# Physicists Explore Human and Artificial Intelligence

**J. Buhmann, R. Divko, H. Ritter and K. Schulten***
Physik-Department
Technische Universität München
D-8046 Garching

## 1. Historical Sketch of Brain Theory

The foundation of modern brain theory[1] is based on the epochal work of the physiologist Sherrington and the anatomist Cajal at the beginning of the twentieth century. Both established the modern view of neural networks as heterogeneous systems composed of single subunits, the neurons. They rejected the theory of Golgi and others that the brain is a continuous net of axons and neurons. Sherrington investigated the electrical firing of neurons and introduced the terminus "synapse" for the connection between the individual neurons. These ideas which drove away the animal ghosts of the continuum theory have been spectacularly confirmed half a century later by electron microscopy photographs of neurons and synapses.

In the forties of this century two mathematicians, McCulloch and Pitts, formulated a mathematical theory which allowed us to describe the behavior of neurons as boolean units. In this theory the complex dynamics of neurons is modeled by a logical element which switches to the "on"-state if enough afferent spikes excite the neuron, and, otherwise, rests in the "off"-state. The afferent excitation is summed up linearly and compared with a threshold value to determine the neural state at the next time step. The theory of McCulloch and Pitts is based on the conviction that information is transmitted by action potentials between the neurons. Only a firing neuron can communicate with other nerve cells. The most prominent principles of the hypothesis that neurons can be described by logical units hold also in the light of modern neurobiology. These fundamental ideas have survived until today and have entered in nearly all theories of neural networks.

Other important biological findings about nerve nets refined the picture of neural networks. Dale stated in the thirties that neurons can produce only one type of neurotransmitter and, therefore, can either inhibit or excite connected neurons, not both. At the same time the unidirectional natures of synapses as connections with a prominent direction was discovered. Twenty years later detailed electron micros-

Institute of Theoretical Physics, University of California, Santa Barbara, CA 93106.

copy photographs proved these findings and revealed the structure of synapses with the synaptic cleft, the presynaptic axon terminals and the postsynaptic dendritic membrane. The presynaptic membrane contains vesicles with neurotransmitters which are ejected into the synaptic cleft and change the electrical properties of the ionic channels in the postsynaptic membrane. Thereby, a postsynaptic electrical potential is induced.

The physiologist Hebb outlined a framework of how the various neurons can act together. He introduced the notion of "neural assembly" as a group of cooperating and densely connected neurons. The neurons of such an assembly should excite each other and should strengthen their mutual synapses under the influence of their activity states. Coincidences of pre- and postsynaptic spikes should be the condition for synaptic growth. Hebb was the first scientist who looked for serious concepts to understand the dynamics of neural networks and who tried to connect the neurons and their dynamics with the behavior of higher vertebrates or men. Hebb's coincidence hypothesis has entered in many modern theories of learning and in our days becomes confirmed by physiological discoveries on the molecular structure of synapses and on the mechanisms of their plasticity.

During the ten years from 1955 to 1965, often characterized as the "golden decade of cybernetics," all these ideas were introduced in various models of neural networks, a prominent one of these proposed by Caianello. He evolved his theory of "Thought Processes and Thinking Machines" on the basis of McCulloch's, Pitt's and Hebb's ideas. At the same time Rosenblatt constructed a layered network model with simple connectivity between the layers to recognize and associate patterns. His "Perceptron," however, showed remarkable limitations later proven by Minsky and Papert. The failure of the "Perceptron" and other problems disappointed the hope that thinking and intelligent behavior of humans could be explained with the help of fast computers and on the basis of cooperating neural assemblies in a few years. The area of "Artificial Intelligence" originally located within theoretical biology separated from the neuronal basis and pursued a more abstract, algorithmic approach.

New impulses, originating from thermodynamics far from equilibrium, from non-linear optics and from non-linear mechanical systems, have initiated a renaissance of brain theory in the late seventies. The experience with simple hydrodynamical systems which show complex formation of structure and internal patterns has suggested that the principles of self-organization, cooperation and competition between units with non-linear dynamics could also be essential features in the neural development and information processing of the brain. This more dynamical view has solved the problems of how the information of the neural wiring is genetically stored and how the local variability of the brain structure and fault tolerance of the brain function can be explained. The formation of neural projections between neural nets, i.e. the retinotopic projection[2] and the building of memories with associative and fault-tolerant properties[3,4] were simulated by non-linear dynamics which involves the electrical activity of neurons as well as changes of the synaptic connectivity.

Hopfield[5] found a concise description of an associative memory, in some respects the "harmonic oscillator" of brain theory. His system of spin-like neurons with dense connectivity obeys a Hamiltonian dynamics with dissipation, i.e. the network settles down in the next local minimum found in the phase space. The structure of the phase space reveals many minima and is equivalent to that of a spin glass, a current research topic of condensed matter physics. The results from spin glass physics have influenced brain theory and initiated new investigations. The emerging technology of massively parallel computers with thousands of processors also stimulates the research in this area.

In our own work we have followed both avenues in the history of brain theory, the avenue of modelling the brain in close agreement to physiological principles, and the avenue of experimenting liberally with digital algorithms which are crude caricatures of the way the brain processes information. In Section 2 of this paper we present results of computer simulations of neural assemblies made up of neuronal units modelled in close analogy to their physiological counterparents. The main new result of our work is that electrical noise apparent in 11 physiological recordings of neural tissue appears to play an important role in higher brain function. In Sections 3 and 4 we investigate algorithms which reproduce brain structure and function, albeit in a way which is hardly reminiscent of the biological system. The algorithms in Section 3 are concerned with the self-organization of information representation (mapping) and the control of motor tasks in the brain and in robots. The algorithms in Section 4 address the problem of optical pattern recognition. The results presented in Sections 2, 3, and 4 have been published previously in the proceedings of the conference "Neural Networks for Computing."[6]

## 2. Autoassociative Neural Network with "Physiological" Neurons

Recently, we simulated the activity and function of neural networks with neuronal units modelled after their physiological counterparents.[7] Neuronal potentials, single neural spikes and their effect on postsynaptic neurons were taken into account. The neural network studied was endowed with plastic synapses. The synaptic modifications were assumed to follow Hebbian rules, i.e. the synaptic strengths increase if the pre- and postsynaptic cells fire a spike synchronously and decrease if there exists no synchronicity between pre- and postsynaptic spikes. The time scale of the synaptic plasticity was that of mental processes, i.e. a tenth of a second, as proposed by von der Malsburg.[8] In this Section we present the model network with deterministic dynamics and we extend our previous study and include random fluctuations of the neural potentials. Such fluctuations can always be observed in electrophysiological recordings.[9] We will demonstrate that random fluctuations of the membrane potentials raise the sensitivity and performance of the neural network. The fluctuations enable the network to react to weak external stimuli which do not affect networks following deterministic dynamics. We argue that fluctuations and noise in the membrane potential are of functional importance in that they trigger the neural firing if a weak receptor input is presented. The noise regulates the level of arousal. It might be an essential feature of the information processing abilities of neuronal networks and not a mere source of disturbance to be suppressed. We will

demonstrate that the neural network investigated here reproduces the computational abilities of formal associative networks.[2-4]

The neural system investigated is composed of a set of interconnected neurons, the membrane potentials of which evolve according to deterministic rules and according to stochastic fluctuations. The connections to sensory organs or to other neural networks are taken into account by a primary set of receptors which send input to the neurons. The receptor-neuron connections form a local, static projection of the activity pattern presented by the receptors as modelled by a one-to-one or a center-surround connectivity. The system is schematically presented in Fig. 1.

## 2.1 Dynamics of the Membrane Potential

The dynamics of the membrane potentials involves two processes, the relaxation of the membrane potential and the neural interaction as determined by the somatic integration rule. Axonal spikes are generated whenever the membrane potential reaches a threshold value. The postsynaptic excitation by presynaptic spikes is described by an exponential activity function with decay time $T_U = 1ms$

$$G_k(\Delta t_k/\tau) = \exp\left(-\frac{\Delta t_k}{\tau}\right) \qquad (2.1)$$

$\Delta t_k = t - t_{0k}$ measures the time that has elapsed since the last spike of neuron $k$ at $t_{0k}$.

The kinetic equations of the membrane potentials $U_i(t)$ which also include the stochastic fluctuations are given by a system of non-linear coupled Langevin equations

$$\frac{dU_i(t)}{dt} = -\frac{U_i(t)}{T_R} + \rho[\Delta t_i]\left(\omega\sigma[A_i(t)] + \frac{\eta}{\sqrt{T_R/2}}\xi(t)\right). \qquad (2.2)$$

The first term in Eq. 2.2 approximates the relaxation of the membrane potential $U_i(t)$ to its resting value $U_0 = 0mV$ within a time interval $T_R = 2.5ms$. The second term in Eq. 2.2 describes the communication of the postsynaptic cell $i$ with the connected neurons and receptors, and adds a Gaussian white noise $\xi(t)$ with the strength $\eta/\sqrt{T_R/2}$. The noise produces a Gaussian distribution of the membrane potential $U_i(t)$ with mean value $U_0 = 0mV$ and variance $\eta = 10mV$. Afferent impinging activities in addition to the noise are integrated to the total postsynaptic excitation $A_i(t)$. The activity of the presynaptic neurons $k$ or receptors $j$ are weighted by the time-dependent synaptic strengths $S_{ik}(t)$ or the static receptor connection strengths $R_{ik}$, respectively,

$$A_i(t) = \sum_k S_{ik}(t)G_k(\Delta t_k/T_U) + \sum_k R_{ik}G_k^R(\Delta t_k^R/T_U). \qquad (2.3)$$
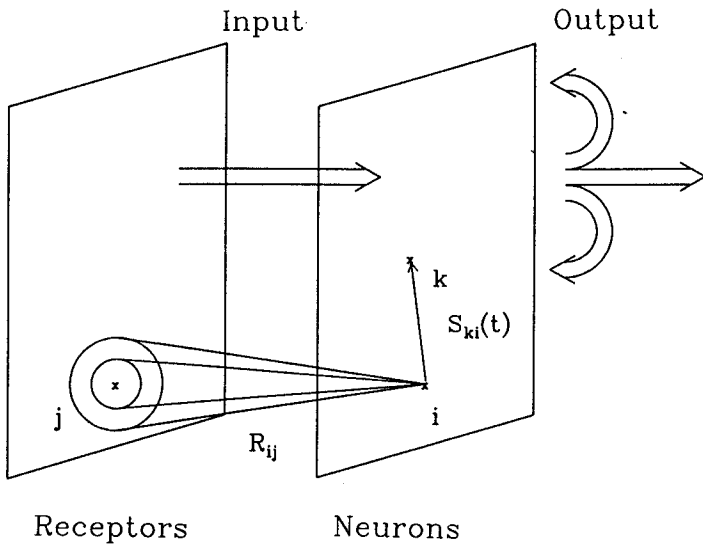
**Figure 1.** Schematic presentation of the neural model investigated: Receptors send spikes to a network of neurons. The resulting activity of the neural network is affected by an activity-dependent alteration of the synapses $S_{ik}(t)$, i.e. the network experiences a feedback as indicated.

