

## Slowly evolving connectivity in recurrent neural networks: I. The extreme dilution regime

This article has been downloaded from IOPscience. Please scroll down to see the full text article.

2004 J. Phys. A: Math. Gen. 37 7653

(<http://iopscience.iop.org/0305-4470/37/31/002>)

View [the table of contents for this issue](#), or go to the [journal homepage](#) for more

Download details:

IP Address: 171.66.16.91

The article was downloaded on 02/06/2010 at 18:29

Please note that [terms and conditions apply](#).

# Slowly evolving connectivity in recurrent neural networks: I. The extreme dilution regime

B Wemmenhove<sup>1</sup>, N S Skantzos<sup>2</sup> and A C C Coolen<sup>3</sup>

<sup>1</sup> Institute for Theoretical Physics, University of Amsterdam, Valckenierstraat 65, 1018 XE Amsterdam, The Netherlands

<sup>2</sup> Departament de Física Fonamental, Facultat de Física, Universitat de Barcelona, 08028 Barcelona, Spain

<sup>3</sup> Department of Mathematics, King's College London, The Strand, London WC2R 2LS, UK

E-mail: wemmenho@science.uva.nl, nikos@ffn.uv.es and tcoolen@math.kcl.ac.uk

Received 24 March 2004

Published 21 July 2004

Online at [stacks.iop.org/JPhysA/37/7653](http://stacks.iop.org/JPhysA/37/7653)

doi:10.1088/0305-4470/37/31/002

## Abstract

We study extremely diluted spin models of neural networks in which the connectivity evolves in time, although adiabatically slowly compared to the neurons, according to stochastic equations which on average aim to reduce frustration. The (fast) neurons and (slow) connectivity variables equilibrate separately, but at different temperatures. Our model is exactly solvable in equilibrium. We obtain phase diagrams upon making the condensed ansatz (i.e. recall of one pattern). These show that, as the connectivity temperature is lowered, the volume of the retrieval phase diverges and the fraction of mis-aligned spins is reduced. Still one always retains a region in the retrieval phase where recall states other than the one corresponding to the 'condensed' pattern are locally stable, so the associative memory character of our model is preserved.

PACS numbers: 75.10.Nr, 05.20.-y, 64.60.Cn

## 1. Introduction

Most statistical mechanical studies of recurrent neural networks have traditionally been concerned with systems in which the dynamical variables are either the neurons (see, e.g., [1–4] or the reviews [5, 6] and references therein), or their interactions (or synapses, see, e.g., [7–10] or the reviews [12–14] and references therein). The first type of processes describe network operation, whereas the second correspond to learning. These areas have by now been investigated quite extensively. In contrast, only a modest number of studies involved coupled dynamical laws for both neurons and interactions [15–22], to reflect the complex dynamical interplay between synapses and neurons found in the real brain. The approach usually adopted

in the latter studies, to obtain analytically solvable models, is the introduction of a hierarchy of adiabatically separated time scales, such that the fast variables (taken to be the neurons) are in equilibrium on the time scales where the slow variables (the interactions, taken to be symmetric) evolve. One can also introduce further levels in the hierarchy by introducing different classes of interactions, each evolving on different characteristic time scales [22]. The resulting formalism involves nested replica theories, with Parisi matrices [23] in which the number of blocks at each level is the ratio of temperatures of subsequent levels in the hierarchy of equilibrating degrees of freedom. Such models can also serve to derive Parisi's replica symmetry breaking (RSB) scheme [24]. In neural network studies, the dynamics of the interactions have usually been governed by Langevin equations in which the deterministic forces are proportional to expectation values of neuronal pair correlations (with the neuron state statistics corresponding to Boltzmann equilibrium, given the instantaneous values of the interactions), potentially biased to reflect the possibility of recall of a pattern. In [18, 19], the interactions were taken to evolve away from an initial state given by Hopfield's [1] interaction matrix, with an extensive number of stored patterns. There it was found that for low interaction temperatures, the network collapsed into an undesirable so-called 'super-ferromagnetic' state, whereas for negative replica dimension (corresponding to anti-Hebbian learning) the storage capacity of the network was found to be enhanced.

All papers dealing with the theory of coupled neuronal and interaction dynamics published so far assumed full connectivity: each neuron interacted with every other neuron, with the magnitude and sign of the interactions evolving in time. Here we propose and study a model of a symmetrically diluted recurrent neural network in which the *connectivity* is allowed to change slowly. On time scales where the neuron variables are in thermodynamic equilibrium, the microscopic realization of the (discrete) dilution variables (reflecting the connectivity) is allowed to evolve slowly and stochastically, driven by forces aiming at a reduction of global frustration, without however changing the actual *values* of the bonds (the latter are frozen, given by Hopfield's [1] recipe). It has been known that one may store information in recurrent neural networks solely by eliminating frustrated bonds, but this has always been done by hand (see, e.g., [25] and references therein). Here the system is allowed to adapt its connectivity autonomously. It should be emphasized that there is an important difference between having dynamic bonds with Hebbian-type dynamical laws, as in [15–17, 20–22], and the present situation of having dynamic connectivity with fixed Hebbian values for active bonds. The former definitions imply irreversible modification or even elimination of stored information, whereas in the present paper, since the *values* assigned to the active bonds are not modified, the slow adaptation is fully reversible (one can always return to random dilution) and all stored information is retained.

The scaling with the system size  $N$  chosen for the average connectivity  $c$  in the system (the average number of bonds per spin) will have a strong influence on the structure of the resulting theory. In this first paper, we consider the so-called 'extreme dilution' regime [26], defined by  $\lim_{N \rightarrow \infty} c^{-1} = \lim_{N \rightarrow \infty} c/N = 0$  (a second paper will be devoted to the finite connectivity regime, where  $c = \mathcal{O}(N^0)$  as  $N \rightarrow \infty$ ). We solve our coupled spin and connectivity dynamics model analytically, in replica-symmetric (RS) ansatz. We find that, as a result of the connectivity adaptation, the network connectivity becomes more ordered to boost retrieval of condensed patterns, as a result of which the system's retrieval phase is enhanced compared to that of the corresponding network with a quenched random connectivity matrix as studied in [26], and that the fraction of 'misaligned spins' is reduced as the temperature of the connectivity variables is lowered. Yet one still retains regions in the phase diagram where the alternative (presently non-condensed) pure retrieval states remain locally stable, so that the system continues to function as an associative memory.

## 2. Definitions

We study diluted Hopfield [1] type recurrent neural networks, with (fast) binary neurons  $\sigma_i \in \{-1, 1\}$  (denoting quiet versus firing states) and  $i = 1, \dots, N$ . The connectivity of the system is defined by connectivity variables  $c_{ij} \in \{0, 1\}$ , with  $c_{ji} = c_{ij}$  and  $c_{ii} = 0$ . Our neurons evolve according to Glauber-type local field alignment at temperature  $T = \beta^{-1}$ , with the fields defined by  $h_i = \sum_j \frac{c_{ij}}{c} \sum_{\mu=1}^p \xi_i^\mu \xi_j^\mu \sigma_j$ , i.e. with Hebbian interactions whenever  $c_{ij} = 1$  (when a bond  $(i, j)$  is present). For frozen connectivity  $\{c_{ij}\}$  our Ising spin neurons would equilibrate to a Boltzmann state characterized by the Hamiltonian

$$H_f(\boldsymbol{\sigma}, \mathbf{c}) = - \sum_{i < j} \frac{c_{ij}}{c} \sum_{\mu=1}^p \xi_i^\mu \xi_j^\mu \sigma_i \sigma_j. \quad (1)$$

Here  $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_N)$  and  $\mathbf{c} = \{c_{ij}\}$ . The  $\{\xi_i^\mu\} \in \{-1, 1\}$  with  $\mu = 1, \dots, p$  represent  $p$  fixed patterns  $\boldsymbol{\xi}^\mu = (\xi_1^\mu, \dots, \xi_N^\mu)$  to be stored and hopefully recalled. Instead of frozen, we now take our connectivity to also evolve in time, albeit on time scales much larger than those of the neuronal relaxation (so the neurons can always be assumed in equilibrium, given the instantaneous connectivity). This slow process is again taken to be of a Glauber type, but at temperature  $\tilde{T} = \tilde{\beta}^{-1}$  and with the connectivity Hamiltonian

$$H_s(\mathbf{c}) = -\frac{1}{\tilde{\beta}} \log Z_f(\mathbf{c}) + \frac{1}{\tilde{\beta}} \log \left( \frac{N}{c} \right) \sum_{i < j} c_{ij} \quad (2)$$

$$Z_f(\mathbf{c}) = \sum_{\boldsymbol{\sigma}} e^{-\beta H_f(\boldsymbol{\sigma}, \mathbf{c})}. \quad (3)$$

The motivation for choosing the latter Hamiltonian is similar to that in [16, 17], in which slow *continuous* variables (interaction strengths) were considered. Our present set-up with the effective Hamiltonian (2) constitutes a translation of this formalism to the present case of discrete slow variables, now governed by Glauber dynamics. In [16, 17], demanding the driving forces in the Langevin equation which describes the slow process to be given by neuron correlations, resulted in gradient descent on an effective Hamiltonian for the slow variables which is similar to (2). In this latter effective Hamiltonian, the first term (equal to the free energy of the fast Hamiltonian) produces neuron correlations with respect to the Boltzmann measure of the fast system. The second term in (2) acts as a chemical potential, ensuring an average number of  $c$  connections per neuron. The pre-factor  $1/\tilde{\beta}$  will be found helpful later.

