

Topology Optimization in Cellular Neural Networks

Varsha Bhambhani and Herbert G. Tanner

Abstract—This paper presents a constrained combinatorial optimization approach to the design of cellular neural networks with sparse connectivity. The method applies to cases where maintaining links between neurons incurs a cost, which could possibly vary between these links. The interconnection topology of the cellular neural network is diluted without significantly degrading its performance, the latter quantified by the average recall probability for the desired patterns engraved into its associative memory. The dilution process selectively removes the links that contribute the least to a metric related to the size of system’s desired memory pattern attraction regions. The metric used here is the magnitude of the network’s nodes’ stability parameters, which have been proposed as a measure for the quality of memorization. We demonstrate by means of an example that this method of network dilution produces cheaper associative memories that in general trade off performance for cost, and in many cases the performance of the diluted network is on par with the original system.

Index Terms—Cellular neural networks, stability parameters, sparse associative memory.

I. INTRODUCTION

A cellular neural network (CNN) is a nonlinear dynamical system, implementing an associative memory, where neuron interactions are limited to the neighboring units only [1]. Research on the design and synthesis of these types of networks has gained increasing attention in recent years because of their topology being relatively sparse, which makes them a promising choice for very-large-scale integration implementations. Several synthesis procedures for designing symmetric/non-symmetric neural network have been proposed [2]–[4], and one of the main directions of current research on the topic is the stability behavior of these systems in the presence of time delays (see for example [5], [6]), which is however beyond the scope of this paper.

Strengths and weaknesses of several network design techniques are discussed in [2]. The analysis results shown in [3] allow one to locate in a systematic manner all equilibrium points of the neural network and to determine the stability properties of the equilibrium points. Different design methods are generalized in [7] in a design procedure for neural networks with sparse connectivity. These results guarantee that the synthesized neural networks have predetermined sparse interconnection structures and store an almost arbitrary set of desired memory patterns as reachable memory vectors. The authors show that a sufficient condition for the existence of such a sparse neural network design is self feedback for every neuron in the network. Critical network issues like robustness and invariance to perturbation have been addressed in [8] where the authors provide upper

bounds for the perturbations of parameters under which desired memories stored in a neural network are preserved. This type of information is of great practical interest during the implementation process of such networks. In [4] the problem of realizing associative memories by designing a CNN that can store given binary vectors with improved performance is formulated as a constrained optimization problem, which can be reduced to a generalized eigenvalue problem (GEVP) [9]. Networks designed using this method exhibit less spurious patterns and higher recall probabilities [4]. Similar interior point methods are also applied to the problem of synchronization in chaotic delayed neural networks [18], [19].

The quality of memorization in associative memories is typically measured by the *average recall probability* [4], [8], which is an approximation of the probability that the network will recall correctly one of the patterns stored in its memory, if presented with a version of this pattern contaminated by noise. This number is defined as the ratio of number of recovered memory patterns (perturbed initial condition vectors which result in same output as the stored memory vector) to the total number of perturbed initial condition vectors. Although this quantification method measures the network’s performance accurately and unambiguously, it can only be applied after the network has been designed and tested on a significant number of test inputs. In other words, it does not allow the designer to *predict* the network’s performance without experimentation.

The concept of stability parameters appears in some earlier work on associative memory networks [10]–[12], as a measure of quality of memorization. Specifically, it has been demonstrated [13], [14] that the sign as well as the magnitude of these numbers are related to the size of the attraction regions of the desired memory patterns. Although the applicability of this concept as a universal measure of memorization quality, (specifically, the use of their size as a direct measure of the absolute sizes of the attraction regions) has long been debated [15], these parameters are generally accepted as a metric for the network’s performance [16].