The sigmoidal function $\sigma[A_i(t)]$ with a linear behavior for small $A_i(t)$ and a saturation value for strong activity prevents potential changes which are unphysiologically large. The total and relative refractory periods are taken into account by the function $\rho[\Delta t_i]$ which suppresses the sensitivity of neuron $i$ to afferent excitation during a total refractory period $T_F = 5ms$. The function also lets the neuron gradually regain its sensitivity to incoming excitation or inhibition during a relative refractory period of $5ms$.

The continuous time evolution of the potential in our model is interrupted if the neuron reaches the threshold $U_T = 30mV$ and fires a spike. Instantaneously the membrane potential is set to a value normally distributed around the refractory potential $U_F = -15mV$. In this event the time of the last spike $t_{0i}$ is updated and the memory function $G_i(\Delta t_i/T_U)$ is set to the value 1. This behavior is represented as follows:

$$\text{if } U_i(t) \geq U_T, \quad \text{then} \quad \begin{cases} t_{0i} = t, \\ U_i(t) \approx U_F, \\ G_i(\Delta t_i/\tau) = 1. \end{cases} \tag{2.4}$$

The reaction of a neuron to a receptor input depends on the coupling constant $\omega$ and the connection strength $R_{ik}$. In the case of strong coupling the excited neuron will always reach the threshold whereas weak coupling causes only small postsynaptic potentials which never reach the threshold. Figure 2 shows the proba-
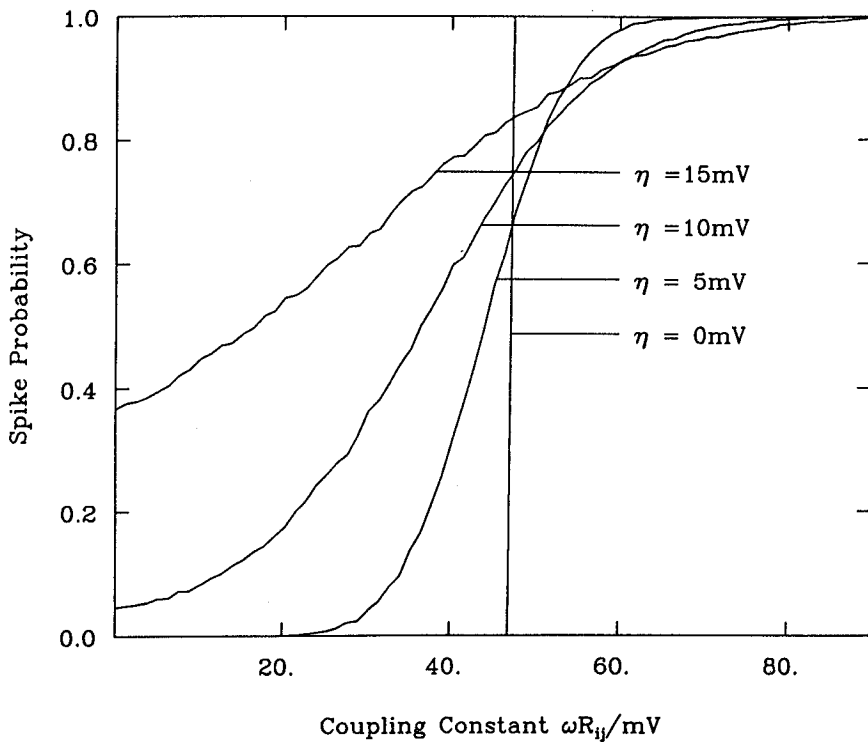
**Figure 2.** The probability to reach the threshold within *5ms* after a receptor spike depends on the coupling between receptors and neurons. The gain of the curve strongly depends on the noise level $\eta$. In our computer simulations we have employed in most cases the value $\eta = 10mV$ and $\omega R_{ik} = 45mV$.

bility that a neuron which received a receptor spike at $t = 0ms$ will fire within *5ms*. This probability is presented as a function of the coupling strength $\omega R_{ik}$ for three different noise levels ($\eta = 0,6,10mV$). Due to the synapse dynamics the mean spike probability of the neuron $\omega \overline{A_i}(t)$ is time-dependent and can be shifted by learning.

## 2.2 Synaptic Plasticity in the Stochastic Neural Network

In our neural network with stochastic firing we introduced a plasticity of the synapses on a time scale of 0.2–0.5s.[2] According to the Hebbian rules the synaptic dynamics were assumed to depend on the synchronicity or asynchronicity of the pre- and postsynaptic spikes. In addition to the Hebbian rules we require for synaptic modifications in the present study that the mean spike frequencies $\overline{v_i}$, $\overline{v_k}$ of both neurons exceed considerably the spontaneous spike rate $v_s \approx 5s^{-1}$. If both neurons satisfy this condition in the case of synchronous firing the synapse can be strengthened. If only the presynaptic neuron fires with a high spike rate the synapse $S_{ik}(t)$ is weakened after each presynaptic spike. Details are described in Ref. 2.

The plasticity of the synapse with the strength $S_{ik}(t)$ connecting neuron $k$ to neuron $i$ is governed by the equation

$$\frac{dS_{ik}}{dt} = \begin{cases} -\dfrac{S_{ik}(t)-S_{ik}(0)}{T_S} + \Omega G_k\!\left(\dfrac{\Delta t_k}{T_M}\right)\kappa(G_i, G_k), & \text{if } S_u \geq |S_{ik}| \geq S_l; \\ -\dfrac{S_{ik}(t)-S_{ik}(0)}{T_S}, & \text{else} \end{cases} \tag{2.5}$$

with

$$\kappa(G_i, G_k) = \begin{cases} 1, & \text{if } G_i > G_k > e^{-1} \ \wedge \ \bar{v}_i \gg v_s \ \wedge \ \bar{v}_k \gg v_s; \\ -1, & \text{if } G_k > e^{-1} > G_i \ \wedge \ \bar{v}_i \ll v_s \ \wedge \ \bar{v}_k \gg v_s; \\ 0, & \text{else.} \end{cases} \tag{2.6}$$

Equation 2.5(a) holds both for excitatory and inhibitory synapses. The first term describes a relaxation process which leads to the gradual loss of stored information. The second term effects a change of the synaptic strength. The influence of this term decays exponentially with the presynaptic activity $G_k(\Delta t_k/T_M)$. The short decay time $T_M = 2.5ms$ guarantees the Hebbian synchronicity condition for synaptic changes. The function $\kappa(G_i, G_k)$ switches between increase of the synaptic strength ($\kappa = 1$), decrease ($\kappa = -1$) and passive relaxation ($\kappa = 0$) of the synapses to the initial value $S_{ik}(0)$. The characteristic time $\Omega^{-1}$ determines the time scale for synaptic modifications. The values assumed for $\Omega^{-1}$ were in the range $0.2 - 0.5s$.

## 2.3 Learning and Association of a Pattern

The neural network presented showed remarkable associative properties in spite of the stochastic fluctuations of the membrane potentials. Starting from a homogeneous structure of synaptic connections with equal numbers of excitatory and inhibitory neurons the network learned a pattern presented by the receptors and associatively reconstructed the original pattern when only incomplete or disturbed patterns were presented.

The simulations of the network were carried out in three different stages. During a first stage which lasted $0.3 - 1.5s$ the neural network had to learn the pattern **brain**, synchronously presented by the receptors with a frequency of $50s^{-1}$. A homogeneous background noise with a spike rate of $10s^{-1}$ was superimposed on the pattern. The coupling constant $\omega R_{ii}$ was set to $45mV$ which effected the firing of about 75 percent of excited neurons. In a second stage lasting $50ms$ the receptors rested quiescent and the electrical activity of the network relaxed to the spontaneous spike rate. During a third stage the receptors presented the test pattern **bra n** which differed from the originally learned pattern by the letter i being left out.

