

# Synthesis of Sigma-Pi Neural Networks by the Breeder Genetic Programming

Byoung-Tak Zhang and Heinz Mühlenbein

*Abstract*— Genetic programming has been successfully applied to evolve computer programs for solving a variety of interesting problems. In the previous work we introduced the breeder genetic programming (BGP) method that has Occam’s razor in its fitness measure to evolve minimal size multilayer perceptrons. In this paper we apply the method to synthesis of sigma-pi neural networks. Unlike perceptron architectures, sigma-pi networks use product units as well as summation units to build higher-order terms. The effectiveness of the method is demonstrated on benchmark problems. Simulation results on noisy data suggest that BGP not only improves the generalization performance, it can also accelerate the convergence speed.

## I. INTRODUCTION

Genetic programming has been successfully used to evolve computer programs for solving many interesting problems in artificial intelligence and artificial life [4, 16, 18, 19]. Similar to usual genetic algorithms (GAs), genetic programming (GP) starts with a population of randomly generated individuals. Each individual is a program that, when executed, is the candidate solution to the problem. These programs are expressed as parse trees, or LISP S-expressions. Fitness proportionate selection and crossover are used to produce increasingly fitter populations of computer programs.

While most GAs use binary strings of fixed size [9], GPs use structured representations of variable length. This is important because it is particularly suited to problems in which the optimal underlying structure must be discovered. One problem with the variable length is that the program size may grow without bound. For example, Kinneer [16] reports that all but a very few of evolved solutions to his sorting problems were so large as to defy any human understanding of them. Tackett [29], in his application of genetic programming to image processing tasks, observes that the size and complexity of trees grows without performance improvement.

In the previous work [35, 36] we introduced the breeder genetic programming (BGP) that employs Occam’s razor in its fitness measure to evolve optimal or minimal size multilayer perceptrons. We apply here the BGP method to synthesize sigma-pi neural networks. Unlike multilayer perceptrons, sigma-pi networks use product units as well as summation units to build higher-order terms. In section

II we illustrate the usefulness of higher-order terms and show how the sigma-pi networks can be used to represent higher-order networks. The representation scheme and the algorithm are described in section III. Simulation results are reported in section IV, followed by concluding remarks in section V.

## II. SIGMA-PI NEURAL NETWORKS

To motivate the approach, we start with a brief description of multilayer neural networks. Multilayer perceptrons are feedforward networks with one or more layers of nodes between the input and output units. These additional layers contain hidden units that are not directly connected to both the input and output units. The input-output relation of the units is given in these networks by weighted sum of inputs

$$u_i = \sum_{j \in R(i)} w_{ij} x_j \quad (1)$$

where  $w_{ij}$  is the connection weight from unit  $j$  to unit  $i$  and  $R(i)$  denotes the receptive field of unit  $i$ . The total input is then transferred to upper layer units by a nonlinear activation function  $f$ , e.g. a threshold function:

$$y_i = f(u_i) = \begin{cases} +1 & \text{if } u_i > \theta_i \\ -1 & \text{otherwise} \end{cases} \quad (2)$$

where  $\theta_i$  is a threshold.

A commonly adopted architecture consists of one hidden layer with full connectivity between neighboring layers. This structure has been very successful for many applications. However, they have some weaknesses:

1. The full connectivity between layers and often only between neighboring layers try to find a prediction over the full input space. This is not necessarily a good strategy. A particular task might not contain a good predictor for the full input space, but might contain functions capable of good prediction on specified regions of input space.
2. They are especially appropriate to approximating additive functions, since they employ linear combinations of inputs. However, the multilayer perceptrons cannot approximate efficiently if there are high order interactions between the inputs. Adding additional hidden layers may help to extend the representational capacity of the network.

The authors are with the German National Research Center for Computer Science (GMD), Schloss Birlinghoven, D-53757 Sankt Augustin, Germany. E-mail: {zhang, muelhen}@gmd.de. **Revised version of this paper appears in Proc. of Int. Conf. on Evolutionary Computation (ICEC-94), IEEE, 1994, pp. 318–323.**

Figure 1: An order 4 neuron

The higher-order terms can be again used as building blocks which are able to capture a high-order correlational structure of the data. In particular, by building a sigma unit which has as input various higher-order terms, we can construct a higher-order network of sigma-pi units:

$$y_i = f_i(u_i) = f_i\left(\sum_k w_k T^{(k)}\right) \quad (8)$$

The problem in using higher-order networks is that the number of terms explodes with the problem size; the number of parameters necessary for specifying an order  $k$  neuron is

$$r_k = \sum_{i=0}^k {}_n C_i, \quad (9)$$

because  $w_{i_1 \dots i_m}^{(m)}$  have  ${}_n C_m$  components. Here  $n$  is the total number of inputs and  ${}_n C_m$  are the binomial coefficients. As an example, an order 4 neuron has  $2^4 = 16$  parameters as shown in Figure 1.