The properties of our system at the largest time scales, where also the connectivity has equilibrated, are characterized by the partition sum of the slow variables:

$$Z_s = \sum_{\mathbf{c}} e^{-\tilde{\beta} H_s(\mathbf{c})} = \sum_{\mathbf{c}} [Z_f(\mathbf{c})]^{\tilde{\beta}/\beta} \exp \left( -\log \left( \frac{N}{c} \right) \sum_{i < j} c_{ij} \right). \quad (4)$$

This sum is interpreted as describing  $n = \tilde{\beta}/\beta$  replicated copies of the fast system, leading to a replica theory with finite replica dimension  $n$ . Minimization of  $H_s(\mathbf{c})$  should give a ‘smart’ arrangement of the connectivity  $\{c_{ij}\}$ , tailored to the realization of the patterns, but constrained to give an average connectivity  $c$ . In the remainder of this paper we calculate phase diagrams and the fraction of mis-aligned spins. Phases are characterized by the values of the replicated overlap and spin-glass-order parameters

$$m_\alpha^\mu = \lim_{N \rightarrow \infty} N^{-1} \sum_i \overline{\xi_i^\mu \sigma_i^\alpha} \quad q_{\alpha\beta} = \lim_{N \rightarrow \infty} N^{-1} \sum_i \overline{\sigma_i^\alpha \sigma_i^\beta}. \quad (5)$$

Here  $\overline{\{\dots\}}$  denotes averaging over all geometries  $\{c_{ij}\}$  and all spin-configurations  $\sigma^\alpha$  in each of the replicas  $\alpha = 1, \dots, n$ , with the Boltzmann measure associated with (4):

$$\overline{G(\{\sigma^\alpha\}, \mathbf{c})} = Z_s^{-1} \sum_{\{\sigma^\alpha\}} \sum_{\mathbf{c}} G(\{\sigma^\alpha\}, \mathbf{c}) \times \exp \left\{ \frac{\beta}{c} \sum_{i < j} c_{ij} \sum_{\mu} \xi_i^\mu \xi_j^\mu \sum_{\alpha=1}^n \sigma_i^\alpha \sigma_j^\alpha - \log \left( \frac{N}{c} \right) \sum_{i < j} c_{ij} \right\}. \quad (6)$$

### 3. Equilibrium analysis

#### 3.1. Calculation of the RS free energy

The thermodynamic properties of the stationary state, with equilibrated connectivity, are derived from the asymptotic free energy per spin  $f = -\lim_{N \rightarrow \infty} (\beta N)^{-1} \log Z_s$ . Upon performing the trace over all geometries in (4) one obtains, with  $\sigma_i = (\sigma_i^1, \dots, \sigma_i^n)$ :

$$Z_s = \sum_{\sigma^1 \dots \sigma^n} \prod_{i < j} \left[ 1 + \frac{c}{N} \exp \left( \frac{\beta}{c} (\xi_i \cdot \xi_j) (\sigma_i \cdot \sigma_j) \right) \right]. \quad (7)$$

In evaluating the free energy we make the usual ‘condensed’ ansatz: only a finite number  $r$  of patterns will be structurally correlated with the system state. The remaining  $\alpha c - r$  patterns can be treated as frozen disorder, over which the free energy may be averaged. For the result we write  $[f]_{\text{dis}}$ . In this paper, we work within the connectivity scaling regime of extreme dilution, where  $\lim_{N \rightarrow \infty} c/N = \lim_{N \rightarrow \infty} c^{-1} = 0$ . Now one obtains

$$-\tilde{\beta} [f]_{\text{dis}} = \lim_{N \rightarrow \infty} \frac{1}{N} \log \sum_{\sigma^1 \dots \sigma^n} \times \exp \left( \frac{\beta}{2N} \sum_{ij} (\sigma_i \cdot \sigma_j) \sum_{\mu \leq r} \xi_i^\mu \xi_j^\mu + \frac{\alpha \beta^2}{4N} \sum_{ij} (\sigma_i \cdot \sigma_j)^2 + \mathcal{O} \left( \frac{N}{c} \right) \right) \quad (8)$$

(modulo irrelevant additive constants). We define the familiar pattern and state overlaps  $m_{\alpha\mu}(\sigma) = N^{-1} \sum_i \xi_i^\alpha \sigma_i$  and  $q_{\alpha\beta}(\{\sigma\}) = N^{-1} \sum_i \sigma_i^\alpha \sigma_i^\beta$ . They are introduced via appropriate  $\delta$ -distributions, so that the spin traces can be carried out. This results in the usual type of steepest descent expression for  $[f]_{\text{dis}}$  (again modulo constants):

$$[f]_{\text{dis}} = \text{extr}_{\{m_{\alpha\mu}, q_{\alpha\beta}\}} \left\{ \frac{\alpha\beta}{4n} \sum_{\alpha \neq \beta} q_{\alpha\beta}^2 + \frac{1}{2n} \sum_{\alpha} \sum_{\mu \leq r} m_{\alpha\mu}^2 - \frac{1}{n\beta} \left\langle \log \sum_{\sigma^1 \dots \sigma^n} \exp \left( \beta \sum_{\alpha} \sum_{\mu \leq r} m_{\alpha\mu} \xi_\mu \sigma_\alpha + \frac{1}{2} \alpha \beta^2 \sum_{\alpha \neq \beta} q_{\alpha\beta} \sigma_\alpha \sigma_\beta \right) \right\rangle_{\xi} \right\} \quad (9)$$

where  $\langle g(\xi) \rangle_{\xi} = 2^{-r} \sum_{\xi \in \{-1, 1\}^r} g(\xi)$ . The parameter  $c$  represents the ensemble averaged connectivity. This follows upon adding suitable generating fields to the slow Hamiltonian:  $H_s(\mathbf{c}) \rightarrow H_s(\mathbf{c}) + \frac{2\lambda}{c} \sum_{i < j} c_{ij}$ . Repeating the above calculation with the added fields shows that  $\lim_{N \rightarrow \infty} \frac{2}{Nc} \sum_{i < j} \overline{c_{ij}} = \lim_{\lambda \rightarrow 0} \frac{\partial [f]_{\text{dis}}}{\partial \lambda} = 1$ , which proves our claim.

We next make the replica-symmetric ansatz for the physical saddle-point:  $m_{\alpha\mu} = m_{\mu}$  for all  $(\alpha, \mu)$  and  $q_{\alpha\beta} = q + \delta_{\alpha\beta}(1 - q)$  for all  $(\alpha, \beta)$ , keeping in mind that the replica dimension

$n$  can take any non-negative value:

$$[f]_{\text{dis}}^{\text{RS}} = \text{extr}_{\{m_\mu, q\}} \left\{ \frac{1}{2} \sum_{\mu \leq r} m_\mu^2 + \frac{1}{4} \alpha \beta [2q + (n-1)q^2] - \frac{1}{n\beta} \left\langle \log \int \text{D}z \cosh^n \left[ \beta \left( \sum_{\mu \leq r} m_\mu \xi_\mu + z \sqrt{\alpha q} \right) \right] \right\rangle_\xi \right\}. \quad (10)$$

Variation of  $\{m_\mu, q\}$  gives the saddle-point equations for our RS order parameters, with the shorthand  $\Xi = \beta \left( \sum_{\mu \leq r} m_\mu \xi_\mu + z \sqrt{\alpha q} \right)$ , which are of the familiar form

$$m_\mu = \left\langle \xi_\mu \frac{\int \text{D}z \tanh(\Xi) \cosh^n(\Xi)}{\int \text{D}z \cosh^n(\Xi)} \right\rangle_\xi \quad (11)$$

$$q = \left\langle \frac{\int \text{D}z \tanh^2(\Xi) \cosh^n(\Xi)}{\int \text{D}z \cosh^n(\Xi)} \right\rangle_\xi. \quad (12)$$

The physical meaning of the RS order parameters is  $m_\mu = \lim_{N \rightarrow \infty} N^{-1} \sum_i \overline{\xi_i^\mu \sigma_i}$  and  $q = \lim_{N \rightarrow \infty} N^{-1} \sum_i \overline{\sigma_i^2}$ , as usual.

### 3.2. Phase transitions and phase diagrams

If one simplifies matters further by assuming only one pattern to be condensed, i.e.  $m_\mu = m \delta_{\mu,1}$ , then equations (11) and (12) reduce to

$$m = \frac{\int \text{D}z \tanh[\beta(m + z \sqrt{\alpha q})] \cosh^n[\beta(m + z \sqrt{\alpha q})]}{\int \text{D}z \cosh^n[\beta(m + z \sqrt{\alpha q})]} \quad (13)$$

$$q = \frac{\int \text{D}z \tanh^2[\beta(m + z \sqrt{\alpha q})] \cosh^n[\beta(m + z \sqrt{\alpha q})]}{\int \text{D}z \cosh^n[\beta(m + z \sqrt{\alpha q})]}. \quad (14)$$

These are recognized to be identical to those of the finite  $n$  model studied in [27] if we re-define the parameters in the latter according to

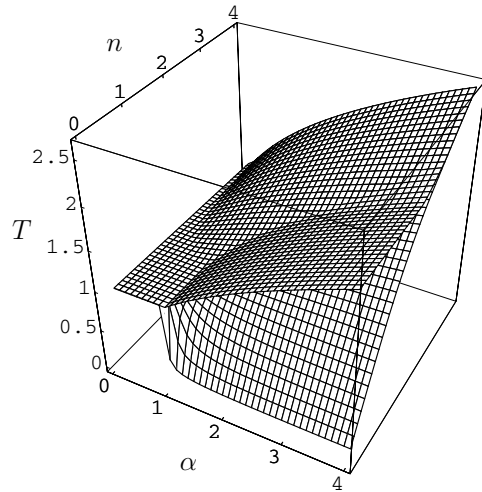
$$J^{(1)} m/k \rightarrow m \quad J^{(2)} q/k \rightarrow \alpha \beta q. \quad (15)$$

This makes sense, since the  $n \rightarrow 0$  limit of our present model (i.e. the symmetrically extremely diluted Hopfield model with quenched random connectivity [26]) is known to map onto the  $n \rightarrow 0$  limit of [27], (i.e. the SK model [28]). Clearly one finds simplified equations for the special dimension values  $n = 1$  (equivalent to having annealed connectivity) and  $n = 2$ , where the Gaussian integrals can be done. For instance, at  $n = 1$  the equation for  $m$  reduces to  $m = \tanh(\beta m)$ , whereas for  $n = 2$  one finds

$$m = \frac{\sinh(2\beta m)}{\cosh(2\beta m) + e^{-2\alpha\beta^2 q}} \quad q = \frac{\cosh(2\beta m) - e^{-2\alpha\beta^2 q}}{\cosh(2\beta m) + e^{-2\alpha\beta^2 q}}. \quad (16)$$

Our RS equations admit three phases: a paramagnetic phase (P) with  $m = q = 0$ , a recall phase (R) where  $m \neq 0$  and  $q > 0$ , and a spin-glass phase (SG) where  $m = 0$  but  $q > 0$ . Since deriving the RS phase transitions has been reduced to appropriate translation of the results found in [27], we will here simply mention the outcome:

- For sufficiently small  $\alpha$  one will find a P  $\rightarrow$  R transition at a finite temperature. For  $\alpha < \alpha_c = \frac{1}{3n-2}$  this transition is second order, and occurs at  $T_R = 1$ .



