

# Analog VLSI Stochastic Perturbative Learning Architectures \*

Gert Cauwenberghs  
Electrical and Computer Engineering  
Johns Hopkins University, Baltimore, MD 21218  
E-mail: [gert@bach.ece.jhu.edu](mailto:gert@bach.ece.jhu.edu)

## Abstract

We present analog VLSI neuromorphic architectures for a general class of learning tasks, which include supervised learning, reinforcement learning, and temporal difference learning. The presented architectures are parallel, cellular, sparse in global interconnects, distributed in representation, and robust to noise and mismatches in the implementation. They use a parallel stochastic perturbation technique to estimate the effect of weight changes on network outputs, rather than calculating derivatives based on a model of the network. This “model-free” technique avoids errors due to mismatches in the physical implementation of the network, and more generally allows to train networks of which the exact characteristics and structure are not known. With additional mechanisms of reinforcement learning, networks of fairly general structure are trained effectively from an arbitrarily supplied reward signal. No prior assumptions are required on the structure of the network nor on the specifics of the desired network response.

Keywords: Neural networks, neuromorphic engineering, reinforcement learning, stochastic approximation

## 1 Introduction

Learning and adaptation are central to the design of neuromorphic VLSI systems that perform robustly in variable and unpredictable environments.

Learning algorithms that are efficiently implemented on general-purpose digital computers do not necessarily map efficiently onto analog VLSI hardware. Even if the learning algorithm supports a parallel and scalable architecture suitable for analog VLSI implementation, inaccuracies in the implementation of the learning functions may significantly affect the performance of the trained system. Learning can only effectively compensate for inaccuracies in the network implementation when their physical sources are contained directly inside the learning feedback loop. Algorithms which assume a particular model for the underlying characteristics of the system being trained perform poorer than algorithms which directly probe the response of the system to external and internal stimuli.

---

\*This work was supported by a National Science Foundation Career Award, and by ARPA/ONR under MURI grant N00014-95-1-0409. Chip fabrication was provided through MOSIS.

A second source of concern in the design of neuromorphic VLSI learning systems has to do with the assumptions made on the particular form of the performance criterion being optimized. In typical physical systems, the learning objectives can not be clearly defined in terms of a target response or desired state of the system. Learning from external discrete rewards, in absence of a well-defined training signal, requires internal mechanisms of credit assignment which make no prior assumptions on the causal relationships of the system and the environment in which it operates. The stereotypical example of a system able to learn from a discrete delayed reward or punishment signal is the pole-balancer trained with reinforcement learning [18].

We use stochastic perturbative algorithms for model-free estimation of gradient information [16] in a general framework that includes reinforcement learning under delayed and discontinuous rewards [17]-[21], suitable for learning in physical systems of which the characteristics nor the optimization objectives are properly defined. Stochastic error-descent architectures for supervised learning [22] and computational primitives of reinforcement learning are combined into an analog VLSI architecture which offers a modular and cellular structure, model-free distributed representation, and robustness to noise and mismatches in the implementation. The combined architecture is applicable to the most general of learning tasks, where an unknown “black-box” dynamical system is adapted using an external “black-box” reinforcement-based delayed and possibly discrete reward signal.

As a proof of principle, we apply the model-free training-free adaptive techniques to blind optimization of a second-order noise-shaping modulator for oversampled data conversion, controlled by a neural classifier. The only evaluative feedback used in training the classifier is a discrete failure signal which indicates when some of the integrators in the modulation loop saturate.

In the following, we review supervised learning and stochastic perturbative techniques, and present a corresponding architecture for analog VLSI implementation. We then cover a generalized form of reinforcement learning, and introduce a stochastic perturbative analog VLSI architecture for reinforcement learning. Neuromorphic implementations in analog VLSI and system examples are also included.

## 2 Supervised Learning

In a metaphorical sense, supervised learning assumes the luxury of a committed “teacher”, who constantly evaluates and corrects the network by continuously feeding it target values for all network outputs. Supervised learning can be reformulated as an optimization task, where the network parameters (weights) are adjusted to minimize the distance between the targets and actual network outputs. Generalization and overtraining are important issues in supervised learning, and are beyond the scope of this paper.

Let  $\mathbf{y}(t)$  be the vector of network outputs with components  $y_i(t)$ , and correspondingly  $\mathbf{y}^{\text{target}}(t)$  be the supplied target output vector. The network contains adjustable parameters (or weights)  $\mathbf{p}$  with components  $p_k$ , and state variables  $\mathbf{x}(t)$  with components  $x_i(t)$  (which may

contain external inputs). Then the task is to minimize the scalar error index

$$\mathcal{E}(\mathbf{p}; t) = \sum_i |y_i^{\text{target}}(t) - y_i(t)|^\nu . \quad (1)$$

in the parameters  $p_i$ , using a distance metric with norm  $\nu > 0$ .

## 2.1 Gradient Descent

Gradient descent is the most common optimization technique for supervised learning in neural networks, which includes the widely used technique of backpropagation (or “dynamic feed-back”) [23] for gradient derivation, applicable to general feedforward multilayered networks.