To avoid the combinatorial explosion, a method is needed to discover and combine useful terms and to eliminate non-essential terms. The following section describes a genetic programming method designed for these purposes.

### III. EVOLVING SIGMA-PI NETS

In order to evolve problem-specific sigma-pi networks we use the breeder genetic programming (BGP) [36]. Similar to GPs, BGP uses tree representation for individuals. BGP uses, however, ranking-based truncation selection as in the breeder genetic algorithm (BGA) [26] instead of fitness-proportionate selection or tournament selection [2]. Another feature of BGP is its fitness function which uses the Occam's razor principle. The truncation selection combined with Occam's razor has been proved useful to balance the accuracy and parsimony of multilayer perceptrons.

The sigma-pi networks are represented as a set of  $m$  trees, where  $m$  is the number of output units. Each tree has an arbitrary number of subtrees. The grammar for describing the genotype is given in Figure 2. Figure 3 shows

Figure 4: *Phenotype of a sigma-pi network. This architecture allows higher-order terms with local receptive fields. Direct connections between non-neighboring layers are also allowed.*

into the mating pool  $\mathcal{B}(g)$ , where  $\tau \in (0, 1]$  is the truncation threshold. Each network in  $\mathcal{B}(g)$  undergoes a local hillclimbing to be described later. This results in revised mating pool  $\mathcal{B}(g)$ . The  $(g+1)$ st generation of size  $M$  is produced by applying crossover operators to randomly chosen parent networks in the mating pool  $\mathcal{B}(g)$ . We use the elitist strategy by replacing the worst fit network  $\mathcal{A}_{worst}(g+1)$  in  $\mathcal{A}(g+1)$  by the best  $\mathcal{A}_{best}(g)$  of  $\mathcal{A}(g)$ . A new population is generated repeatedly until an acceptable solution is found or the variance of the fitness  $V(g)$  falls below a specified limit value  $V_{min}$ :

$$V(g) = \frac{1}{M} \sum_{i=1}^M (F_i(g) - \bar{F}(g))^2 \leq V_{min} \quad (12)$$

where  $\bar{F}(g)$  is the average fitness of the individuals in  $\mathcal{A}(g)$ . The algorithm also stops if it reaches the maximum number of generations,  $g_{max}$ .

The weights of a network are trained by a local hillclimbing method to each selected parent. Many hillclimbing methods are possible. We use a simple next-ascent hillclimbing procedure: given an individual  $B_i$ , a better individual  $B_i^{new}$  is sought by repeatedly applying the mutation

```

procedure BGP( $M, \tau, V_{min}, g_{max}$ )
  population size: int  $M$ 
  truncation rate: real  $\tau$ 
  variance: real  $V_{min}, V(g)$ 
  generation: int  $g_{max}, g$ 
  fitness values: real  $F_i$ 
  population: array  $\mathcal{A} = (A_1, A_2, \dots, A_M)$ 
  mating pool: array  $\mathcal{B} = (B_1, B_2, \dots, B_{\tau \cdot M})$ 

   $g \leftarrow 0$ 
   $\mathcal{A}(0) \leftarrow \text{Initialize}(M)$ 
  while ( $V(g) > V_{min}$ ) do
     $F_i(g) \leftarrow \text{Evaluate}(A_i(g)) \quad \forall i \in \{1, \dots, M\}$ 
    if (solution found or  $g \geq g_{max}$ ) stop
     $\mathcal{B}(g) \leftarrow \text{Select}(\mathcal{A}(g), \mathcal{F}(g), \tau)$ 
     $B_i(g) \leftarrow \text{Hillclimb}(B_i(g)) \quad \forall i \in \{1, \dots, \tau M\}$ 
     $\mathcal{A}(g+1) \leftarrow \text{Mate}(\mathcal{B}(g), M)$ 
     $\mathcal{A}_{worst}(g+1) \leftarrow \mathcal{A}_{best}(g)$ 
     $g \leftarrow g+1$ 
  endwhile
endprocedure

```

Figure 5: The top-level structure of the BGP algorithm

operator to the weights until there is no weight configuration found having better fitness in each sweep through the individual. The sequence of mutation is defined as the depth-first search order.

