

A Non-Binary Associative Memory with Exponential Pattern Retrieval Capacity and Iterative Learning

K. Raj Kumar, Amir Hesam Salavati, and Amin Shokrollahi

Laboratoire d'algorithmique (ALGO)

Ecole Polytechnique Fédérale de Lausanne (EPFL), 1015 Lausanne, Switzerland

E-mail: {raj.kumar,hesam.salavati,amin.shokrollahi}@epfl.ch

Abstract—We consider the problem of neural association for a network of non-binary neurons. Here, the task is to first memorize a set of given patterns using a network of neurons whose states assume values from a finite number of non-negative integer levels. Later, the same network should be able to recall a previously memorized pattern from its noisy version. Prior works in this area consider storing a finite number of *purely random* patterns, and have shown that the pattern retrieval capacities (maximum number of patterns that can be memorized) scale only linearly with the number of neurons in the network.

In our formulation of the problem, we concentrate on exploiting redundancy and internal structure of the patterns in order to improve the pattern retrieval capacity. We show that if the given patterns have a suitable structure, i.e. come from a sub-space of the set of all possible patterns, we can have exponential pattern retrieval capacities in terms of the length of the patterns. Another example is when the patterns have strong minor components. We will use this minor components (or the basis vectors of the patterns null space) to both increase the pattern retrieval capacity and error correction capabilities.

An iterative algorithm is proposed for the learning phase. In addition, two simple neural update algorithms are presented for the recall phase and using analytical results and simulations, we show that the suggested methods can tolerate a fair amount of errors in the input.

I. INTRODUCTION

Neural associative memory is a particular class of neural memory capable of memorizing (learning) a number of given patterns and recall them later in presence of noise, i.e. retrieve the correct memorized pattern from a given noisy version. Since 1982 and due to the seminal work of Hopfield [4], various artificial neural networks have been designed to mimic the task of the neural associative memory. See for instance [13], [14], [10], [12], [18].

In essence, the neural associative memory is very similar to the one faced in communication systems where the goal is to reliably transmit and efficiently decode a set of patterns (so called codewords) over a noisy channel. In both cases, we have to learn a set of patterns (codebook), and retrieve the correct one from a noisy version. More interestingly, the techniques used to implement an artificial neural associative memory looks very similar to some of the methods used in codes on graphs to design reliable and efficient codes.

However, despite the similarity in the task and techniques employed in both problems, there is a huge gap in terms

of another important criterion: the efficiency. Using binary codewords of length n , all modern codes on graphs are capable of reliably transmitting 2^{rn} codewords of length n over a noisy channel (and recover the transmitted codeword from the corrupted received pattern), with $0 < r < 1$ being the code rate [20]. In many cases, one can also obtain the optimal r , i.e. the largest possible value that permits the almost sure recovery of transmitted codewords from the corrupted received version.

However, using current neural networks of size n to memorize a set of *randomly* chosen patterns, the maximum number of patterns that can be reliably memorized scales linearly in n [11], [13]. In other words, a large body of the current works on artificial neural associative memories concentrate on the ability of the network to memorize any set of patterns, randomly chosen from the pool of all possibilities, and being able to retrieve the correct memorized pattern from a noisy version (e.g., [4], [13], [14], [10]). Although this gives the network a sense of generality, it severely limits its efficiency since in many cases, it is possible to have two patterns that look very similar to each other. Now distinguishing these two patterns from a corrupted input version would be really difficult and one has to limit the number of patterns that can be memorized in order to maintain the reliability of the answers in the recall phase. This is in sharp contrast to coding techniques where one designs codewords such that they have the maximum possible similarity distance from each other, i.e. the codewords are selected such that they do not look like each other at all. Furthermore, add the above issue to the simple nature neurons, due to which they can only perform rather simple operations, to conclude that most of the techniques used in coding theory can not be implemented by neural networks.

Therefore, the above two reasons can be among the possible explanation about the huge difference in the pattern retrieval capacity of neural networks and that of coding techniques.

In this paper, we focus on increasing the pattern retrieval capacity of neural associative memories. We propose a neural network which is able to memorize and reliably recall an exponential number of pattern with length n , if the patterns come from a subspace with dimension $k < n$. In other words, we show that if the patterns in the training set has some internal *structure* and *redundancy*, the proposed model exploits this structure in order to increase the number of

patterns that can be memorized as well as eliminating noise during the recall phase. In [2], we explained some preliminary results in which two efficient recall algorithms were proposed for the case where the neural graph was expander. Here, we extend the previous results to general sparse neural graphs as well as proposing a simple learning algorithm to capture the internal structure of the patterns in order to be used later in the recall phase.

The remainder of this paper is organized as follows: In section II, we will discuss the neural model used in this paper and formally define the problem. We explain the proposed learning algorithm in section III. Section IV is dedicated to the recall algorithm and analytically investigating its performance in retrieving corrupted patterns. In section V we address the pattern retrieval capacity and show it can be exponential in n . Simulation results are discussed in section VI. Finally, section VII concludes the paper and discusses future research topics.

II. PROBLEM FORMULATION AND THE NEURAL MODEL

A. The Model

We choose to work with non-binary neurons, i.e. neurons whose state is an integer from a finite set of integer values $\mathcal{S} = \{0, 1, \dots, S - 1\}$. A natural way of interpreting this model is to think of the integer states as the short-term firing rate of neurons. In other words, the state of a neuron in this model indicates the number of spikes fired by the neuron in a fixed short time interval.

Like in other neural networks, neurons can do only simple operations. We consider neurons that can do *linear summation* over the input and possibly apply a *non-linear function* (such as thresholding), to produce the output. More specifically, neuron x updates its state based on the states of its neighbors $\{s_i\}_{i=1}^n$ as follows:

- 1) It computes the weighted sum $h = \sum_{i=1}^n w_i s_i$, where w_i denotes the weight of the input link from s_i .
- 2) It updates its state as $x = f(h)$, where $f : \mathbb{R} \rightarrow \mathcal{S}$ is a possibly non-linear function from the field of real numbers \mathbb{R} to \mathcal{S} .

B. The Problem

We assume to be given C vectors of length n with integer-valued entries belonging to \mathcal{S} . In order to have a neural associative memory with large pattern retrieval capacity, we concentrate on learning the set of patterns that have some internal structure. More specifically, we assume these patterns to come from a subspace with dimension $k \leq n$. Note that if $k = n$, then we are back to the original associative memory addressed by Hopfield and others [4], [13], [14], [10], [12], [18]. However, if $k < n$, we will end up with much larger pattern retrieval capacities, as will be shown later.

We would like to memorize these patterns by finding a set of non-zero vectors $w_1, \dots, w_m \in \mathbb{R}^n$ that are orthogonal to the set of given patterns. Furthermore, we are interested in rather sparse vectors. Denoting pattern μ by x^μ and focusing

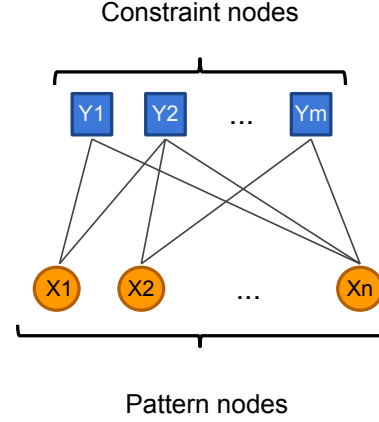


Fig. 1. A bipartite graph that represents the constraints on the training set.

on one such vector w , we can formulate the problem as:

$$\min \sum_{\mu=1}^C |x^\mu \cdot w|^2 \quad (1a)$$

subject to

$$\|w\|_0 \leq q \quad (1b)$$

$$\|w\|_2^2 \geq \epsilon \quad (1c)$$

where \cdot represents the inner-product, $q \in \mathbb{N}$ determines the degree of sparsity and $\epsilon \in \mathbb{R}^+$ prevents the all-zero solution.

By repeating the above problem for different vectors w , we will have $m = n - k$ sparse vectors who form the basis for the null space of the patterns which would like to learn. Therefore, the inherent structure of the patterns are captured in the obtained null-basis, denoted by the matrix $W \in \mathbb{R}^{m \times n}$, which can be interpreted as the adjacency matrix of a neural graph. The overall network model is shown in Figure 1.

In the recall phase, the neural network should retrieve the correct memorized pattern from a possibly corrupted version. In this case, the states of the pattern neurons x_1, x_2, \dots, x_n are initialized with the given input pattern. Here, we assume that the noise is integer valued and additive¹. Therefore, assuming the input to the network is a corrupted version of pattern x^μ , the state of the pattern nodes are $x = x^\mu + z$, where z is the noise. Now the neural network should use the given states together with the fact that $Wx^\mu = 0$ to retrieve the pattern x^μ . Any algorithm designed for this purpose should be simple enough to be implemented by neurons. Therefore, we must find a simple algorithm capable of eliminating noise using only neural operation.

Before presenting our solution, we briefly overview the literature and explain the relations between the previous works and the presented approach in this paper.

¹It must be mentioned that neural states below 0 and above $S - 1$ will be set to 0 and $S - 1$, respectively.