Figure 3a shows the activity of the network at the beginning of the learning phase. At $t = 180ms$ the receptors corresponding to the pattern **brain** had just fired. Within $3ms$, 75 percent of the excited neurons reach the threshold and fire. The other neurons are only gradually excited and fail to fire. The network reaction to a
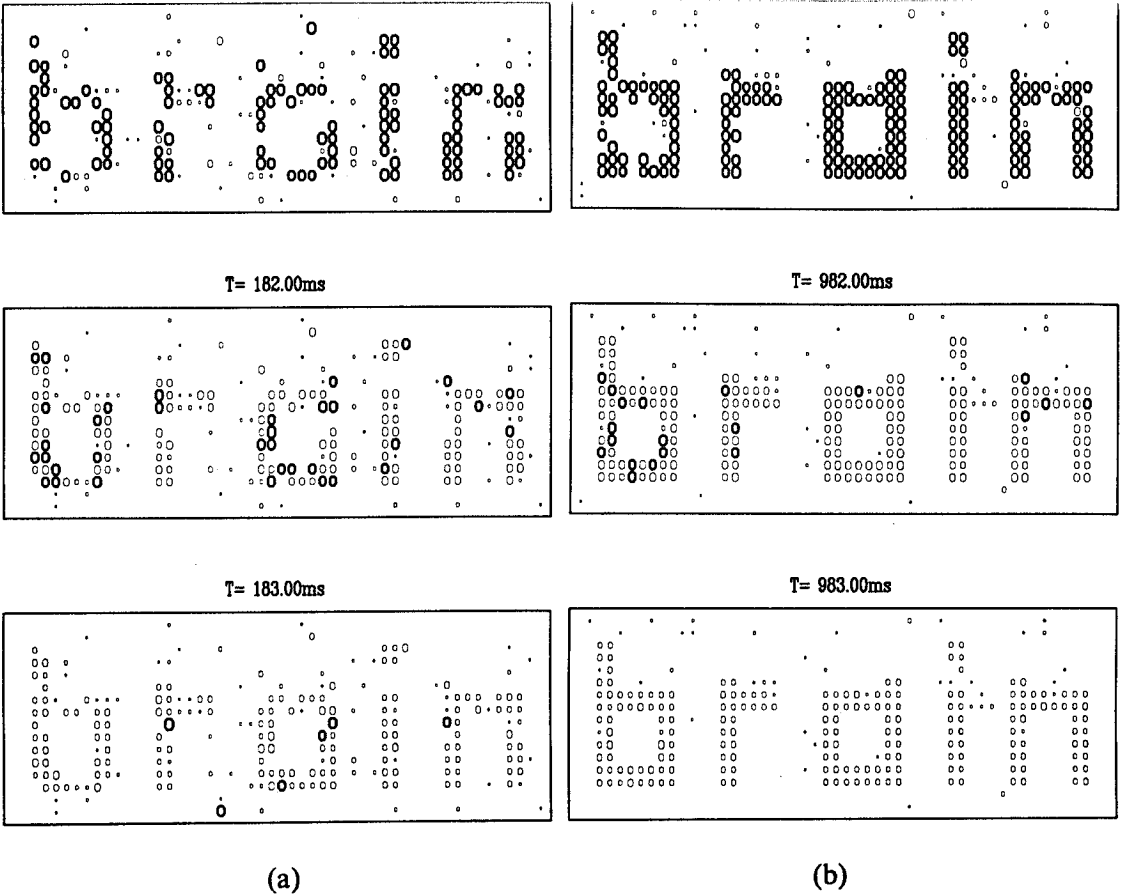
T= 182.00ms                                      T= 982.00ms





T= 183.00ms                                      T= 983.00ms





(a)                                              (b)

**Figure 3.** The activity function $G_i(\Delta t_i/T_M)$ is shown for an untrained (a) and an instructed (b) network. At the beginning of the learning session (t = 181ms, 182ms, 183ms) 75 percent of the excited neurons fire after a receptor spike. At the end of the learning stage (t = 981ms, 982ms, 983ms) nearly all excited neurons have synchronized their firing behavior and reach the threshold.

receptor input at the end of the learning stage is shown in Fig. 3b. Due to the acquired excitatory synaptic connections between neurons receiving input directly from the pattern **brain** (pattern neurons) the assembly reacts more synchronously and the fault level, given by the number of pattern neurons which fail to fire, nearly vanishes.

The success of the learning session is documented in Fig. 4. The incomplete test pattern **bra n** is associatively restored by the network. The neurons representing the missing letter **i** react with a delay time of $1 - 3ms$, i.e. they fire nearly synchronously with the neurons excited by the test pattern.

The synchronization of the neural activity and the associative abilities of the network can be understood on account of the synaptic structure acquired during the learning session. Figures 5a and 5b show the afferent synapses of neuron (37,4) [presented by a star] after the training. All the neurons representing the pattern **brain** have developed saturated excitatory or inhibitory synapses to the reference neuron. During the association task the excitatory synapses saturated at a strength
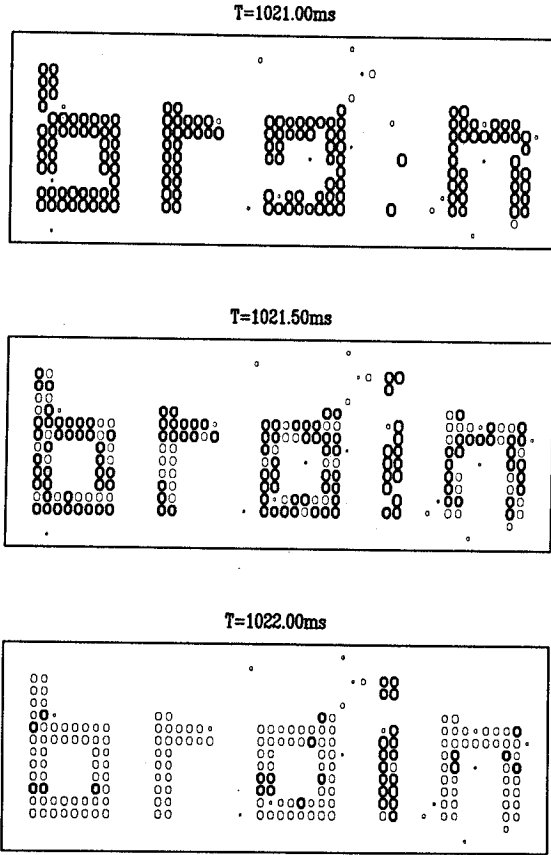
**Figure 4.** Network activity during the association task: The network associates the missing letter i by excitatory interaction within 2 milliseconds (t = 1021ms, 1021.5ms, 1022ms).

value $S_u$ support the firing of the reference cell, whereas the inhibitory synapses saturated at $-S_l$ do not prevent the reference cell from firing. Afferent synapses of the reference cell (37,4) coming from a background neuron rest at the initial synaptic strength.
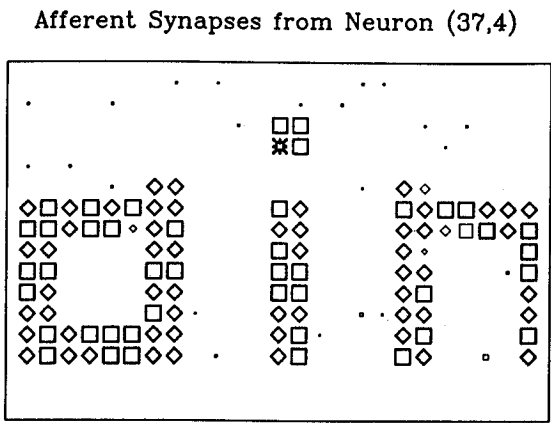


Afferent Synapses from Neuron (37,4)

**Figure 5.** The size of the squares and the diamonds encodes the changes $S_{ik}(t)-S_{ik}(0)$ which the excitatory and inhibitory synapses acquired during the learning session, respectively.

Due to fluctuations of the membrane potential which raise the sensitivity of the neurons the network can also learn a pattern which at any given time is only partially presented by the receptors. At each time interval the invisible fraction of the pattern (50 percent of the receptors) is chosen randomly. The uninstructed network has to learn the total pattern from the detected spike coincidences. The evolution of the synapses is demonstrated for the case of the afferent synapses of neuron (37,4) which represents the dot on the letter $i$. During the learning stage which lasts $3.7s$ the network has built up a synaptic structure which contains the information of the whole pattern (Fig. 6). This simulation demonstrates that the synchronization of all pattern receptors at any given time is not a necessary condition for learning.

## 2.4 Conclusion

We have presented a model neural network with a high level of endogenous noise acting on the cellular potentials. This noise, which is inherent in all biological neurons, does not destroy the abilities of the network to learn and associatively reconstruct patterns. On the contrary, the noise controls the level of arousal and makes the network capable to react to a weak receptor input otherwise neglected. We argue that noise has a functional importance in neural systems. The explicit simulation of single spikes allows us to test the influence of single neural events which are averaged over by mean spike rate models.[4] In addition, the nonspecific influence of large neural nets (neural activity bath) on small neural assemblies can also be studied by stochastic dynamics.

On the basis of the Hebbian rules which detect synchronicities between pre- and postsynaptic spikes, a second condition for synaptic changes is introduced to protect the synaptic structure against destruction by spontaneous activity. The mean spike rates $\bar{v}$ of the pre- and postsynaptic neurons have to exceed considerably the spontaneous spike rate $v_s$ for an increase of the synaptic strengths. For a decrease of the synaptic strengths the postsynaptic spike rate must be considerably below $v_s$. With this modified rule the network can also learn highly noisy patterns and patterns which are presented by a partially asynchronous receptor activity.

Afferent Synapses from Neuron (37,4)          Afferent Synapses from Neuron (37,4)
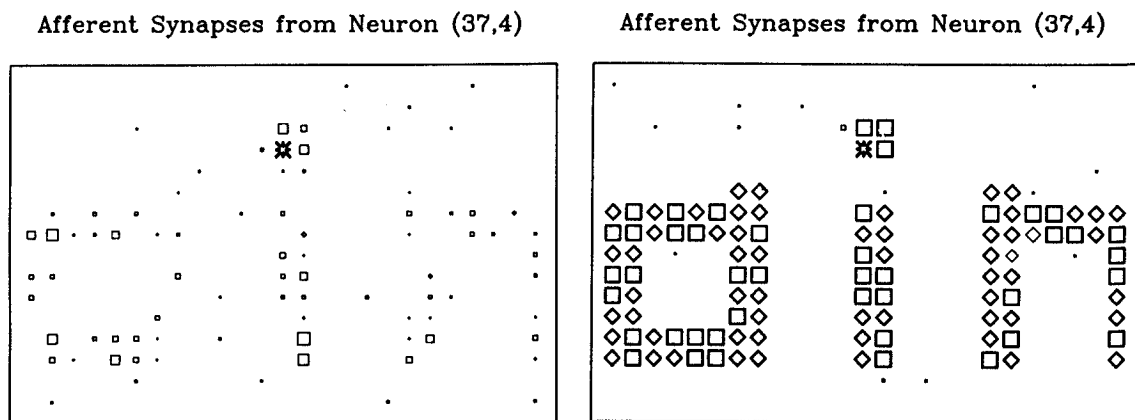


**Figure 6.** Evolution of the afferent synapses of neuron (37,4) for the time $t = 1s$ (a) and $t = 3.6s$ (b) during the learning stage. 50 percent of the pattern **brain** is invisible.

