change state asynchronously and there are no transmission delays across connections, the network's stable states are those states α that minimize a certain energy measure $E(\alpha)$. Let w_{ij} denote the weight of the connection between the i th and j th units; let θ_i denote the threshold of the i th unit; and let s_i^α denote the state (0 or 1) of the i th unit when the network as a whole is in state α . Then the energy of a state is the sum of the active units' thresholds minus the sum of the weights of connections between pairs of active units:

$$E(\alpha) = \sum_i s_i^\alpha \theta_i - \sum_{i < j} s_i^\alpha s_j^\alpha w_{ij}$$

This energy measure derives from an analogy Hopfield draws with spin glasses in physics, which operate under the same sorts of constraints as the neural networks he was studying. The stable states of these networks are called *local energy minima* because energy cannot be lowered any further by an individual unit's flipping state. Hopfield showed that networks that meet his constraints will settle into an energy minimum from any starting state because each state change either leaves the energy unchanged or reduces it; thus the energy decreases monotonically as the network moves from its initial state to a stable state. In general, however, the particular minimum energy state the network will end up in cannot be predicted from the starting state, and there is no guarantee that it will be a global minimum⁸ rather than a local one.

6.2. Matching as Parallel Constraint Satisfaction

The argument that a valid rule match corresponds to a minimum energy state, in fact, to a global energy minimum, is based on reformulating the match as a constraint satisfaction problem. Weighted connections between units cause them to impose constraints on each other and the energy of a state is a measure of how much it violates the constraints. So a minimum energy state is one in which as many constraints are satisfied as possible. The following sorts of constraints are present:

- Due to their high thresholds, clause units cannot become active unless their corresponding working memory units are active.
- Due to mutual inhibition, only about 28 clause units can be active simultaneously in each space, which is just enough to represent one triple.
- Rule and bind units influence the clause units. A triple can remain active in C1 or C2 space only if it is supported by a population of rule and bind units, i.e., it must match some rule's left hand side and contain the symbol voted for by the active bind clique in its first position.

⁸A global minimum is a state whose energy is less than or equal to the energy of all other states the network could be in.

- Active clause units excite the rule and the bind units with which they are compatible. For example, C1 units whose receptive field includes (F A B) will try to turn on any rule units whose first clause is (=x A B), and any bind units that support the variable value F.
- Rule space is organized as a winner take all network. Rule units excite others that vote for the same rule and inhibit those that vote for different rules.
- Bind units form a coarse coded winner take all network. They excite other units that vote for the same symbol (or symbols, if they have more than one in common), and inhibit units that vote for different symbols.

Considered individually, the C1, C2, rule, and bind spaces have many equivalent stable states. For instance, if bind space wasn't connected to clause spaces that are influenced by working memory, its 25 stable states would be completely equivalent. Rule space has as many stable states as there are rules; if rule space wasn't connected to the clause spaces then its stable states would also be equivalent. But considered together, the various spaces interact with each other so that the only way all their constraints can be satisfied — thus putting the network into a global energy minimum — is for the C1 and C2 spaces to settle into representations of triples that are in fact present in working memory and match one of the rules, while rule space settles into a state where that particular rule is the winner, and bind space settles into a state where the active symbol is the one that appears in the first position of both the triples in C1 and C2 spaces.

Constraint satisfaction in a Hopfield net is not a foolproof match technique because it is possible for the network to get stuck in a local energy minimum that does not represent a valid match. This occurs when a winner-take-all space, either rule or bind, settles so deeply into an undesirable stable state (all the units of one incorrect clique on, the remaining units off) that the other spaces cannot dislodge it.

In practice, the Hopfield net version of DCPS had no trouble finding the global energy minimum when the answer to the match problem was clear. However, in more difficult cases where there were many elements in working memory, many similar rules, or many partial matches possible but only one correct one, the network would often get stuck in a local minimum. In order to improve the chances of settling into the global minimum, DCPS was converted to a Boltzmann machine.

6.3. Boltzmann Machines

A Boltzmann machine (Fahlman, Hinton & Sejnowski, 1983; Ackley, Hinton & Sejnowski, 1985) is a Hopfield network whose units behave stochastically as a function of their *energy gap*. A unit's energy gap is the amount by which its activation exceeds its threshold. The energy gap of the *i*th unit when the network as a whole is in state α , written $\Delta E_i(\alpha)$, is defined as:

$$\Delta E_i(\alpha) = \left(\sum_j s_j^\alpha w_{ij} \right) - \theta_i$$

While the deterministic units of a Hopfield network turn on whenever their energy gap is positive, i.e., whenever their input exceeds their threshold, in a Boltzmann machine a unit's energy gap determines only the probability that it will turn on, in accordance with the Boltzmann distribution. Let $p_i(\alpha)$ denote the probability that the i th unit is on when the network as a whole is in state α . This probability is given by the formula

$$p_i(\alpha) = \frac{1}{1 + e^{-\Delta E_i(\alpha)/T}}$$

The parameter T in the above equation is called the *temperature*. At very high temperatures units behave almost randomly, i.e., the probability that a unit will turn on is approximately 0.5. (It is slightly above 0.5 for units with large positive energy gaps, slightly below 0.5 for units with large negative energy gaps.) On the other hand, when the temperature is close to zero the behavior of the units is almost deterministic, i.e., the Boltzmann machine acts like a Hopfield net (see figure 10). At moderate temperatures units tend to turn on when their energy gaps are positive, but they have a small probability of turning on even if their energy gap is negative, and a small probability of turning off even if their energy gap is positive. So at moderate temperatures a Boltzmann machine will occasionally move uphill in energy space, although the trend is still to move downhill. The higher the temperature the more likely an uphill move will be made.