Reported approaches on the design of CNNs (without time delays) start with, and build on some given network topology. The choice of this starting point is justified from the fact that typically the physical platform on which the network is implemented is fixed. In this paper, we raise the question of how could the design process be different if the network topology was also part of the design. It could be the case that implementing and maintaining certain neuron connections might be more expensive than others, and then one may be forced to strike a balance between connectivity and performance. Simple stability analysis of

The authors are with the Cooperative Robots Laboratory, Department of Mechanical Engineering, University of Delaware, Newark, Delaware. Their work is supported by NSF-IIS award #0822845.

neural networks implementing associative memories suggests similar links between the network of connections and the stability of the dynamical system as in consensus networks. In the case of neural networks, however, the topologies used are typically standard: either complete graphs (in Hopfield networks) or grid-like structures (in CNNs). Yet, available design methodologies do not place specific conditions on the network's structure. In addition, the cost of communication (delayed or otherwise) between neurons is commonly ignored. But if neurons communicate at a nontrivial cost, there is benefit in designing networks that perform just as well, but are less expensive.

This paper presents an approach to the optimization of the network topology of CNNs in which communication links between neurons may incur variable cost. We selectively "trim" network links in an effort to trade network performance for smaller communication cost. Based on existing design tools, we perform combinatorial optimization on a portion of a CNN that includes the *most expensive* interconnection links. This produces a sparser CNN, the performance of which can be comparable to the original network. This performance is quantified in terms of the network recall probability, and in the proposed optimization algorithm approach is captured by the neural network's stability parameters.

II. PRELIMINARIES

A two-dimensional continuous time zero-input CNN consisting of $n = MN$ cells in a $M \times N$ array, as first introduced by Chua and Yang [1], can be described by

$$\dot{x}_{ij} = -x_{ij} + \sum_{(k,l) \in N_r(i,j)} W_{ij,kl} y_{kl} + d_{ij}, \quad y_{ij} = \text{sat}(x_{ij}),$$

where $1 \leq i \leq M$, $1 \leq j \leq N$, and $\text{sat}(x_{ij}) \triangleq \frac{1}{2}(|x_{ij} + 1| - |x_{ij} - 1|)$. Here x_{ij} and y_{ij} are the state and output of the (i, j) th cell respectively, and $N_r(i, j)$ is an r -neighborhood of the (i, j) th cell defined as $N_r(i, j) \triangleq \{(k, l) : \max\{|k - i|, |l - j|\} \leq r\}$, for $1 \leq i \leq M, 1 \leq j \leq N$. This system can be expressed as

$$\dot{x} = -x + T \text{sat}(x) + b, \quad (1a)$$

$$y = \text{sat}(x) \quad (1b)$$

where $x = [x_{11}, x_{12}, \dots, x_{21}, \dots, x_{MN}]^T \in \mathbb{R}^n$ is the stack vector of all neuron states, $y = [y_{11}, y_{12}, \dots, y_{21}, \dots, y_{MN}]^T \in H^n$ is the output vector (H^n is the n -dimensional hypercube $[-1, +1]^n$), $T = [T_{ij}] \in \mathbb{R}^{n \times n}$ is the network connection weight matrix, $b \in \mathbb{R}^n$ is the network's bias vector and for vector arguments, the saturation function is defined elementwise.

Let B^n denote the set of bipolar vectors in H^n , namely those whose elements are either $+1$ or -1 . For $i = 1, \dots, n$, the initial condition vectors of (1) should always satisfy $|x_i(0)| \leq 1$. The interconnection topology of the network can be described by an adjacency matrix S , and thus a weight T_{ij} is non-zero only if $S_{ij} = 1$. Vector $\alpha \in H^n$ is a *memory vector* for (1) if the latter has an asymptotically stable equilibrium point $\beta \in \mathbb{R}^n$ such that $\alpha = \text{sat}(\beta)$ [1].

The synthesis problem for a CNN can be stated as follows:

Problem 1 (Synthesis): Given a CNN (1), implemented on network expressed by an adjacency matrix S , along with the set of desired bipolar memory vectors $\alpha^1, \dots, \alpha^m \in B^n$, determine the network connection weights T_{ij} and bias parameters b_i so that the obtained neural network can store all desired memory patterns as reachable memory vectors.