# 3. Topology Conserving Mappings in the Brain

It has been known for a long time that the brain has a modular structure. This is especially conspicuous in the neocortex, which is the brain structure most predominant in man, and which forms essentially a folded, large, two-dimensional neural sheet. Different cognitive abilities, such as touch, audition, vision, and motor capabilities reside in precisely circumscribed regions of this sheet, the so-called brain areas. These areas can be subdivided further, corresponding to different subtasks of each modality. Each area consists of a two-dimensional arrangement of groups of neuronal cells, called cortical columns, with each group devoted to the processing of a small part of the overall information impinging from input pathways upon its area. For the processing of information in the brain, neural activity is mapped between external sensory regions and cortical areas, or between two cortical areas. It is a very important organizational feature of the brain that these mappings are continuous, i.e. that neighboring neuronal groups are connected to neighboring cortical sites of the cortical or sensory regions to which they are mapped. As a consequence, neuronal wiring in the brain preserves the relation of neighborhood, i.e. the neuronal wiring establishes mappings which are topology conserving.

Due to the huge amount of neurons to be connected, the cortical "wiring diagram" cannot be prespecified in detail genetically but must self-organize during the ontogeny and maturation of the brain. In fact, it has been observed that the evolution of such mappings is shaped by sensory experiences. This has been revealed, for example, by neurophysiological experiments which show that the somatosensory map is plastic even through adulthood and can adapt to a changing sensory environment, e.g. brought about by changing "calibration" of the sensors or sensory injury and loss.[10-12]

To understand the principles inherent in the evolution of such mappings, we have studied the formation of the connections between the touch receptors distributed over the skin of a hand and the somatosensory cortex, using a mathematical model originally put forward by Kohonen.[10,11] In this model, the somatosensory cortex is represented as an array of vectors $u(x)$, each vector $u(x)$ corresponding to a neuronal group situated at location $x$ in the cortex. The value of $u(x)$ specifies the position of the receptive field of the group, which is the region of the skin the group is connected to. Initially the connections are completely unordered, i.e. each group is connected to a sensor at a random position in the skin. Each touch stimulus, delivered to the hand at a location $y$, is assumed to create a local region of excitation in the neuronal sheet, which is centered around that particular neuronal group whose receptive field lies closest to the location of the stimulus, i.e. whose vector $u(x)$ matches $y$ best. All neuronal groups in this activated region are now supposed to readjust their connectivity to get their receptive fields closer to the location $y$ of the stimulus, which is mathematically realized by shifting their vectors $u(x)$ towards the stimulus location $y$ by an amount decreasing with increasing distance from the center of the excited region. It turns out that the cumulative effect of such local adjustments, caused by a sufficiently long lasting train of touch stimuli, is capable of establishing a completely ordered mapping between the
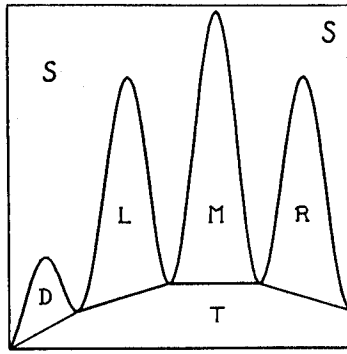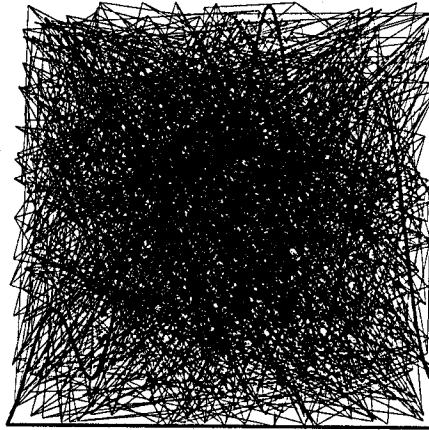
**Figure 7.** Model of the hand surface used in the simulation.

neuronal groups and the touch receptors. This is illustrated in Figs. 7 through 12 in a simulation for the case of the mapping between the surface of a (model-) hand (shown in Fig. 7) and an array of 30×30 vectors, representing a cortical region of 900 neuronal groups. In Fig. 8 we see the initial state of the connectivity, depicted in two different ways. In the upper diagram the array of neuronal groups is projected onto the hand surface via the vectors $u(x)$: each array element $x$ is shown



```
ML..L..TTTLL..L..LR.LMT......M
.T.LRR..R.R.MMD.T.M.R..R.RML..
....RTT..M.M..TRR..L...MMT..R.
....RLL......R..LM.T.LD.T.T.RTM
.RT.TM.....T.RDL.L..M.L.R...DL
T.LLR..T.M...L.....T.TR...LT.M
......T..RM...T.T.T..MM.T.RR...
..........M.RM.TT...R....LR.R.
.M..TR...LT.RL.RR..TR.T.M..MR.
.MTTR..TM.T....T...TT..DDM.RLR
RTRTTD......R.DTL.T..T.....T.M
LTTM.........T...M...L.M.M.R
MT.T...T...L...M......RTT...
.T...M..M...T....T...MTT...T
.R.D..LRMM..R........M.T.TML...
ML..R.R.R.....M..M......M.M.T.R
MT.L..LRL.RT....L.TL.TTTR.T.L.
.....T.T.T...MR..T....TLD..T.T
..L..L..DT...LM...L.TL...TM.RD
MMTML..L...MLTMM.........T..
...M..LT.M...T...TMTT..L...TLT
MMLL..T.RRR.M.M.R.L..R.T..MRT.
...T..R...TM.MRRR.L..R.LL...LM
........R.T....LL..MLM....T...
.LMT.L...L...R..TT.....MMT.TT.
..M...T.R..TR.RL....L.T.R..L..
.....MRT...RMTRLT.ML..MTLMLL.
M..MMDR..MTM..R..L.T.....L..T..
....RTMMM....TMRTMM.....MT...
..M.RM.........R..M.RM....M.RT
```

**Figure 8.** Initial state of the connectivity between hand surface and cortical array.

```
MMMMM..LLLLLTTTTTTTTTTTTTTTTTTM
MMMMMM.LLLLLLTTTTTTTTTTTTTTTTTT
MMMMM..LLLLLLTTTTTTTTTTTTTTTTTT
MMMM..LLLLLLLTTTTTTTTTTTTTTTTTT
MMMM..LLLLLLL.TTTTTTTTTTTTTTTTT
MMM..LLLLLLL..TTTTTTTTTTTTTTTTT
MM...LLLLLLL...OTTTTTTTTTTTTTTT
MM..LLLLLLLL.....TTTTTTTTTTTTTT
M...LLLLLLLL.....TTTTTTTTTTTTTT
M...LLLLLLLLL....LLLLTTTTTTTTTT
M..LLLLLLLLLLLLLLLLLLLTTTTTTTTT
M..LLLLLLLLLLLLLLLLLLLLTTTTTTTT
M...LLLLLLLLLLLLLLLLL.MTTTTTTTT
M...LLLLLLLLLLLLLLLL.MMMTTTTTTT
MM..LLLLLLLLLLLLLLL..MMMMTTTTTT
MM...LLLLLLLLLLLLLL.MMMMMTTTTTT
MMM...LLLLLLLLLLLL.MMMMMTTTTTTT
MMM....LLLLLLLLLL...MMMMMTTTTTT
MMMM.....LLLL....MMMMMMTTTTTTT
MMMMM.........MMMMMMMMMTTTTTT
MMMMMMM...MMMMMMMMMMMMMTTTTTTT
MMMMMMMMMMMMMMMMMM.......TTTTTTT
MMMMMMMMMMMMMM......RRRRRTTTTTT
MMM...MMMMMM...RRRRRRRRRRTTTTT
.........MMM...RRRRRRRRRRRTTTTT
........MM...RRRRRRRRRRRRTTTTT
.......MM..RRRRRRRRRRRRRTTTTT
..LLLL....MM..RRRRRRRRRRRRTTTTT
LLLLLLL...MM..RRRRRRRRRRRRTTTTT
LLLLLLLL...M..RRRRRRRRRRRRTTT
```
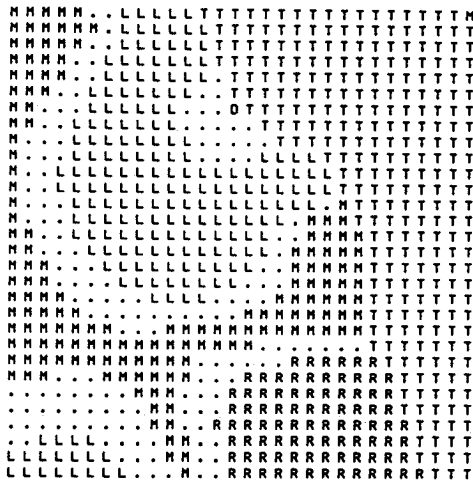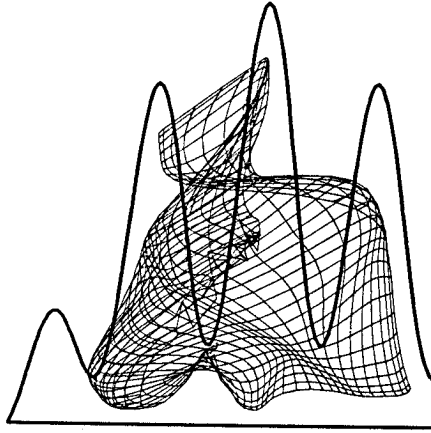
Figure 9. Connectivity between hand surface and cortical array after 500 adaptation steps due to model stimuli.