The mutation operation is performed by replacing the value of a node,  $c_i$ , of the tree by another, i.e. by finding the class  $C_k$  of  $c_i$  and replacing  $c_i$  by another member  $c_j, j \neq i$  in the set  $C_k$ . Here the class  $C_k$  must first be found because not every value (node) can be mutated to arbitrary values. For example, a weight value must be drawn from the set  $\{+1, -1\}$ . The thresholds are mutated the same way as the weights. The index for the input units can be mutated by another input index. We also allow a nonterminal symbol  $S$  to be mutated by a  $P$  and vice versa, i.e. changing the type of neural units. This flexibility ensures that multilayer perceptrons can also be evolved from sigma-pi networks.

The crossover operator adapts the size, depth and receptive field shape of the network architecture. The crossover operation starts with choosing randomly two parents,  $B_i$  and  $B_j$ , from the mating pool  $\mathcal{B}(g)$ . The nodes in the tree are numbered according to the depth-first search order and crossover sites  $c_i$  and  $c_j$  are chosen at random with the following conditions:  $1 \leq c_i \leq \text{Size}(B_i)$  and  $1 \leq c_j \leq \text{Size}(B_j)$ . The length  $\text{Size}(B_k)$  of an individual  $B_k$  is defined as the total number of units and weights. The subtrees of two parent individuals,  $B_i$  and  $B_j$ , are exchanged at the given crossover points to form two offspring  $B'_i$  and  $B'_j$ . The label of the nodes,  $c_i$  and  $c_j$ , must belong to the same class, i.e. either both  $U$ -type or both  $W$ -type nodes. The number of arguments of each operator plays no role because the syntactically correct subtree under the node is completely replaced by another syntactically correct subtree.

We use the following equation as the fitness measure

$$F(D|W, A) = \frac{E(D|W, A)}{m \cdot N} + \frac{C(W|A)}{N \cdot C_{max}}. \quad (13)$$

A theoretical background behind this fitness function is discussed in [36]. The first term expresses the accuracy penalty caused by the error for the training set:

$$E(D|W, A) = \sum_{i=1}^N \sum_{j=1}^m (y_{ij} - o_j(x_i; W, A))^2. \quad (14)$$

Here  $y_{ij}$  denotes the  $j$ th component of the  $i$ th desired output vector  $y_i$ , and  $o_j(x_i; W, A)$  denotes the  $j$ th actual output of the network with the architecture  $A$  and the set of weights  $W$  for the  $i$ th training input vector  $x_i$ .

The second term in the fitness function expresses the complexity penalty of the network, often called Occam's razor. The complexity is defined as

$$C(W|A) = W(A) + 10 \cdot L(A) + U(A), \quad (15)$$

where  $W(A) = \sum_{k=1}^K w_k^2$  is the number of weights in the network for binary weights.  $L(A)$  and  $U(A)$  denote the number of layers and units, respectively. The  $L(A)$  term penalizes a deep architecture which requires a large execution time after training. The  $U(A)$  term penalizes a large number of units whose realization is more expensive than weights. Notice that the  $L(A)$  term is multiplied by 10 to penalize it more strongly than others. Depending on applications, one may weight three terms differently.  $C_{max}$  is a normalization factor used for the complexity term to be between 0 and 1. In all experiments we set  $C_{max} = 1000$ , assuming the problems can be solved by  $C(W|A) \leq 1000$ .

Notice in equation (13) that the complexity term  $C(W|A)$  is divided by  $N$ , the number of training examples, to have the complexity term play a minor role in determining the total fitness value of the network. This ensures a small network be preferred to a large network only if both of them achieve a comparable performance. Otherwise, the evolution may not lead to a solution by preferring smaller networks which lack the capacity to learn the training set.

#### IV. SIMULATION RESULTS

The method was tested on the parity problem. We performed two kinds of experiments separately. In the first, we are interested to know whether the use of product units is effective and, if yes, to what extent. In these experiments, noise-free examples are used. For the second series of experiments, we use noisy data. The generalization performance and the learning speed of different strategies are compared to study the effect of Occam's razor for the construction of sigma-pi networks.

##### A. Clean Data

The accuracy and convergence speed of sigma-pi networks are studied. We also measured the network complexity in terms of the number of layers, units and connections.