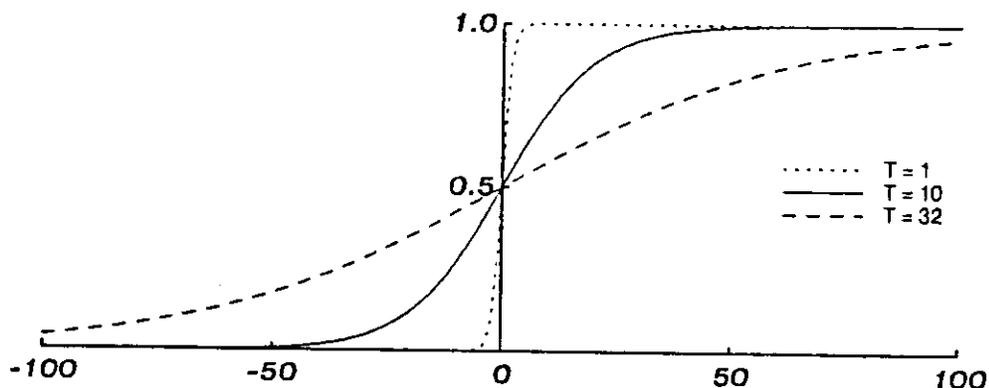


Figure 10: Graph of the Boltzmann equation for three different temperature values. This sigmoid curve shows the probability p_i that unit i will be active as a function of its energy gap ΔE_i .

If a Boltzmann machine starts out at high temperature and is very gradually cooled to a temperature

close to zero, it is likely to end up in a state that is a global energy minimum. The probability that this will happen can be brought arbitrarily close to 1.0 by lowering the temperature sufficiently slowly (Geman & Geman, 1984). This stochastic search technique, which is known as simulated annealing (Kirkpatrick *et al.*, 1983), has been applied with good results to optimization problems unconnected with neural networks, and has also been applied to a variety of problems in low level vision (Marroquin, 1985).

6.4. Matching by Simulated Annealing

The ability to move uphill in energy space allows the Boltzmann version of DCPS to escape local energy minima as it searches for the global minimum. In practice, we have not had to use a genuine annealing search in order to get acceptable performance from the network. When we ran the network at zero temperature, it got trapped in poor local minima, but we discovered that this could be avoided by running at three distinct temperatures. Figure 11 shows the temperature schedule used in the current version of the model.

-
1. **Initialize:** turn off all rule, bind, and clause units.
 2. **Randomize:** run for 2 cycles at temperature 300. This temperature is high enough to ensure that all units which have any chance of being part of the solution have a reasonable chance of turning on, but it is low enough that completely irrelevant units are unlikely to be on.
 3. **Match:** run for up to 10 cycles at temperature 32; stop if the energy is negative after any cycle.
 4. **Cleanup:** run for 4 cycles at a temperature which is effectively zero. (We actually used 0.1 to avoid dividing by zero.)
 5. **Rebias:** raise the threshold of all clause, rule, and bind units by 50.
 6. **Verify:** run for 5 cycles at temperature of effectively zero.

Figure 11: The temperature schedule used in the Boltzmann machine version of DCPS.

The network is initialized for matching by turning off all rule, bind, and clause units, leaving it in a zero energy state. Next its state is "randomized" by running it at a relatively high temperature of 300 for two cycles.⁹ As figure 10 shows, units behave fairly randomly at this temperature, but they are still more likely to be active if their energy gap is positive than negative. At this temperature we have

⁹A cycle is N random updates, where N is the number of units in the network. Although the updating of units is done randomly, on average each unit will get one chance to update its state during each cycle.

observed that the units that support the correct match and units that support partial matches are the ones that are on most often; units unrelated to a legal match become active less frequently. With so many units on, the energy of the network becomes quite high; with six rules (240 rule units) it varies between 8,000 and 12,000. See figure 12.

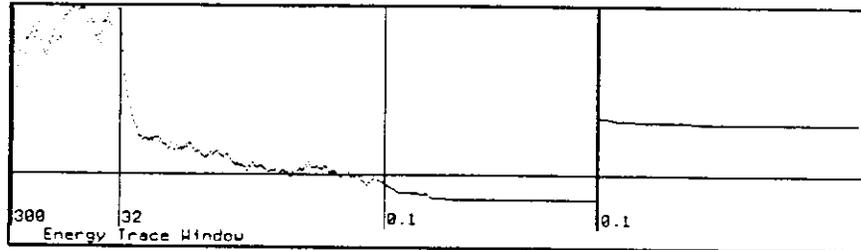


Figure 12: A graph of the energy level as the network follows the temperature schedule of figure 11. Thin vertical divisions mark temperature changes, with the new temperature shown at the bottom of the graph. A thick division marks the point where thresholds are raised in the rebiasing step.

The real matching work is performed in the next step of the schedule, at a temperature of 32. The precipitous drop in temperature from 300 to 32 is more suggestive of quenching than annealing but has no adverse effect on the match. The continued activity of rule, bind, and clause units now depends more strongly on support received from other units, but the network retains enough flexibility at this moderate temperature to explore various match possibilities rather than sink into the nearest local minimum. Cliques for a particular rule in rule space or symbol in bind space may become very active, fade away, and become active again. Triples may materialize in the clause spaces, be partially replaced by other triples, and then perhaps return. The energy of the network rises and falls, but the general trend is decreasing. Once the energy falls below zero¹⁰ the system is deep enough into a local minimum that it is unlikely to get out, so we move on to the cleanup step of the temperature schedule. In this step the network is run at a very low temperature, 0.1. Only units with positive energy gaps will remain active at this temperature. The result is that the clause spaces are left with roughly 28 units on, rule and bind spaces each have one clique active (40 units on), and the network is indicating as clearly as possible what it thinks the correct match should be.

¹⁰This value is approximate and was determined empirically.