as a point at location $u(x)$ in the hand surface and the projected positions of elements which are nearest neighbors in the array are connected by lines. The lower diagram shows the array itself. For each element a letter indicates the position of the skin receptor it is linked to, with the letters D,L,M,R,T referring to the regions given in Fig. 7. Dots indicate elements so far unconnected to receptors in the hand (for mathematical convenience these elements are given values of the vector $u$ lying in the region S not belonging to the hand). Both diagrams show that the initial connectivity is completely unordered. The following pictures show the gradual evolution of an ordered mapping under a train of touch stimuli scattered randomly over the hand surface. Figure 9 shows partial order already after 500 touch stimuli, whereas in Fig. 10, after 20,000 touch stimuli, a very regular connectivity has evolved. If finger M is amputated, i.e. in the continuation of the algorithm, no further touch stimuli originate from region M (see Fig. 11), the connections between the cortex and the receptors are redistributed such that the neuronal groups which were connected to the removed region M finally get devoted to the adjacent areas L, R and T (Fig. 12). This behavior is in good qualitative agreement with physiological experiments, e.g. carried out in the brains of monkeys.[13-15]
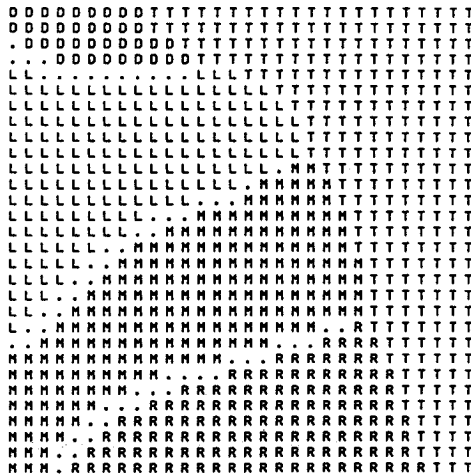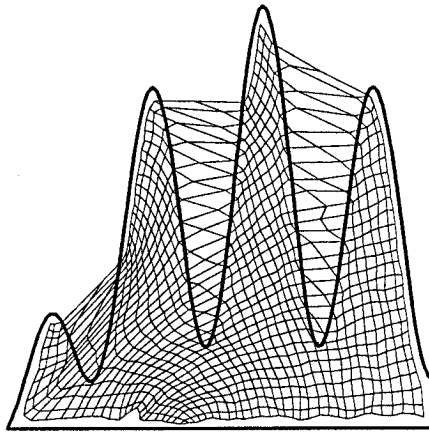
**Figure 10.** Ordered connectivity between hand surface and cortical array after 20,000 adaptation steps.

There is neurophysiological evidence that the role of such mappings need not be restricted to the mediation between sensory input and a cortical target field, but that they can be involved as well in the generation of output. The presence of topology conserving mappings of variables concerned with motor responses, e.g. the mapping of the next movement vector for a saccade of the eyes to a location of increased activity in the superior colliculus,[15,16] suggests that topology conserving mappings may also subserve the acquisition and execution of motor tasks. In learning a new motor capability, initially the movement has to be generated fully consciously and step-by-step until it becomes increasingly automatic with further practice. This process of learning could be brought about by an internal teacher, who teaches, by a series of consciously generated movement instances, a topology conserving mapping the correct input-output relation for the task to be performed. As a consequence of the gradual formation of the map, the brain can take over the part of the teacher to an ever larger extent, rendering the execution of the task more and more automatic until, finally, the internal teacher becomes completely dispensable and thus free to dedicate itself to other tasks. We studied this process by the model task of learning how to balance a pole against gravity.[17] The motion of the pole was simulated in discrete time steps and monitored by an array of

```
D D D D D D D D T T T T T T T T T T T T T T T T T T T T T T
D D D D D D D D T T T T T T T T T T T T T T T T T T T T T T
. D D D D D D D D D T T T T T T T T T T T T T T T T T T T T
. . . D D D D D D D D T T T T T T T T T T T T T T T T T T T
L L . . . . . . . . . . L L L T T T T T T T T T T T T T T T
L L L L L L L L L L L L L L L T T T T T T T T T T T T T T T
L L L L L L L L L L L L L L L L T T T T T T T T T T T T T T
L L L L L L L L L L L L L L L L L T T T T T T T T T T T T T
L L L L L L L L L L L L L L L L L L T T T T T T T T T T T T
L L L L L L L L L L L L L L L L L L L T T T T T T T T T T T
L L L L L L L L L L L L L L L L L . . . . T T T T T T T T T
L L L L L L L L L L L L L L L . . . . . . . T T T T T T T T
L L L L L L L L L L L L L . . . . . . . . . . T T T T T T T
L L L L L L L L L . . . . . . . . . . . . . . T T T T T T T
L L L L L L L L . . . . . . . . . . . . . . . . T T T T T T
L L L L L . . . . . . . . . . . . . . . . . . . T T T T T T
L L L L . . . . . . . . . . . . . . . . . . . . . T T T T T
L L L . . . . . . . . . . . . . . . . . . . . . . T T T T T
L L . . . . . . . . . . . . . . . . . . . . . . . T T T T T
L . . . . . . . . . . . . . . . . . . . . . . R T T T T T T
. . . . . . . . . . . . . . . . . . . . R R R T T T T T T T
. . . . . . . . . . . . . . . . . . R R R R R R T T T T T T
. . . . . . . . . . . . . . . R R R R R R R R R R T T T T T
. . . . . . . . . . . . . R R R R R R R R R R R R R T T T T
. . . . . . . . . . R R R R R R R R R R R R R R R R T T T T
. . . . . . . . . R R R R R R R R R R R R R R R R R T T T T
. . . . . . . R R R R R R R R R R R R R R R R R R R R T T T
. . . . . . R R R R R R R R R R R R R R R R R R R R R T T T
. . . . . R R R R R R R R R R R R R R R R R R R R R R T T T
```
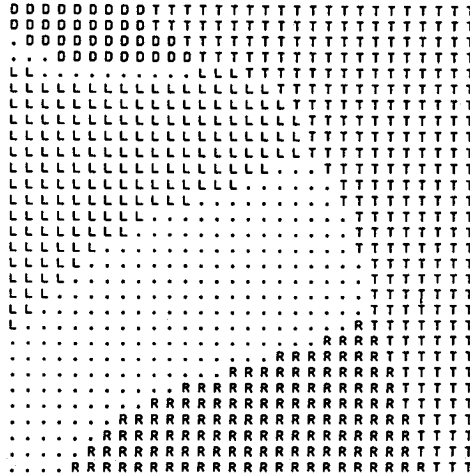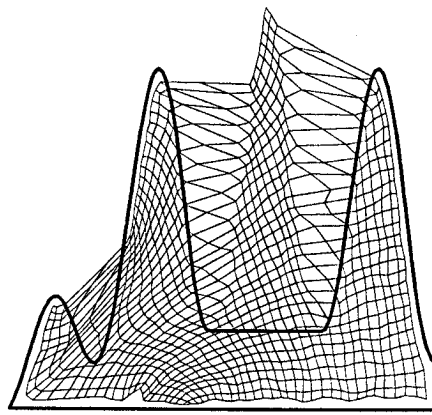
**Figure 11.** State of Fig. 4 after removal of region M, i.e. the middle finger.

vectors $u(x)$. Initially a teacher generated a sequence of balancing forces capable of balancing the pole. As in the previous example, each vector codes for a certain state of the pole, here represented by two successive pole inclinations, separated by one time step of the simulation. In addition to this "sensory part," a third component codes for an "action": here the action specifies the value of a force, which is to be applied over one time step, whenever the state of the pole matches the "sensory part" of the vector. The "sensory input" to the array is now generated by the motion of the pole and consists of a sequence of pairs of successive pole inclinations. At each time step the unit $x^*$, whose "sensory part" matches the state of the pole most closely, is taken to be the center of the region of local adjustment of the values $u(x)$ in the array. Such an adjustment consists of refining the "sensory parts" of the vectors of all units lying within a small neighborhood of unit $x^*$. These units are matched more closely to the latest pair of successive pole inclinations and at the same time their "action part" is altered towards the force value supplied by the teacher. In the course of time the balancing force delivered by the teacher at each time step is replaced gradually by the output evoked by the "action part" of the array element residing in the center of the region of adjustment. Finally the contribution of the teacher vanishes and the control is exerted by the array alone. Figure 13 shows the gradual improvement of the balancing capability
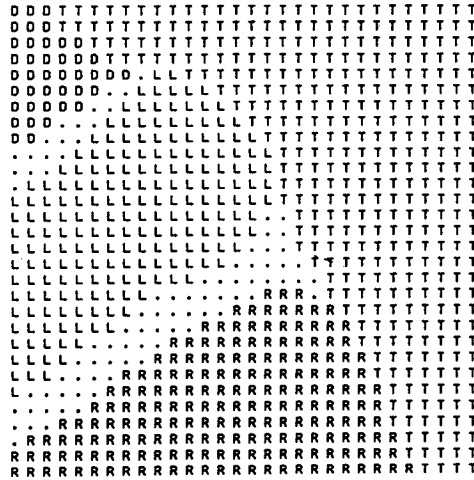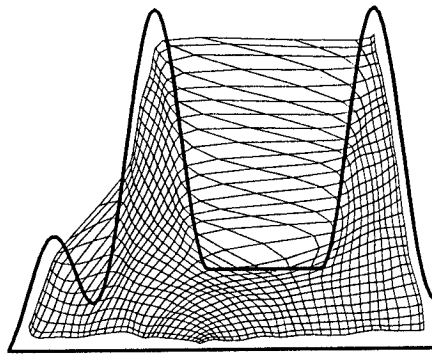
```
D D D T T T T T T T T T T T T T T T T T T T T T T T T T
D D D T T T T T T T T T T T T T T T T T T T T T T T T T
D D D D D T T T T T T T T T T T T T T T T T T T T T T T
D D D D D D T T T T T T T T T T T T T T T T T T T T T T
D D D D D D D D . L L T T T T T T T T T T T T T T T T T
D D D D D D . . L L L L T T T T T T T T T T T T T T T T
D D D D D . . L L L L L L T T T T T T T T T T T T T T T
D D D . . . L L L L L L L L T T T T T T T T T T T T T T
D D . . . L L L L L L L L L L T T T T T T T T T T T T T
. . . . L L L L L L L L L L L L T T T T T T T T T T T T
. . . L L L L L L L L L L L L L L T T T T T T T T T T T
. L L L L L L L L L L L L L L L L T T T T T T T T T T T
L L L L L L L L L L L L L L L L L T T T T T T T T T T T
L L L L L L L L L L L L L L L L . . T T T T T T T T T T
L L L L L L L L L L L L L L L . . . T T T T T T T T T T
L L L L L L L L L L L L L L L . . . T T T T T T T T T T
L L L L L L L L L L L L L L . . . . . . T T T T T T T T T
L L L L L L L L L L L L L . . . . . . . T T T T T T T T
L L L L L L L L L L . . . . . . . R R R . T T T T T T T
L L L L L L L L L . . . . . R R R R R R T T T T T T T T
L L L L L L L . . . . R R R R R R R R R R T T T T T T T
L L L L L . . . . R R R R R R R R R R R T T T T T T T T
L L L L . . . . R R R R R R R R R R R R R T T T T T T T
L L L . . . R R R R R R R R R R R R R R R R T T T T T T
L . . . . R R R R R R R R R R R R R R R R R T T T T T T
. . . . . R R R R R R R R R R R R R R R R R T T T T T T
. . . R R R R R R R R R R R R R R R R R R R T T T T T T
. R R R R R R R R R R R R R R R R R R R R R R T T T T T
R R R R R R R R R R R R R R R R R R R R R R R T T T T T
R R R R R R R R R R R R R R R R R R R R R R R T T T T T
```