C. Related Works

Designing a neural associative memory has been an active field for the past three decades. Hopfield was the first to design an artificial neural associative memory in his seminal work back in 1982 [4]. The so-called Hopfield network is inspired by Hebbian learning [citehebb] and is composed of binary-valued (± 1) neurons, which together are able to memorize a certain number of patterns. The pattern retrieval capacity of a Hopfield network of n neurons was derived later by Amit et al. [6] and shown to be $0.13n$, under vanishing bit error probability requirement. Later, McEliece et al. proved that the capacity of Hopfield networks under vanishing block error probability requirement is $O(n/\log(n))$ [11].

While the connectivity graph of a Hopfield network is a complete graph, Komlos et al. [9] extended the work of McEliece to sparse neural graphs. Their results are of particular interest as physiological data is also in favor of sparsely interconnected neural networks. Komlos and Paturi have considered a network in which each neuron is connected to d other neurons, i.e., a d -regular network. Assuming that the network graph satisfies certain connectivity measures, they prove that it is possible to store a linear number of *random* patterns (in terms of d) with vanishing bit error probability or $C = O(d/\log n)$ random patterns with vanishing block error probability. Furthermore, they show that in spite of the capacity reduction, the error correction capability remains the same as the network can still tolerate a number of errors linear in n .

It is also known that the capacity of neural associative memories could be enhanced if the patterns are *sparse* in the sense that at any time instant many of the neurons are silent [7]. However, even these schemes fail when required to correct a fair amount of erroneous bits as the information retrieval is not better compared to that of normal networks.

In addition to neural networks capable of learning patterns gradually, in [13], the authors calculate the weight matrix offline (as opposed to gradual learning) using the pseudo-inverse rule [7] which in return help them improve the capacity of a Hopfield network to $n/2$, under vanishing block error probability condition, while being able to correct *one bit* of error in the recall phase. Although this was a significant improvement to the $n/\log(n)$ scaling of the pattern retrieval capacity in [11], it comes at the price of much higher computational complexity and the lack of the ability to learn patterns gradually online.

Extension of associative memories to non-binary neural models has also been explored in the past. Hopfield addressed the case of continuous neurons and showed that similar to the binary case, neurons with states between -1 and 1 can memorize a set of random patterns, albeit with less capacity [5]. In [14] the authors investigated a multi-state complex-valued neural associative memories for which the estimated capacity is $C < 0.15n$. Under the same model but using a different learning method, Muezzinoglu et al. [10] showed that the capacity can be increased to $C = n$. However the complex-

ity of the weight computation mechanism is prohibitive. To overcome this drawback, a Modified Gradient Descent learning Rule (MGDR) was devised in [15], with similar results in terms of capacity.

Given that even very complex offline learning methods can not improve the capacity of binary or multi-state Hopfield networks, a line of recent work has made considerable efforts to exploit the inherent structure of the patterns in order to increase both capacity and error correction capabilities. Such methods either make use of higher order correlations of patterns or focus merely on those patterns that have some sort of redundancy. As a result, they differ from previous methods in which the network was designed to be able to memorize any random set of patterns. Pioneering this prospect, Berrou and Gripon [18] achieved considerable improvements in the pattern retrieval capacity of Hopfield networks, by utilizing Walsh-Hadamard sequences. The only slight downside to the proposed method is the use of a decoder based on the winner-take-all approach which requires a separate neural stage, increasing the complexity of the overall method. Using low correlation sequences has also been considered in [12], where the authors introduced two novel mechanisms of neural association that employ binary neurons to memorize patterns with low correlation properties. The network itself is very similar to that of Hopfield, with a slightly modified weighting rule. Therefore, similar to a Hopfield network, the complexity of the learning phase is small. However, the authors failed to increase the pattern retrieval capacity beyond n and it was shown that the pattern retrieval capacity of the proposed model is $C = n$, while being able to correct a fair number of erroneous input bits.

In contrast to the pairwise correlation of the Hopfield model [4], Peretto et al. [17] deployed *higher order* neural models: the state of the neurons not only depends on the state of their neighbors, but also on the correlation among them. Under this model, they showed that the storage capacity of a higher-order Hopfield network can be improved to $C = O(n^{p-2})$, where p is the degree of correlation considered. The main drawback of this model is again the huge computational complexity required in the learning phase, as one has to keep track of $O(n^{p-2})$ neural links and their weights in the learning phase.

To address this difficulty while being able to capture higher-order correlations, The present authors introduced a bipartite neural graph inspired from iterative coding theory [2]. Under the assumptions that the bipartite graph is known, sparse, and expander, the proposed algorithm increased the pattern retrieval capacity to $C = O(a^n)$, for some $a > 1$, closing the gap between the pattern retrieval capacities achieved in neural networks and that of coding techniques. The main drawbacks in the proposed approach were the lack of a learning algorithm as well as the expansion assumption on the neural graph.. The sparsity criterion on the other hand, as it was mentioned earlier, is biologically more meaningful.

In this paper, we focus on extending the results described in [2] in several directions: first, we will suggest a learning algorithm, to find the neural connectivity matrix from the

patterns in the training set. Secondly, we provide an analysis of the proposed error correcting algorithm in the recall phase and investigates its performance as a function of input noise and network model. Finally, we discuss some variants of the error correcting method which achieve better performance marks in practice.

It worths mentioning that an extension of this approach to a multi-level neural network is considered in [29]. There, the novel structure enables better error correction. However, the learning algorithm lacks the ability to learn the patterns one by one and requires the patterns to be presented all at the same time in the form of a big matrix. In contrats, the learning algorithm proposed in this paper is capable of learning patterns gradually as they are presented to the network one by one.

Another important point to note is that learning linear constraints by a neural network is hardly a new topic as one can learn a matrix orthogonal to a set of patterns in the training set (i.e., $Wx^\mu = 0$) using simple neural learning rules (we refer the interested readers to [3] and [16]). However, to the best of our knowledge, finding such a matrix subject to the sparsity constraints has not been investigated before. This problem can also be regarded as an instance of compressed sensing [21], in which the measurement matrix is given by the big patterns matrix $\mathcal{X}_{C \times n}$ and the set of measurements are the constraints we look to satisfy, denoted by the tall vector b , which for simplicity reasons we assume to be all zero. Thus, we are interested in finding a sparse vector w such that $\mathcal{X}w = 0$. Nevertheless, many decoders proposed in this area are very complicated and cannot be implemented by a neural network using simple neuron operations. Some exceptions are [1] and [19] which are closely related to the learning algorithm proposed in this paper.

III. LEARNING PHASE

Since the patterns are assumed to be coming from a subspace in the n -dimensional space, we adapt the algorithm proposed by Oja and Karhunen [28] to learn the null-space basis of the subspace defined by the patterns. In fact, a very similar algorithm is also used in [3] for the same purpose. However, since we need the basis vectors to be sparse (due to requirements of the algorithm used in the recall phase), we add an additional term to make the final result of the learning algorithm rather sparse.

Another difference with the proposed method and that of [3] is that in [3] the authors design their method such that the resulted dual vectors form an orthogonal set. Although one can easily extend our suggested method to such a case as well, we find this requirement unnecessary in our case. This gives us the additional advantage to make the algorithm *parallel* and *adaptive*. Parallel in the sense that we can design an algorithm to learn one constraint and repeat it several times in order to find all of the constraints with high probability. And adaptive in the sense that we can determine the number of constraints on-the-go, i.e. start by learning just a few constraints. If needed (for instance due to bad performance in the recall phase), one can easily identify additional constraints. This increases

the flexibility of the algorithm and provides a nice trade-off between the time spent on learning and the performance in the recall phase.

A. Overview of the proposed algorithm

In order to develop a simple iterative algorithm, we formulate the problem in an optimization framework. The problem to find a constraint vector w is given by equation (1). However, instead of tackling problem (1) directly, we make a slight modification and consider the following optimization problem:

$$\min \sum_{\mu=1}^C |x^\mu \cdot w|^2 + \beta g(w). \quad (2a)$$

subject to:

$$\|w\|_2 = 1 \quad (2b)$$

In the above problem, we have replaced the constraint $\|w\|_0 \leq q$ with a penalty function $g(w)$ since $\|\cdot\|_0$ is not easy to handle analytically. The function $g(w)$ is chosen such that it favors sparsity. For instance one can pick $g(w)$ to be $\|\cdot\|_1$, which leads to ℓ_1 -norm penalty and is widely used in compressed sensing applications [1], [19]. In this paper, we consider the function

$$g(w) = \sum_{i=1}^n \tanh(\sigma w_i^2)$$

where σ is chosen appropriately. Intuitively, $\tanh(\sigma w_i^2)$ approximates $|\text{sgn}(w_i)|$ in ℓ_0 -norm. Therefore, the larger σ is, the closer $g(w)$ will be to $\|\cdot\|_0$. By calculating the derivative of the objective function, and by considering the update due to each randomly picked pattern x^μ , we will get the following iterative algorithm:

$$y(t) = x(t) \cdot w(t) \quad (3a)$$

$$\tilde{w}(t+1) = w(t) - \alpha_t (y(t)x(t) + \beta \Gamma(w(t))) \quad (3b)$$

$$w(t+1) = \frac{\tilde{w}(t+1)}{\|\tilde{w}(t+1)\|} \quad (3c)$$

In the above equations, t is the iteration index, $x(t)$ is the sample pattern chosen at iteration t uniformly at random from the patterns in the training set \mathcal{X} , and α_t is a small positive constant. Finally, $\Gamma(w) : \mathcal{R}^n \rightarrow \mathcal{R}^n = \nabla g(w)$ is the gradient of the penalty term for non-sparse solutions.