**Figure 1.** RS phase diagram in the space of control parameters. We show the critical temperature(s) as surface(s) over the  $(\alpha, n)$  plane. The high-temperature phase is paramagnetic (P). At low temperature, we find the retrieval phase R. For  $(\alpha, n)$  values with two critical temperatures, the latter define the boundaries of a spin-glass phase SG. The P  $\rightarrow$  R transitions are second order for  $\alpha < \frac{1}{3n-2}$ , and first order elsewhere. The P  $\rightarrow$  SG transitions are second order for  $n < 2$  and first order elsewhere. For large  $\alpha$  the SG  $\rightarrow$  R transitions become second order, but for small  $\alpha$  they are first order.

- For larger  $\alpha$ , lowering the temperature will lead first to a P  $\rightarrow$  SG transition, followed at some yet lower temperature by a SG  $\rightarrow$  R transition<sup>4</sup>. For  $n \leq 2$  the P  $\rightarrow$  SG transition is second order, and occurs for  $T_{SG} = \sqrt{\alpha}$ .
- The SG  $\rightarrow$  R transition is second order for  $\alpha \rightarrow \infty$ , where its transition temperature tends to  $T_c = n$ , but first order for sufficiently small  $\alpha$ .
- The effects of increasing the replica dimension  $n$  are (i) a reduction of the size in the phase diagram of the SG phase, and (ii) a change of the orders of the P  $\rightarrow$  R and P  $\rightarrow$  SG transitions, from second order (for small  $n$ ) to first order (for large  $n$ ).

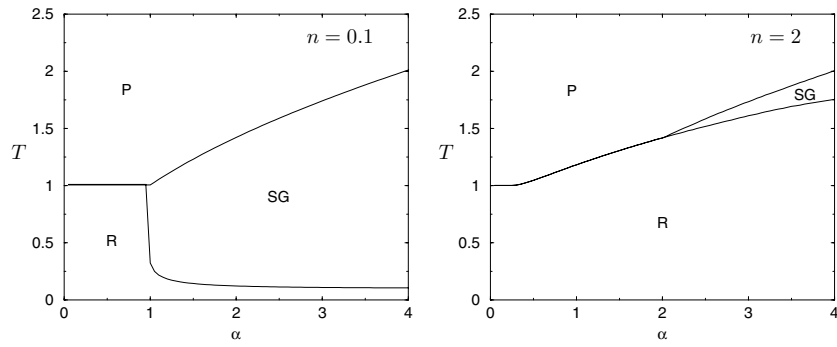
Numerical solution of equations (13) and (14) leads to the RS phase diagram drawn in figure 1. Figure 2 shows intersections of this diagram in the planes of constant replica dimension  $n = 0.1$  and  $n = 2$ . All transitions discussed and drawn above refer to bifurcations of locally stable solutions, since for recurrent neural networks the time scales where thermodynamic stability would become an issue are in practice never reached.

We finally turn to replica symmetry breaking. The location in our phase diagram of second-order RSB phase transition follows upon the inspection of the eigenvalues of the Hessian. Since our model can be mapped onto the nonzero- $n$  SK-model [27], we can read off the eigenvalues from [29]. The dangerous eigenvalue  $\lambda_{RSB}$  is the one associated with the so-called replicon mode:

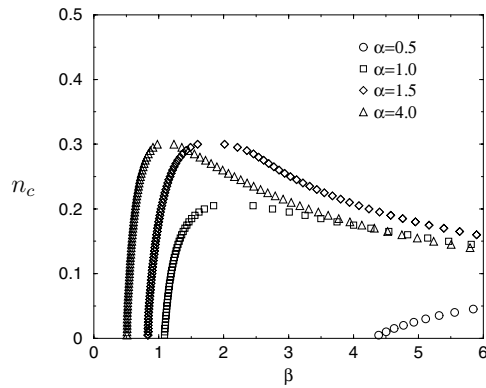
$$\lambda_{RSB} = \alpha\beta^2[1 - \alpha\beta^2[1 - 2q + h(m, q)]] \quad (17)$$

$$h(m, q) = \frac{\int Dz \tanh^4[\beta(m + z\sqrt{\alpha q})] \cosh^n[\beta(m + z\sqrt{\alpha q})]}{\int Dz \cosh^n[\beta(m + z\sqrt{\alpha q})]}. \quad (18)$$

<sup>4</sup> For small values of  $n$ , the latter SG  $\rightarrow$  R transition is expected to disappear when replica symmetry breaking is taken into account.



**Figure 2.** Intersections of the phase diagram shown in figure 1, at  $n = 0.1$  (left) and  $n = 2$  (right). We have paramagnetic (P), recall (R) and spin-glass (SG) phases. We note that there is no critical value for  $\alpha$  above which recall is no longer possible. Instead, the  $SG \rightarrow R$  transition line will approach the line  $T = n$  for large  $\alpha$ . Since RSB phenomena appear to be confined to  $n < 1$  (see below), this is not an artefact of the RS assumption. For large  $n$ , all phase transitions ultimately become first order.



**Figure 3.** Location of the AT instability  $n_c$ , shown as a function of the inverse temperature  $\beta = T^{-1}$ , and for a number of different storage ratios  $\alpha$ . Replica symmetry breaking is seen to be limited to the values of  $n$  below 0.32. We also note the non-monotonic dependence on the temperature of the critical dimension  $n_c$  for fixed  $\alpha$ .

Replica symmetry is stable only if  $\lambda_{RSB} > 0$ . For each combination  $(\alpha, T)$ , one finds a critical value  $n_c(\alpha, T)$  (the AT line) below which replica symmetry is unstable. Examples of the behaviour of  $n_c(\alpha, T)$  are shown in figure 3. Replica symmetry breaking is found to occur only for  $n < 0.32$ .

Compared to diluted neural network models with static random connectivity, the main effect of introducing dynamic connectivity (with the present Glauber dynamics, aimed at reducing frustration) is to reduce the spin-glass phase in favour of the recall phase. The connectivity adjusts itself autonomously in order to retrieve the condensed pattern optimally, to such an extent that for sufficiently low temperature there is no upper limit on the storage ratio (provided we do not leave the ‘extreme dilution’ scaling regime  $\lim_{N \rightarrow \infty} c/N = \lim_{N \rightarrow \infty} c^{-1} = 0$ ). This then raises the question of whether the other (non-condensed) patterns can be retrieved at all after the connectivity has been tailored to the recall of one specific condensed pattern. This is investigated in section 5.



#### 4. Fraction of misaligned spins

We expect the observed improvement of retrieval performance due to the slow connectivity dynamics to be reflected in a reduction with increasing  $n$  of the fraction of frustrated bonds in the system. To verify this we calculate a different but similar quantity: the fraction  $\phi$  of misaligned spins, i.e. those where  $\sigma_i$  and local field  $h_i$  have opposite sign:

$$\phi = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_i \left[ \left\langle \theta \left[ -\sigma_i \sum_j \frac{c_{ij}}{c} \sum_{\mu} \xi_i^{\mu} \xi_j^{\mu} \sigma_j \right] \right\rangle_{\text{dis}} \right]. \quad (19)$$

To calculate this object one could introduce further replicas, but here we follow an alternative route: we solve our model first for finite  $c$ , in which case joint replicated spin–field distributions (in terms of which  $\phi$  can be expressed) become the natural order parameters, followed by taking the limit  $c \rightarrow \infty$ .