In general terms, gradient descent minimizes the scalar performance index  $\mathcal{E}$  by specifying incremental updates in the parameter vector  $\mathbf{p}$  according to the error gradient  $\nabla_{\mathbf{p}}\mathcal{E}$ :

$$\mathbf{p}(t + 1) = \mathbf{p}(t) - \eta \nabla_{\mathbf{p}}\mathcal{E}(t) . \quad (2)$$

One significant problem with gradient descent and its variants for on-line supervised learning is the complexity of calculating the error gradient components  $\partial\mathcal{E}/\partial p_k$  from a model of the system. This is especially so for complex systems involving internal dynamics in the state variables  $x_j(t)$ :

$$\frac{\partial\mathcal{E}}{\partial p_k} = \sum_{i,j} \frac{\partial\mathcal{E}(t)}{\partial y_i} \cdot \frac{\partial y_i(t)}{\partial x_j} \cdot \frac{\partial x_j(t)}{\partial p_k} \quad (3)$$

where derivation of the dependencies  $\partial x_j/\partial p_k$  over time constitutes a significant amount of computation that typically scales super-linearly with the dimension of the network [22]. Furthermore, the derivation of the gradient in (3) assumes accurate knowledge of the model of the network ( $\mathbf{y}(t)$  as a function of  $\mathbf{x}(t)$ , and recurrence relations in the state variables  $\mathbf{x}(t)$ ). Accurate model knowledge cannot be assumed for analog VLSI neural hardware, due to mismatches in the physical implementation which can not be predicted at the time of fabrication. Finally, often a model for the system being optimized may not be readily available, or may be too complicated for practical (real-time) evaluation. In such cases, a black-box approach to optimization is more effective in every regard. This motivates the use of the well-known technique of stochastic approximation [24] for blind optimization in analog VLSI systems. We apply this technique to supervised learning as well as to more advanced models of “reinforcement” learning under discrete delayed rewards. The connection between stochastic approximation techniques and principles of neuromorphic engineering will be illustrated further below, in contrast with gradient descent.

## 2.2 Stochastic Approximation Techniques

Stochastic approximation algorithms [24] have long been known as effective tools for constrained and unconstrained optimization under noisy observations of system variables [25]. Applied to on-line minimization of an error index  $\mathcal{E}(\mathbf{p})$ , the algorithms avoid the computational burden of gradient estimation by directly observing the dependence of the index  $\mathcal{E}$  on randomly

applied perturbations in the parameter values. Variants on the Kiefer-Wolfowitz algorithm for stochastic approximation [24], essentially similar to random-direction finite-difference gradient descent, have been formulated for blind adaptive control [26], neural networks [16],[27] and the implementation of learning functions in VLSI hardware [28],[29],[22],[30].

The broader class of neural network learning algorithms under this category exhibit the desirable property that the functional form of the parameter updates is “model-free”, *i.e.*, independent of the model specifics of the network or system under optimization. The model-free techniques for on-line supervised learning are directly applicable to almost any observable system with deterministic, slowly varying, and possibly unknown characteristics. Parallel implementation of the stochastic approximation algorithms results into efficient and modular learning architectures that map well onto VLSI hardware. Since those algorithms use only direct function evaluations and no derivative information, they are functionally simple, and their implementation is independent of the structure of the system under optimization. They exhibit robust convergence properties in the presence of noise in the system and model mismatches in the implementation.

A brief description of the stochastic error-descent algorithm follows below, as introduced in [22] for efficient supervised learning in analog VLSI. The integrated analog VLSI continuous-time learning system used in [31], [32] forms the basis for the architectures outlined in the sections that follow.

### 2.3 Stochastic Supervised Learning

Let  $\mathcal{E}(\mathbf{p})$  be the error functional to be minimized, with  $\mathcal{E}$  a scalar deterministic function in the parameter (or weight) vector  $\mathbf{p}$  with components  $p_i$ . The stochastic algorithm specifies incremental updates in the parameters  $p_i$  as with gradient descent (2), although using a stochastic approximation to the true gradient

$$\frac{\partial \mathcal{E}(t)^{\text{est}}}{\partial p_i} = \pi_i(t) \cdot \hat{\mathcal{E}}(t) \quad (4)$$

where the differentially perturbed error

$$\hat{\mathcal{E}}(t) = \frac{1}{2\sigma^2} (\mathcal{E}(\mathbf{p}(t) + \boldsymbol{\pi}(t)) - \mathcal{E}(\mathbf{p}(t) - \boldsymbol{\pi}(t))) \quad (5)$$

is obtained from two direct observations of  $\mathcal{E}$  under complementary activation of a parallel random perturbation vector  $\boldsymbol{\pi}(t)$  with components  $\pi_i(t)$  onto the parameter vector  $\mathbf{p}(t)$ . The perturbation components  $\pi_i(t)$  are fixed in amplitude and random in sign,  $\pi_i(t) = \pm\sigma$  with equal probabilities for both polarities. The algorithm essentially performs gradient descent in random directions in the parameter space, as defined by the position of the perturbation vector.

As with exact gradient descent, iteration of the updates using (4) converges in the close proximity of a (local) minimum of  $\mathcal{E}$ , provided the perturbation amplitude  $\sigma$  is sufficiently small. The rate of convergence is necessarily slower than gradient descent, since every observation (5)

only reveals scalar information about the gradient vector in one dimension. However, the amount of computation required to compute the gradient at every update may outweigh the higher convergence rate offered by gradient descent, depending on the model complexity of the system under optimization. When applied to on-line supervised learning in recurrent dynamical systems, the stochastic algorithm provides a net computational efficiency rivaling that of exact gradient descent. Computational efficiency is defined in terms of the total number of operations required to converge, *i.e.*, reach a certain level of  $\mathcal{E}$ . A formal derivation of the convergence properties is presented in [22].