**Figure 12.** Readaptation of connectivity after 50,000 further iterations, subsequent to state in Fig. 11.

of the array at three different stages of learning. For both the "sensory" and the "action" parts of all array elements, random entries were taken as starting values. After learning for 100 seconds of simulated time, the array can balance the pole moderately well (bottom diagram of Fig. 13). The performance has improved after 500 seconds (middle diagram) until, after 1000 seconds (low diagram), only very little fluctuation around the unstable vertical pole position remains during a balancing trial.

In the course of the learning process the array has achieved two things simultaneously: First, it has established a mapping between inputs and appropriate responses by having learned the relation between the state of the pole and a necessary force to keep the pole balanced. Second, it has distributed the values of the "sensory parts" of its vectors $u$ such that they populate only that region of the state space of the pole, which is actually visited by the motion.

However, this kind of learning still requires a teacher. Where does the teacher get its instructions from? Even motions such as that of a simple invertible pendulum may require preplanning if there are restrictions upon the available control force. For instance, if the initial state of the pendulum is the stable resting state, it may be desired to turn it into the inverted unstable equilibrium position by applying a
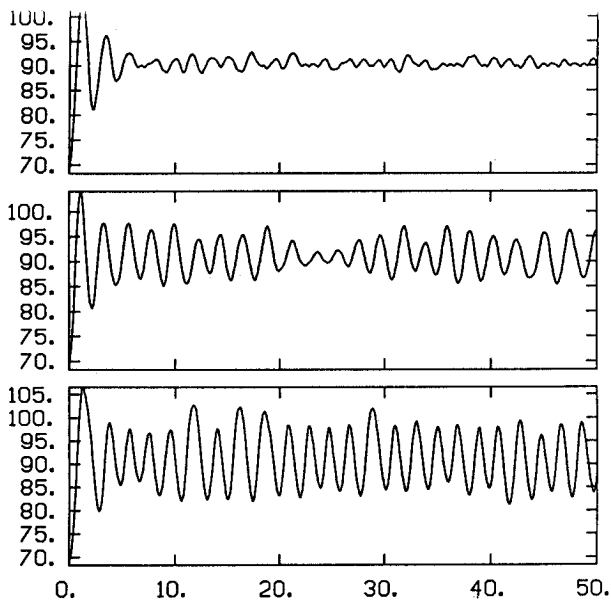
**Figure 13.** Balancing Capability of Array after different Periods of Learning. Each diagram shows the time evolution of pole inclination $\theta$ (vertical) during the first 50 seconds under the control of the array after releasing the pole from a resting position with $\theta = 70°$. *Lower diagram:* after 100 sec of training; *Middle diagram:* after 500 sec of training; *Top diagram:* after 1000 sec of training.

suitable torque at its pivot. In this situation the direct approach of simply turning it up may fail if the admissible torque for the task is too weak. Instead, the pendulum first has to be swung back and forth several times before the weak torque suffices to complete the motion. To plan a motion trajectory for tasks of this kind, we have considered a formulation of the task as a path search problem in the phase space of the system.[18] The solution of this problem can then be achieved using a physical analogy from the realm of diffusion. The computation can be performed in a fully parallel and "neuronal plausible" manner. In order to apply this method, the phase space of the system is discretized into a lattice of nodes, and possible transitions between nodes are represented as directed links. The presence of a link depends on the equation of motion and the available constraints on the control force. The given initial state and the desired final state of the motion can then be represented by a starting node A and an end node B in the lattice. The problem to find a trajectory between these two states is now transformed to the search for a lattice path, the path consisting of a series of oriented links connecting A and B. This latter search problem is solved by considering B as the source of a fictitious substance, which spreads over the lattice by diffusion. The steady state concentration of this substance can be calculated by a simple relaxation algorithm. The path sought is found by starting at A and following the steepest gradient of the concentration until B is reached.

The application of the method to the example mentioned above, i.e. moving a pendulum from the stable to the unstable position, is presented in Figs. 14 through 16. Figure 14 shows the discretized portion of the phase space of the pendulum, with the links indicating transitions possible with a maximum torque of 0.5 (a torque value of 1 corresponds to holding the pendulum fixed in horizontal position). If the
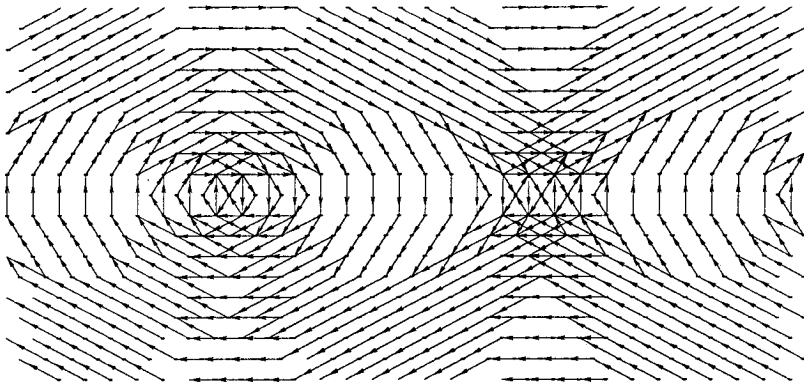
**Figure 14.** Discretized portion of the phase space for the pendulum described in the text. The angle $\theta$ is measured with respect to the vertical axis.

task is to bring the pendulum from the downward resting position into the upward resting state, a torque of 0.5 is not sufficient to achieve this via a direct path. Instead, the pendulum first has to be swung through a few oscillatory cycles to accumulate kinetic energy. The trajectory found by the above method, therefore, first spirals outward from its left starting point until it reaches its destination (Fig. 15). The resulting inclination angle versus time of the pendulum is shown in Fig. 16.

Although, in its present formulation, the method consists of a "neural-like" computation, it is unlikely that this particular algorithm may be realized in this manner in actual brains. Nevertheless it is important to explore ways to formulate tasks commonly solved by our nervous system in a manner amenable to highly parallel, neural-like computation, as only by experimenting within a broad spectrum of different algorithmic alternatives can we hope to finally come closer to the actual working principles of the brain.



**Figure 15.** Trajectory found by the algorithm. Starting point A is the left focus of the spiral, corresponding to the pendulum resting downward ($\theta = 180°$, $\dot{\theta} = 0$). Target point B is the right end point of the trajectory, corresponding to the upward position of the pendulum ($\theta = 0°$, $\dot{\theta} = 0$).
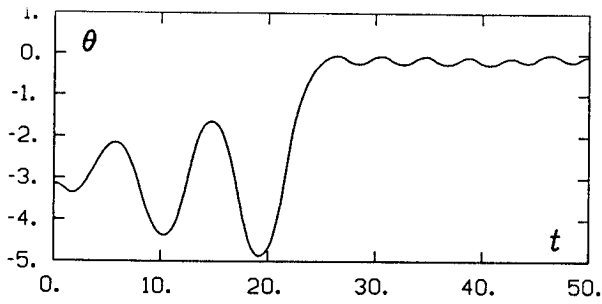
**Figure 16.** Time evolution of the angle $\theta$ of the pendulum corresponding to the simulation of the trajectory in Fig. 15.

# 4. Stochastic Spin Models for Pattern Recognition

## 4.1 Introduction

In this Section we exploit, for the recognition of patterns, the properties of physical spin systems to assume long range order and, thereby, to establish a global interpretation of patterns. For this purpose we chose spins which can take a discrete set of values to code for local features of the patterns to be processed (feature spins). A Hamiltonian is chosen for the system which entails a field contribution and interactions between the feature spins. The field incorporates the information on the input pattern. The spin-spin interaction represents *a priori* knowledge on relationships between features, e.g. continuity properties. The Hamiltonian is chosen such that the ground state of the feature spin system corresponds to the best global interpretation of the pattern. The ground state is reached in the course of local stochastic dynamics, this process being simulated by the method of Monte Carlo annealing.[19] Our study is related to work presented in in Refs. 20 and 21.