6.5. Detecting Failed Matches

There are two ways in which the match can fail. The simplest is when the network fails to settle into any energy minimum at all. In this case very few of the units will have positive energy gaps, so when the temperature drops to 0.1 they will eventually all turn off. The more difficult case to detect is when the network has settled into a local energy minimum representing a partial match. The energy of a partial match is moderately negative, typically around -2500. When the temperature drops to zero the network settles to the very bottom of the energy minimum and stays there.

All correct matches have energies below a certain value, which distinguishes them from partial matches. However, in connectionist models it is better if the behavior of individual units does not depend on measuring global properties of the network such as energy. To detect failed matches without measuring energy directly we use a technique called rebiasing. After the network has run for four cycles at a temperature of 0.1, in the cleanup phase the thresholds of all rule, bind, and clause units are raised by a value of 50, or equivalently, an inhibitory bias of -50 is applied to each unit. This has the effect of reshaping the energy landscape as shown schematically¹¹ in figure 13. The correct match is still a deep energy minimum, but it is much narrower and its absolute energy is now considerably higher than zero. More importantly, a partial match that was a local minimum before is now located on a slope that leads down to the zero energy state with all units turned off. After rebiasing, the network is run for five more cycles at a temperature of 0.1. If units remain active at the end of this step, the network is indicating a correct match. If a partial match was found, units will gradually turn off as a result of rebiasing, causing the energy to drop to zero as shown in figure 14.

One might wonder why the thresholds of the rule, bind, and clause units were not originally set at the higher level, eliminating the need for rebiasing. This would make the energy minima too narrow, making them more difficult for the search to find. Also, after rebiasing the energy of the correct match state becomes moderately positive. At high temperatures the network could find a better state simply by turning all its units off. When rebiasing is delayed until a low temperature has been reached, the network remains trapped in the state (now with positive energy, but still a local minimum) it was in if it managed to find the correct match.

We have also considered the possibility of more flexible temperature schedules for coping with failed matches. After running for 10 cycles in the match phase at a temperature of 32, if the energy is not low enough for the network to have settled into the global minimum, it is probably in a state indicating a partially valid match. Either rule space has settled onto the right rule but bind space

¹¹The true energy landscape is not continuous, and nor can it be represented by a two-dimensional graph. It is an assignment of real values to the corners of an N-dimensional boolean hypercube representing the states of the network, where N is the number of units.

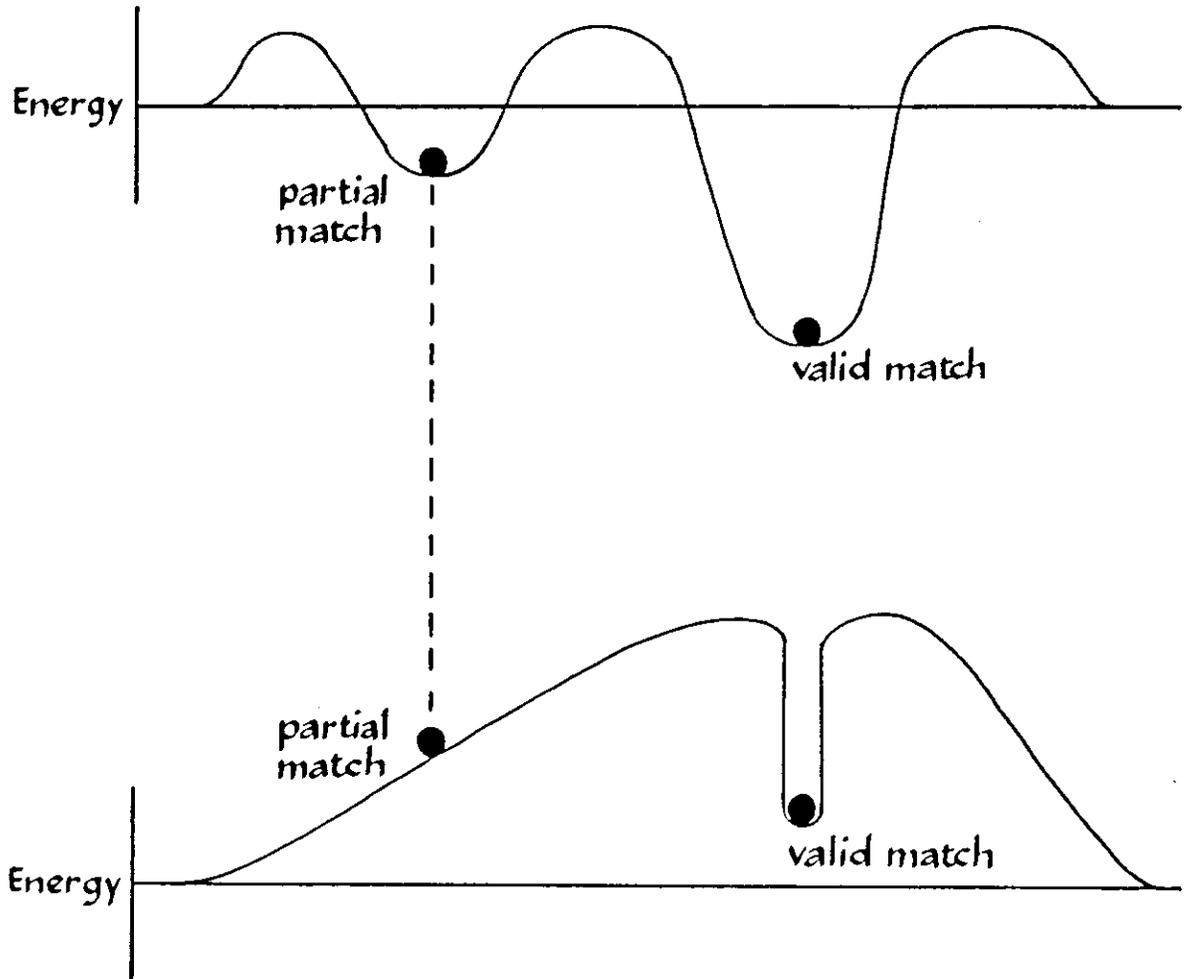


Figure 13: Effect of rebiasing on the energy landscape. The global energy minimum becomes a deep but narrower energy minimum. States representing local energy minima end up on a slope leading down to a zero energy state.