##### 4.1. Calculation of the joint spin–field distribution

To do so we have to adapt and generalize the calculation in [30] by first introducing the  $2^p$  so-called sub-lattices [31], with  $\xi_i = (\xi_i^1, \dots, \xi_i^p)$ :

$$I_{\xi} = \{i | \xi_i = \xi\} \quad p_{\xi} = |I_{\xi}|/N. \quad (20)$$

Since  $c$  is assumed finite, so is  $p = \alpha c$ . We write  $\sum_{\xi} p_{\xi} \Phi(\xi) = \langle \Phi(\xi) \rangle_{\xi}$ ; for randomly drawn patterns  $\lim_{N \rightarrow \infty} p_{\xi} = 2^{-p}$ . For finite  $c$  our analysis will start to resemble that in [32]. In each sublattice, we may define a joint distribution for replicated spins and fields, and (with a modest amount of foresight) conjugate fields:

$$P_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) = |I_{\xi}|^{-1} \sum_{i \in I_{\xi}} \delta_{\sigma, \sigma_i} \delta[\mathbf{h} - \mathbf{h}_i(\{\sigma\})] \delta[\hat{\mathbf{h}} - \hat{\mathbf{h}}_i(\{\sigma\})] \quad (21)$$

where  $\sigma \in \{-1, 1\}^n$ ,  $\mathbf{h}, \hat{\mathbf{h}} \in \mathbb{R}^n$  and  $h_i^{\alpha}(\{\sigma\}) = \sum_j \frac{c_{ij}}{c} (\xi_i \cdot \xi_j) \sigma_j^{\alpha}$ . In evaluating the free energy per spin we write the fast Hamiltonian in terms of replicated fields and introduce (21) by inserting suitable integrals over  $\delta$ -functions. This is done first only for discrete values of  $\mathbf{h}$ , with the continuum limit (converting integrals into path integrals) to be taken after the thermodynamic limit. We abbreviate  $\{dP d\hat{P}\} = \prod_{\xi, \sigma, \mathbf{h}, \hat{\mathbf{h}}} [dP_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) d\hat{P}_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}})]$  and find

$$\begin{aligned} -\tilde{\beta} f &= \lim_{N \rightarrow \infty} \frac{1}{N} \log \sum_{\mathbf{c}} \sum_{\sigma^1 \dots \sigma^n} \exp \left( \frac{\beta}{2} \sum_i \sigma_i \cdot \mathbf{h}_i(\{\sigma\}) - \log \left( \frac{N}{c} \right) \sum_{i < j} c_{ij} \right) \\ &= \lim_{N \rightarrow \infty} \frac{1}{N} \log \int \{dP d\hat{P}\} \exp \left( N \left\langle \sum_{\sigma} \int d\mathbf{h} d\hat{\mathbf{h}} P_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) \right. \right. \\ &\quad \left. \left. \times \left[ i \hat{P}_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) + \frac{1}{2} \beta (\sigma \cdot \mathbf{h}) \right] \right\rangle_{\xi} \right) \int \prod_i \left[ \frac{d\mathbf{h}_i d\hat{\mathbf{h}}_i}{(2\pi)^n} e^{i\hat{\mathbf{h}}_i \cdot \mathbf{h}_i} \right] \\ &\quad \times \sum_{\sigma^1 \dots \sigma^n} \exp \left( -i \sum_{\xi} \sum_{i \in I_{\xi}} \hat{P}_{\xi}(\sigma_i, \mathbf{h}_i, \hat{\mathbf{h}}_i) \right) \\ &\quad \times \prod_{i < j} \left[ 1 + \frac{c}{N} \exp \left( -\frac{i}{c} (\xi_i \cdot \xi_j) [(\hat{\mathbf{h}}_i \cdot \sigma_j) + (\hat{\mathbf{h}}_j \cdot \sigma_i)] \right) \right] \end{aligned}$$

$$\begin{aligned}
&= \lim_{N \rightarrow \infty} \frac{1}{N} \log \int \{dP d\hat{P}\} \exp \left( N \left\langle \sum_{\sigma} \int d\mathbf{h} d\hat{\mathbf{h}} P_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) \right. \right. \\
&\quad \times \left. \left. \left[ i\hat{P}_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) + \frac{1}{2}\beta(\sigma \cdot \mathbf{h}) \right] \right\rangle_{\xi} \right) \\
&\quad \times \exp \left( \frac{c}{2} N \left\langle \left\langle \sum_{\sigma\sigma'} \int d\mathbf{h} d\mathbf{h}' d\hat{\mathbf{h}} d\hat{\mathbf{h}}' P_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) P_{\xi'}(\sigma', \mathbf{h}', \hat{\mathbf{h}}') \right. \right. \right. \\
&\quad \times \left. \left. \left. \exp \left( -\frac{i}{c}(\xi \cdot \xi')[(\hat{\mathbf{h}} \cdot \sigma') + (\hat{\mathbf{h}}' \cdot \sigma)] \right) \right\rangle_{\xi'} \right\rangle_{\xi} \right) \\
&\quad \times \int \prod_i \left[ \frac{d\mathbf{h}_i d\hat{\mathbf{h}}_i}{(2\pi)^n} \exp(i\hat{\mathbf{h}}_i \cdot \mathbf{h}_i) \right] \sum_{\sigma^1 \dots \sigma^n} \exp \left( -i \sum_{\xi} \sum_{i \in I_{\xi}} \hat{P}_{\xi}(\sigma_i, \mathbf{h}_i, \hat{\mathbf{h}}_i) \right) \\
&= \lim_{N \rightarrow \infty} \frac{1}{N} \log \int \{dP d\hat{P}\} \exp \left( N \left\langle \sum_{\sigma} \int d\mathbf{h} d\hat{\mathbf{h}} P_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) \right. \right. \\
&\quad \times \left. \left. \left[ i\hat{P}_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) + \frac{1}{2}\beta(\sigma \cdot \mathbf{h}) \right] \right\rangle_{\xi} \right) \\
&\quad \times \exp \left( \frac{c}{2} N \left\langle \left\langle \sum_{\sigma\sigma'} \int d\mathbf{h} d\mathbf{h}' d\hat{\mathbf{h}} d\hat{\mathbf{h}}' P_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) P_{\xi'}(\sigma', \mathbf{h}', \hat{\mathbf{h}}') \right. \right. \right. \\
&\quad \times \left. \left. \left. \exp \left( -\frac{i}{c}(\xi \cdot \xi')[(\hat{\mathbf{h}} \cdot \sigma') + (\hat{\mathbf{h}}' \cdot \sigma)] \right) \right\rangle_{\xi'} \right\rangle_{\xi} \right) \\
&\quad \times \exp \left( N \left\langle \log \left\{ \int \left[ \frac{d\mathbf{h} d\hat{\mathbf{h}}}{(2\pi)^n} \exp(i\hat{\mathbf{h}} \cdot \mathbf{h}) \right] \sum_{\sigma \in \{-1,1\}^n} \exp(-i\hat{P}_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}})) \right\} \right\rangle_{\xi} \right) \\
&= \text{extr}_{\{P, \hat{P}\}} \left\{ \left\langle \sum_{\sigma} \int d\mathbf{h} d\hat{\mathbf{h}} P_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) \left[ i\hat{P}_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) + \frac{1}{2}\beta(\sigma \cdot \mathbf{h}) \right] \right\rangle_{\xi} \right. \\
&\quad + \frac{1}{2}c \left\langle \left\langle \sum_{\sigma\sigma'} \int d\mathbf{h} d\mathbf{h}' d\hat{\mathbf{h}} d\hat{\mathbf{h}}' P_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}}) P_{\xi'}(\sigma', \mathbf{h}', \hat{\mathbf{h}}') \right. \right. \\
&\quad \times \left. \left. \left. \exp \left( -\frac{i}{c}(\xi \cdot \xi')[(\hat{\mathbf{h}} \cdot \sigma') + (\hat{\mathbf{h}}' \cdot \sigma)] \right) \right\rangle_{\xi'} \right\rangle_{\xi} \right. \\
&\quad \left. + \left\langle \log \left\{ \int \left[ \frac{d\mathbf{h} d\hat{\mathbf{h}}}{(2\pi)^n} \exp(i\hat{\mathbf{h}} \cdot \mathbf{h}) \right] \sum_{\sigma \in \{-1,1\}^n} \exp(-i\hat{P}_{\xi}(\sigma, \mathbf{h}, \hat{\mathbf{h}})) \right\} \right\rangle_{\xi} \right\}. \quad (22)
\end{aligned}$$

Extremization with respect to  $P_\xi(\sigma, \mathbf{h}, \hat{\mathbf{h}})$  and  $\hat{P}_\xi(\sigma, \mathbf{h}, \hat{\mathbf{h}})$  gives the following two saddle-point equations:

$$\hat{P}_\xi(\sigma, \mathbf{h}, \hat{\mathbf{h}}) = ic \left\langle \sum_{\sigma'} \int d\hat{\mathbf{h}}' P_{\xi'}(\sigma', \hat{\mathbf{h}}') \times \exp \left( -\frac{i}{c} (\boldsymbol{\xi} \cdot \boldsymbol{\xi}') [(\hat{\mathbf{h}} \cdot \boldsymbol{\sigma}') + (\hat{\mathbf{h}}' \cdot \boldsymbol{\sigma})] \right) \right\rangle_{\xi'} + \frac{1}{2} i\beta(\boldsymbol{\sigma} \cdot \mathbf{h}) \quad (23)$$

$$P_\xi(\sigma, \mathbf{h}, \hat{\mathbf{h}}) = \frac{\exp(i\hat{\mathbf{h}} \cdot \mathbf{h} - i\hat{P}_\xi(\sigma, \mathbf{h}, \hat{\mathbf{h}}))}{\int d\mathbf{h}' d\hat{\mathbf{h}}' \exp(i\hat{\mathbf{h}}' \cdot \mathbf{h}') \sum_{\sigma' \in \{-1, 1\}^n} \exp(-i\hat{P}_\xi(\sigma', \mathbf{h}', \hat{\mathbf{h}}'))}. \quad (24)$$