## 2.4 Supervised Learning Architecture

While alternative optimization techniques based on higher-order extensions on gradient descent will certainly offer superior convergence rates, the above stochastic method achieves its relative efficiency at a much reduced complexity of implementation. The only global operations required are the evaluations of the error function in (5), which are obtained from direct observations on the system under complementary activation of the perturbation vector. The operations needed to generate and apply the random perturbations, and to perform the parameter update increments, are strictly local and identical for each of the parameter components. The functional diagram of the local parameter processing cell, embedded in the system under optimization, is shown in Figure 1. The complementary perturbations and the corresponding error observations are performed in two separate phases on the same system, rather than concurrently on separate replications of the system. The sequential activation of the complementary perturbations and corresponding evaluations of  $\mathcal{E}$  are synchronized and coordinated with a three-phase clock:

$$\begin{aligned} \phi^0 & : \mathcal{E}(\mathbf{p}, t) \\ \phi^+ & : \mathcal{E}(\mathbf{p} + \boldsymbol{\pi}, t) \\ \phi^- & : \mathcal{E}(\mathbf{p} - \boldsymbol{\pi}, t) . \end{aligned} \tag{6}$$

This is represented schematically in Figure 1 by a modulation signal  $\phi(t)$ , taking values  $\{-1, 0, 1\}$ . The extra phase  $\phi^0$  ( $\phi(t) = 0$ ) is not strictly needed to compute (5)—it is useful otherwise, *e.g.* to compute finite difference estimates of second order derivatives for dynamic optimization of the learning rate  $\eta(t)$ .

The local operations are further simplified owing to the binary nature of the perturbations, reducing the multiplication in (4) to an exclusive-or logical operation, and the modulation by  $\phi(t)$  to binary multiplexing. Besides efficiency of implementation, this has a beneficial effect on the overall accuracy of the implemented learning system, as will be explained in the context of VLSI circuit implementation below.

## 2.5 Supervised Learning in Dynamical Systems

In the above, it was assumed that the error functional  $\mathcal{E}(\mathbf{p})$  is directly observable on the system by applying the parameter values  $p_i$ . In the context of on-line supervised learning in dynamical

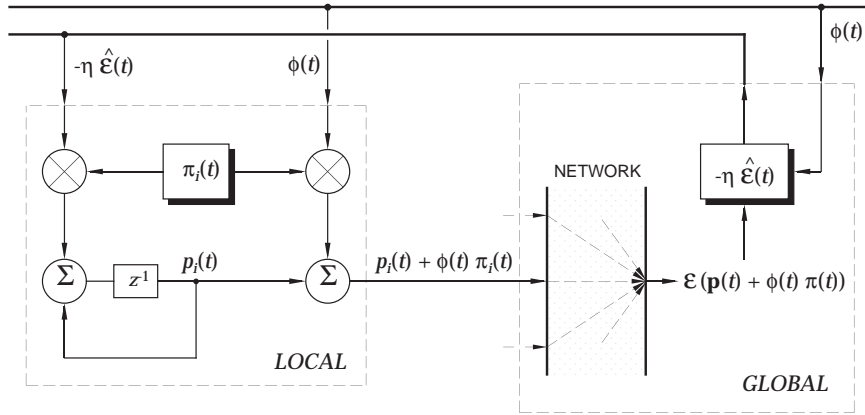


Figure 1: Architecture implementing stochastic error-descent supervised learning. The learning cell is locally embedded in the network. The differential error index is evaluated globally, and communicated to all cells.

systems, the error functional takes the form of the average distance of norm  $\nu$  between the output and target signals over a moving time window,

$$\mathcal{E}(\mathbf{p}; t_i, t_f) = \int_{t_i}^{t_f} \sum_i |y_i^{\text{target}}(t') - y_i(t')|^\nu dt', \quad (7)$$

which implicitly depends on the training sequence  $\mathbf{y}^{\text{target}}(t)$  and on initial conditions on the internal state variables of the system. An on-line implementation prohibits simultaneous observation of the error measure (7) under different instances of the parameter vector  $\mathbf{p}$ , as would be required to evaluate (5) for construction of the parameter updates. However, when the training signals are periodic and the interval  $T = t_f - t_i$  spans an integer multiple of periods, the measure (7) under constant parameter values is approximately invariant to time. In that case, the two error observations needed in (5) can be performed in sequence, under complementary piecewise constant activation of the perturbation vector.

In the context of on-line supervised learning in dynamical systems, the requirement of periodicity on the training signals is a limitation of the stochastic error-descent algorithm. Next, this requirement will be relaxed, along with some more stringent assumptions on the nature of supervised learning. In particular, a target training signal will no longer be necessary. Instead, learning is based on an external reward signal that conveys only partial and delayed information about the performance of the network.

### 3 Learning From Delayed and Discrete Rewards

Supervised learning methods rely on a continuous training signal that gives constant feedback about the direction in which to steer the response of the network to improve its performance. This continuous signal is available in the form of target values  $\mathbf{y}^{\text{target}}(t)$  for the network outputs  $\mathbf{y}(t)$ . More widely applicable but also more challenging are learning tasks where target outputs

or other continuous teacher feedback are not available, but instead only a non-steady, delayed reward (or punishment) signal is available to evaluate the *quality* of the outputs (or actions) produced by the network. The difficulty lies in assigning proper credit for the reward (or punishment) to actions that were produced by the network in the past, and adapting the network accordingly in such a way to *reinforce* the network actions leading to reward (and conversely, avoid those leading to punishment).