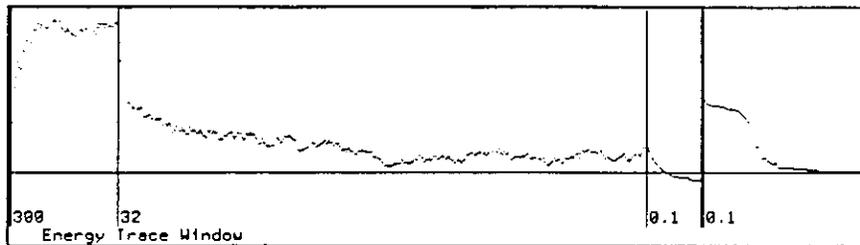


Figure 14: Detection of a partial match by rebiasing. Energy drops to zero as units turn off after their thresholds are raised.

picked the wrong symbol, or else the reverse has occurred. To recover, we could run for a few cycles at a slightly higher temperature, around 40, to kick the network out of its local minimum, and then enter the match phase again.

7. Rule Firing

After a rule has matched successfully it must be fired, which means performing its right hand side actions that update working memory. The ability of rules to update a persistent symbol structure whose contents determines the next rule that will match is what enables DCPS to exhibit interesting sequential behavior. We first consider the problem of right hand side update actions that are variable-free, and then move on to the general case where variables may appear in any position of a triple.

7.1. Variable-Free Actions

The right hand sides of rules are implemented in DCPS as globally gated connections from rule units back to working memory units. The gate is closed during the match process so that rule units cannot affect working memory at all. During the rule firing portion of the production system's recognize-act cycle the gate is opened briefly; at this time, rule units that excite or inhibit working memory units can cause them to change their state. In the absence of outside stimuli, working memory units have a built in hysteresis property that causes them to retain their current state. When the gate is closed prior to the next match cycle, working memory will be frozen in its updated state.

Consider the rule units that implement Rule-1 on page 11. This rule adds the triples (G A B) and (P D Q) to working memory and deletes (F C D). The units that implement Rule-1 will have gated excitatory connections to (G A B) and (P D Q) receptors, and gated inhibitory connections to (F C D) receptors. The hysteresis levels of working memory units are set so that no one rule unit can force them to change state; instead, the concerted action of several units is required. This is another feature of the model that contributes to its tolerance for unreliability in individual

components: if a few random rule units fire spontaneously, they will have no effect on working memory.

Architectures with one-way and/or gated connections admittedly violate the definitions of a Hopfield network or a Boltzmann machine. DCPS requires these types of connections in order to produce sequential behavior; without them it would simply settle into an energy minimum and stay there. Fortunately these special connections only come into play during the rule firing phase of the recognize-act cycle. In the rule matching phase the network is equivalent to one that is a pure Hopfield net/Boltzmann machine, because all the functioning connections are bidirectional and there are no gates opening or closing. The theoretical results of Hopfield and of Hinton and Sejnowski are therefore applicable to DCPS.¹²

7.2. Actions Requiring Instantiated Variable Values

To instantiate variable values into right hand side actions requires a cooperative effort between rule and bind spaces. Consider the $+(G =x P)$ action in Rule-2 on page 15. The Rule-2 units would collectively make excitatory connections to all working memory units that are receptors of $(G =x P)$ for any value of $=x$. On average there will be $6^2/25^2 \times 2000$ or roughly 115 such working memory units. However, these connections are individually gated by bind unit cliques: Connections (synapses) from rule units to working memory units are only effective if the connection itself receives some excitatory stimulation from bind units (see figure 15). This is equivalent to saying that the input to working memory units is the conjunction of the activity coming from a rule unit and a clique of bind units.¹³ Thus, if the network has settled into a state where the F clique is the winner in bind space, only connections from rule units to units that are receptors for $(G F P)$ will be enabled. Each such connection from Rule space back to working memory must be stimulated by several bind units in order to be effective; this is necessary because individual bind units vote for three different symbols; only the collection as a whole votes for a unique symbol. The requirement for support from multiple bind units also makes the network resistant to noise that could occur during rule firing due to randomly malfunctioning bind units.

Gated connections are also needed to allow actions to delete items from memory, because bind units by themselves have no way to tell whether the value they represent is needed for an add action

¹²A similar "equivalent network" argument can be made for the use of regulatory units, even though those cells exert their influence during the match phase.

¹³The use of gated (or conjunctive) connections may appear to violate the normal ground rules of connectionist modeling. It is not difficult, however, to find biological structures that exhibit the crucial property of gated connections: A local non-linear interaction between two synapses. Poggio and Torre (1978) have shown that such interactions can be expected to occur in the dendrites of cortical neurons, and Kandel and Schwartz (1982) have demonstrated the importance, in the sea slug *Aplysia*, of presynaptic facilitation, which is a different way of achieving local, non-linear synaptic interactions.