Insertion of (23) into (24) gives a saddle-point equation in terms of  $P$  only:

$$P_\xi(\sigma, \mathbf{h}, \hat{\mathbf{h}}) = Z_\xi^{-1} \exp \left( i\hat{\mathbf{h}} \cdot \mathbf{h} + \frac{1}{2} \beta(\boldsymbol{\sigma} \cdot \mathbf{h}) + c \left\langle \sum_{\sigma'} \int d\hat{\mathbf{h}}' P_{\xi'}(\sigma', \hat{\mathbf{h}}') \exp \left( -\frac{i}{c} (\boldsymbol{\xi} \cdot \boldsymbol{\xi}') [(\hat{\mathbf{h}} \cdot \boldsymbol{\sigma}') + (\hat{\mathbf{h}}' \cdot \boldsymbol{\sigma})] \right) \right\rangle_{\xi'} \right) \quad (25)$$

with  $Z_\xi$  denoting a normalization constant. According to (21), the physical meaning of the saddle-point is

$$P_\xi(\sigma, \mathbf{h}, \hat{\mathbf{h}}) = \lim_{N \rightarrow \infty} \frac{1}{|I_\xi|} \sum_{i \in I_\xi} \overline{\langle \delta_{\sigma, \sigma_i} \delta[\mathbf{h} - \mathbf{h}_i(\{\sigma\})] \delta[\hat{\mathbf{h}} - \hat{\mathbf{h}}_i(\{\sigma\})] \rangle}. \quad (26)$$

We next make the one pattern condensed ansatz (this is not essential for being able to proceed, but will simplify and compactify our derivation significantly), which here implies  $P_\xi(\sigma, \mathbf{h}, \hat{\mathbf{h}}) = P_{\xi_1}(\sigma, \mathbf{h}, \hat{\mathbf{h}})$ , and we send  $c \rightarrow \infty$ . As a result  $(\boldsymbol{\xi} \cdot \boldsymbol{\xi}')/\sqrt{c} = \xi_1 \xi'_1/\sqrt{c} + \sqrt{\alpha} z$  where  $z$  is a zero-average unit-variance Gaussian variable, and

$$P_\xi(\sigma, \mathbf{h}, \hat{\mathbf{h}}) = Z_\xi^{-1} \exp \left( i\hat{\mathbf{h}} \cdot \mathbf{h} + \frac{1}{2} \beta(\boldsymbol{\sigma} \cdot \mathbf{h}) - \left\langle \sum_{\sigma'} \int d\hat{\mathbf{h}}' P_{\xi'}(\sigma', \hat{\mathbf{h}}') \times \left[ i(\xi \xi') [(\hat{\mathbf{h}} \cdot \boldsymbol{\sigma}') + (\hat{\mathbf{h}}' \cdot \boldsymbol{\sigma})] + \frac{\alpha}{2} [(\hat{\mathbf{h}} \cdot \boldsymbol{\sigma}') + (\hat{\mathbf{h}}' \cdot \boldsymbol{\sigma})]^2 \right] \right\rangle_{\xi'} \right). \quad (27)$$

In the right-hand side of (27) we are seen to need only the following moments of our distributions (which include the previously encountered  $\{m_\alpha, q_{\alpha\beta}\}$ ):

$$\begin{aligned} m_\alpha &= \left\langle \xi \sum_{\sigma} \int d\hat{\mathbf{h}} P_\xi(\sigma, \hat{\mathbf{h}}) \sigma_\alpha \right\rangle_\xi & q_{\alpha\beta} &= \left\langle \sum_{\sigma} \int d\hat{\mathbf{h}} P_\xi(\sigma, \hat{\mathbf{h}}) \sigma_\alpha \sigma_\beta \right\rangle_\xi \\ k_\alpha &= i \left\langle \xi \sum_{\sigma} \int d\hat{\mathbf{h}} P_\xi(\sigma, \hat{\mathbf{h}}) \hat{h}_\alpha \right\rangle_\xi & L_{\alpha\beta} &= \left\langle \sum_{\sigma} \int d\hat{\mathbf{h}} P_\xi(\sigma, \hat{\mathbf{h}}) \hat{h}_\alpha \hat{h}_\beta \right\rangle_\xi \\ K_{\alpha\beta} &= i \left\langle \sum_{\sigma} \int d\hat{\mathbf{h}} P_\xi(\sigma, \hat{\mathbf{h}}) \sigma_\alpha \hat{h}_\beta \right\rangle_\xi. \end{aligned}$$

Integration by parts over the fields in (27) shows that  $k_\alpha = -\frac{1}{2}\beta m_\alpha$ ,  $K_{\alpha\beta} = -\frac{1}{2}\beta q_{\alpha\beta}$  and  $L_{\alpha\beta} = -\frac{1}{4}\beta^2 q_{\alpha\beta}$ . The replicated joint spin–field distributions can now be written as

$$P_\xi(\boldsymbol{\sigma}, \mathbf{h}) = \frac{\exp\left(\xi\beta\mathbf{m}\cdot\boldsymbol{\sigma} + \frac{1}{2}\alpha\beta^2\boldsymbol{\sigma}\cdot\mathbf{q}\boldsymbol{\sigma} - \frac{1}{2\alpha}(\mathbf{h} - \xi\mathbf{m} - \alpha\beta\mathbf{q}\boldsymbol{\sigma})\mathbf{q}^{-1}(\mathbf{h} - \xi\mathbf{m} - \alpha\beta\mathbf{q}\boldsymbol{\sigma})\right)}{\sum_{\boldsymbol{\sigma}'} \int d\mathbf{h}' \exp\left(\xi\beta\mathbf{m}\cdot\boldsymbol{\sigma}' + \frac{1}{2}\alpha\beta^2\boldsymbol{\sigma}'\cdot\mathbf{q}\boldsymbol{\sigma}' - \frac{1}{2\alpha}(\mathbf{h}' - \xi\mathbf{m} - \alpha\beta\mathbf{q}\boldsymbol{\sigma}')\mathbf{q}^{-1}(\mathbf{h}' - \xi\mathbf{m} - \alpha\beta\mathbf{q}\boldsymbol{\sigma}')\right)} \quad (28)$$

with  $\mathbf{m} = \{m_\alpha\}$  and  $\mathbf{q} = \{q_{\alpha\beta}\}$ . The latter obey the following familiar closed equations which in the RS ansatz lead one back to (13), as they should:

$$m_\alpha = \left\langle \xi \frac{\sum_{\boldsymbol{\sigma}} \sigma_\alpha \exp\left(\xi\beta\mathbf{m}\cdot\boldsymbol{\sigma} + \frac{1}{2}\alpha\beta^2\boldsymbol{\sigma}\cdot\mathbf{q}\boldsymbol{\sigma}\right)}{\sum_{\boldsymbol{\sigma}} \exp\left(\xi\beta\mathbf{m}\cdot\boldsymbol{\sigma} + \frac{1}{2}\alpha\beta^2\boldsymbol{\sigma}\cdot\mathbf{q}\boldsymbol{\sigma}\right)} \right\rangle_\xi \quad (29)$$

$$q_{\alpha\beta} = \left\langle \frac{\sum_{\boldsymbol{\sigma}} \sigma_\alpha \sigma_\beta \exp\left(\xi\beta\mathbf{m}\cdot\boldsymbol{\sigma} + \frac{1}{2}\alpha\beta^2\boldsymbol{\sigma}\cdot\mathbf{q}\boldsymbol{\sigma}\right)}{\sum_{\boldsymbol{\sigma}} \exp\left(\xi\beta\mathbf{m}\cdot\boldsymbol{\sigma} + \frac{1}{2}\alpha\beta^2\boldsymbol{\sigma}\cdot\mathbf{q}\boldsymbol{\sigma}\right)} \right\rangle_\xi. \quad (30)$$

#### 4.2. Fraction of mis-aligned spins in the RS ansatz

The fraction  $\phi$  defined in (19) can be written as  $\phi = \langle \phi_\xi \rangle_\xi$ , where the sublattice fractions  $\phi_\xi$  are expressed in terms of the replicated distributions (28) in the following way:

$$\begin{aligned} \phi_\xi &= \frac{1}{2} - \frac{1}{2} \sum_{\boldsymbol{\sigma}} \int d\mathbf{h} P_\xi(\boldsymbol{\sigma}, \mathbf{h}) \frac{1}{n} \sum_{\gamma} \sigma_\gamma \operatorname{sgn}[h_\gamma] \\ &= \frac{1}{2} - \frac{1}{2n} \sum_{\gamma} \left\{ \left( \sum_{\boldsymbol{\sigma}} \sigma_\gamma \exp\left(\beta\mathbf{m}\cdot\boldsymbol{\sigma} + \frac{1}{2}\alpha\beta^2\boldsymbol{\sigma}\cdot\mathbf{q}\boldsymbol{\sigma}\right) \int d\mathbf{x} \right. \right. \\ &\quad \times \operatorname{sgn}\left[ m_\gamma + \alpha\beta \sum_{\beta} q_{\gamma\beta} \sigma_\beta + \sqrt{\alpha} x_\gamma \right] \exp\left(-\frac{1}{2}\mathbf{x}\cdot\mathbf{q}^{-1}\mathbf{x}\right) \\ &\quad \left. \left. \times \left( \sum_{\boldsymbol{\sigma}} \exp\left(\beta\mathbf{m}\cdot\boldsymbol{\sigma} + \frac{1}{2}\alpha\beta^2\boldsymbol{\sigma}\cdot\mathbf{q}\boldsymbol{\sigma}\right) \int d\mathbf{x} \exp\left(-\frac{1}{2}\mathbf{x}\cdot\mathbf{q}^{-1}\mathbf{x}\right) \right)^{-1} \right\}. \quad (31) \end{aligned}$$