### 3.1 Reinforcement Learning Algorithms

Several iterative approaches to dynamic programming have been applied to solve the credit assignment problem for training a neural network with delayed rewards [17]-[21]. They all invoke an “adaptive critic element” which is trained along with the network to predict the future reward signal from the present state of the network. We define “reinforcement learning” essentially as given in [18], which includes as special cases “time difference learning” or TD( $\lambda$ ) [19], and, to some extent, Q-learning [20] and “advanced heuristic dynamic programming” [21]. The equations are listed below in general form to clarify the similarity with the above supervised perturbative learning techniques. It will then be shown how the above architectures are extended to allow learning from delayed and/or impulsive rewards.

Let  $r(t)$  be the discrete delayed reward signal for state vector  $\mathbf{x}(t)$  of the system (components  $x_j(t)$ ).  $r(t)$  is zero when no signal is available, and is negative for a punishment. Let  $y(t)$  be the (scalar) output of the network in response to an input (or state)  $\mathbf{x}(t)$ , and  $q(t)$  the predicted future reward (or “value function”) associated with state  $\mathbf{x}(t)$  as produced by the adaptive critic element. The action taken by the system is determined by the polarity of the network output,  $\text{sign}(y(t))$ . For example, in the pole balancing experiment of [18],  $y(t)$  is hard-limited and controls the direction of the fixed-amplitude force exerted on the moving cart. Finally, let  $\mathbf{w}$  and  $\mathbf{v}$  (components  $w_i$  and  $v_i$ ) be the weights of the network and the adaptive critic element, respectively. Then the weight updates are given by

$$\begin{aligned} \Delta w_i(t) &= w_i(t+1) - w_i(t) \\ &= \alpha \hat{r}(t) \cdot e_i(t) \end{aligned} \tag{8}$$

$$\begin{aligned} \Delta v_i(t) &= v_i(t+1) - v_i(t) \\ &= \beta \hat{r}(t) \cdot d_i(t) \end{aligned} \tag{9}$$

where the “eligibility” functions  $e_i(t)$  and  $d_i(t)$  are updated as

$$e_i(t+1) = \delta e_i(t) + (1 - \delta) \text{sign}(y(t)) \frac{\partial y(t)}{\partial w_i} \tag{10}$$

$$d_i(t+1) = \lambda d_i(t) + (1 - \lambda) \frac{\partial q(t)}{\partial v_i}$$

and the reinforcement  $\hat{r}(t)$  is given by

$$\hat{r}(t) = r(t) + \gamma q(t) - q(t-1) . \tag{11}$$

The parameters  $\delta$  and  $\lambda$  define the time span of credit assigned by  $e_i(t)$  and  $d_i(t)$  to actions in the past, weighting recent actions stronger than past actions:

$$\begin{aligned} e(t) &= (1 - \delta) \sum_{t'=-\infty}^{t-1} \delta^{t-t'-1} \text{sign}(y(t')) \frac{\partial y(t')}{\partial w_i} \\ d(t) &= (1 - \lambda) \sum_{t'=-\infty}^{t-1} \lambda^{t-t'-1} \frac{\partial q(t')}{\partial v_i}. \end{aligned} \quad (12)$$

Similarly, the parameter  $\gamma$  defines the time span for the prediction of future reward by  $q(t)$ , at convergence:

$$q(t) \approx \sum_{t'=t+1}^{\infty} \gamma^{t'-t-1} r(t'). \quad (13)$$

For  $\gamma = 1$  and  $y \equiv q$ , the equations reduce to TD( $\lambda$ ). Convergence of TD( $\lambda$ ) with probability one has been proven in the general case of linear networks of the form  $q = \sum v_i x_i$  [33].

Learning algorithms of this type are neuromorphic in the sense that they emulate classical (pavlovian) conditioning in pattern association as found in biological systems [6] and their mathematical and cognitive models [34],[7]. Furthermore, as shown below, the algorithms lend themselves to analog VLSI implementation in a parallel distributed architecture which, unlike more complicated gradient-based schemes, resembles the general structure and connectivity of biological neural systems.

### 3.2 Reinforcement Learning Architecture

While reinforcement learning does not perform gradient descent of a (known) error functional, the eligibility functions  $e_i(t)$  and  $d_i(t)$  used in the weight updates are constructed from derivatives of output functions to the weights. The eligibility functions in equation (10) can be explicitly restated as (low-pass filtered) gradients of an error function

$$\mathcal{E}(t) = |y(t)| + q(t) \quad (14)$$

with

$$\begin{aligned} e_i(t+1) &= \delta e_i(t) + (1 - \delta) \frac{\partial \mathcal{E}(t)}{\partial w_i} \\ d_i(t+1) &= \lambda d_i(t) + (1 - \lambda) \frac{\partial \mathcal{E}(t)}{\partial v_i}. \end{aligned} \quad (15)$$

Rather than calculating the gradients in (15) from the network model, we can again apply stochastic perturbative techniques to estimate the gradients from direct evaluations on the network. Doing so, all properties of robustness, scalability and modularity that apply to stochastic error descent supervised learning apply here as well. As in (4), stochastic approximation estimates of the gradient components in (15) are

$$\frac{\partial \mathcal{E}(t)^{\text{est}}}{\partial w_i} = \omega_i(t) \hat{\mathcal{E}}(t) \quad (16)$$



$$\frac{\partial \mathcal{E}(t)^{\text{est}}}{\partial v_i} = v_i(t) \hat{\mathcal{E}}(t)$$

where the differential perturbed error

$$\hat{\mathcal{E}}(t) = \frac{1}{2\sigma^2} (\mathcal{E}(\mathbf{w} + \boldsymbol{\omega}, \mathbf{v} + \mathbf{v}, t) - \mathcal{E}(\mathbf{w} - \boldsymbol{\omega}, \mathbf{v} - \mathbf{v}, t)) \quad (17)$$

is obtained from two-sided parallel random perturbation  $\mathbf{w} \pm \boldsymbol{\omega}$  simultaneous with  $\mathbf{v} \pm \mathbf{v}$  ( $|\omega_i| = |v_i| = \sigma$ ).