Remark 1. *The i^{th} entry of the function $\Gamma(w(t)) = \nabla g(w(t))$ was driven to be $f_i(w(t)) = \partial g(w(t))/\partial w_i(t) = 2\sigma_i w_i(t)(1 - \tanh^2(\sigma_i w_i(t)^2))$. This function has the interesting property that for very small values of $w_i(t)$, $f_i(w(t)) \simeq 2\sigma_i w_i(t)$ and for relatively larger values of $w_i(t)$, we get $f_i(w(t)) \simeq 0$. Therefore, by proper choice of β , equation (3b) suppresses small entries of $w(t)$ by pushing them towards zero. In other words, this function favors sparser results.*

Following the same approach as [28] and assuming α_t to be small enough such that equation (3c) can be expanded as

Algorithm 1 Iterative Learning

Input: Set of patterns $x^\mu \in \mathcal{X}$ with $\mu = 1, \dots, C$, stopping point ε .

Output: w

while $\sum_\mu |x^\mu \cdot w(t)|^2 > \varepsilon$ **do**
 Choose $x(t)$ at random from patterns in \mathcal{X}
 Compute $y(t) = x(t) \cdot w(t)$
 Update $w(t+1) = w(t) - \alpha_t y(t) \left(x(t) - \frac{y(t)w(t)}{\|w(t)\|^2} + \beta \Gamma(w(t)) \right)$.
 $t \leftarrow t + 1$.

end while

powers of α_t . In this case, we can approximate the algorithm (3) with the following simpler version:

$$y(t) = x(t) \cdot w(t) \quad (4a)$$

$$w(t+1) = w(t) - \alpha_t \left(y(t) \left(x(t) - \frac{y(t)w(t)}{\|w(t)\|^2} \right) + \beta \Gamma(w(t)) \right) \quad (4b)$$

In the above approximation, we also omitted the term $\alpha_t \beta (w(t) \cdot \Gamma(w(t))) w(t)$ since $w(t) \cdot \Gamma(w(t)) \simeq \|\Gamma(w(t))\|^2$ and $\|\Gamma(w(t))\|^2$ tends to zero as σ_t grows with t . In words, $y(t)$ is the projection of $x(t)$ on the basis vector $w(t)$. If for a given data vector $x(t)$, $y(t)$ is equal to zero, namely, the data is orthogonal to the current weight vector $w(t)$, then according to equation (4b) the weight vector will not be updated. However, if the data vector $x(t)$ has some projection over $w(t)$ then the weight vector is updated towards the direction to reduce this projection. The overall learning algorithm for one constraint node is given by algorithm 1.

Since we are interested in finding m basis vectors, we have to do the above procedure m times in parallel.

Remark 2. *Although we are interested in finding a sparse graph, note that too much sparseness is not desired. Because we are going to use the feedback sent by the constraint nodes to eliminate input noise at pattern nodes. Now if the graph is too sparse, the number of feedbacks received by each pattern node is too small to be reliable. Therefore, we must adjust the penalty coefficient β such that resulting neural graph is rather small. In the section on experimental results, we compare the error correction performance for different choices of β .*

B. Convergence analysis

In order to prove that the algorithm 1 converges to the proper solution, we use results from statistical learning. More specifically, we benefit from the convergence of Stochastic Gradient Descent (SGD) algorithms [25]. To prove the convergence, let $E(w) = \sum_\mu |x^\mu \cdot w|^2$ be the cost function. Furthermore, let $A = \mathbb{E}\{x^T x\}$ be the correlation patterns for the patterns in the training set. Therefore, due to uniformity assumption for the patterns in the training set, one can rewrite $E(w) = w^T A w$, where we have omitted the normalizing $1/C$ constant for simplicity. finally, denote $A_\mu = x^m u(x^\mu)^T$. Now consider the following assumptions:

- A1. $\|A\|_2 \leq \Upsilon < \infty$ and $\sup_\mu \|A_\mu\|_2 = \|x^\mu\|^2 \leq \zeta < \infty$.
A2. $\alpha_t \geq 0$, $\sum \alpha_t = \infty$ and $\sum \alpha_t^2 < \infty$.

The following lemma proves the convergence of the learning algorithm 1 to a local minimum w^* .

Lemma 1. *Let assumption A1 and A2 above hold. Then, the algorithm 1 converges to a local minimum for which $\nabla C(w) = 0$.*

Proof: To prove the lemma, we use the convergence results in [25] and show that the necessary assumptions for the convergence holds for the proposed algorithm. To start, assumption (i) holds trivially as the cost function is three-times differentiable, with continuous derivatives. Furthermore, $C(w) \geq 0$. Assumption (ii) holds because of our choice of the step size α_t , as mentioned in the lemma description.

Assumption (iii) ensures that the vector w could not escape by becoming larger and larger. Due to the constraint $\|w\| = 1$, this assumption holds as well. Assumption (iv) holds as well because:

$$\begin{aligned} \mathbb{E}_\mu (2A_\mu w + \beta \Gamma(w))^2 &= 4w^T \mathbb{E}_\mu (A_\mu^2) w + \beta^2 \|\Gamma(w)\|^2 + 4\beta w^T \mathbb{E}_\mu (A_\mu \Gamma(w)) \\ &\leq 4\|w\|^2 \zeta^2 + \beta^2 \|w\|^2 + 4\beta \Upsilon \|w\|^2 \\ &= \|w\|^2 (4\zeta^2 + 4\beta \Upsilon + \beta^2) \end{aligned}$$

Finally, assumption (v) holds because:

$$\begin{aligned} \|2A_\mu w + \beta \Gamma(w)\|^2 &= 4w^T A_\mu^2 w + \beta^2 \|\Gamma(w)\|^2 + 4\beta w^T A_\mu \Gamma(w) \\ &\leq \|w\|^2 (4\zeta^2 + 4\beta \zeta + \beta^2) \end{aligned} \quad (6)$$

Therefore, $\exists E > D$ such that as long as $|w|^2 < E$:

$$\sup_{\|w\|^2 < E} \|2A_\mu w + \beta \Gamma(w)\|^2 \leq E(2\zeta + \beta)^2 = \text{constant} \quad (7)$$

Since all necessary assumptions hold for the learning algorithm 1, it converges to a local minimum where $\nabla C(w) = 0$. ■

Next, we prove the desired result, i.e. the fact that in the local minimum, the resulting weight vector is orthogonal to the patterns, i.e. $Aw = 0$.

Theorem 2. *In the local minimum where $\nabla C(w^*) = 0$, the optimal vector w^* is orthogonal to the patterns in the training set, i.e. $Aw^* = 0$.*

Proof: Since $\nabla C(w^*) = 2Aw^* + \beta \Gamma(w^*) = 0$, we have:

$$w^* \cdot \nabla C(w^*) = 2(w^*)^T A w^* + \beta w^* \cdot \Gamma(w^*) \quad (8)$$

The first term is always greater than or equal to zero. Now as for the second term, we have that $|\Gamma(w_i)| \leq |w_i|$ and $\text{sgn}(w_i) = \text{sgn}(\Gamma(w_i))$, where w_i is the i^{th} entry of w . Therefore, $0 \leq w^* \cdot \Gamma(w^*) \leq \|w^*\|^2$. Therefore, both terms on the right hand side of (8) are greater than or equal to zero. And since the left hand side is known to be equal to zero, we conclude that $(w^*)^T A w^* = 0$ and $\Gamma(w^*) = 0$. The former means $(w^*)^T A w^* = \sum_\mu (w^* \cdot x^\mu)^2 = 0$. Therefore, we must have $w^* \cdot x^\mu = 0$, for all $\mu = 1, \dots, C$. Which simply means the vector w^* is orthogonal to all the patterns in the training set. ■

Remark 3. Note that the above theorem only proves that the obtained vector is orthogonal to the data set and says nothing about its degree of sparsity. The reason is that there is no guarantee that the dual basis of a subspace be sparse. The introduction of the penalty function $g(w)$ in problem (4a) only encourages sparsity by suppressing the small entries of w , i.e. shifting them towards zero if they are really small or leaving them intact if they are rather large. And from the fact that $\Gamma(w^*) = 0$, we know this is true as the entries in w^* are either large or zero, i.e. there are no small entries. Our experimental results in section VI show that in fact this strategy works perfectly and the learning algorithm results in sparse solutions.

C. Avoiding the all-zero solution

Although in problem (2) we have the constraint $\|w\|_2 = 1$ to make sure that the algorithm does not converge to the trivial solution $w = \mathbf{0}$, due to approximations we made when developing the optimization algorithm, we should make sure to choose the parameters such that the all-zero solution is still avoided.