We see that the  $\phi_1 = \phi_{-1}$ . At this stage we make the RS ansatz, putting  $m_\alpha = m$  and  $q_{\alpha\beta} = q + \delta_{\alpha\beta}(1 - q)$ , which results in

$$\begin{aligned} \phi_{\text{RS}} &= \frac{1}{2} - \frac{1}{2} \left\{ \left( \int Dz \sum_{\boldsymbol{\sigma}} \sigma_1 \exp\left(\beta \sum_{\alpha} \sigma_\alpha (m + z\sqrt{\alpha}q)\right) \int Dx \operatorname{sgn}\left[ m + \alpha\beta \left( \sigma_1 + q \sum_{\beta>1} \sigma_\beta \right) \right. \right. \right. \\ &\quad \left. \left. \left. + \sqrt{\alpha}x \right] \right) \left( \int Dz \sum_{\boldsymbol{\sigma}} \exp\left(\beta \sum_{\alpha} \sigma_\alpha (m + z\sqrt{\alpha}q)\right) \right)^{-1} \right\} \\ &= \frac{1}{2} - \frac{1}{2} \left\{ \left( \int Dx Dy Dz \sum_{\boldsymbol{\sigma}} \sigma_1 \exp\left(\beta \sum_{\alpha} \sigma_\alpha (m + z\sqrt{\alpha}q) + (x - iy)\sqrt{\alpha}\beta \right. \right. \right. \\ &\quad \left. \left. \left. \times \left( \sigma_1 + q \sum_{\alpha>1} \sigma_\alpha \right) \right) \operatorname{sgn}[m + \sqrt{\alpha}x] \right) \left( \int Dz [2 \cosh[\beta(m + z\sqrt{\alpha}q)]]^n \right)^{-1} \right\}. \end{aligned}$$

We carry out the spin summations over  $\sigma_\alpha$  with  $\alpha > 1$ . A shift in the complex plane for the variable  $z$  in the numerator,  $z \rightarrow z - \sqrt{q}(x - iy)$ , followed by integration over  $y$  and some simple manipulations, converts this expression into

$$\begin{aligned}
\phi_{\text{RS}} &= \frac{1}{2} - \frac{1}{4} \left\{ \left( \int \text{D}x \text{D}z \operatorname{sgn} \left[ z\sqrt{q} + x\sqrt{1-q} + \beta\sqrt{\alpha}(1-q) + \frac{m}{\sqrt{\alpha}} \right] \right. \right. \\
&\quad \times \exp(\beta(z\sqrt{\alpha q} + m)) \cosh^{n-1}[\beta(z\sqrt{\alpha q} + m)] \left. \left( \int \text{D}z \cosh^n[\beta(z\sqrt{\alpha q} + m)] \right)^{-1} \right\} \\
&\quad - \frac{1}{4} \left\{ \left( \int \text{D}x \text{D}z \operatorname{sgn} \left[ z\sqrt{q} + x\sqrt{1-q} + \beta\sqrt{\alpha}(1-q) - \frac{m}{\sqrt{\alpha}} \right] \exp(\beta(z\sqrt{\alpha q} - m)) \right. \right. \\
&\quad \times \cosh^{n-1}[\beta(z\sqrt{\alpha q} - m)] \left. \left( \int \text{D}z \cosh^n[\beta(z\sqrt{\alpha q} - m)] \right)^{-1} \right\} \\
&= \frac{1}{2} - \frac{1}{4} \left\{ \left( \int \text{D}z \operatorname{Erf} \left[ \frac{z\sqrt{\alpha q} + \beta\alpha(1-q) + m}{\sqrt{2\alpha(1-q)}} \right] \exp(\beta(z\sqrt{\alpha q} + m)) \cosh^{n-1} \right. \right. \\
&\quad \times [\beta(z\sqrt{\alpha q} + m)] \left. \left( \int \text{D}z \cosh^n[\beta(z\sqrt{\alpha q} + m)] \right)^{-1} \right\} \\
&\quad - \frac{1}{4} \left\{ \left( \int \text{D}z \operatorname{Erf} \left[ \frac{z\sqrt{\alpha q} + \beta\alpha(1-q) - m}{\sqrt{2\alpha(1-q)}} \right] \exp(\beta(z\sqrt{\alpha q} - m)) \cosh^{n-1} \right. \right. \\
&\quad \times [\beta(z\sqrt{\alpha q} - m)] \left. \left( \int \text{D}z \cosh^n[\beta(z\sqrt{\alpha q} - m)] \right)^{-1} \right\}. \tag{32}
\end{aligned}$$

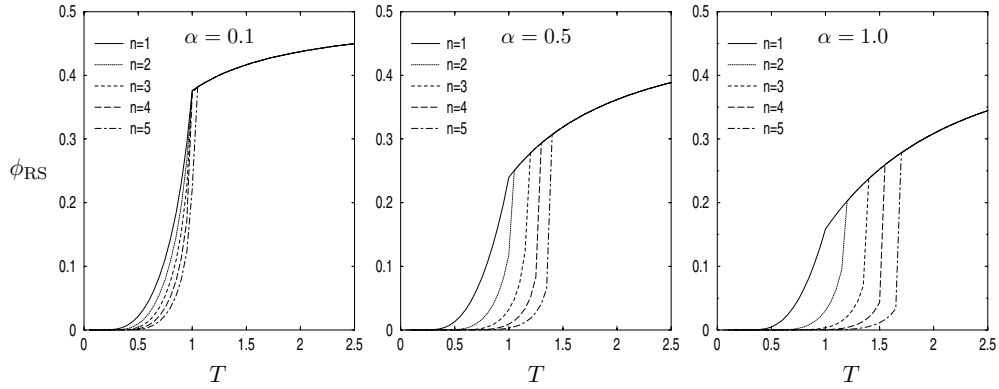
In the paramagnetic state, where  $m = q = 0$ , this simplifies further to

$$\phi_{\text{RS}} = \frac{1}{2} - \frac{1}{2} \operatorname{Erf} \left[ \beta \sqrt{\frac{\alpha}{2}} \right]. \tag{33}$$

In the recall and spin-glass states the evaluation of (32) requires the (numerical) solution of the RS order parameters  $\{m, q\}$  from (13). Examples of the resulting curves as functions of temperature are shown in figure 4, for  $\alpha \in \{0.1, 0.5, 1.0\}$  and  $n \in \{1, 2, 3, 4, 5\}$ . We note that for these values of  $n$ , the replica symmetry should be stable. The fraction of misaligned spins is seen to decrease with increasing  $n$  (i.e., with decreasing connectivity temperature), as expected. This effect becomes more pronounced for larger  $\alpha$ , where the amount of frustration to be reduced by the connectivity dynamics should indeed be largest. The first-order phase transitions (see subsection 3.2) induce discontinuities in  $\phi$  at the critical temperature, clearly recognizable for  $\alpha = 0.5, 1$  and  $n = 2, 3, 4, 5$ , and just visible for  $\alpha = 0.1$  and  $n = 5$ . In the paramagnetic phase (large  $T$ ) we see that  $\phi_{\text{RS}}$  is independent of  $n$ , in accordance with (33).

### 4.3. Comparison with numerical simulations

In order to perform numerical simulations, we need an explicit stochastic dynamical equation for updating of the network connectivity variables  $\mathbf{c} = \{c_{ij}\}$ , which must approach the appropriate Boltzmann equilibrium state characterized by the slow Hamiltonian (2). Here we used a Glauber-type Markov process, where candidate bonds  $c_{ij}$  are drawn randomly at



**Figure 4.** The fraction  $\phi_{\text{RS}}$  of misaligned spins as a function of temperature, in the RS ansatz, for integer values of  $n$  between 1 and 5, and  $\alpha = 0.1$ ,  $\alpha = 0.5$  and  $\alpha = 1$  respectively, as a function of temperature. The degree of alignment of spins with their local fields increases with  $n$  outside the paramagnetic phase. In the paramagnetic phase, there is no dependence on  $n$ . One sees clearly the effect of the first-order phase transitions, at  $\alpha = 0.1$  only for  $n = 5$ , and for  $\alpha = 0.5$  and  $\alpha = 1$  for all  $n \geq 2$ , appearing as discontinuities in  $\phi$  at the critical temperature.

each iteration step and then flipped with probability  $W[F_{ij}\mathbf{c}; \mathbf{c}]$ , where  $F_{ij}$  denotes the bond switch operator defined by  $F_{ij}c_{ij} = 1 - c_{ij}$ ,  $F_{ij}c_{k\ell} = c_{k\ell}$  if  $(i, j) \neq (k, \ell)$ :

$$W[F_{ij}\mathbf{c}; \mathbf{c}] = \frac{1}{2} \left\{ 1 - \tanh \left( \frac{\tilde{\beta}}{2} [H_s(F_{ij}\mathbf{c}) - H_s(\mathbf{c})] \right) \right\}. \quad (34)$$

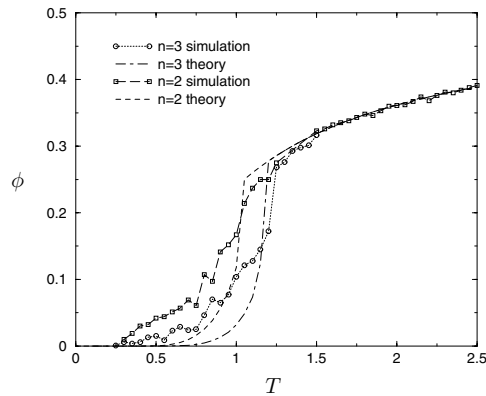
Detailed balance is built in. Upon inserting the slow Hamiltonian (2) and using the scaling property  $\lim_{N \rightarrow \infty} c/N = 0$  of our present extreme dilution regime, one can for large  $N$  rewrite our transition probabilities as

$$W[F_{ij}\mathbf{c}; \mathbf{c}] = \frac{1}{2} \left\{ 1 - \tanh \left[ \frac{1}{2} (2c_{ij} - 1) \left[ \log \left( \frac{c}{N} \right) + \frac{\tilde{\beta}}{c} \sum_{\mu} \xi_i^{\mu} \xi_j^{\mu} \langle \sigma_i \sigma_j \rangle \right] \right] \right\} \quad (35)$$

where, as before,  $\langle \dots \rangle$  indicates an equilibrium average for the fast system, in Boltzmann equilibrium with Hamiltonian (1), for a given connectivity matrix  $\mathbf{c}$ .