The adjacency matrix S which determines which T_{ij} , typically encodes a lattice structure and can be selected arbitrarily as long as certain conditions are satisfied [1], [3], [4]: If $\alpha \in B^n$ and $\beta = T\alpha + b$ are such that

$$\alpha_i \beta_i = \alpha_i \left(\sum_{j=1}^n T_{ij} \alpha_j + b_i \right) > 1, \quad \forall i = 1, \dots, n, \quad (2)$$

then (α, β) is a pair of a memory vector and an asymptotically stable equilibrium point of system (1). Also if $\alpha \in B^n$ and $\beta = T\alpha + b$ are such that for any $i = 1, \dots, n$, $\alpha_i \beta_i = \alpha_i \left(\sum_{j=1}^n T_{ij} \alpha_j + b_i \right) < 1$, then $\alpha \in B^n$ cannot be a memory vector. System (1) is globally stable if T is symmetric.

In addition to aforementioned stability criteria, specific stability and robustness properties for these networks are established in terms of the elements of the network's weight matrix T and bias vector b [1], [3], [4]:

Let $\alpha \in B^n$ be a memory vector of system (1) and let $k \geq 1$ be an integer. If $\tilde{T} = T - I_n$ and b satisfy

$$\alpha_i \left(\sum_{j=1}^n (\tilde{T}_{ij} \alpha_j + b_i) \right) > 2(k-1) \max_{1 \leq j \leq n} |\tilde{T}_{ij}|, \quad (3)$$

for $i = 1, \dots, n$, then any binary vector $\alpha^* \in B^n$ such that $1 \leq h(\alpha^*, \alpha) \leq k$ ($h(\alpha^*, \alpha) \triangleq \sum_i |\alpha_i^* - \alpha_i|$ denotes the Hamming distance) has the following properties:

- 1) α^* is not a memory vector (asymptotically stable equilibrium point for (1)).
- 2) if $x(0) = \alpha^*$ and $\alpha_i^* \neq \alpha_i$, then $x_i(t)$ converges to α_i .

Satisfying (3) with large k increases both the attractivity and the robustness of the stored memory vector $\alpha \in B^n$ and decrease the probability of existence of spurious patterns in vertices near α [4].

Problem 1 can be formulated as a GEVP [4]. For $i, j = 1, \dots, n$, and $k = 1, \dots, m$,

$$\min(-\delta), \quad \text{s.t.} \quad (4a)$$

$$(-\delta) \text{diag}[2q_1, \dots, 2q_n] - \text{diag}[-p_1, \dots, -p_n] > 0 \quad (4b)$$

$$\alpha_i^{(k)} \left(\sum_{j=1}^n \tilde{T}_{ij} \alpha_j^{(k)} + b_i \right) - p_i > 0, \quad (4c)$$

$$q_i - \tilde{T}_{ij} > 0, \quad (4d)$$

$$\tilde{T}_{ij} + q_i > 0, \quad (4e)$$

$$\tilde{T}_{ii} = 0, \quad (4f)$$

$$\tilde{T}_{ij} = \tilde{T}_{ij}^T = \tilde{T}_{ij}|_S \quad (4g)$$

$$L < q_i < U, \quad (4h)$$

where p_i and q_i for $i = 1, \dots, n$ are additional slack

variables used to cast the design problem as a linear matrix inequality (LMI) [17], and L and U are the lower and upper bounds for the design variables in the GEVP.

In this paper we quantify the quality of memorization in the network in terms of the magnitude of the stability parameters $K_{i\mu}$ defined as:

$$K_{i\mu} = \alpha_i^{(\mu)} h_i = \frac{\alpha_i^{(\mu)} \sum_j c_{ij} \alpha_j^{(\mu)}}{\|c_i\|_2}, \quad (5)$$

where $c_{ij} = T_{ij}$, $\|c_i\| = \sqrt{\sum_j c_{ij}^2}$, and superscript $(\cdot)^{(\mu)}$ indexes the set of desired memory vectors. The stability parameters $K_{i\mu}$, (one for each pair of neuron node and memory vector) are numbers which have been proposed as a measure of quality of memorization and it has been hypothesized that they are linked to the size of the attraction regions of the neural network [12], [15], [16]. The stability parameters and the degree of symmetry of connection matrix control the system dynamics and the domain sizes. Positive values of stability parameters indicate stability of a pattern.