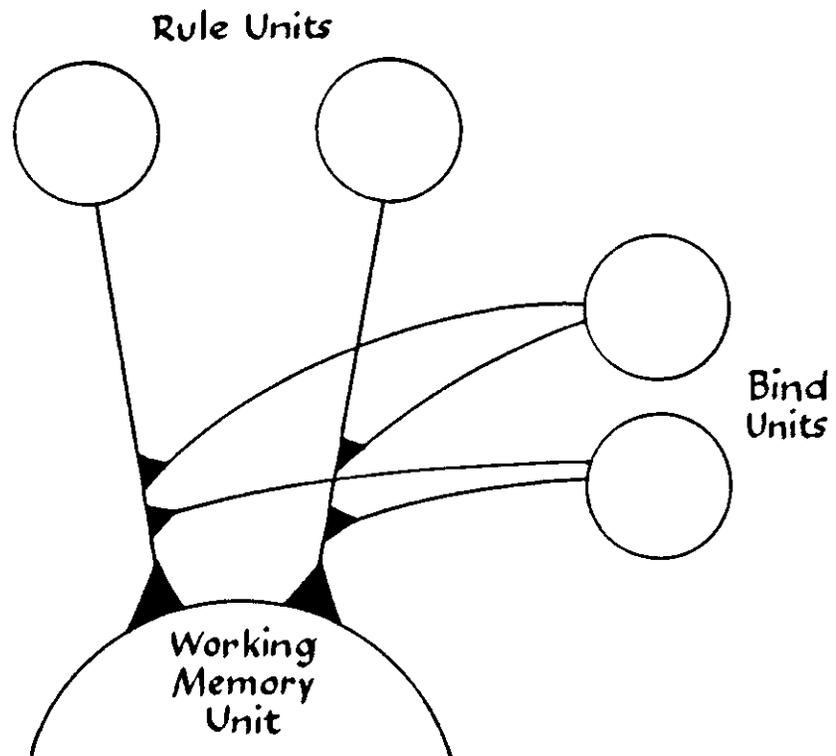


Figure 15: Right hand side actions involving variables are gated by excitatory connections from bind units onto the synapses that rule units make with working memory units.

or a delete action. In the case of delete actions such as $-(=x \ R \ =x)$, the connections from rule units to working memory are inhibitory, but the bind units' effects are always excitatory. By using gated connections, we allow the bind units to select the inhibitory connections that will be allowed to influence working memory.

7.3. Functions on Variable Values

Instead of instantiating the exact value of a variable into right hand side actions, we can instantiate some function of that value. The function will be "computed" by the gating pattern that bind units apply to the rule units' connections to working memory. For example, consider the increment and decrement functions. We will use $>x$ and $<x$ in right hand side actions to denote values one greater and one less than the value to which the variable is bound, e.g., if the variable is bound to F, then $<x$ appearing in a right hand side action would be instantiated as E and $>x$ as G. Modular arithmetic should be used so that every symbol has a successor and predecessor: the successor of Y is A, and the predecessor of A is Y.

Figure 16 shows how the increment function could be used to sequentially step through a series of working memory elements by bumping a counter. The left hand side pattern $(=x \ R \ R)$ refers to the counter value, which is maintained as a triple in working memory and incremented by a right hand side action. On successive firings, rule Seq-1 will step through the triples $(A \ R \ R)$, $(B \ R \ R)$, $(C \ R \ R)$, etc., and leave behind another trail of triples $(A \ B \ A)$, $(B \ C \ B)$, and so on.

Rule:

Seq-1: $(=x \ R \ R) \ (=x \ R \ R) \ \rightarrow \ -(=x \ R \ R) \ \ +(>x \ R \ R) \ \ +(>x \ >x \ =x)$

Initial contents of working memory:

$(A \ R \ R)$

Figure 16: Use of the right hand side increment function to step a counter.

The implementation of the increment and decrement functions is straightforward. In actions that don't compute functions of the bound variable, such as $-(=x \ R \ R)$, the rule units make connections to all working memory units that could match the action, and each connection is gated by bind units of the appropriate type. For example, a connection from a rule unit to an $(A \ R \ R)$ unit would be gated by a set of bind units that vote for A, while a connection to a $(B \ R \ R)$ unit would be gated by bind units that vote for B. To compute a function on the right hand side, the connections that implement the right hand side action are simply gated by bind units specified by the function. Thus, in rule Seq-1, the increment function that appears in $+(>x \ R \ R)$ can be implemented by using bind units that vote for A to gate $(B \ R \ R)$ connections, bind units that vote for B to gate $(C \ R \ R)$ connections, and so on. Any mapping from symbols to symbols can be computed in this way.

8. Experimental Results

8.1. Measured Performance

DCPS has run a six-rule loop overnight through more than one thousand rule firings without error. Working memory contained two triples at a time, and each rule firing involved one addition and one deletion. In the current version of the model, a rule match takes about ninety seconds on a Symbolics 3600 running Common Lisp. Part of this time is spent updating a graphic display as each unit changes state, so that the network's progress can be monitored during the match.

The capacity of the coarse coded working memory of DCPS depends in part on the number of units

used (2000 in our experiments), and also on the similarity of the items that are stored. With a 25 symbol alphabet there are 25 maximally dissimilar triples; an example is the set (A A A), (B B B), through (Y Y Y). This entire set can be stored in DCPS' working memory without losing the ability to distinguish between present and absent triples. (As external observers measuring the memory's capacity, we used 75% activation as the dividing line.) On the other hand, only 5 to 6 elements from a maximally similar set, such as {(A A A) (A A B) ... (A A Y)} can be stored before local blurring begins to interfere with the accuracy of recall. When randomly-generated triples were stored, the measured capacity of the memory was 20 triples on average, varying from a low of 12 to a high of 29.

The number of rules the system can represent appears to be limited only by synergistic effects and by the number of possible partial matches during each search. The largest production system we have run to date, which used a slightly modified version of DCPS as part of a parse tree manipulation task, had 17 rules (Touretzky, 1986b).

The matching portion of an annealing typically involves 6 probes of each unit, where a probe consists of computing the unit's energy gap, deciding whether or not it will change state, and notifying its neighbors if its state does change. Failed matches are detected after 10 probes, when the cleanup portion of the temperature schedule is begun.

8.2. Difficult Match Cases

Early in the development of DCPS we adopted the simplifying assumption that match problems would always have unique answers, so that only one rule and one variable binding could constitute a valid match. This allowed us to avoid the issue of conflict resolution (Brownston *et al.*, 1985), which, although interesting, is not central to our enterprise. But even with this simplifying assumption some match problems are more demanding than others, and situations can be contrived in which DCPS has difficulty finding the correct solution. Two such situations are discussed below.