A side benefit of the low-pass filtering of the gradient in (15) is an improvement of the stochastic gradient estimate (16) through averaging. As with stochastic error descent supervised learning, averaging reduces the variance of the gradient estimate and produces learning increments that are less stochastic in nature, although this is not essential for convergence of the learning process [22].

Figure 2 shows the block diagram of a reinforcement learning cell and an adaptive critic cell, with stochastic perturbative estimation of the gradient according to (16).  $\text{LP}_\delta$  and  $\text{LP}_\lambda$  denote first-order low-pass filters (15) with time constants determined by  $\delta$  and  $\lambda$ . Other than the low-pass filtering and the global multiplicative factor  $\hat{r}(t)$ , the architecture is identical to that of stochastic error descent learning in Figure 1. As before, the estimation of  $\hat{\mathcal{E}}(t)$  does not require separate instances of perturbed and non-perturbed networks shown in Figure 2, and can be computed sequentially by evaluating the output of the network and adaptive critic in three phases for every cycle of  $t$ :

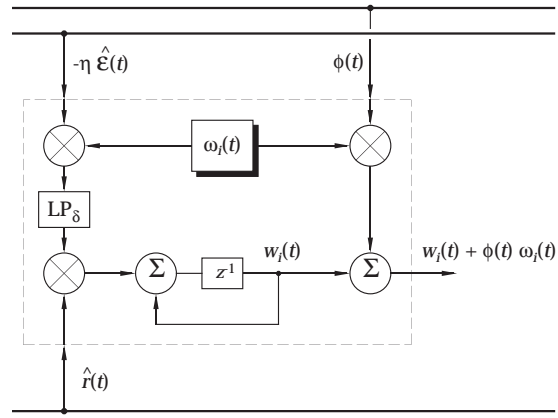
$$\begin{aligned} \phi^0 & : \mathcal{E}(\mathbf{w}, \mathbf{v}, t) \\ \phi^+ & : \mathcal{E}(\mathbf{w} + \boldsymbol{\omega}, \mathbf{v} + \mathbf{v}, t) \\ \phi^- & : \mathcal{E}(\mathbf{w} - \boldsymbol{\omega}, \mathbf{v} - \mathbf{v}, t) . \end{aligned} \quad (18)$$

In systems with a continuous-time output response, we assume that the time lag between consecutive observations of the three phases of  $\mathcal{E}$  is not an issue, which amounts to choosing an appropriate sampling rate for  $t$  in relation to the bandwidth of the system.

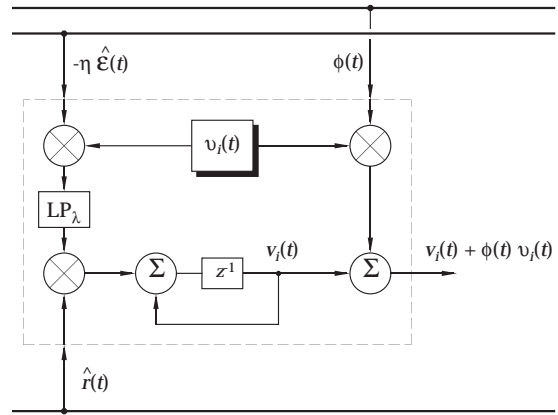
Similarities between the above cellular architectures for supervised learning and reinforcement learning are apparent: both correlate local perturbation values  $\pi_i$ ,  $\omega_i$  or  $v_i$  with a global scalar index  $\hat{\mathcal{E}}$  that encodes the differential effect of the perturbations on the output, and both incrementally update the weights  $p_i$ ,  $w_i$  or  $v_i$  accordingly. The main difference in reinforcement learning is the additional gating of the correlate product with a global reinforcement signal  $\hat{r}$  after temporal filtering. For many applications, the extra overhead that this implies in hardware resources is more than compensated by the utility of the reward-based credit assignment mechanism, which does not require an external teacher. An example is given below in the case of oversampled A/D conversion.

### 3.3 System Example: Stable Higher-Order Noise-Shaping Modulation

We evaluated both exact gradient and stochastic perturbative embodiments of the reinforcement learning algorithms on an adaptive neural classifier, controlling a higher-order noise-



(a)



(b)

Figure 2: General architecture implementing reinforcement learning using stochastic gradient approximation. (a) Reinforcement learning cell. (b) Adaptive critic cell.

shaping modulator used for oversampled A/D data conversion [37]. The order- $n$  modulator comprises a cascade of  $n$  integrators  $x_i(t)$  operating on the difference between the analog input  $u(t)$  and the binary modulated output  $y(t)$ :

$$\begin{aligned} x_0(t+1) &= x_0(t) + a(u(t) - y(t)) \\ x_i(t+1) &= x_i(t) + a x_{i-1}(t), \quad i = 1, \dots, n-1 \end{aligned} \quad (19)$$

where  $a = 0.5$ . The control objective is to choose the binary sequence  $y(t)$  such as to keep the excursion of the integration variables within bounds,  $|x_i(t)| < x_{\text{sat}}$  [36].