Spin systems are characterized by a set of values for the spin variable $S_{i,j}$, a lattice on which they are defined, and by an interaction energy. In the two-dimensional Ising model the spins take the values $\pm1$ and the interaction energy is defined by the Hamiltonian

$$E = -J \sum_{<(i,j),(k,l)>} S_{i,j} S_{k,l} - \sum_{(i,j)} H_{i,j} S_{i,j} \qquad (4.1)$$

where the brackets indicate summation over nearest neighbors. In the ferromagnetic case $(J > 0)$ the first term, the exchange interaction, gives a negative contribution if neighboring spins point in the same direction. This term creates a tendency for an alignment of all the spins. The second term describes the interaction of the spins with a local magnetic field $H_{i,j}$ tending to align the spins locally with the field. The regularizing effect of the exchange interaction will be utilized in the following to solve pattern recognition tasks under the constraint that pattern features are expected to vary continuously.

For the purpose of picture processing, the spins are chosen to code local features of a pattern. Examples for attributes coded by such feature spins are intensities, disparities between corresponding points in a stereogram or edges of different directions. Several different types of feature spins, interacting with each other and with external fields, may be needed to solve a specific pattern recognition problem.

At finite temperatures the feature spin system shows fluctuations like its physical counterpart. Certain values of a feature spin at a certain lattice point are more probable than others. One may consider the value of a feature spin as the hypothesis that the picture has a certain local attribute at this point.

At high temperatures all hypotheses are equally probable. After carefully cooling down to low temperatures (simulated annealing[19]) the fluctuations eventually disappear. At zero temperature the feature spins take definite values indicating the final global hypothesis about a pattern.

The final hypothesis, i.e. the ground state, achieved by the system after cooling down to low temperatures depends on the interaction among the feature spins as well as on the interaction with the external field. The interaction among the feature spins contains an *a priori* global knowledge on relationships to be expected to hold between the features of a pattern. Correct interpretations of patterns must meet certain constraints, e.g. the constraints of continuity, which have to be realized by the feature spin configurations in the final hypothesis. Such configurations can be achieved by a properly chosen interaction between the feature spins. For example, a Potts model type interaction between intensity spins yields a smooth change of brightness. The external field serves to communicate the pattern to be processed to the system of feature spins. Examples for pattern attributes coded by the external fields are local brightness or edges of various directions.

## 4.3 Stereo Vision

Whereas a certain degree of depth vision can be obtained from perspective distortion or from hidden parts of a scene, full stereo vision is a result of binocular perception. The projections of an object in both eyes differ slightly from each other. This difference (disparity) allows the reconstruction of the three-dimensional information. Figure 17 shows an image pair of dot patterns appearing completely random when viewed monocularly. But when viewed one through each eye the two pictures fuse showing a three-dimensional structure (square hovering over the ground). The absence of higher level structures in the patterns shows that disparity alone can be used to obtain three-dimensional information from a stereogram.

To obtain stereo information the disparity of corresponding points in the two retina projections of an image must be determined. The problem is to assign correspondences between points of the two pictures. This is a difficult task because of the so-called "false target problem"[22] occurring in its extreme in Julesz patterns.[23]

**Figure 17.** Julesz pattern[23] with 50% black dots. This random-dot stereogram of 50×50 pixels is generated by copying the right image from the left one, shifting a square-shaped region of 30×30 pixels slightly to the left and filling the gap caused by the shift with a new random pattern.

Every black pixel could correspond to every other black pixel. To restrict all possible combinations of points from both pictures to physically plausible correspondences the following matching conditions must hold:

- Compatibility: Black dots can only match black dots and vice versa.

- Continuity: The physical feature disparity varies smoothly almost everywhere over the image.

- Uniqueness: Except in rare cases each point from one image can match only one point from the other image.



**Figure 18.** Ground state of the disparity spin system corresponding to the Julesz pattern in Fig. 17.
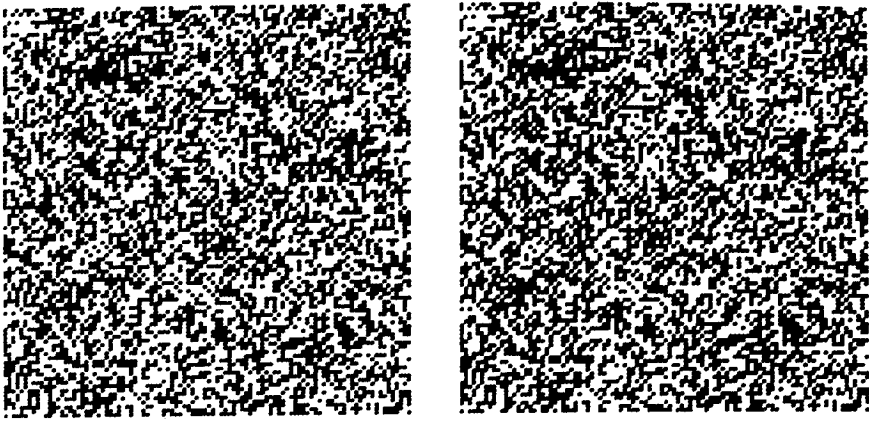
**Figure 19.** Julesz pattern of an eight level pyramid.

The following spin model is designed to find the correspondences between the pixels of both pictures of a random-dot stereogram and to measure the disparities. This information will be contained in the ground state of the spin system.

### 4.4 A Spin Model for Stereo Vision

The feature coded by a spin is the disparity of corresponding pixels. A disparity spin with a value $S_{i,j} \in \{0, \pm 1, ..., \pm N\}$ at lattice site $(i,j)$ stands for the hypothesis of a correspondence between the pixel $(i,j)$ in the right picture of the stereogram and the pixel $(i, j + S_{i,j})$ in the left picture shifted $S_{i,j}$ units to the right. Both pixels are assumed to correspond to the same original point of an object and to have the disparity $S_{i,j}$.

The Hamiltonian of the disparity spin system is split into two contributions
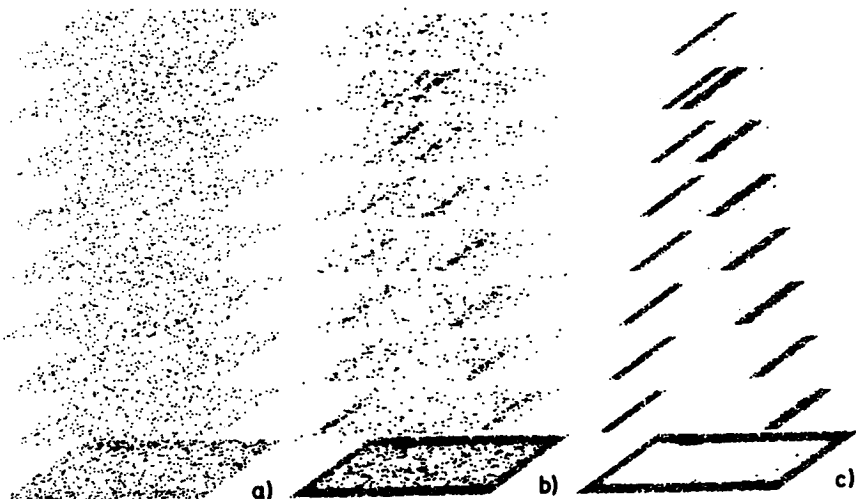


**Figure 20.** Behavior of the disparity spin system for the Julesz pattern in Fig. 19.

These terms reflect the continuity and the compatibility conditions, respectively. The distance of an observer to a point on the surface of an object is a smoothly varying property. To achieve a corresponding property for the values of the disparity spin, a Potts model interaction is chosen

$$E_{exchange} = -J \sum_{<(i,j),(k,l)>} F(S_{i,j}, S_{k,l}) ,$$

$$F(S_{i,j}, S_{k,l}) = \begin{cases} 1, & \text{if } S_{i,j} = S_{k,l} \\ q, & \text{if } S_{i,j} = S_{k,l} \pm 1; \ q < 1. \\ 0, & \text{else} \end{cases} \quad (4.3)$$

If the value of a disparity spin is $S_{i,j}$ and if this hypothesis is correct, the pixel $(i,j)$ in the right picture and the pixel $(i, j + S_{i,j})$ in the left picture have identical surroundings. If the disparity hypothesis $S_{i,j}$ is wrong the surroundings may be completely different. Therefore, the comparison of the neighborhoods of the two points assumed to correspond to each other indicates a possibly correct match.

Whereas many features can (and for real pictures must) be used for comparison, we restrict ourselves in the present application to the most simple choice and compare the pixel intensities in a square shaped region only. Comparison is established by the following energy contribution

$$E_{field} = \sum_{(i,j)} G(S_{i,j}) ,$$

$$G(S_{i,j}) = G_0 \sum_{k=i-w}^{k=i+w} \sum_{l=j-h}^{l=j+h} \left| H_{k,l+S_{i,j}}^{left} - H_{k,l}^{right} \right| . \quad (4.4)$$

Here $H_{i,j}^{left}$ and $H_{i,j}^{right}$ denote the intensity of the pixel $(i,j)$ in the left and in the right picture of the stereogram and $G_0 = [(2h + 1)(2w + 1)]^{-1}$ is a normalization constant. Correct disparity spin configurations are characterized by low energy contributions.

For the input pattern shown in Fig. 17 the disparity field obtained is presented in Fig. 18. Starting from the temperature $T = 1.5$ the annealing process was stopped at a low temperature $T = 0.01$. For the interaction parameters in Eq. 4.3 we have assumed the values $J = 2$ and $q = 0$ and for the maximal disparity the value $N = 5$.

A more complicated stereogram containing the three-dimensional information of an eight level pyramid is shown in Fig. 19. The solutions of the disparity spin system

($N = 10$) with interaction parameters $J = 2$ and $q = 0.2$ are shown for three temperatures in Fig. 20.