The problem addressed in this paper is the following variant of Problem 1:

Problem 2 (Topology optimization): Given (1), implemented on network expressed by a *weighted* graph with adjacency matrix \hat{S} , with the set of desired bipolar memory vectors $\alpha^1, \dots, \alpha^m \in B^n$, determine the connection weights T_{ij} and bias parameters b_i of a subnetwork of S , so that this subnetwork stores all desired memory patterns as reachable memory vectors, and recalls them (almost) as well as the complete network.

In this paper, the performance metric is formulated using the values of the neural network's stability parameter values. The design process includes an optimization stage, in which the high cost edges that contribute the least to the engraving of the desired patterns on the network's memory, are selectively trimmed.

III. DILUTION OF NETWORK CONNECTIVITY

A CNN can be viewed as a collection of sub-networks that are linked to each other through communication links. We assume that communication incurs a cost, and in a particular scenario of interest, information flow across sub-networks is more expensive compared to communication within each subnetwork. In this case, one may be interested in minimizing communication cost, while maintaining the functionality and performance of the whole network above a certain threshold.

In this paper, reducing communication cost is achieved by trimming those expensive neuron links which do not contribute significantly to the ability of the network to retrieve its memorized information. The communication costs are captured by the weights of the weighted adjacency matrix \hat{S} . Given the (unweighted) adjacency matrix S of the CNN, along with the set of desired bipolar memory vectors $\alpha^1, \dots, \alpha^m \in B^n$, the design process begins by determining the network parameters T_{ij} and b_i through the solution of the GEVP (4). The resulting network maximizes the recall

probability of the patterns it has been designed for, without considering the cost of using the different network links. The next step is to dilute the connectivity of S , by judiciously removing some of the high cost links identified in \hat{S} , in a way so that the quality of memorization is not severely affected. Balancing performance against communication cost is achieved through (combinatorial) optimization over the network links, subject to the stability constraints stated in Section II.

The optimization process is iterative. For a given (intermediate) topology S , the algorithm determines the neural network weights T_{ij} for $i, j = 1 \dots, n$ and biases b , and based on the selected patterns $\alpha_i^{(\mu)}$ for $\mu = 1, \dots, m$ to be memorized, an $n \times m$ stability parameter matrix is formed by repeated application of (5): $K = [K_{i\mu}]_{i=1, \dots, n; \mu=1, \dots, m}$. For the prescribed patterns $\alpha^1 \dots \alpha^m$ to be stored effectively in the network's memory, all stability parameters $K_{i\mu}$ must be nonnegative and as large as possible. For the network topology encoded in S , the value of the following objective function is evaluated:

$$\bar{K} = \sum_{i=1}^n \sum_{\mu=1}^m K_{i\mu}. \quad (6)$$

Thus \bar{K} , being the sum of all nodes' stability parameters for all chosen memory vectors, quantifies the collective ability of the network to recall all desired memories.

Remark 1: Several performance metrics have been explored as alternatives to (6), based on different norms of the stability parameter matrix K , such as the minimum row (or column) sums, the (absolute) minimum element of K , etc. When the average recall probability of each design was evaluated, it was determined that the sum of all stability parameters captured more accurately the ability of the network to recall memory patterns.

On the other hand, a naive implementation of a "branch-and-bound" approach, where links are divided into "promising" and "not-promising" for deletion groups according to their associated \bar{K} value, and only the "promising" possibilities are explored in the subsequent steps, will generally fail. This is because due to the combinatorial nature of the problem, an edge whose sole deletion has an adverse effect on the stability parameters may even improve the value of \bar{K} when combined with additional edge removals.

For a given cost threshold ν , and performance threshold κ , the high-cost edges that are candidates for deletion are identified in the residual adjacency matrix

$$R = \frac{1}{2} \left(\text{sign}(\hat{S} - \nu S) + S \right),$$

where the sign function is evaluated element-wise on the matrix argument. For every nonzero (i, j) element in R , we remove the associated (i, j) edge from S . If (3) is satisfied then we store the resulting value of \bar{K} . The high-cost edge associated with the highest stored \bar{K} value is marked for deletion, and the step is repeated for the topology S' , \hat{S}' in which neurons i and j are not linked.