To this end, denote $w'(t) = w(t) - \alpha_t y(t)(x(t) - \frac{y(t)w(t)}{\|w(t)\|^2})$ and consider the following inequalities:

$$\begin{aligned} \|w(t)\|^2 &= \|w(t) - \alpha_t y(t)(x(t) - \frac{y(t)w(t)}{\|w(t)\|^2}) - \alpha_t \beta \Gamma(w(t))\|^2 \\ &= \|w'(t)\|^2 + \alpha_t^2 \beta^2 \|\Gamma(w(t))\|^2 - 2\alpha_t \beta \Gamma^T(w(t))w'(t) \\ &\geq \|w'(t)\|^2 - 2\alpha_t \beta \Gamma^T(w(t))w'(t) \end{aligned}$$

Now in order to have $\|w(t)\|^2 > 0$, we must have that $2\alpha_t \beta \|\Gamma(w(t))\|^T w'(t) \leq \|w'(t)\|^2$. Given that, $|\Gamma^T(w(t))w'(t)| \leq \|w'(t)\| \|\Gamma(w(t))\|$, it is therefore sufficient to have $2\alpha_t \beta \|\Gamma(w(t))\| \leq \|w'(t)\|$. On the other hand, we have:

$$\begin{aligned} \|w'(t)\|^2 &= \|w(t)\|^2 + \alpha_t^2 y(t)^2 \|x(t) - \frac{y(t)w(t)}{\|w(t)\|^2}\|^2 \\ &\geq \|w(t)\|^2 \end{aligned} \quad (10)$$

As a result, in order to have $\|w(t)\|^2 \geq 0$, it is sufficient to have $2\beta_t \|\Gamma(w(t))\| \leq \|w(t)\|$. Finally, since we have $\Gamma(w(t)) \leq w(t)$ (entriewise), we know that $\|\Gamma(w(t))\| \leq \|w(t)\|$. Therefore, having $2\alpha_t \beta < 1 \leq \|w(t)\|/\|\Gamma(w(t))\|$ ensures $\|w(t)\| > 0$.

Remark 4. Interestingly, the above choice for the function $\Gamma(w)$ looks very similar to the soft thresholding function (11) by [1] to perform iterative compressed sensing. The authors show that their choice of the sparsity function is very competitive in the sense that one can not get much better results by choosing other thresholding functions. However, one main difference between their work and that of ours is that we enforce the sparsity as a penalty in equation (3b) while they apply the soft thresholding function in equation (11) to the whole w , i.e. if the updated value of w is larger than a

threshold, it is left intact while it will be put to zero otherwise.

$$\eta_t(x) = \begin{cases} x - \theta_t & \text{if } x > \theta_t; \\ x + \theta_t & \text{if } x < -\theta_t \\ 0 & \text{otherwise.} \end{cases} \quad (11)$$

where θ_t is the threshold at iteration t .

D. Making the Algorithm Parallel

In order to find m constraints, we need to repeat algorithm 1 m times. Fortunately, we can repeat this process in parallel, which speeds up the algorithm and is more meaningful from a biological point of view as each constraint neuron can act independently of other neighbors. Although doing the algorithm in parallel may result in redundant constraints once in a while, our experimental results show that starting from different random initial points, the algorithm converges to different distinct constraints most of the time. Besides, as long as we have enough distinct constraints, the recall algorithm in the next section works just fine and there is no need to learn all the distinct basis vectors of null space defined by the training patterns. Therefore, we will use the parallel version to have a faster algorithm in the end.

IV. RECALL PHASE

In the recall phase, we are going to exploit the facts that the connectivity matrix of the neural graph is sparse and orthogonal to the memorized patterns. Therefore, given a noisy version of the learned patterns, we can use the feedback at the constraint neurons in Fig. 1 to eliminate noise. More specifically, the linear input sum to the constraint neurons is given by the vector $W \cdot (x^\mu + z) = W \cdot x^\mu + W \cdot z = W \cdot z$, with z being the integer-valued input noise (biologically speaking, the noise can be interpreted as a neuron skipping some spikes or firing more spikes than it should). We will use sparsity of the neural graph to eliminate z iteratively, as explained in the following.

A. The Recall Algorithms

The proposed algorithm for the recall phase comprises a series of forward and backward iterations. Two different methods are suggested in this paper, which slightly differ from each other in the way pattern neurons are updated. The first one is based on the Winner-Take-All approach (WTA) and is given by algorithm 2. In this version, only the pattern node that receives the highest amount of feedback updates its state while the other pattern neurons maintain their current states. The winner-take-all circuitry can be easily added to the neural model shown in Figure 1 using any of the classic WTA methods [7].

The second approach, given by algorithm 3, is much simpler and in every iteration, each pattern neuron decides locally to update its current state or not. More specifically, if the amount of feedback received by a pattern neuron exceeds a threshold, the neuron updates its state and remain intact otherwise.

²Note that in practice, we replace the condition $h_i = 0$ and $h_i > 0$ with $|h_i| < \varepsilon$ and $h_i > \varepsilon$ for some small positive number ε .

Algorithm 2 Recall Algorithm: Winner-Take-All

Input: Connectivity matrix W , iteration t_{\max} **Output:** x_1, x_2, \dots, x_n

- 1: **for** $t = 1 \rightarrow t_{\max}$ **do**
- 2: *Forward iteration:* Calculate the weighted input sum $h_i = \sum_{j=1}^n W_{ij}x_j$, for each constraint neuron y_i and set:

$$y_i = \begin{cases} 1, & h_i < 0 \\ 0, & h_i = 0 \\ -1, & \text{otherwise}^2 \end{cases}.$$

- 3: *Backward iteration:* Each neuron x_j with degree d_j computes

$$g_j = \frac{\sum_{i=1}^m W_{ij}y_i}{d_j}.$$

- 4: Find

$$j^* = \arg \max_j |g_j|.$$

- 5: Update the state of winner j^* : If $g_{j^*} \neq 0$, then set $x_{j^*} = x_{j^*} + \text{sign}(g_{j^*})$.
 - 6: $t \leftarrow t + 1$
 - 7: **end for**
-

Algorithm 3 Recall Algorithm: Bit-Flipping

Input: Connectivity matrix W , threshold φ , iteration t_{\max} **Output:** x_1, x_2, \dots, x_n

- 1: **for** $t = 1 \rightarrow t_{\max}$ **do**
- 2: *Forward iteration:* Calculate the weighted input sum $h_i = \sum_{j=1}^n W_{ij}x_j$, for each neuron y_i and set:

$$y_i = \begin{cases} 1, & h_i < 0 \\ 0, & h_i = 0 \\ -1, & \text{otherwise}^2 \end{cases}.$$

- 3: *Backward iteration:* Each neuron x_j with degree d_j computes

$$g_j = \frac{\sum_{i=1}^m W_{ij}y_i}{d_j}.$$

- 4: Update the state of each pattern neuron j according to $x_j = x_j + \text{sgn}(g_j)$ only if $|g_j| > \varphi$.
 - 5: $t \leftarrow t + 1$
 - 6: **end for**
-

It is worthwhile to mention that the bit-flipping algorithm is very similar to the bit-flipping algorithm of Sipser and Spielman [24] and a similar approach in [23] for compressive sensing methods.

Remark 5. *To give the reader some insight about why the neural graph should be sparse in order for the above algorithms to work, consider the backward iteration of both algorithms: they are based on counting the fraction of received input feedbacks from the neighbors of a pattern node. In the*

extreme case, if the neural graph is complete, then a single noisy pattern node results in the violation of all constraint nodes in the forward iteration. As a result, in the backward iteration all the pattern nodes receive feedback from their neighbors and it is impossible to tell which of the pattern nodes is the noisy one.

However, if the graph is sparse, a single noisy pattern node only makes some of the constraints unsatisfied. Consequently, in the recall phase only the nodes which share the neighborhood of the noisy node receive input feedbacks. And the fraction of the received feedbacks would be much larger for the original noisy node. Therefore, by merely looking at the fraction of received feedback from the constraint nodes, one can identify the noisy pattern nodes with high probability as long as the graph is sparse and input noise is fairly bounded.

B. Performance analysis

In order to find out analytical estimations on the recall probability of error as a function of input noise, we assume the connectivity graph W to be sparse. If W is also an expander-graph, we get better performance guarantees. However, finding an expander graph by means of an iterative learning method, assuming its existence, is very difficult. Therefore, in this section we limit ourselves to the case of sparse graphs and the analysis for the expander graphs can be found in appendix B and in [2].

To start the analysis, we consider an ensemble of random neural graphs with a given right-degree distribution and investigate the average error performance. In other words, W is assumed to be randomly constructed in the sense that in each construction round, a pattern node is selected and based on its degree, it connects randomly to constraint nodes. The degree-distribution is determined by the result of the learning algorithm so we do not aim to optimize the degree distribution here to get a better error performance.

Finally, we assume that the errors do not cancel each other out in the constraint nodes. This is in fact a realistic assumption because the neural graph is weighted, with weights belonging to the real field, and the noise values are integers.

Now we take the following steps in order to bound the probability of making an error *in one iteration* for the bit-flipping algorithm:

- 1) We start by dividing the errors in two types: a correct node updates itself mistakenly (P_{e_1}) and a noisy node update itself in the wrong direction (P_{e_2}).
- 2) We find an explicit relationship for the average of P_{e_1} over all nodes as a function of the number of constraint neurons affected by the noisy pattern nodes.
- 3) We then find an upper *bound* on the average of P_{e_2} as a function of the number of constraints affected by the noisy pattern nodes.
- 4) Having found the average bit error probability, we can estimate the block error probability for the recall algorithm.

Note that if we choose the bit-flipping update threshold $\varphi = 1$,

roughly speaking, we will have the winner-take-all algorithm.² Therefore, it is sufficient to analyze the bit-flipping method.

Let \mathcal{E}_t denote the set of erroneous pattern nodes at iteration t , and $\mathcal{N}(\mathcal{E}_t)$ be the set of constraint nodes that are connected to the nodes in \mathcal{E}_t , i.e. these are the constraint nodes that have at least one neighbor in \mathcal{E}_t . In addition, let $\mathcal{N}^c(\mathcal{E}_t)$ denote the other constraint nodes that do not have any connection to any node in \mathcal{E}_t . Finally, Let \mathcal{C}_t be the set of correct pattern nodes.