For the adaptive classifier, we specify a one-hidden-layer neural network, with inputs  $x_i(t)$  and output  $y(t)$ . The network has  $m$  hidden units, a  $\tanh(\cdot)$  sigmoidal nonlinearity in the hidden layer, and a linear output layer. For the simulations we set  $n = 2$  and  $m = 5$ . The case  $n = 2$  is equivalent to the single pole-balancing problem [18].

The only evaluative feedback signal used during learning is a failure signal which indicates when one or more integration variables saturate,  $|x_i(t)| \geq x_{\text{sat}}$ . In particular, the signal  $r(t)$  counts the number of integrators in saturation:

$$r(t) = -b \sum_i H(|x_i(t)| - x_{\text{sat}}) \quad (20)$$

where  $b = 10$ , and where  $H(\cdot)$  denotes a step function ( $H(x) = 1$  if  $x > 0$  and 0 otherwise). The adaptive critic  $q(t)$  is implemented with a neural network of identical structure as for  $y(t)$ . The learning parameters in (8), (10) and (11) are  $\delta = 0.8$ ,  $\lambda = 0.7$ ,  $\gamma = 0.9$ ,  $\alpha = 0.05$  and  $\beta = 0.001$ . These values are consistent with [18], adapted to accommodate for differences in the time scale of the dynamics (19). The perturbation strength in the stochastic version is  $\sigma = 0.01$ .

Figure 3 shows the learning performance for several trials of both versions of reinforcement learning, using exact and stochastic gradient estimates. During learning, the input sequence  $u(t)$  is random, uniform in the range  $-0.5 \dots 0.5$ . Initially, and every time failure occurs ( $r(t) < 0$ ), the integration variables  $x_i(t)$  and eligibilities  $e_k(t)$  and  $d_k(t)$  are reset to zero. Qualitative differences observed between the exact and stochastic versions in Figure 3 are minor. Further, in all but one of the 20 cases tried, learning has completed (*i.e.*, consequent failure is not observed in finite time) in fewer than 20 consecutive trial-and-error iterations. Notice that a non-zero  $r(t)$  is only generated at failure, *i.e.*, less than 20 times, and no other external evaluative feedback is needed for learning.

Figure 4 quantifies the effect of stochastic perturbative estimation of the gradients (15) on the quality of reinforcement learning. The correlation index  $c(t)$  measures the degree of conformity in the eligibilities (both  $e_i(t)$  and  $d_i(t)$ ) between stochastic and exact versions of reinforcement learning. Correlation is expressed as usual on a scale from  $-1$  to  $1$ , with  $c(t) = 1$  indicating perfect coincidence. While  $c(t)$  is considerably less than 1 in all cases,  $c(t) > 0$  about 95 % of the time, meaning that *on average* the sign of the parameter updates (8) for exact and stochastic versions are consistent in at least 95 % of the cases. The scatterplot  $c(t)$  vs.  $\hat{r}(t)$  also illustrates how the adaptive critic produces a positive reinforcement  $\hat{r}(t)$  in most

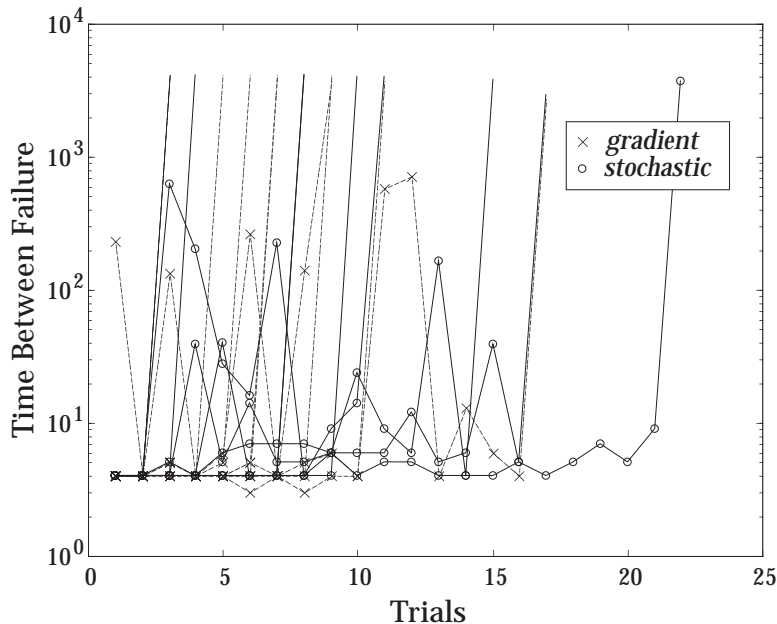


Figure 3: Simulated performance of stochastic perturbative (o) and gradient-based (x) reinforcement learning in a second-order noise-shaping modulator. Time between failures for consecutive trials from zero initial conditions.

of the cases, even though the “reward” signal  $r(t)$  is never positive by construction. Positive reinforcement  $\hat{r}(t)$  under idle conditions of  $r(t)$  is desirable for stability. Notice that the failure-driven punishment points (where  $r(t) < 0$ ) are off-scale of the graph and strongly negative.