The process for evaluating the edges which are candidate

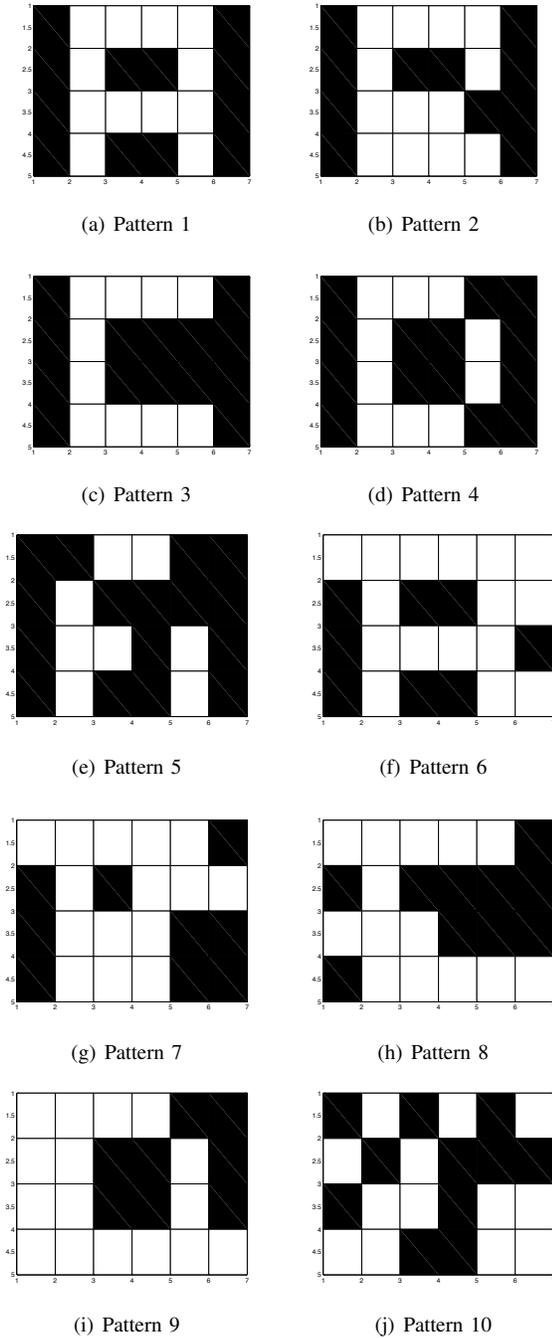


Fig. 2. Ten patterns to be stored as memory vectors

(shown in Fig. 3). It can be seen that the original network topology has same recall probability graph as the sparser topology derived with selective deletion of costly edges. However with a relatively large number of deleted edges, there is a noticeable decrease in network performance. The objective of optimization is thus to balance communication cost versus average recall probability.

Figure 4 shows the effect on recall probability of the CNN parameterized by the number of edges deleted. The vertical axis represents the average of the recall probabilities for

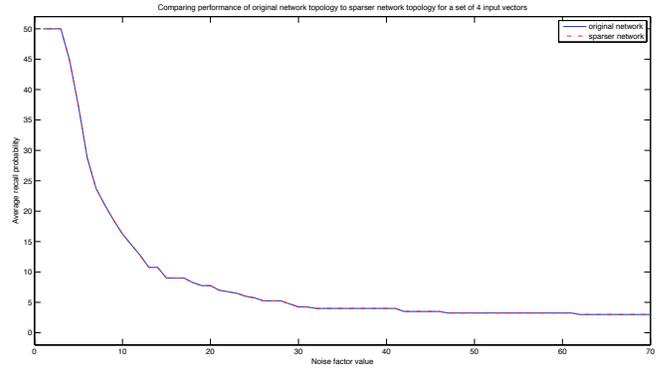


Fig. 3. Network performance vs. selective deletion of edges for the original and the optimized network. The curves are initially indistinguishable.