In the simplest match cases there are no partial matches to worry about; the triples in working memory that do not match the winning rule do not match any of the other rules either. In more complex cases several feasible-looking matches exist with relatively low energy states; the system is forced to search among them to find the lowest one. This involves calling up different triples in the clause spaces for each possibility. As the number of partial matches increases DCPS becomes more likely to settle into a local minimum representing a partially successful match rather than finding the lowest energy state associated with the one correct match. Figure 17 shows a set of rules and working memory elements that produce this behavior. In theory, annealing long enough and slowly enough would solve the problem, since the correct match is always a deeper energy minimum than any partial match.

Rules:

Comb-1: (=x A A) (=x B B) --> ...

Comb-2: (=x C C) (=x D D) --> ...

Comb-3: (=x E E) (=x F F) --> ...

Contents of working memory:

(J A A) (K B B)
 (K C C) (J D D)
 (M E E) (M F F)

Figure 17: A match situation in which combinatorial complexity hinders the search for a valid match.

In this match scenario there are six triples in working memory; the clause spaces must select from among the 36 possible ways to form a pair of triples the one combination that produces a correct match. What makes this problem difficult is the fact that four pairs of triples have fairly deep energy minima representing almost-successful partial matches. See table 2. In these partial matches, either both clauses on the left hand side of rule Comb-1 or Comb-2 are satisfied but the variable binding constraint is not, or else only one of the left hand side clauses is satisfied but the variable binding constraint is met because both clause spaces support the same bind clique (J or K.) The source of the combinatorial confusion is the fact that all three rules and all three bind cliques are capable of getting full support from the clause spaces, so it's difficult to choose among them; what differentiates partial from complete matches is the fact that rule and bind space can't *both* get full support except when the rule is Comb-3 and the variable =x is bound to M.

DCPS does not search a combinatorial space by sequentially enumerating the possibilities. The partial representations of competing triples coexist simultaneously in the clause spaces, while rule and bind winner-take-all spaces host similar competitions. The stochastic nature of the Boltzmann machine causes some competitors in a space to fade out, and possibly fade back in again, until the network as a whole settles deeply enough into an energy minimum that a clear winner emerges in each space.

Figure 18 illustrates another contrived case where it is difficult for DCPS to conclude the match correctly. (M J J) is present in working memory but none of the rules Syn-1 through Syn-4 can match, due to their second clause. While all rules compete with each other as a result of being in a

<u>Degree of Match</u>	<u>Triple in Clause 1</u>	<u>Triple in Clause 2</u>	<u>Rule Supported</u>	<u>Binding Supported</u>
Partial	(J A A)	(K B B)	Comb-1	half J, half K
Partial	(J A A)	(J D D)	half Comb-1, half Comb-2	J
Partial	(K C C)	(J D D)	Comb-2	half J, half K
Partial	(K C C)	(K B B)	half Comb-1, half Comb-2	K
Complete	(M E E)	(M F F)	Comb-3	M

Table 2: The four partial matches generated by the rules in figure 17 have fairly deep energy minima, but there is a global minimum, representing the one complete match, in which all constraints are satisfied.

winner-take-all network, the Syn rules also *help* each other by supporting (M J J) as the first clause. This unwanted synergy, which occurs whenever failing rules have related left hand sides, interferes with the search for the correct match. In order to find this match, the lone Anti rule must override the four Syn rules and get the pattern for (M R R) into C1 space. The more Syn rules there are to support (M J J), the harder this will be.

Rules:

Syn-1: (=x J J) (=x A A) --> ...

Syn-2: (=x J J) (=x B B) --> ...

Syn-3: (=x J J) (=x C C) --> ...

Syn-4: (=x J J) (=x D D) --> ...

Anti: (=x R R) (=x S S) --> ...

Contents of working memory:

(M J J)
(M R R)
(M S S)

Figure 18: A match situation in which synergistic action between four rules that generate partial matches can prevent the system from finding the correct match.

9. Discussion

9.1. Alternative implementations of working memory

There are two broad approaches to implementing a working memory in a connectionist network. The obvious method, which we use here, is to set aside a separate group of units whose activity encodes the current contents of working memory. A less obvious alternative is to use temporary modifications of connection strengths to make it easier to recreate patterns of activity that have recently occurred. The advantage of this second method is that it does not require any extra units to act as a memory, and the memory is automatically content-addressable — recent patterns can be reconstructed from any sufficiently large subpattern. A particularly simple version of the second method is to implement working memory by temporarily lowering thresholds. In DCPS1, for example, the only effect of the units in the working memory space is to provide additional input to units in the clause1 and clause2 spaces, so we could remove the working memory units and exactly mimic their effects by temporary reductions of the thresholds of units in the clause spaces. This would also get rid of all the one-to-one connections between the working memory and clause spaces.

One disadvantage of using thresholds instead of units is that each time a new item is inserted (or deleted) it is necessary to lower (or raise) thresholds in both clause spaces, because there is no way of knowing in advance whether the item will subsequently match the first or the second clause of a rule.

Some important properties of the working memory are broadly independent of whether it is implemented as activity levels, temporary threshold changes, or temporary weight changes. Because the working memory for each item is distributed over many units, thresholds, or weights, there will be interference if more than a few items are stored at once, and the interference will be greater as the items become more similar. This is a necessary consequence of using distributed representations to allow many more possible items than there are storage sites. We interpret the well-known limitations of human short-term memory as an indication that it too may involve the use of distributed representations.

9.2. Multiple interacting distributed representations

In the introduction we alluded to a problem that arises when there are interactions between several groups of units that each use distributed representations. Each unit takes part in the representation of many different items and its causal effects on units in other spaces must reflect this fact. This means that a unit in one space will generally provide excitatory input to a great many units in another space, and so there is a danger that the activation within each space will become more and more diffuse as time progresses. In DCPS, the tendency for activation to become more diffuse is

counteracted by using lateral inhibition within the spaces. This suppresses units that are only supported by a small fraction of the units in other spaces and concentrates the activation on units which receive multiple excitatory input.