At the high temperature $T = 0.8$ the fluctuations of the disparity spin values are large. This is demonstrated by a snapshot of the dynamics shown in Fig. 20a. Figure 20b illustrates that at the intermediate temperature $T = 0.3$ the system still fluctuates; however, the disparity field already indicates the presence of different disparity planes. Figure 20c shows that at the low temperature $T = 0.1$ the fluctuations almost disappeared and that the disparity spin system achieves the correct interpretation of the Julesz pattern, an eight level pyramid.

## 4.5 A Spin Model for Picture Restoration

Restoration of noisy pictures can be simplified if expected relations between picture attributes are known. As an example, we consider a chessboard-like pattern as input for a picture restoring system. There are several *a priori* qualities present in such a pattern: the intensity in a square is constant, at a square's border are straight edges in a vertical or a horizontal direction, and edges are continuous. A system of feature spins instructed with this knowledge can restore noisy chessboard patterns. Such system entails three kinds of feature spins

- Intensity spins which take the values $\pm 1$ for black and white colors, respectively.

- Horizontal edge spins which take the values $+1$ for intensity changes from white to black, the value $-1$ from black to white and the value $0$ in the case of an absence of any edges.

- Vertical edge spins which follow corresponding rules.

The intensity spins are defined on a square lattice. The edge spins are located between neighboring intensity spins.

The Hamiltonian of the picture restoring spin system can be written

$$E_{total} = \begin{cases} E_{i-ifield} + E_{h-hfield} + E_{v-vfield} \\ + E_{i-i} + E_{h-h} + E_{v-v} + E_{i-h} + E_{i-v} \end{cases} \qquad (4.5)$$

where the indices $i,h,v$ refer to intensity, horizontal and vertical edge spins, respectively, with the corresponding fields *ifield, hfield, vfield*. To implement the continuity condition for the intensity spins $I_{i,j} \in \{-1, +1\}$ an Ising-like interaction is assumed. To implement the continuity property of edges, an attractive interaction in the proper direction is employed for the horizontal edge spins $H_{i,j} \in \{-1, 0, +1\}$ as well as for the vertical edge spins $V_{i,j} \in \{-1, 0, +1\}$.

$$E_{i-i} = -J_i \sum_{<(i,j),(k,l)>} I_{i,j} I_{k,l} \qquad (4.6)$$

**Table 1.** Compatibility condition table $T(i,j,k)$ where $i,j$ denote pairs of intensity spins and $k$ denotes the edge spin in between.

| i | j | k | T |
|---|---|---|---|
| 1 | 1 | 0 | 1 |
| −1 | −1 | 0 | 1 |
| 1 | −1 | −1 | 1 |
| −1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 0 |
| −1 | −1 | 1 | 0 |
| 1 | 1 | −1 | 0 |
| −1 | −1 | −1 | 0 |
| 1 | −1 | 0 | 0 |
| −1 | 1 | 0 | 0 |
| 1 | −1 | 1 | 0 |
| −1 | 1 | 1 | 0 |
| 1 | −1 | −1 | 0 |
| −1 | 1 | −1 | 0 |

$$E_{h-h} = -J_h \sum_{(i,j)} \delta(H_{i,j+1}, H_{i,j}) \, ;$$

$$E_{v-v} = -J_v \sum_{(i,j)} \delta(V_{i+1,j}, V_{i,j}) \, . \qquad (4.7)$$

To obtain compatibility between the hypotheses of edge spins and intensity spins an interaction energy favoring consistent configurations is added
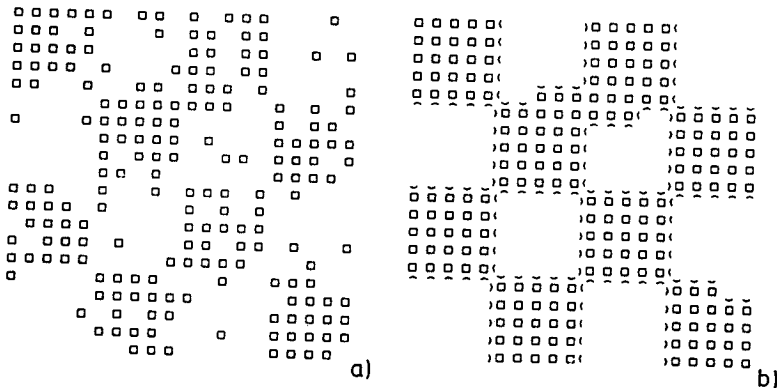


**Figure 21.** Input pattern and restored chessboard pattern for the feature spin system described by Eqs. 4.5 through 4.10.

$$E_{i-h} = -J_{i-h} \sum_{(i,j)} T(I_{i,j}, I_{i+1,j}, H_{i,j}) ;$$

$$E_{i-v} = -J_{i-v} \sum_{(i,j)} T(I_{i,j}, I_{i,j+1}, V_{i,j}) . \qquad (4.8)$$

Here T contains the compatibility condition listed in Table I.

The pattern to be processed is coded as a field $F_{i,j} \in \{ +1, -1 \}$ corresponding to black and white pixels at position $(i, j)$. The interaction between intensity spins and the field is chosen like in the Ising model

$$E_{i-ifield} = -J_{ifield} \sum_{(i,j)} I_{i,j} F_{i,j}. \qquad (4.9)$$

The field for edge spins codes intensity changes. The interaction between edge spins and the corresponding fields is

$$E_{h-hfield} = -J_{hfield} \sum_{(i,j)} \delta(H_{i,j}, (F_{i+1,j}-F_{i,j})/2) ;$$

$$E_{v-vfield} = -J_{vfield} \sum_{(i,j)} \delta(V_{i,j}, (F_{i,j+1}-F_{i,j})/2) . \qquad (4.10)$$

As input for the picture restoration spin system we chose a distorted chessboard pattern, measuring 20×20 pixels. Twenty percent of the pixels were randomly reversed from black to white and vice versa. The resulting pattern is presented in Fig. 21a. The aim of the restoration process is to find the chessboard closest to this picture. The interaction parameters assumed for the restoration are $J_i = 1$, $J_h = J_v = 4.5$, $J_{i-h} = J_{i-v} = 4.5$, $J_{ifield} = 1$, $J_{hfield} = J_{vfield} = 3$. The temperature was lowered in 12 steps from an initial value of $T = 8$ to the final value of $T = 0.05$. Figure 21b shows the result: the chessboard pattern has been restored to a very large degree.

## Acknowledgments

## References

1. G. Palm, A. Aertsen, Brain Theory, Springer, Berlin-Heidelberg, New York (1986); H. Gardner, The Mind's New Science, Basic Books, New York (1985).

2.  C. von der Malsburg, Self-Organization of Orientation Sensitive Cells in the Striate Cortex, Kybernetik, **14**, 85 (1973).

3.  L.N. Cooper, Proceedings of the Nobel Symposium on Collective Properties of Physical Systems, ed. Lundqvist, Academic, New York, p. 252 (1973).

4.  T. Kohonen, Self-Organization and Associative Memory, Springer, Berlin-Heidelberg, New York, p. 128 (1984).

5.  J.J. Hopfield, Neural Networks and Physical Systems with Emergent Collective Computational Abilities, Proc. Natl. Acad. Sci. USA, **79**, 2554 (1982).

6.  J.J. Hopfield, Neural Networks for Computing, American Institute of Physics Publication, in press.

7.  J. Buhmann, and K. Schulten, Associative Recognition and Storage in a Model Network of Physiological Neurons, Biol. Cybern., **54**, 319 (1986).

8.  C. von der Malsburg, The Correlation Theory of Brain Function, Internal Report 81/2, Dept. Neurobiologie, MPI f. Biophysikalische Chemie, Göttingen (1981).

9.  M. Abeles, Local Cortical Circuits, Springer, Berlin-Heidelberg, New York (1982).

10. T. Kohonen, Self-Organized Formation of Topologically Correct Feature Maps, Biol. Cybern., **43**, 59 (1982).

11. T. Kohonen, Analysis of a Simple Self-Organizing Process, Biol. Cybern., **44**, 135 (1982).

12. H. Ritter and K. Schulten, On the Stationary State of Kohonen's Self-Organizing Sensory Mapping, Biol. Cybern., **54**, 99 (1986).

13. J. Fox, The Brain's Dynamic Way of Keeping in Touch, Science, **225**, 820 (1984).

14. J.H. Kaas, M.M. Merzenich and H.P. Killackey, The reorganization of somatosensory cortex following peripheral nerve damage in adult and developing mammals, Ann. Rev. Neuroscience, **6**, 325 (1983).

15. E.R. Kandel and J.H. Schwartz, Principles of Neural Science, Elsevier, New York (1985).

16. B.E. Stein, H.P. Clamann and S.J. Goldberg, Superior Colliculus: Control of Eye Movements in Neonatal Kittens, Nature, **??**, ?? (1980).

17. H. Ritter and K. Schulten, Topology Conserving Mappings for Learning Motor Tasks, Proceedings of the 1st Conference on Neural Networks and Computing, Utah, in press (1986).

18. H. Ritter and K. Schulten, Planning a Dynamic Trajectory as Path Finding in Phase Space, Wopplot Proceedings, Munich, submitted 1986.

19. S. Kirkpatrick, C.D. Gelatt and M.P. Vecchi, Optimization by Simulated Annealing, Science, **220**, 671 (1983).

20. S. Geman and D. Geman, IEEE Tran. Pattern Anal. Machine Intell., Vol. PAMI-6, 721 (1984).

21. P. Kienker, T.J. Sejnowski, G.E. Hinton and L.E. Schuhmacher, Separating Figure from Ground with a Parallel Network, in press (1986).

22. D. Marr, Vision, Freeman (1982).

23. B. Julesz, Foundations of Cyclopean Perception, Univ. Chicago Press (1971).