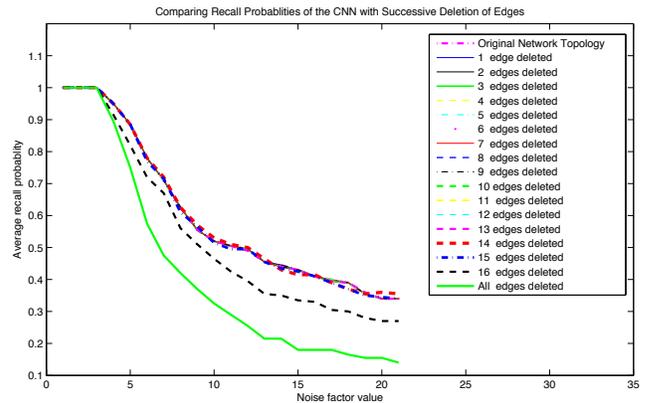


Fig. 4. Network performance, parameterized by number of edges removed.

the four given initial condition vectors and the horizontal axis represents the noise factor or the amplification factor by which the set of noise is multiplied. There is always a gradual decrease in performance as the amount of noise injected in the pattern increases. The recall probability of the original network design evolves with the amount of noise very similar to that of the sparser network produced through the topology optimization process. The additional curves correspond to even sparser networks, and demonstrate that reducing the connectivity further, has an observable impact on network performance. Also if the new \bar{K} value is close to the original \bar{K} , performance is still comparable. But when \bar{K} drops significantly compared to its original value, we see a notable performance degradation.

Next we consider all ten input patterns shown in Fig. 2 to test the ability of the network to store additional patterns with lean interconnection topologies. A comparison is made between the original network model and the network with the topology determined by the optimization algorithm (shown in Fig. 5) for a set of all ten input memory vectors. Simulation results again show that the original network topology has same recall probability graph as the sparser topology derived with selective deletion of costly edges. This proves the validity of the above proposed method in presence of larger set of input data to be stored. It is also observed that the

number of input vectors to be stored is increased, the recall probability of the network decreases gradually and it is able to tolerate less noise (Fig. 6). This is to be expected, since more attractors now share the same area of the state space; the region of attraction for each one of these stable equilibria is reduced.

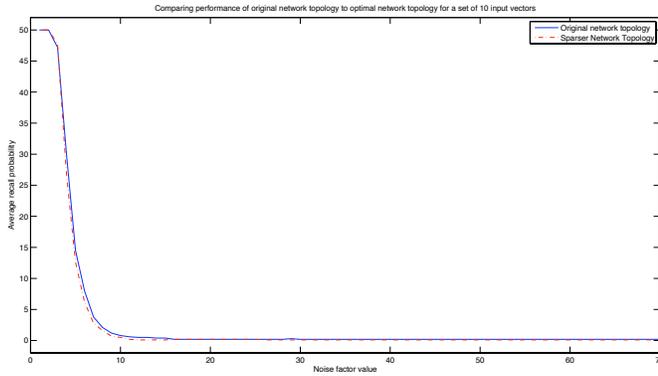


Fig. 5. Network performance, original to sparser comparison for 10 memory vectors.

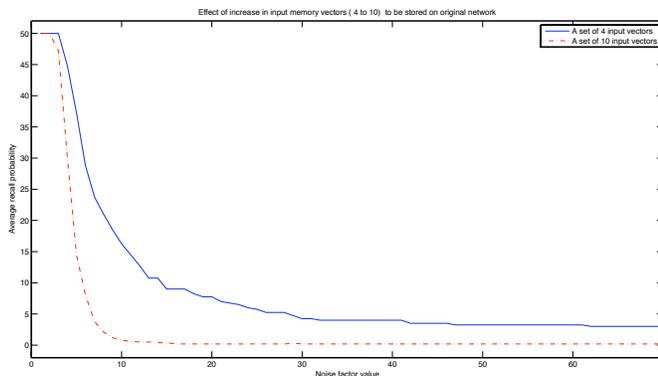


Fig. 6. Original network performance for different set of input vectors.