The results are summarized in Table 1. For these experiments we used noise-free data consisting of  $2^n$  examples for problem size  $n$ . We performed 10 runs for each problem size. The population was initialized for every individual to contain sigma and pi units with 50% probability each. The depth of initialized network was limited to 3. The truncation rate was 40%. For each problem size  $n$ , we used the population size  $M = 100n/2$  and the maximum generation  $g_{max} = 10n$ . For example,  $M = 400$  and  $g_{max} = 80$  for  $n = 8$ .

n	arch	layr	unit	conn	accr	gen
2	S	2.2	4.0	12.9	95.0	12.5
	S/P	2.0	2.2	6.5	100.0	2.9
4	S	4.9	17.9	55.0	81.9	40.0
	S/P	2.4	3.6	14.8	100.0	8.5
6	S	4.4	30.2	161.4	89.5	46.2
	S/P	2.9	8.9	40.2	98.4	35.8
8	S	5.9	65.6	414.1	86.4	63.4
	S/P	3.6	21.0	106.8	98.1	56.2

Table 1: Results for parity problems

To test the effectiveness of pi units we also run 10 experiments using only sigma units. The experiments were the same as before except that all units are initialized as sigma units. As the results show, the additional use of pi units consistently improved the performance in accuracy as well as in complexity reduction.

### B. Noisy Data

In the second set of experiments we used the parity problem of input size 8. A total of 128 correct examples were generated randomly to get a training set and then noise was inserted to this data by randomly changing the output value with 5% probability. This means, on average, 6 or 7 examples out of 128 have false outputs. The generalization performance of the best solution in each generation was tested by the complete data set of  $2^8 = 256$  noise-free examples.

Figure 6 shows a typical evolution of the training and generalization error of the best fit sigma-pi networks. In spite of the noise, a good correspondence is observed between learning and generalization performance. Figure 7 shows the corresponding evolution of the complexity of the best fit network in each generation. Notice that the change of network performance is closely related with the change of its complexity.

The performance of the BGP with the fitness function (13) was compared with a method that uses just the error term as the fitness measure, i.e.

$$F(D|W, A) = \frac{E(D|W, A)}{m \cdot N}. \quad (16)$$

Both methods used the same noisy data of the 8-parity problem. For each method, 10 runs were executed until the 80th generation to observe the training and generalization performance of the solutions.

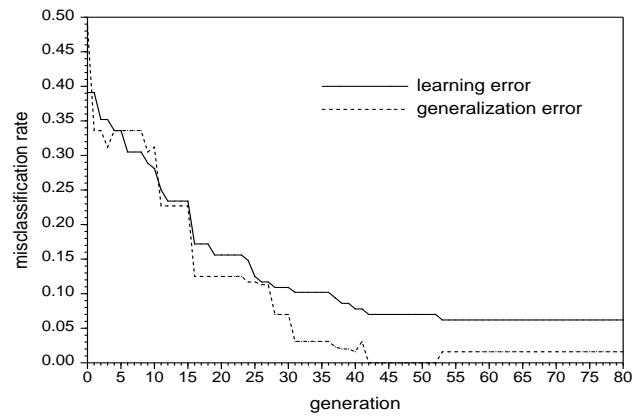


Figure 6: The evolution of the network performance for noisy data of the 8-input parity function. Also shown is the generalization performance on the complete test set of noise-free examples.

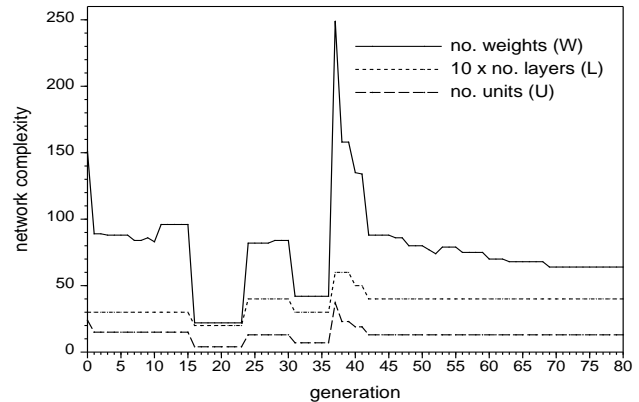


Figure 7: The evolution of the network size for noisy data of the 8-input parity function.