Based on the error correcting algorithm and the above notations, at a given iteration two types of error are possible:

- 1) A node $x \in \mathcal{C}_t$ decides to update its value. The probability of this phenomenon is denoted by P_{e_1} .
- 2) A node $x \in \mathcal{E}_t$ updates its value in the wrong direction.

Let P_{e_2} denote the probability of error for this type.

1) *Error probability - type 1:* To begin, let P_1^x be the probability that a node $x \in \mathcal{C}_t$ with degree d_x makes a correct decision and *does not* updates its state. We have:

$$P_1^x = \Pr\left\{\frac{|\mathcal{N}(\mathcal{E}_t) \cap \mathcal{N}(x)|}{d_x} < \varphi\right\} \quad (12)$$

where $\mathcal{N}(x)$ is the neighborhood of x . Assuming random construction of the graph and relatively large graph sizes, one can approximate P_1^x by

$$P_1^x = \sum_{i=0}^{\lceil \varphi d_x \rceil - 1} \binom{d_x}{i} \left(\frac{S_t}{m}\right)^i \left(1 - \frac{S_t}{m}\right)^{d_x - i} \quad (13)$$

Where $S_t = \mathbb{E}(|\mathcal{N}(\mathcal{E}_t)|)$ is the average neighborhood size of \mathcal{E}_t .

As a result of the above equations, we have:

$$P_{e_1} = 1 - \mathbb{E}_{d_x}(P_1^x) \quad (14)$$

2) *Error probability - type 2:* A node $x \in \mathcal{E}_t$ makes a wrong decision if the net input sum it receives has a different sign than the sign of noise it experiences. Instead of finding an exact relation, we bound this probability by the probability that the neuron x shares at least half of its neighbors with other neurons, i.e. $P_{e_2} \leq \Pr\left\{\frac{|\mathcal{N}(\mathcal{E}_t^*) \cap \mathcal{N}(x)|}{d_x} \geq 1/2\right\}$, where $\mathcal{E}_t^* = \mathcal{E}_t \setminus x$. Letting $P_2^x = 1 - \Pr\left\{\frac{|\mathcal{N}(\mathcal{E}_t^*) \cap \mathcal{N}(x)|}{d_x} \geq 1/2\right\}$, we will have:

$$P_2^x = \sum_{i=0}^{\lfloor d_x/2 \rfloor} \binom{d_x}{i} \left(\frac{S_t}{m}\right)^i \left(1 - \frac{S_t}{m}\right)^{d_x - i} \quad (15)$$

Therefore, we will have:

$$P_{e_2} \leq 1 - \mathbb{E}_{d_x}(P_2^x) \quad (16)$$

²It must be mentioned that choosing $\varphi = 1$ does not yield the WTA algorithm exactly because in the original WTA, only one node is updated in each round. However, in this version with $\varphi = 1$, all nodes that receive feedback from all their neighbors are updated. Nevertheless, the performance of the both algorithms is rather similar.

Combining equations (14) and (14), the symbol error probability at iteration would be

$$\begin{aligned} P_b(t) &= \Pr\{x \in \mathcal{C}_t\}P_{e_1} + \Pr\{x \in \mathcal{E}_t\}P_{e_2} \\ &= \frac{n - |\mathcal{E}_t|}{n}P_{e_1} + \frac{|\mathcal{E}_t|}{n}P_{e_2} \\ &= 1 - \frac{n - |\mathcal{E}_t|}{n}\bar{P}_1^x - \frac{|\mathcal{E}_t|}{n}\bar{P}_2^x \end{aligned} \quad (17)$$

where $\bar{P}_i^x = \mathbb{E}_{d_x}\{P_i^x\}$.

And finally, the average block error rate is given by the probability of at least one pattern node x makes a mistake. Therefore:

$$P_e(t) = 1 - (1 - P_b(t))^n \quad (18)$$

Equation (18) gives the probability of making a mistake in iteration t . Therefore, we can bound the overall probability of error as

$$P_E \leq P_e(0) \quad (19)$$

In other words, if we made a mistake in the first round, we assume everything will go wrong till the end. Obviously, this is not necessarily true and in practice one might do much better. In fact simulation results confirm this conjecture. However, as the initial number of noisy nodes grow, the above bound becomes tight.

Using the following lemma we will have the average neighborhood size which is widely used in deriving the probabilities of error above.

Lemma 3. *The average neighborhood size S_t in iteration t is given by:*

$$S_t = m \left(1 - \left(1 - \frac{\bar{d}}{m}\right)^{m c E_t}\right) \quad (20)$$

Proof: The proof is given in appendix C. ■

Remark 6. *Although the randomness assumption may not be correct due to the fact that the resulting graphs may not correspond to the dual subspace of the given training set, over all ensembles of random training sets this assumption might hold. Our simulation results in the following confirm this assumption as well. Nevertheless, for a particular realization of the neural graph, one might do better or worse than average.*

C. Some Practical Modifications

D. Simultaneous Learn and Recall

So far, the algorithm yields a sparse matrix which satisfies the set of linear equations $Wx = 0$, for all vectors x in the training data set. However, since we have to go over an exponentially many such vectors (and for each vector update m vectors of length n), the whole learning process is extremely slow even for rather small values of k . Therefore, instead of going over all possible patterns in the subspace, we pick a portion of them at random, perform the learning phase and when finished, go directly to the recall phase.

Now during the recall phase, two cases may occur: either we are given a noisy version of a pattern we have memorized

before, or we will have a noisy version of a new pattern. One can distinguish these two cases based on the values returned by constraint nodes. In the former case, the deviation from the desired constraint values is rather small if the amount of input noise is limited to some extent. Therefore, we can carry on with the usual error correcting algorithm in the recall phase.

In the latter case, where we are given a noisy version of a new pattern, the deviation from the constraints are usually high, indicating the existence of a pattern not learned before. In this case, we perform a learning phase. Note that since the pattern we are learning is usually noisy, we might have to repeat this process several times for a pattern to cancel out learning noise. Obviously, there is a trade off between accuracy and speed here, but this scenario makes sense biologically as well. Because we do not learn everything at once but rather do the process in a gradual manner. In the next section, we compare this approach with the traditional training in terms of speed and recall performance.

V. PATTERN RETRIEVAL CAPACITY

It is interesting to see that the number of patterns C does not have any effect in the learning or recall algorithm. Because as long as the patterns come from a subspace, the learning algorithm will yield a matrix which is orthogonal to all of the patterns in the training set. And in the recall phase, all we deal with is Wz , with z being the noise which is independent of the patterns.

Therefore, in order to show that the pattern retrieval capacity is exponential with n , all we need to show is that there exists a valid training set \mathcal{X} with C patterns of length n for which $C \propto a^n$, for some $a > 1$. By valid we mean the patterns should come from a subspace with dimension $k < n$ and the entries in the patterns should be non-negative integers. The next theorem proves the required result.

Theorem 4. *Let \mathcal{X} be a $C \times n$ matrix which comprise of C vectors of length n with non-negative integers entries between 0 and $S - 1$. Then, there exists an \mathcal{X} for which $C \propto a^n$, with $a > 1$, such that $\text{rank}(\mathcal{X}) = k < n$.*

Proof: The proof is based on construction: we construct a data set \mathcal{X} with the required properties. To start, consider a matrix $G \in \mathbb{R}^{k \times n}$ with $\text{rank } k = n$. Assume $k = rn$, with $0 < r < 1$. Let the entries of G be non-negative integers, between 0 and $\gamma - 1$. We start constructing the patterns in the data set as follows: consider a random vector $u^\mu \in \mathbb{R}^k$ with integer-valued-entries between 0 and $v - 1$. We set the pattern $x^\mu \in \mathcal{X}$ to be $x^\mu = u^\mu \cdot G$, if all the entries of x^μ are between 0 and $S - 1$. Obviously, since both u^μ and G have only non-negative entries, all entries in x^μ are non-negative. Therefore, it is the $S - 1$ upper bound that we have to worry about.

The j^{th} entry in x^μ is equal to $x_j^\mu = u^\mu \cdot g_j$, where g_j is the j^{th} column of G . Suppose g_j has d_j non-zero elements. Then, we have:

$$x_j^\mu = u^\mu \cdot g_j \leq d_j(\gamma - 1)(v - 1)$$

Therefore, denoting $d^* = \max_j d_j$, we could choose γ , v and d^* such that

$$S - 1 \geq d^*(\gamma - 1)(v - 1) \quad (21)$$

to ensure all entries of x^μ are less than S . In this case, we will have $v^k = v^{rn}$ patterns belonging to \mathcal{X} . Therefore, $C \geq v^{rn}$, which would be an exponential number in n for $v \geq 2$.