V. CONCLUSIONS

This paper suggests a topology optimization approach to cellular neural network design, as a method for realizing associative memories using sparser (thus more efficient, or less expensive) networks. The optimization criterion utilizes the stability parameters of the network as a quantitative measure of its ability to recall patterns contaminated with different levels of noise. These stability parameters have been associated with the size of the attraction regions of the neural network. Simulations and comparisons show that a sparser topology can be achieved without significantly degrading the

performance of the network, by selectively deleting those weights from the optimized network which contribute the least (as measured by the arithmetic sum of the stability parameters) to ability of the network to recall the desired patterns.

REFERENCES

- [1] L. O. Chua and L. Yang, "Cellular neural networks: theory and applications," *IEEE Transactions on Circuits and Systems*, vol. 35, pp. 1257–1290, 1988.
- [2] A. N. Michel and J. A. Farrell, "Associative memories via artificial neural networks," *IEEE Control Systems Magazine*, vol. 1990, pp. 6–17, 1988.
- [3] D. Liu and A. N. Michel, "Cellular neural networks for associative memories," *IEEE Transactions on Circuits and Systems*, vol. 40, pp. 119–121, 1993.
- [4] J. Park and Y. Park, "An optimization approach to design of cellular neural networks," *International Journal of Systems Science*, vol. 31, pp. 1585–1591, 2000.
- [5] S. Arik and V. Tavsanoglu, "Global asymptotic stability analysis of bidirectional associative memory neural networks with constant time delays," *Neurocomputing*, vol. 68, pp. 161–176, 2005.
- [6] J. H. Park, "Robust stability of bidirectional associative memory neural networks with time delays," *Physics Letters A*, vol. 349, no. 6, p. 494499, 2008.
- [7] D. Liu and A. N. Michel, "Sparsely interconnected neural networks for associative memories with applications to cellular neural networks," *IEEE Transactions on Circuits and Systems -II, Analog and Digital Signal Processing*, vol. 41, pp. 295–307, 1994.
- [8] —, "Robustness analysis and design of a class of neural networks with sparse interconnecting structure," *Neurocomputing*, vol. 12, pp. 59–76, 1996.
- [9] Y. Nesterov and A. Nemirovskii, *Interior Point Polynomial Algorithms in Convex Programming*. Studies in Applied and Numerical Mathematics, 1994, vol. 13.
- [10] S.-I. Amari, "Characteristics of randomly connected threshold-element networks and network systems," *Proceedings of the IEEE*, vol. 59, no. 1, pp. 35–47, 1971.
- [11] W. Krauth and M. Mézard, "Learning algorithms with optimal stability in neural networks," *Journal of Physics A: Mathematical and General*, vol. 20, no. 11, pp. L745–L752, 1987.
- [12] E. Gardner, "The space of interactions in neural network models," *Journal of Physics A: Mathematical and General*, vol. 21, no. 1, pp. 257–270, 1988.
- [13] B. M. Forrest, "Content-addressability and learning in neural networks," *Journal of Physics A: Mathematical and General*, vol. 21, no. 1, pp. 245–255, 1988.
- [14] T. B. Kepler and L. Abbott, "Domains of attraction in neural networks," *Journal de Physique*, vol. 49, no. 10, pp. 1657–1662, 1988.
- [15] A. C. C. Coolen, "On the relation between stability parameters and sizes of domains of attraction in attractor neural networks," *Europhysics Letters*, vol. 16, pp. 73–78, 1991.
- [16] K. E. Kurten, "Adaptive architectures for hebbian network models," *Journal de Physique I*, vol. 2, pp. 615–624, 1992.
- [17] S. Boyd, L. Elghaoui, E. Feron, and V. Balakrishnan, *Linear Matrix Inequalities in Systems and Control Theory*. Studies in Applied and Numerical Mathematics, 1994, vol. 15.
- [18] M. Liu, "Optimal exponential synchronization of general chaotic delayed neural network: An lmi approach," *Neural Networks*, vol. 22, pp. 949–957, 2009.
- [19] D. L. T. Ma, H. Zhang and Z. Wang, "A novel lmi approach to global impulsive exponential synchronization of chaotic delayed neural networks," *Proceedings of 48th IEEE Conference of Decision and Control*, pp. 626–631, 2009.