Table 2 shows the average network size found at the 80th generation. The corresponding performance and learning time are shown in the table. The learning time is measured in millions of evaluations of arithmetic operations associated with calculating activation values of neural units. The results show that applying Occam's razor achieves significantly better performance for this problem. Without Occam's razor the network size increased to an arbitrarily large size, which makes it difficult to find a useful building block to combine. Another advantage of using Occam's ra-

	F = E	F = E + C
number of layers	$7.8 \pm 0.1$	$3.5 \pm 0.2$
number of units	$65.4 \pm 1.8$	$19.0 \pm 2.7$
number of weights	$391.8 \pm 23.1$	$103.8 \pm 19.2$
learning accuracy	$75.1 \pm 2.5$	$89.3 \pm 1.8$
generalization	$61.2 \pm 0.8$	$89.9 \pm 3.8$
num. evaluations	$2630.4 \pm 136.1$	$1432.0 \pm 112.2$

Table 2: Performance with and without Occam's razor

zor is the accelerated convergence. In these experiments, the proposed fitness function decreased the network size by approximately four times and the speed-up factor of learning was two.

## V. CONCLUDING REMARKS

Necessity and usefulness of higher-order neural networks have been well-known. However the explosively increasing number of terms has hampered the design and training of higher-order networks. The present work shows the potential effectiveness of genetic programming to handle this problem. In particular, we show how the sigma-pi neural networks can represent the higher-order terms and how BGP can be extended to synthesize problem-specific sigma-pi networks. We also confirm the usefulness of Occam's razor in genetic programming for the improvement of generalization and convergence speed. Experimental results show that sigma-pi networks solve parity-like problems significantly better than perceptron networks. In another set of experiments, we found that the method can also evolve usual multilayer perceptrons with slightly more costs when the product units are unnecessary for solving the problem.

## ACKNOWLEDGEMENT

This research was supported in part by the Real-World Computing Program under the project SIFOGA.

## REFERENCES

- [1] S. Amari, "Dualistic geometry of the manifold of higher-order neurons," *Neural Networks*, vol. 4, pp. 443-451, 1991.
- [2] P. J. Angeline and J. B. Pollack, "Competitive environments evolve better solutions for complex tasks," in *Proc. Fifth Int. Conf. Genetic Algorithms (ICGA-93)*, Morgan Kaufmann, 1993, pp. 264-270.
- [3] T. Bäck and H.-P. Schwefel, "An overview of evolutionary algorithms for parameter optimization," *Evolutionary Computation*, vol. 1, pp. 1-23, 1993.
- [4] R. J. Collins, *Studies in Artificial Evolution*, Ph.D. thesis, University of California, Los Angeles, 1992.
- [5] R. Durbin and D. E. Rumelhart, "Product units: A computationally powerful and biologically plausible extension to back-propagation networks," *Neural Computation*, vol. 1, pp. 133-142, 1989.
- [6] J. A. Feldman and D. H. Ballard, "Connectionist models and their properties," *Cognitive Science*, vol. 6, pp. 205-254, 1982.
- [7] J. Franke, "On the functional classifier," *Proc. First Int. Conf. Document Analysis and Recognition*, 1991.
- [8] C. L. Giles and T. Maxwell, "Learning, invariance, and generalization in high-order neural networks," *Applied Optics*, vol. 26, no. 23, pp. 4972-4978, 1987.
- [9] D. E. Goldberg, *Genetic Algorithms in Search, Optimization & Machine Learning*. Addison Wesley, 1989.
- [10] F. Gruau, "Genetic synthesis of Boolean neural networks with a cell rewriting developmental process," Tech. Rep., Laboratoire de l'Informatique du Parallélisme, 1992.
- [11] S. A. Harp, T. Samad, and A. Guha, "Towards the genetic synthesis of neural networks," in *Proc. Third Int. Conf. Genetic Algorithms (ICGA-89)*, D. Schaffer, Ed. Morgan Kaufmann, pp. 360-369.
- [12] J. H. Holland, *Adaptation in Natural and Artificial Systems*, Ann Arbor: The University of Michigan Press, 1975.