We tried reinforcement learning on higher-order modulators,  $n = 3$  and higher. Both exact and stochastic versions were successful for  $n = 3$  in the majority of cases, but failed to converge for  $n = 4$  with the same parameter settings. On itself, this is not surprising since higher-order delta-sigma modulators tend to become increasingly prone to unstabilities and sensitive to small changes in parameters with increasing order  $n$ , which is why they are almost never used in practice [37]. It is possible that more advanced reinforcement learning techniques such as “Advanced Heuristic Dynamic Programming” (AHDP) [21] would succeed to converge for orders  $n > 3$ . AHDP offers improved learning efficiency using a more advanced, gradient-based adaptive critic element for prediction of reward, although it is not clear at present how to map the algorithm efficiently onto analog VLSI.

The above stochastic perturbative architectures for both supervised and reinforcement learning support common “neuromorphs” and corresponding analog VLSI implementations. Neuromorphs of learning in analog VLSI are the subject of next section.

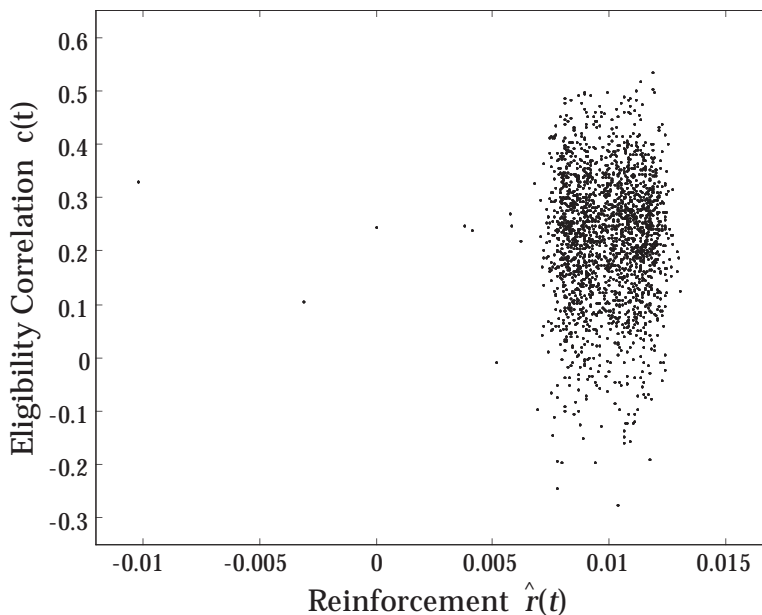


Figure 4: Effect of stochastic perturbative gradient estimation on the reinforcement learning.  $c(t)$  quantifies the degree of conformity in the eligibilities  $e_i(t)$  and  $d_i(t)$  between exact and stochastic versions.

## 4 Neuromorphic Analog VLSI Learning

### 4.1 Adaptive Circuits

The model-free nature of the stochastic perturbative learning algorithms does not impose any particular conditions on the implementation of computational functions required for learning. By far the most critical element in limiting learning performance is the quality of the parameter update increments and decrements, in particular the correctness of their polarity. Relative fluctuations in amplitude of the learning updates do not affect the learning process to first order, since their effect is equivalent to relative fluctuations in the learning rate. On the other hand, errors in the polarity of the learning updates might adversely affect learning performance even at small update amplitudes.

A binary controlled charge-pump adaptive element is described next. Offering precise control of the update polarity, this circuit element provides the primitives for learning as well as memory in the analog VLSI systems covered further below.

#### 4.1.1 Charge-pump adaptive element

Figure 5 shows the circuit diagram of a charge-pump adaptive element implementing a volatile synapse. The circuit is a simplified version of the charge pump used in [14] and [32]. When enabled by  $EN_n$  and  $EN_p$  (at GND and  $V_{dd}$  potentials, respectively), the circuit generates an incremental update of which the polarity is determined by  $\overline{POL}$ . The amplitude of the

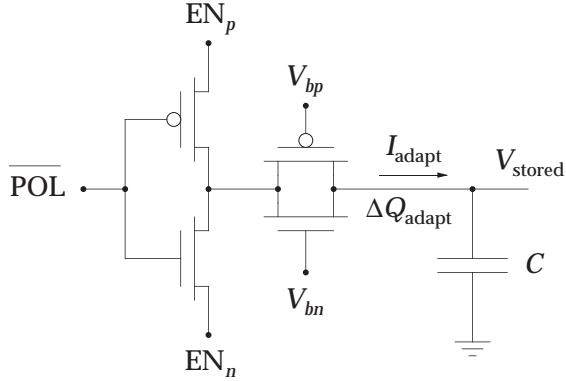


Figure 5: Charge-pump adaptive element implementing a volatile synapse.

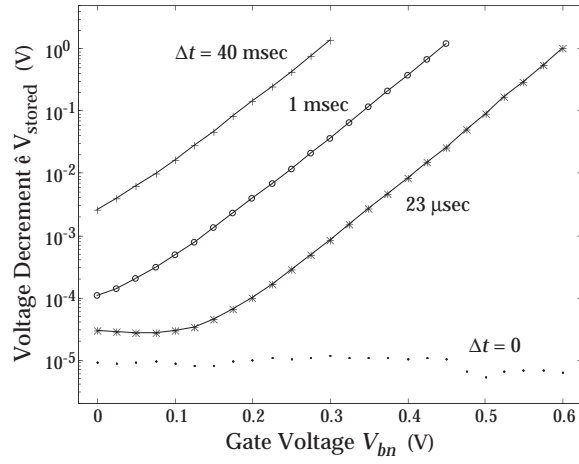
current supplying the incremental update is determined by gate voltages  $V_{bn}$  and  $V_{bp}$ , biased deep in subthreshold to allow fine (sub-fC) increments if needed. The increment amplitude is also determined by the duration of the enabled current, controlled by the timing of  $EN_n$  and  $EN_p$ . When both  $EN_n$  and  $EN_p$  are set midway between GND and  $V_{dd}$ , the current output is disabled. Notice that the switch-off transient is (virtually) free of clock feedthrough charge injection, because the current-supplying transistors are switched from their source terminals, with the gate terminals being kept at constant voltage [14].