As an example, if G is selected to be a sparse 200×400 matrix with 0/1 entries (i.e. $\gamma = 2$) and $d^* = 10$, and u is also chosen to be a vector with 0/1 elements (i.e. $v = 2$), then it is sufficient to have $S \geq 11$, i.e. the maximum firing rate of neurons should be 11 to have a pattern retrieval capacity of $C = 2^{rn}$. ■

Remark 7. *Note that the inequality (21) was obtained for the worst-case scenario and in fact is very loose. Therefore, even if it does not hold, we will still be able to memorize a very large number of patterns since a big portion of the generated vectors x^μ will have entries less than S . These vectors correspond the message vectors u^μ who are "sparse" as well, i.e. do not have all entries greater than zero. The number of such vectors is a polynomial in n , the degree of which depends on the number of non-zero entries in u^μ .*

VI. SIMULATION RESULTS

A. Simulation Scenario

We have simulated the proposed learning and recall algorithms for three different network sizes $n = 100, 200, 400$, with $k = n/2$ for all cases. For each case, we considered three different set-ups with different values for β and θ in the learning algorithm 1, and different φ for the bit-flipping recall algorithm 3.

In all cases, we generated 50 training set at random using the approach explained in the proof of theorem 4, i.e. we generated a generator matrix G at random with 0/1 entries and $d^* = 10$ (??? verify!). We also used 0/1 generating message words u and put $S = 11$ to ensure the validity of the generated training set.

However, since in this setup we will have 2^n patterns to memorize, doing a simulation over all of them would take forever. Therefore, we have selected a random sample sub-set \mathcal{X} each time with size $C = 10^5$ for each of the 50 generated sets and used this subset as the training set.

For each set-up, we performed the learning algorithm over these 50 instances, and investigated the average sparsity of the learned constraints over the ensemble of 50 instances. As explained earlier, all the constraints for each network were learned in parallel, i.e. to obtain $m = n - k$ constraints, we executed algorithm 1 from random initial points each time.

As for the recall algorithms, the error correcting performance was assessed for each set-up, averaged over the 50 ensembles. The empirical results are compared to the theoretical bounds derived in section IV-B as well.

The recall algorithm is first executed just over the patterns in the training set, i.e. the patterns that already have been learned, to assess the true error correction performance of the algorithm

and compare it with the theoretical bounds. Afterwards, and to investigate the performance of the simultaneous learn and recall scenario, explained in section IV-C, the recall algorithm is executed over not just the training set but over all of the possible 2^n patterns. In this case, in response to a noisy pattern if the network decides it has learned the pattern before, it attempts to eliminate the noise. Otherwise, it tries to learn the pattern and will not do error correction. The performance of these two cases are compared to each other as well.

B. Learning Phase Results

In the learning algorithm, we pick each pattern from the training set and adjust the weights accordingly. Once we have gone over all the patterns, we repeat this operation again several times to make sure that update for one pattern does not adversely affect the other learned patterns. Let t denote the number of times we go over the patterns in the training set. Then we set $\alpha_t \propto 1/t$ to ensure the conditions of theorem 1 is satisfied. Interestingly, all of the constraints converged in at most two learning iterations for all different set-ups. Therefore, the learning is very fast in this case.

Figure 2 illustrates the percentage of variable nodes with the specified sparsity measure defined as $\rho = \kappa/n$, where κ is the number of non-zero elements. From the figure one notices that as n increases, the weight vectors become sparser.

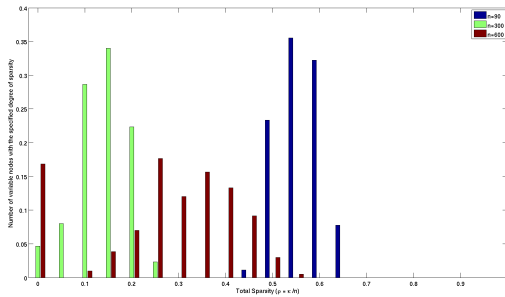


Fig. 2. The percentage of variable nodes with the specified sparsity measure and different values of network sizes. The sparsity measure is defined as $\rho = \kappa/n$, where κ is the number of non-zero elements.

C. Recall Phase Results

For the recall phase, in each trial we pick a pattern randomly from the training set, corrupt a given number of its symbols with ± 1 noise and use the suggested algorithm to correct the errors. A pattern error is declared if the output does not match the correct pattern.

Figure 3 illustrates the pattern error rates for different network sizes. Note that the results are in close match with the theoretical upper bound derived in section IV-B. Note that since we stop the learning after 99% of the patterns had learned, it is natural to see some recall errors even for 1 initial erroneous node.

Figure 4 shows the recall error rates for different network sizes. Note that a recall error is declared if the network decides

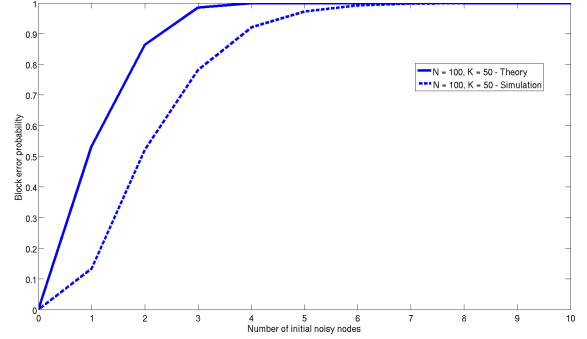


Fig. 3. Pattern error rate against the initial number of erroneous nodes and comparison with theoretical upper bounds

that it has learned the given pattern before and makes an attempt to correct the errors but fails to succeed.

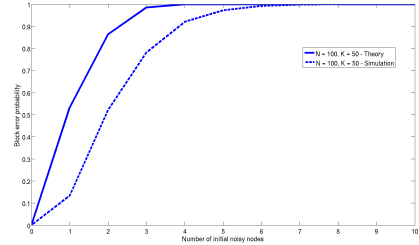


Fig. 4. Pattern error rate against the initial number of erroneous nodes and comparison with theoretical upper bounds

If the network mistakes between an already learned pattern and a new one, a learning decision error is declared. Table VI-C indicates the fraction of incorrect learning decisions (both false positives and false negatives) of the algorithm.

TABLE I
SIMULATION PARAMETERS

Parameter	φ	t_{\max}	ε
Value	0.8	$20\ z\ _0$	0.01

VII. FUTURE WORKS

In order to extend the multi-level neural network, we must first find a way to generate patterns that belong to a subspace with dimensions $nL - m_g$, where m_g lies within the inside of bounds $L(n - k) < m_g < nL - k$. This will give us a way to investigate the trade off between the maximum number of memorizable patterns and the degree of error correction possible.

Furthermore, so far we have assumed that the second level enforces constraints in the same space. However, it is possible that the second level imposes a set of constraints in a totally different space. For this purpose, we need a mapping from one

space into another. A good example is the written language. While they are local constraints on the spelling of the words, there are some constraints enforced by the grammar or the overall meaning of a sentence. The latter constraints are not on the space of letters but rather the space of grammar or meaning. Therefore, in order to for instance to correct an error in the word *_at*, we can replace *_* with either *W*, to get *hat*, or *c* to get *cat*. Without any other clue, we can not find the correct answer. However, let's say we have the sentence "The *_at* ran away". Then from the constraints in the space of meaning we know that the subject must be an animal or a person. Therefore, we can return *cat* as the correct answer. Finding a proper mapping is the subject of our future work.

ACKNOWLEDGMENT

The authors would like to thank Prof. Amin Shokrollahi for helpful comments and discussions. This work was supported by Grant 228021-ECCSciEng of the European Research Council.

APPENDIX A EXPANDER GRAPHS

This section contains the definitions and the necessary background on expander graphs.

Definition 1. A regular (d_p, d_c, n, m) bipartite graph W is a bipartite graph between n pattern nodes of degree d_p and m constraint nodes of degree d_c .

Definition 2. An $(\alpha n, \beta d_p)$ -expander is a (d_p, d_c, n, m) bipartite graph such that for any subset \mathcal{P} of pattern nodes with $|\mathcal{P}| < \alpha n$ we have $|\mathcal{N}(\mathcal{P})| > \beta d_p |\mathcal{P}|$ where $\mathcal{N}(\mathcal{P})$ is the set of neighbors of \mathcal{P} among the constraint nodes.

The following result from [24] shows the existence of families of expander graphs with parameter values that are relevant to us.

Theorem 5. [24] Let W be a randomly chosen (d_p, d_c) -regular bipartite graph between n d_p -regular vertices and $m = (d_p/d_c)$ d_c -regular vertices. Then for all $0 < \alpha < 1$, with high probability, all sets of αn d_p -regular vertices in W have at least

$$n \left(\frac{d_p}{d_c} (1 - (1 - \alpha)^{d_c}) - \sqrt{\frac{2d_c \alpha h(\alpha)}{\log_2 e}} \right)$$

neighbors, where $h(\cdot)$ is the binary entropy function.

The following result from [26] shows the existence of families of expander graphs with parameter values that are relevant to us.

Theorem 6. Let d_c, d_p, m, n be integers, and let $\beta < 1 - 1/d_p$. There exists a small $\alpha > 0$ such that if W is a (d_p, d_c, n, m) bipartite graph chosen uniformly at random from the ensemble of such bipartite graphs, then W is an $(\alpha n, \beta d_p)$ -expander with probability $1 - o(1)$, where $o(1)$ is a term going to zero as n goes to infinity.

APPENDIX B

ANALYSIS OF THE RECALL ALGORITHMS FOR EXPANDER GRAPHS

A. Analysis of the Winner-Take-All Algorithm

We prove the error correction capability of the winner-take-all algorithm in two steps: first we show that in each iteration, only pattern neurons that are corrupted by noise will be chosen by the winner-take-all strategy to update their state. Then, we prove that the update is in the right direction, i.e. toward removing noise from the neurons.