- [13] H. Iba, T. Kurita, H. de Garis, and T. Sato, "System identification using structured genetic algorithm," in *Proc. Fifth Int. Conf. Genetic Algorithms (ICGA-93)*, S. Forrest, Ed. Morgan Kaufmann, 1993, pp. 279-286.
- [14] A. G. Ivakhnenko, "Polynomial theory of complex systems," *IEEE Trans. Systems, Man, and Cybernetics*, vol. 1, pp. 364-378, 1971.
- [15] H. Kargupta and R. E. Smith, "System identification with evolving polynomial networks," in *Proc. Fourth Int. Conf. Genetic Algorithms (ICGA-91)*, R. Belew and L. Booker, Ed. Morgan Kaufmann, 1991, pp. 370-376.
- [16] K. E. Kinnear Jr., "Generality and difficulty in genetic programming: Evolving a sort," in *Proc. Fifth Int. Conf. Genetic Algorithms (ICGA-93)*, Morgan Kaufmann, 1993, pp. 287-294.
- [17] H. Kitano, "Designing neural networks using genetic algorithms with graph generation system," *Complex Systems*, vol. 4, pp. 461-476, 1990.
- [18] J. R. Koza, *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, 1992.
- [19] J. R. Koza and J. P. Rice, *Genetic Programming: The Movie*. MIT Press, 1992.
- [20] Z. Michalewicz, *Genetic Algorithms + Data Structures = Evolutionary Programs*. Berlin: Springer, 1992.
- [21] G. F. Miller, P. M. Todd, and S. U. Hegde, "Designing neural networks using genetic algorithms," in *Proc. Third Int. Conf. Genetic Algorithms (ICGA-89)*, Morgan Kaufmann, 1989, pp. 379-384.
- [22] M. Minsky and S. Papert, *Perceptrons: An Introduction to Computational Geometry*. MIT Press, 1969, 1988.
- [23] H. Mühlenbein, "Evolution in time and space—The parallel genetic algorithm," in *Foundations of Genetic Algorithms*, G. Rawlins, Ed. Morgan Kaufmann, 1991, pp. 316-338.
- [24] H. Mühlenbein, "Parallel genetic algorithms in combinatorial optimization," in *Computer Science and Operations Research*, G. Balci, R. Sharda, and S. A. Zenios, Ed. Pergamon, 1992, pp. 441-456.
- [25] H. Mühlenbein, "Evolutionary algorithms: Theory and applications," in *Local Search in Combinatorial Optimization*, E. H. L. Aarts and J. K. Lenstra, Ed. Wiley, 1993.
- [26] H. Mühlenbein and D. Schierkamp-Voosen, "Predictive models for the breeder genetic algorithm I: Continuous parameter optimization," *Evolutionary Computation*, vol. 1, no. 1, pp. 25-49, 1993.
- [27] D. E. Rumelhart, G. E. Hinton, and J. L. McClelland, "A general framework for parallel distributed processing," in *Parallel Distributed Processing*, Vol. 1, D. E. Rumelhart and J. L. McClelland, Ed. MIT Press, 1986, pp. 45-76.
- [28] R. Sorkin, "A quantitative Occam's razor," *International Journal of Theoretical Physics*, vol. 22, pp. 1091-1104, 1983.
- [29] W. A. Tackett, "Genetic programming for feature discovery and image discrimination," in *Proceedings of the Fifth International Conference on Genetic Algorithms (ICGA-93)*, S. Forrest, Ed. Morgan Kaufmann, 1993, pp. 303-309.
- [30] D. Whitley, T. Starkweather, and C. Bogart, "Genetic algorithms and neural networks: Optimizing connections and connectivity," *Parallel Computing*, vol. 14, pp. 347-361, 1990.
- [31] R. Zembowicz and J. M. Żytkow, "Discovery of equations: Experimental evaluation of convergence," in *Proc. of AAAI-92*, Menlo Park, CA: AAAI Press, 1992, pp. 70-75.
- [32] B. T. Zhang, *Learning by Genetic Neural Evolution*. Ph.D. thesis in German, ISBN 3-929037-16-5, Sankt Augustin: Infix-Verlag, 1992. Also available as Informatik Berichte No. 93, Institut für Informatik I, Universität Bonn, D-53117 Bonn, 1992.
- [33] B. T. Zhang, "Accelerated learning by active example selection," to appear in *International Journal of Neural Systems*, 1993.
- [34] B. T. Zhang, "Teaching neural networks by genetic exploration," *GMD-Arbeitspapiere*, No. 805, German National Research Center for Computer Science, 1993.
- [35] B. T. Zhang and H. Mühlenbein, "Genetic programming of minimal neural nets using Occam's razor," in *Proc. Fifth Int. Conf. Genetic Algorithms (ICGA-93)*, S. Forrest, Ed. Morgan Kaufmann, 1993, pp. 342-349.
- [36] B. T. Zhang and H. Mühlenbein, "Evolving optimal neural networks using genetic algorithms with Occam's razor," to appear in *Complex Systems*, 1993.