In the present type of systems with multiple adiabatically separated time scales and nested equilibrations, the verification of theoretical results by numerical simulations is known to be a highly demanding task. Even without the evolving connectivity, full equilibration of the spins requires relaxation times which diverge with  $N$  faster than polynomially. If on top of this one aims to also approach a connectivity equilibrium, which involves  $\mathcal{O}(N^2)$  stochastic degrees of freedom, the system sizes accessible in practice for numerical experimentation are small. Thus profound finite-size effects are unavoidable. It turns out that, when simulating the process (35), the connectivity equilibration times are indeed extremely long, especially close to phase transitions. This limits our ambitions regarding size, with the standard CPU resources at our disposal, to the order of  $N \sim 10^2$  spins. Since in our chosen scaling regime of extreme dilution, we have to simultaneously minimize  $c^{-1}$  and  $cN^{-1}$ , we have in our numerical experiments chosen  $c = \sqrt{N}$ .

Different macroscopic quantities could in principle be used for testing our theory against experiments. The advantage of observables such as  $m$  and  $\phi$  is that they can be measured



**Figure 5.** Comparison between simulation measurements (all with  $N = 200$ ) and RS theoretical predictions for the fraction of misaligned spins  $\phi = N^{-1} \sum_i \theta[-\sigma_i h_i]$  (where  $h_i$  is the local field at site  $i$ ), as functions of temperature. The data shown refer to  $\alpha = 0.5$  with  $n = 2$  (simulations: connected squares; theory: dashed lines) or  $n = 3$  (simulations: connected circles; theory: dotted-dashed lines). Because of the need to equilibrate two nested disordered processes, conventional computer resources limit experimentation to modest values of  $N$ . In spite of the resulting finite size effects, the graph does show satisfactory qualitative agreement between theory and experiment.

instantaneously, in contrast with the spin-glass order parameter  $q$ . Here we have opted for the fraction of misaligned spins  $\phi$ . The results are shown in figure 5, where we observe qualitative agreement between theory and the simulations. The deviations observed in such experiments are found to decrease with the increasing system size  $N$ , albeit slowly.

## 5. Stability of non-condensed retrieval states

The significant enlargement of the retrieval phase caused by our connectivity adaptation (see, e.g., figures 1 and 2) could have a drawback that retrieval of patterns other than the condensed one becomes impossible. Here we address the question of whether the present ‘tailoring’ of the connectivity variables  $\{c_{ij}\}$  to one condensed state will leave a finite basin of attraction for the non-condensed patterns, or whether recalling the latter requires a rewiring of the system (e.g., by temporarily raising the temperature  $\tilde{T}$ ) to undo the established connectivity. For large  $\alpha$  most retrieval states must be unstable for any given connectivity in the extreme dilution scaling regime, since Gardner-type optimal capacity calculations for diluted networks predict a finite storage capacity [33].

To answer our question we will study a second (fast) spin system of  $N$  spins  $\tau = \{\tau_i\}$ , governed again by the fast Hamiltonian (1), with patterns and connectivity identical to that of the first. In particular, the connectivity statistics are again given by

$$P(\mathbf{c}) = Z_s^{-1} e^{-\tilde{\beta} H_s(\mathbf{c})}. \quad (36)$$

The slow Hamiltonian  $H_s(\mathbf{c})$  continues to be defined in terms of the original spins  $\sigma$ , assumed in a condensed state characterized by a finite overlap with the first pattern, and will therefore be tailored towards the recall of that particular pattern. By studying in the  $\tau$  system the properties of states which are condensed in patterns two or higher, we gain access to the stability of non-condensed retrieval states in the original  $\sigma$  system.

The connectivity-averaged free energy per spin of our new system is calculated by using the replica trick in its conventional form, i.e. via

$$\begin{aligned}
 [f_\tau] &= - \lim_{\hat{n} \rightarrow 0} \lim_{N \rightarrow \infty} \frac{1}{\beta \hat{n} N} \log \left\{ \sum_{\mathbf{c}} P(\mathbf{c}) \left[ \sum_{\tau} \exp(-\beta H_f(\tau, \mathbf{c})) \right]^{\hat{n}} \right\} \\
 &= - \lim_{\hat{n} \rightarrow 0} \lim_{N \rightarrow \infty} \frac{1}{\beta \hat{n} N} \log \left\{ Z_s^{-1} \sum_{\{\sigma^\alpha\}} \sum_{\{\tau^\gamma\}} \sum_{\mathbf{c}} \exp \left( - \log \left( \frac{N}{c} \right) \sum_{i < j} c_{ij} \right) \right. \\
 &\quad \left. \times \exp \left( \frac{\beta}{c} \sum_{i < j} c_{ij} (\xi_i \cdot \xi_j) \left[ \sum_{\alpha=1}^n \sigma_i^\alpha \sigma_j^\alpha + \sum_{\gamma=1}^{\hat{n}} \tau_i^\gamma \tau_j^\gamma \right] \right) \right\}. \tag{37}
 \end{aligned}$$

The next stages of analysis are sufficiently similar to those followed earlier to justify limiting ourselves to giving the final result in the RS approximation. If again we assume at most  $r$  patterns to be condensed we find

$$\begin{aligned}
 [f_\tau]^{\text{RS}} &= \text{extr}_{\{\hat{m}_\mu, \hat{q}, a\}} \left\{ \frac{1}{2} \sum_{\mu \leq r} \hat{m}_\mu^2 - \frac{1}{4} \alpha \beta (\hat{q} - 1)^2 + \frac{1}{2} \alpha \beta n a^2 - \frac{1}{\beta} \log 2 \right. \\
 &\quad \left. - \frac{1}{\beta} \left\langle \frac{\int \text{Dy Dz} \cosh^n(\Xi_1) \log \cosh(\Xi_2)}{\int \text{Dy} \cosh^n(\Xi_1)} \right\rangle_{\xi} \right\} \tag{38}
 \end{aligned}$$

$$\Xi_1 = \beta (\mathbf{m} \cdot \xi + y \sqrt{\alpha q}) \tag{39}$$

$$\Xi_2 = \beta \left( \hat{\mathbf{m}} \cdot \xi + \frac{a}{q} \sqrt{\alpha q} [y + z \sqrt{\hat{q} q / a^2 - 1}] \right). \tag{40}$$

In addition to the previously encountered order parameters  $\{\mathbf{m}, q\}$ , which relate to the  $\sigma$  system (and continue to be defined as the solution of the earlier saddle-point equations), we now have new order parameters  $\{\hat{\mathbf{m}}, \hat{q}, a\}$ , whose physical meaning is found to be

$$\hat{m}_\mu = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_i \overline{\langle \xi_i^\mu \tau_i \rangle} \quad \hat{q} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_i \overline{\langle \tau_i \rangle^2} \quad a = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_i \overline{\langle \sigma_i \tau_i \rangle}.$$

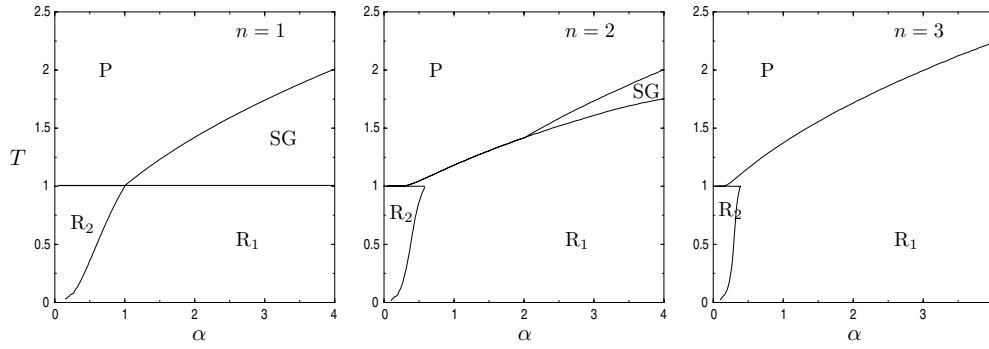
The new order parameters are to be solved from the saddle-point equations

$$\begin{aligned}
 \hat{m}_\mu &= \left\langle \xi_\mu \frac{\int \text{Dy Dz} \tanh(\Xi_2) \cosh^n(\Xi_1)}{\int \text{Dy} \cosh^n(\Xi_1)} \right\rangle_{\xi} \\
 \hat{q} &= \left\langle \frac{\int \text{Dy Dz} \tanh^2(\Xi_2) \cosh^n(\Xi_1)}{\int \text{Dy} \cosh^n(\Xi_1)} \right\rangle_{\xi} \\
 a &= \left\langle \frac{\int \text{Dy Dz} \tanh(\Xi_1) \tanh(\Xi_2) \cosh^n(\Xi_1)}{\int \text{Dy} \cosh^n(\Xi_1)} \right\rangle_{\xi}. \tag{41}
 \end{aligned}$$

It can be shown that solutions of these equations will obey  $a^2 \leq \hat{q} q$  (to be expected in view of the square root in  $\Xi_2$ ).