Lemma 8. If the constraint matrix W is an $(\alpha n, \beta d_p)$ expander and the original number of erroneous neurons are less than $\epsilon_{min} = \lfloor \frac{\beta}{1-\beta} \rfloor$, then in each iteration of the winner-take-all algorithm only the corrupted pattern nodes update their value and the other nodes remain intact. For $\beta = 3/4$, the algorithm will always pick the correct node if we have two or fewer erroneous nodes.

Proof: For simplicity, we restrict our attention to the case $\beta = 3/4$. If we have only one node x_i in error, it is obvious that the corresponding node will always be the winner of the winner-take-all algorithm unless there exists another node that has the same set of neighbors as x_i . However, this is impossible as because of the expansion properties, the neighborhood of these two nodes must have at least $2\beta d_p$ members which for $\beta = 3/4$ is equal to $3d_p/2$. As a result, no two nodes can have the same neighborhood and the winner will always be the correct node.

In the case where there are two erroneous nodes, say x_i and x_j , let Q be the set $\{x_i, x_j\}$ and $\mathcal{N}(Q)$ be the corresponding neighborhood on the constraint nodes side. Furthermore, assume x_i and x_j share $d_{p'}$ of their neighbors so that $|\mathcal{N}(Q)| = 2d_p - d_{p'}$. First of all note that because of the expansion properties and for $\beta = 3/4$:

$$|\mathcal{N}(Q)| = 2d_p - d_{p'} > 2\beta d_p \Rightarrow d_{p'} < d_p/2.$$

Now we have to show that there are no nodes other than x_i and x_j that can be the winner of the winner-take-all algorithm. To this end, note that only those nodes that are connected to $N(Q)$ will receive some feedback and can hope to be the winner of the process. So let's consider such a node x_ℓ that is connected to d_{p_ℓ} of the nodes in $N(Q)$. Let Q' be $Q \cup \{x_\ell\}$ and $\mathcal{N}(Q')$ be the corresponding neighborhood. Because of the expansion properties we have $|\mathcal{N}(Q')| = d_p - d_{p_\ell} + |\mathcal{N}(Q)| > 3\beta d_p$. Thus:

$$d_{p_\ell} < d_p + |\mathcal{N}(Q)| - 3\beta d_p = 3d_p(1 - \beta) - d_{p'}.$$

Now, note that the nodes x_i and x_j will receive some feedback from at least $d_p - d_{p'}$ edges because those are the edges that are uniquely connected to them and noise from the other erroneous nodes cannot cancel them out. Since $d_p - d_{p'} > 3d_p(1 - \beta) - d_{p'}$, for $\beta = 3/4$, we conclude that $d_p - d_{p'} > d_{p_\ell}$ which proves that no node outside Q can be picked during the winner-take-all algorithm as long as $|Q| \leq 2$ for $\beta = 3/4$. ■

In the next lemma, we show that the state of erroneous neurons is updated in the direction of reducing the noise.

Lemma 9. *If the constraint matrix W is an $(\alpha n, \beta d_p)$ expander and the original number of erroneous neurons is less than $e_{min} = \lfloor \frac{\beta}{1-\beta} \rfloor$, then in each iteration of the winner-take-all algorithm the winner is updated toward reducing the noise.*

Proof: As before, we only focus on the case $\beta = 3/4$. When there is only one erroneous node, it is obvious that all its neighbors agree on the direction of update and the node reduces the amount of noise by one unit.

If there are two nodes x_i and x_j in error, since the number of their shared neighbors is less than $d_p/2$ (as we proved in the last lemma), more than half of their neighbors agree on the direction of update. Therefore, whoever the winner is will be updated to reduce the amount of noise by one unit. ■

The following theorem sums up the results of the previous lemmas to show that the winner-take-all algorithm is guaranteed to perform error correction.

Theorem 7. *If the constraint matrix W is an $(\alpha n, \beta d_p)$ expander, then the winner-take-all algorithm is guaranteed to correct at least $e_{min} = \lfloor \frac{\beta}{1-\beta} \rfloor$ positions in error, irrespective of the magnitudes of the errors.*

Proof: The proof is immediate from Lemmas 8 and 9. ■

B. Analysis of the Bit-Flipping Algorithm

Roughly speaking, one would expect the bit-flipping algorithm to be sub-optimal in comparison to the winner-take-all strategy, since the pattern neurons need to make independent decisions, and are not allowed to cooperate amongst themselves. In this subsection, we show that despite this restriction, the bit-flipping algorithm is capable of error correction; the suboptimality in comparison to the winner-take-all algorithm can be quantified in terms of a larger expansion factor β being required for the graph.

Theorem 8. *If the constraint matrix W is an $(\alpha n, \beta d_p)$ expander with $\beta > \frac{4}{5}$, then the bit-flipping algorithm with $\gamma = \frac{3}{5}$ is guaranteed to correct at least two positions in error, irrespective of the magnitudes of the errors.*

Proof: As in the proof for the winner-take-all case, we will show our result in two steps: first, by showing that for a suitable choice of the bit-flipping threshold γ , that only the positions in error are updated in each iteration, and that this update is towards reducing the effect of the noise.

a) Case 1: First consider the case that only one pattern node x_i is in error. Let x_j be any other pattern node, for some $j \neq i$. Let x_i and x_j have $d_{p'}$ neighbours in common. As argued in the proof of Lemma 8, we have that

$$d_{p'} < 2d_p(1 - \beta). \quad (22)$$

Hence for $\beta = \frac{4}{5}$, x_i receives non-zero feedback from at least $\frac{3}{5}d_p$ constraint nodes, while x_j receives non-zero feedback from at most $\frac{2}{5}d_p$ constraint nodes. In this case, it is clear that

setting $\gamma = \frac{3}{5}$ will guarantee that only the node in error will be updated, and that the direction of this update is towards reducing the noise.

b) Case 2: Now suppose that two distinct nodes x_i and x_j are in error. Let $Q = \{x_i, x_j\}$, and let x_i and x_j share $d_{p'}$ common neighbours. If the noise corrupting these two pattern nodes z_i and z_j are such that $\text{sign}(z_i) = \text{sign}(z_j)$, then both x_i and x_j receive $-\text{sign}(z_i)$ along all d_p edges that they are connected to during the backward iteration. Now suppose that $\text{sign}(z_i) \neq \text{sign}(z_j)$. If $|z_i| = |z_j|$, then x_i (x_j) receives non-zero feedback only from the $d_p - d_{p'}$ edges in $\mathcal{N}(\{x_i\}) \setminus Q$ (resp. $\mathcal{N}(\{x_j\}) \setminus Q$) during the backward iteration, and the feedback is such that the noise is reduced at the end of the update. If on the other hand $|z_i| > |z_j|$, then x_i receives $-\text{sign}(z_i)$ along all d_p of it's incoming links during the backward iteration. However, x_j receives the correct $-\text{sign}(z_j)$ feedback along only the $d_p - d_{p'}$ edges from $\mathcal{N}(\{x_j\}) \setminus Q$, and receives incorrect feedback of $\text{sign}(z_j)$ along the $d_{p'}$ edges from Q .

From the above and from (22), we conclude that when two pattern nodes are in error, at least one of the erroneous pattern nodes receives correct feedback along $d_p - 2d_p(1 - \beta)$ or more edges; and, in the event of a node x_j (say) receiving incorrect feedback along some of it's incoming links, we have that $g_j = (-\text{sign}(z_j)) \frac{d_p - 2d_{p'}}{d_p}$. For $\beta = \frac{4}{5}$, we have from (22) that $\text{sign}(g_j) = -\text{sign}(z_j)$ since $d_p - 2d_{p'} > 0$; hence irrespective of the value of γ , it is not possible that the node x_j is updated such that the noise magnitude is increased.

Let us now examine what happens to a node x_ℓ that is different from the two erroneous nodes x_i, x_j . Suppose that x_ℓ is connected to d_{p_ℓ} nodes in $\mathcal{N}(Q)$. From the proof of Lemma 8, we know that

$$\begin{aligned} d_{p_\ell} &< 3d_p(1 - \beta) - d_{p'} \\ &\leq 3d_p(1 - \beta). \end{aligned}$$

Hence x_ℓ receives at most $3d_p(1 - \beta)$ non-zero messages during the backward iteration.

For $\beta > \frac{4}{5}$, we have that $d_p - 2d_p(1 - \beta) > 3d_p(1 - \beta)$. Hence by setting $\beta = \frac{4}{5}$ and $\gamma = [d_p - 2d_p(1 - \beta)]/d_p = \frac{3}{5}$, it is clear from the above discussion that we have ensured the following in the case of two erroneous pattern nodes:

- At least one erroneous pattern node is updated in each iteration of our algorithm such that the noise magnitude is reduced.
- No erroneous pattern node can be updated such that the noise magnitude is increased.
- No pattern node other than the erroneous pattern nodes is updated. ■

C. Choice of Parameters

In order to put together the results of the previous two subsections and obtain a neural associative scheme that stores an exponential number of patterns and is capable of error

correction, we need to carefully choose the various relevant parameters. We summarize some design principles below.

- From Theorems 6 and ??, the choice of β depends on d_p , according to $\frac{1}{2} + \frac{1}{4d_p} < \beta < 1 - \frac{1}{d_p}$.
- Choose $d_c, S, S', r \triangleq m/n$ such that $S^n > (d_c S)^{rn}$ and $S' > d_c S$, so that Theorem 4 yields an exponential number of patterns.
- For a fixed α , n has to be chosen large enough so that an $(\alpha n, \beta d_p)$ expander exists according to Theorem 6, and so that $\alpha n/2 > e_{min} = \lfloor \frac{\beta}{1-\beta} \rfloor$.