Winner take all networks, bind spaces, clause spaces (or pullout networks), and coarse coded symbol representations are generally useful bits of machinery that have been profitably incorporated into other connectionist models. Touretzky (1986a) describes a system for manipulating recursive data structures, called BoltzCONS, that was assembled by rearranging the components of DCPS. BoltzCONS has only one pullout network instead of two, but it has three independent bind spaces.

The representation of rules in DCPS is only "semi-distributed." Although rules are represented by collections of units, each unit is associated with a single rule, rather than being coarse coded. Sharing units between similar rules is counterproductive in this architecture, because rules with similar left hand sides may have totally dissimilar or even directly opposed right hand sides. Consider the two rules Sim-1 and Sim-2 below: one tries to add the triple (H H H) and ones tries to delete it. The rule units common to Sim-1 and Sim-2, which should be in the majority because the rules are so similar, would have both excitatory and inhibitory connections to (H H H) working memory units. Thus, the majority of the rule units would have no action at all. More sophisticated versions of DCPS, which we are presently considering, may be able to exploit similarity among rules by segregating left hand side and right hand side operations into different collections of units.

Sim-1: (=x A B) (=x C D) --> +(H H H)

Sim-2 (=x B B) (=c C D) --> -(H H H)

9.3. Similarity and generalization

One automatic consequence of using distributed representations is that similar items tend to have similar effects. This is a helpful effect if the particular distributed patterns that are used impose a similarity metric that reflects the important distinctions in the domain. If, for example, "cheese" and "chalk" have rather different representations but "cheese" and "cheddar" have rather similar representations, a connectionist network will tend to make sensible generalizations (Hinton *et al.*, 1986). There have been many demonstrations of this effect when the experimenter chooses the distributed representations (Hinton, 1981c; Rumelhart & McClelland, 1986). More recently, Hinton (1986) has described a network that can construct the appropriate distributed representations for itself, so the generalizations cannot be said to have been determined by the experimenter.

DCPS does not currently make any use of similarity between triples or between rules, and it therefore fails to make good use of the properties that a connectionist implementation could provide. We view DCPS as only the first step in the development of connectionist symbol manipulation

architectures. Future advances should lead to models which make better use of the powerful constraint satisfaction and generalization abilities of connectionist networks. Such models would be more than mere implementations of conventional symbol processing ideas because the connectionist substrate would provide important computational properties that are not available in standard implementations.

9.4. Seriality and variable binding

DCPS is implemented in a massively parallel network and yet it is unable to bind the variables in more than one rule at a time. It can perform a parallel search over rules that contain variables to discover which rule fits the contents of working memory best and during this search it considers many different rules and many different variable bindings in parallel, but it is unable to represent particular *conjunctions* of rules and variable bindings. Its only method of representing such a conjunction is by settling on a *single* rule and a *single* binding of each variable. This means that it is using simultaneity to represent the binding, and simultaneity cannot be used for representing several different bindings at once.

Many different variable bindings could be explicitly represented at the same time if we dedicated a separate unit to each possible conjunction of a rule and a variable binding, but this is equivalent to eliminating variables altogether by having many different, variable-free versions of each rule. Newell (1980) has advanced the idea that variable binding may be one of the things that forces people to be sequential processors, and DCPS corroborates this view. By separating the rule space from the bind space we achieve great economies in the the number of units required, but the cost is that the only way to explicitly represent which binding goes with which rule is to settle on one bound rule at a time.

Acknowledgements

This work was sponsored by a grant from the System Development Foundation, and by National Science Foundation grants IST-8516330 and IST-8520359. We thank Scott Fahlman, Jay McClelland, David Rumelhart, and Terry Sejnowski for helpful discussions, and Bruce Krulwich for contributing to the continued development of the simulation.

Appendix A. Model Parameters

DCPS is one of the largest connectionist models built to date. Tables A-1 through A-3 give the number of units in each space and the types, numbers, and weights of their connections. In these tables, thresholds are expressed as connections with weight $-\theta$ to a "true unit" whose state is always 1.

Clause Spaces: 2000 units each

Source of Connections	Number of Connections	Weight per Connection
Working memory unit	1	+900
Other clause units	1999	-2 (<i>mutual inhibition</i>)
Rule units	avg. 7 per rule	+5
Bind units	avg. 40	+10
True unit	1	-939 (<i>threshold</i>)

Table A-1: Parameters of clause units.

Rule Space: 40 units per rule

Source of Connections	Number of Connections	Weight per Connection
C1 clause units	40	+5
C2 clause units	40	+5
Sibling rule units	39	+2
Rival rule units	40 per rival rule	-2
WM units (gated)	40 per RHS action	<i>n/a</i>
True unit	1	-69 (<i>threshold</i>)

Table A-2: Parameters of rule units.

Bind Space: 333 coarse coded units. 3 symbols per unit; 40 units per symbol.

Source of Connections	Number of Connections	Weight per Connection
C1 clause units	240	+ 10
C2 clause units	240	+ 10
Sibling bind units	avg. 107	+ 2
Rival bind units	avg. 225	-2
True unit	1	-119 (<i>threshold</i>)

Table A-3: Parameters of bind units.

Appendix B. Generating Receptive Fields for Working Memory Units.

In our simulation, each triple in working memory is represented by activity in about 28 units. We initially chose the receptive fields of working memory units at random in the obvious way: Six different random letters are chosen for the first position, six for the second, and six for the third. Unfortunately this introduces large sampling errors. Triples represented by as few as 20 or as many as 36 active units are quite common. This can make it hard to distinguish between triples that are present but have few units to represent them and triples that are absent but have accidental activation in some of their many units. If the expected number of active units per triple was much larger than 28, the law of large numbers would eliminate this problem, but in our simulation we used a heuristic method for making the number of units per triple be more uniform.