We now adopt a condensed ansatz which corresponds to the  $\tau$  system being in a condensed state which differs from that of the  $\sigma$  system (where the latter drives the connectivity evolution):  $m_\mu = m \delta_{\mu 1}, \hat{m}_\mu = \hat{m} \delta_{\mu 2}$ . Solutions of this type must have  $a = 0$ , which is reasonable





**Figure 6.** Cross-sections for fixed  $n$  of the expanded phase diagram, in which the previous retrieval phase R has been separated into two sub-regions: R<sub>1</sub> defines the phase where only the nominated pattern can be recalled to which the connectivity has been tailored, and R<sub>2</sub> defines the phase where, in spite of the biased connectivity, all stored patterns can still be recovered. From left to right:  $n = 1, 2, 3$ .

considering that any finite correlation between the  $\tau$  and  $\sigma$  systems makes  $\hat{m}_1 = 0$  highly improbable. For  $a = 0$  our two systems decouple, with the equations for  $\hat{m}$  and  $\hat{q}$  reducing to

$$\hat{m} = \int Dz \tanh(\beta[\hat{m} + z\sqrt{\alpha\hat{q}}]) \quad (42)$$

$$\hat{q} = \int Dz \tanh^2(\beta[\hat{m} + z\sqrt{\alpha\hat{q}}]). \quad (43)$$

These are exactly the RS equations of the model with frozen random dilution [26]. The  $\hat{m}_1 = a = 0$  solutions of our saddle-point equations could be unstable against perturbations in  $\hat{m}_1$  and  $a$ . In the paramagnetic phase, an expansion of the free energy up to second order in the order parameters gives

$$[f_\tau]^{\text{RS}} = \text{extr}_{\{\hat{m}, \hat{q}, a\}} \left\{ \frac{1}{2}(1 - \beta) \sum_\mu \hat{m}_\mu^2 - \frac{1}{4}\alpha\beta\hat{q}^2 + \frac{1}{2}\alpha\beta a^2 + \text{higher orders} \right\}$$

indicating that the physical solution of the saddle-point equations is the one which minimizes the free energy with respect to  $\hat{\mathbf{m}}$  and  $a$ , and maximizes it with respect to  $\hat{q}$ , as is usual in the limit  $\hat{n} \rightarrow 0$  [23]. Expansion around  $\hat{m}_1 = a = 0$ , with nonzero  $\hat{m}_2$  and  $\hat{q}$ , reveals that a second-order instability in  $a$  occurs at the temperature

$$T_c = \sqrt{\alpha(1 - \hat{q})[1 + (n - 1)q]}. \quad (44)$$

Below  $T_c$  the  $\tau$  system will be captured in the  $\hat{m}_\mu = \hat{m}\delta_{\mu 1}$  state, with  $\hat{m}_1 = m$  (so retrieval of states other than that to which the connectivity has adapted is impossible), whereas above  $T_c$  the  $\tau$  system can be in a locally stable condensed state different from that in which the  $\sigma$  system is found. The free energy of the  $\hat{m}_\mu = m\delta_{\mu 1}$  state is, however, always lower than that of other retrieval states, at any temperature. This implies that in the latter states the  $\tau$  system can be at most locally stable. The line (44) has been calculated in the RS ansatz; this seems reasonable, since in the model of [26] RSB does not occur for  $n \geq 1$ . In figure 6, we show for  $n \in \{1, 2, 3\}$  the line (44) which separates the previous retrieval phase R into two sub-regions, one R<sub>1</sub> where only recall of one single pattern is possible (the one to which the connectivity is tailored), and a second region R<sub>2</sub> where, in spite of the biased connectivity, multiple patterns can be recalled. As expected, increasing  $n$  (i.e. reducing the connectivity temperature, so the ‘tailoring’ of the connectivity becomes more effective) reduces the size of the R<sub>2</sub> region.

## 6. Conclusion

We have studied extremely diluted recurrent neural networks in which the connectivity is allowed to evolve on time scales which are adiabatically slower than the equilibration time of the (fast) neurons. In contrast to earlier studies, the actual *values* of the bonds remain frozen (they are here given by Hopfield's [1] recipe) and only the connectivity is dynamic, which implies that the slow adaptation is reversible and will not wipe out any stored information. Our motivation was to investigate whether, by having a connectivity dynamics which aims to reduce frustration, the information retrieval properties of the system can be improved. As in earlier models with slow bond dynamics [16–22] the equilibrium properties of our model are described by a replica theory with nonzero replica dimension  $n$ , where  $n = \tilde{\beta}/\beta$  is the ratio between the temperature of the (fast) neurons and the temperature of the (slow) connectivity.

We have calculated phase diagrams, reflecting the stationary state of the slowest stochastic system (i.e. the connectivity). They reveal a boosting of the retrieval phase, compared to the frozen connectivity case, as soon as  $n > 0$ . In fact, for nonzero  $n$  the storage capacity diverges at low temperatures, as long as  $p \ll N$ . This at first sight somewhat surprising result is explained by the observation that, in tailoring the connectivity to the recall of a single condensed pattern, the system sacrifices the recall quality of an infinite number of non-nominated patterns. RSB effects are as always confined to small values of  $n$  (below approx. 0.32). In order to measure the expected reduction in frustration as a result of the connectivity dynamics, we have calculated the fraction of mis-aligned spins (where spin and local field are of the opposite sign). This fraction is indeed found to decrease with decreasing temperature of the connectivity. In order to examine in which region of the phase diagram retrieval states other than the condensed pattern are still locally stable, we studied a pair of identical diluted networks, both with the same Boltzmann-type connectivity distribution. The connectivity is tailored to reduce frustration in only the first of the two copies. This allows one to study scenarios corresponding to the recall of patterns (in the second copy) which are not the one to which the connectivity is adapted. Such recall is seen to be possible, but only in a sub-region of the recall phase, whose size decreases with increasing  $n$ .

It would appear an interesting question to examine to what extent the properties of our model with slowly evolving connectivity persist in more (biologically) realistic scenarios, e.g. when the average number of connections  $c$  per neuron remains finite when  $N \rightarrow \infty$ . Such studies will involve order-parameter functions (see e.g. [32, 34]), and require finite  $n$  generalizations of finite connectivity replica theory.

## Acknowledgments

BW and NS acknowledge financial support from the FOM Foundation (Fundamenteel Onderzoek der Materie, the Netherlands) and the Ministerio de Educación, Cultura y Deporte (Spain, grant SB2002-0107).

## References

- [1] Hopfield J J 1982 *Proc. Natl Acad. Sci. USA* **79** 2554
- [2] Amit D J, Gutfreund H and Sompolinsky H 1985 *Phys. Rev. A* **32** 1007
- [3] Amit D J, Gutfreund H and Sompolinsky H 1985 *Phys. Rev. Lett.* **55** 1530
- [4] Derrida B, Gardner E and Zippelius A 1987 *Europhys. Lett.* **4** 167
- [5] Coolen A C C 2001 *Handbook of Biological Physics* vol 4 ed F Moss and S Gielen (Amsterdam: Elsevier) p 531

- 
- [6] Coolen A C C 2001 *Handbook of Biological Physics* vol 4 ed F Moss and S Gielen (Amsterdam: Elsevier) p 597
  - [7] Gardner E 1988 *J. Phys. A: Math. Gen.* **21** 257
  - [8] Hertz J A, Krogh A and Thorgersson G I 1989 *J. Phys. A: Math. Gen.* **22** 2133
  - [9] Biehl M and Schwarze H 1992 *Europhys. Lett.* **20** 733
  - [10] Kinouchi O and Caticha N 1992 *J. Phys. A: Math. Gen.* **25** 6243
  - [11] Heimel J A F and Coolen A C C 2001 *J. Phys. A: Math. Gen.* **34** 9009
  - [12] Kinzel W and Oppen M 1991 *Physics of Neural Networks I* (Berlin: Springer)
  - [13] Watkin T L H, Rau A and Biehl M 1993 *Rev. Mod. Phys.* **65** 499
  - [14] Mace C W H and Coolen A C C 1998 *Stat. Comput.* **8** 55
  - [15] Shinomoto S 1987 *J. Phys. A: Math. Gen.* **20** L1305
  - [16] Penney R W, Coolen A C C and Sherrington D 1993 *J. Phys. A: Math. Gen.* **26** 3681
  - [17] Coolen A C C, Penney R W and Sherrington D 1993 *Phys. Rev. B* **48**,16 116
  - [18] Feldman D E and Dotsenko V S 1994 *J. Phys. A: Math. Gen.* **27** 4401
  - [19] Dotsenko V, Franz S and Mézard M 1994 *J. Phys. A: Math. Gen.* **27** 2351
  - [20] Jongen G, Bollé D and Coolen A C C 1998 *J. Phys. A: Math. Gen.* **31** L737
  - [21] Jongen G, Anemuller J, Bollé D, Coolen A C C and Perez-Vicente C 2001 *J. Phys. A: Math. Gen.* **34** 3957
  - [22] Uezu T and Coolen A C C 2002 *J. Phys. A: Math. Gen.* **35** 2761
  - [23] Mézard M, Parisi G and Virasoro M A 1987 *Spin-Glass Theory and Beyond* (Singapore: World Scientific)
  - [24] van Mourik J and Coolen A C C 2001 *J. Phys. A: Math. Gen.* **34** L111
  - [25] Jonker H J J and Coolen A C C 1993 *J. Phys. A: Math. Gen.* **26** 563
  - [26] Watkin T L H and Sherrington D 1991 *Europhys. Lett.* **14** 791
  - [27] Sherrington D 1980 *J. Phys. A: Math. Gen.* **13** 637
  - [28] Sherrington D and Kirkpatrick S 1975 *Phys. Rev. Lett.* **35** 1792
  - [29] De Almeida J R L and Thouless D J 1978 *J. Phys. A: Math. Gen.* **11** 983
  - [30] Thomsen M, Thorpe M F, Choy T C, Sherrington D and Sommers H J 1986 *Phys. Rev. B* **33** 1931
  - [31] van Hemmen J L and Kühn R 1986 *Phys. Rev. Lett.* **57** 913
  - [32] Wemmenhove B and Coolen A C C 2003 *J. Phys. A: Math. Gen.* **36** 9617
  - [33] Shim G M, Kim D and Choi M Y 1993 *J. Phys. A: Math. Gen.* **26** 3741
  - [34] Perez-Castillo I and Skantzos N S 2003 *Preprint cond-mat/0309655*