Once we choose a judicious set of parameters according to the above requirements, we have a neural associative memory that is guaranteed to recall an exponential number of patterns even if the input is corrupted by errors in two coordinates. Our simulation results will reveal that a greater number of errors can be corrected in practice.

APPENDIX C

AVERAGE NEIGHBORHOOD SIZE

Now it is time to obtain an expression for $S_e = \mathbb{E}(|\mathcal{N}(\mathcal{E}_t)|)$. To do so, we assume the following procedure for constructing a right-irregular bipartite graph:

- In each iteration, we pick a variable node x with a degree randomly determined according to the given degree distribution.
- Based on the given degree d_x , we pick d_x constraint nodes uniformly at random *with replacement* and connect x to the constraint node.
- We repeat this process n times, until all variable nodes are connected.

Note that the assumption that we do the process with replacement is made to simplify the analysis. This assumption becomes more exact as n grows.

With some abuse of notations, let S_e denote the size of the neighborhood of \mathcal{E} when the size of \mathcal{E} is equal to e . We write S_e recursively in terms of e as follows:

$$\begin{aligned} S_{e+1} &= \mathbb{E}_{d_x} \left\{ \sum_{j=0}^{d_x} \binom{d_x}{j} \left(\frac{S_e}{m}\right)^{d_x-j} \left(1 - \frac{S_e}{m}\right)^j (S_e + j) \right\} \\ &= \mathbb{E}_{d_x} \{ S_e + d_x(1 - S_e/m) \} \end{aligned} \quad (6)$$

nonumber

$$= S_e + \bar{d}(1 - S_e/m)$$

Where $\bar{d} = \mathbb{E}_{d_x} \{d_x\}$ is the average degree of the pattern nodes. In words, the first line calculates the average growth of the neighborhood when a new variable node is added to the graph. Noting that $S_1 = \bar{d}$, one obtains:

$$S_t = m \left(1 - \left(1 - \frac{\bar{d}}{m} \right)^{|mcE_t|} \right) \quad (26)$$

In order to verify the correctness of the above analysis, we have performed some simulation the results of which are given in this section.

Figure 5 illustrates the average neighborhood size in each iteration for a randomly chosen degree distribution with $n =$

400 and $m = 200$ being the number of pattern and constraint nodes, respectively. We generated 100 random graphs and the dashed line shows the average neighborhood size over these graphs. It is obvious that the theoretical value approximates the simulation results rather exactly.

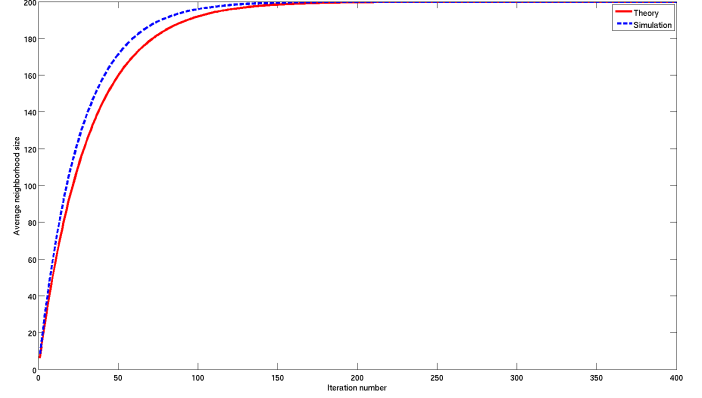


Fig. 5. The theoretical estimation and simulation results for the neighborhood size of an irregular graph with a given degree-distribution for $n = 400$, $m = 200$ and over 2000 random graphs.

REFERENCES

- [1] D. L. Donoho, A. Maleki, A. Montanari, *Message passing algorithms for compressed sensing*, Proc. Nat. Acad. Sci., Vol. 106, 2009, pp. 18914-18919.
- [2] K.R. Kumar, A.H. Salavati and A. Shokrollahi, *Exponential pattern retrieval capacity with non-binary associative memory*, Proc. IEEE Information Theory Workshop, 2011.
- [3] L. Xu, A. Krzyzak, E. Oja, Neural nets for dual subspace pattern recognition method, Int. J. Neur. Syst., Vol. 2, No. 3, 1991, pp. 169-184.
- [4] J. J. Hopfield, *Neural networks and physical systems with emergent collective computational abilities*, Proc. Natl. Acad. Sci., Vol. 79, 1982, pp. 2554-2558.
- [5] J. J. Hopfield, *Neurons with graded response have collective computational properties like those of two-state neurons*, Proc. Natl. Acad. Sci., Vol. 81, No. 10, 1984, pp. 3088 - 3092.
- [6] D. Amit, H. Gutfreund, H. Sompolinsky, *Storing infinite numbers of patterns in a spin-glass model of neural networks*, Physic. Rev. Lett., Vol. 55, 1985, pp. 1530-1533.
- [7] J. Hertz, A. Krogh, R. G. Palmer, *Introduction to the theory of neural computation*, USA: Addison-Wesley, 1991.
- [8] D. O. Hebb, *The organization of behavior*, New York: Wiley & Sons, 1949.
- [9] J. Komlos, R. Paturi, *Effect of connectivity in an associative memory model*, J. Computer and System Sciences, 1993, pp. 350-373.
- [10] M. K. Muezzinoglu, C. Guzelis, J. M. Zurada, *A new design method for the complex-valued multistate Hopfield associative memory*, IEEE Trans. Neur. Net., Vol. 14, No. 4, 2003, pp. 891-899.
- [11] R. McEliece, E. Posner, E. Rodemich, S. Venkatesh, *The capacity of the Hopfield associative memory*, IEEE Trans. Inf. Theory, Jul. 1987.
- [12] A. H. Salavati, K. R. Kumar, W. Gerstner, A. Shokrollahi, *Neural Pre-coding Increases the Pattern Retrieval Capacity of Hopfield and Bidirectional Associative Memories*, To be presented at the IEEE Intl. Symp. Inform. Theory (ISIT - 11), St. Petersburg, Aug 2011.
- [13] S. S. Venkatesh, D. Psaltis, *Linear and logarithmic capacities in associative neural networks*, IEEE Trans. Inf. Theory, Vol. 35, No. 3, 1989, pp. 558-568.
- [14] S. Jankowski, A. Lozowski, J.M., Zurada, *Complex-valued multistate neural associative memory*, IEEE Tran. Neur. Net., Vol. 1, No. 6, 1996, pp. 1491-1496.

- [15] D. L. Lee, *Improvements of complex-valued Hopfield associative memory by using generalized projection rules*, IEEE Tran. Neur. Net., Vol. 12, No. 2, 2001, pp. 439-443.
- [16] E. Oja, T. Kohonen, *The subspace learning algorithm as a formalism for pattern recognition and neural networks*, Neural Networks, Vol. 1, 1988, pp. 277-284.
- [17] P. Peretto, J. J. Niez, *Long term memory storage capacity of multiconnected neural networks*, Biological Cybernetics, Vol. 54, No. 1, 1986, pp. 53-63.
- [18] V. Gripon, C. Berrou, *Sparse neural networks with large learning diversity*, IEEE Trans. on Neural Networks, Vol. 22, No. 7, 2011, pp. 10871096.
- [19] J. Tropp, J. S. J. Wright, *Computational methods for sparse solution of linear inverse problems*, Proc. IEEE, Vol. 98, No. 6, 2010, pp. 948-958.
- [20] T. Richardson and R. Urbanke. *Modern Coding Theory*. Cambridge University Press, 2008.
- [21] E. Cands, T. Tao, *Near optimal signal recovery from random projections: Universal encoding strategies?*, IEEE Trans. on Information Theory, Vol. 52, No. 12, 2006, pp. 5406 - 5425.
- [22] R. Gold, *Optimal binary sequences for spread spectrum multiplexing*, IEEE Trans. Inf. Theory, Vol. 13, No. 4, 1967, pp. 619-621.
- [23] W. Xu, B. Hassibi, *Efficient compressive sensing with deterministic guarantees using expander graphs*, Proc. Inf. Theo. Workshop (ITW), 2007, pp. 414-419.
- [24] M. Sipsper and D. Spielman, *Expander codes*, IEEE Trans. Inform. Theory, Vol. 42, No. 6, pp. 1710-1722, 1996.
- [25] L. Bottou, *Online algorithms and stochastic approximations*, in David Saad, editor, *Online Learning and Neural Networks*. Cambridge University Press, Cambridge, UK, 1998.
- [26] D. Burshtein and G. Miller, *Expander graph arguments for message passing algorithms*, IEEE Trans. Inform. Theory, pp. 782-790, Feb. 2001.
- [27] H. J. Kushner, G. G. Yin, *Stochastic approximation algorithms and applications*, Springer Verlag, 1997.
- [28] E. Oja, J. Karhunen, *On stochastic approximation of eigenvectors and eigenvalues of the expectation of a random matrix*, J. Math. Analysis and Applications, Vol. 106, No. 1, 1985, pp 6984.
- [29] A. H. Salavati, A. Karbasi, *Multi-Level Error-Resilient Neural Networks with Learning*, Submitted to the IEEE Int. Sym. Inf. Theory (ISIT 2012), 2012.