We started with a set of receptive fields that were chosen so that every letter occurred equally often in each of the three positions. We then considered all possible triples, and recorded how many units encoded each triple. We defined a cost function which was the sum (over all possible triples) of the square of the difference between $6^3/25^3 \times 2000 = 27.65$ and the number of units encoding the triple. This measure is minimized when the number of units per triple is as uniform as possible. We performed gradient descent in this cost function by selecting moves which reduced the cost function but preserved the number of times a letter occurred in each position. A candidate move consisted of taking the receptive fields of two units and swapping two letters in corresponding positions. If for example, two letters from the second position are swapped, the two receptive fields

```
((A B C D E F) (G H I J K L) (M N O P Q R))
((T U V W X Y) (P Q R S T U) (A C E G I K))
```

might become

```
((A B C D E F) (G H R J K L) (M N O P Q R))
((T U V W X Y) (P Q I S T U) (A C E G I K))
```

Candidate moves were selected at random, and were accepted whenever they reduced the cost function or left it unaltered. This was continued until no more improvements were encountered. We considered using simulated annealing to improve the solution, but simple gradient descent was already rather slow and it gave an adequate solution. The standard deviation was reduced from 4.9 to 1.5.

References

- Ackley, D. H., Hinton, G. E., & Sejnowski, T. J. A learning algorithm for Boltzmann machines. *Cognitive Science*, 1985, 9, 147-169.
- Ballard, D. H. & Hayes, P. J. Parallel logical inference. *Proceedings of the Sixth Annual Conference of the Cognitive Science Society*, pp. 114-123. Boulder, Colorado. June, 1984.
- Ballard, D. H. Parallel logical inference and energy minimization. Technical report TR 142, Computer Science Department, University of Rochester, Rochester, NY. March, 1986.
- Barnden, J. A. On short-term information processing in connectionist theories. *Cognition and Brain Theory*, 1984, 7, 25-59.
- Brownston, L., Farrell, R., Kant, E., and Martin, N. *Programming Expert Systems in OPS5*. New York: Addison-Wesley, 1985.
- Derthick, M. A. & Plaut, D. C. Is distributed connectionism compatible with the physical symbol system hypothesis? *Proceedings of the Eighth Annual Conference of the Cognitive Science Society*, 639-644. Amherst, MA. August, 1986.
- Fahlman, S. E. Hinton, G. E. & Sejnowski, T. J. Massively parallel architectures for AI: Netl, Thistle, and Boltzmann Machines. In *Proceedings of the National Conference on Artificial Intelligence*. Washington D.C.: August 1983.
- Feldman, J. A. & Ballard, D. H. Connectionist models and their properties. *Cognitive Science*, 1982, 6, 205-254.
- Hinton, G. E. A parallel computation that assigns canonical object-based frames of reference. In *Proceedings of the Seventh International Joint Conference on Artificial Intelligence*, Vol 2, 683-685. Vancouver BC, Canada. August, 1981a.
- Hinton, G. E. Shape representation in parallel systems. In *Proceedings of the Seventh International Joint Conference on Artificial Intelligence*, Vol 2, 1088-1096. Vancouver BC, Canada. August, 1981b.
- Hinton, G. E. Implementing semantic networks in parallel hardware. In G. E. Hinton & J. A. Anderson (Eds.) *Parallel Models of Associative Memory*. Hillsdale, NJ: Erlbaum, 1981c.
- Hinton, G. E. Learning distributed representations of concepts. In *Proceedings of the Eighth Annual Conference of the Cognitive Science Society*, pp. 1-12. Amherst, Mass., August, 1986.
- Hinton, G. E., McClelland, J. M. & Rumelhart, D. E. Distributed representations. In D. E. Rumelhart and J. L. McClelland (Eds.), *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*. Volume 1. Cambridge, MA: Bradford Books, 1986.
- Hopfield, J. J. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences USA*, 1982, 79 pp. 2554-2558.
- Kandel, E., & Schwartz, J. Molecular biology of memory: Modulation of transmitter release. 1982, *Science*, 218, 433-443.

Kirkpatrick, S. Gelatt, C. D. & Vecchi, M. P. Optimization by simulated annealing. *Science*, 1983, 220, 671-680.

Marroquin, J. L., Probabilistic solution of inverse problems. Technical Report AI-TR-860, MIT Artificial Intelligence Laboratory, Cambridge Mass, 1985.

Mozer, M. C. The perception of multiple objects: a parallel, distributed processing approach. Unpublished thesis proposal, Institute for Cognitive Science, University of California, San Diego. La Jolla, CA: 1984.

Newell, A. Harpy, production systems and human cognition. In Cole, R. (ed), *Perception and Production of Fluent Speech*, pp. 289-395. Hillsdale NJ: Erlbaum, 1980.

Poggio, T. & Torre, V. A new approach to synaptic interactions. In R. Heim and G. Palm (Eds.) *Approaches to Complex Systems*. Berlin: Springer. 1978.

Rolls, E. T. Neurons in the cortex of the temporal lobe and in the amygdala of the monkey with responses selective to faces. *Human Neurobiology*, 1984, 3, 209-222.

Rumelhart, D. E. & McClelland, J. L. (Eds.), *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*. Volume 1. Cambridge, MA: Bradford Books, 1986.

Touretzky, D. S. BoltzCONS: reconciling connectionism with the recursive nature of stacks and trees. *Proceedings of the Eighth Annual Conference of the Cognitive Science Society*, 522-530. Amherst, Mass., August, 1986a.

Touretzky, D. S. Representing and transforming recursive objects in a neural network, or "Trees do grow on Boltzmann machines." *Proceedings of the 1986 IEEE International Conference on Systems, Man, and Cybernetics*, pp. 12-16. Atlanta, GA, October, 1986b.