Measurements on a charge pump with  $C = 0.5$  pF fabricated in a  $2 \mu\text{m}$  CMOS process are shown in Figure 6. Under pulsed activation of  $EN_n$  and  $EN_p$ , the resulting voltage increments and decrements are recorded as a function of the gate bias voltages  $V_{bn}$  and  $V_{bp}$ , for both polarities of POL, and for three different values of the pulse width  $\Delta t$  ( $23 \mu\text{sec}$ ,  $1 \text{ msec}$  and  $40 \text{ msec}$ ). In all tests, the pulse period extends 2 msec beyond the pulse width. The exponential subthreshold characteristics are evident from Figure 6, with increments and decrements spanning four orders of magnitude in amplitude. The lower limit is mainly determined by junction diode leakage currents, as shown in Figure 6 (a) for  $\Delta t = 0$  ( $0.01 \text{ mV}$  per 2 msec interval at room temperature). This is more than adequate to accommodate learning over a typical range of learning rates. Also, the binary control POL of the polarity of the update is effective for increments and decrements down to  $0.05 \text{ mV}$  in amplitude, corresponding to charge transfers of only a few hundred electrons.

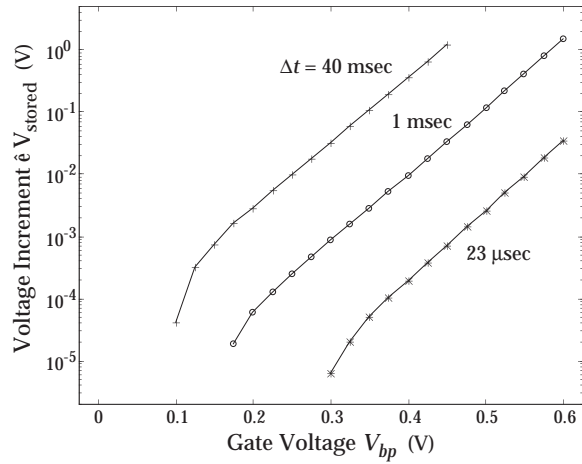
#### 4.1.2 Analog storage

Because of the volatile nature of the adaptive element used, a dynamic refresh mechanism is required if long-term local storage of the weight values after learning is desired. A robust and efficient self-contained mechanism that does not require external storage is “partial incremental refresh” [14]

$$p_i(t+1) = p_i(t) - \delta Q(p_i(t)) \quad . \quad (21)$$



(a)



(b)

Figure 6: Measured characteristics of charge-pump adaptive element in  $2 \mu\text{m}$  CMOS with  $C = 0.5 \text{ pF}$ . (a) n-type decrements ( $\text{POL} = 0$ ); (b) p-t increments ( $\text{POL} = 1$ ).

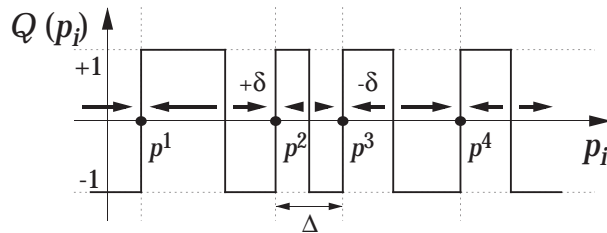


Figure 7: Example illustrating binary quantization  $Q$  and partial incremental refresh.

obtained from binary quantization  $Q$  of the parameter value. Stable discrete states of the analog dynamic memory under periodic actication of (21) are located at the positive transitions of  $Q$ , illustrated in Figure 7. Long-term stability and robustness to noise and errors in the quantization requires that the separation between neighboring discrete states  $\Delta$  be much larger than the amplitude of the parameter updates  $\delta$ , which in turn needs to exceed the spontaneous drift in the parameter value due to leakage between consecutive refresh cycles [14].

Partial incremental refresh can be directly implemented using the adaptive element in Figure 8 by driving  $\overline{\text{POL}}$  with a binary function of the weight value [15]. As in [15], the binary quantization function can be multiplexed over an array of storage cells, and can be implemented by retaining the LSB from A/D/A conversion [41] of the value to be stored. Experimental observation of quantization and refresh in a fabricated 128-element array of memory cells has confirmed stable retention of analog storage at 8-bit effective resolution over a time interval exceeding  $10^9$  refresh cycles (several days) [15].

A non-volatile equivalent of the charge-pump adaptive element in Figure 5, which does not require dynamic refresh, is described in [9]. Correspondingly, a non-volatile learning cell performing stochastic error descent can be obtained by substitution of the core adaptive element in Figure 8 below, and more intricate volatile and non-volatile circuits implementing stochastic reinforcement learning can be derived from extensions on Figure 8 and [9]. The non-volatile solution is especially attractive if long-term storage is a more pressing concern than speed of adaptation and flexibility of programming.

#### 4.1.3 Stochastic perturbative learning cell

The circuit schematic of a learning cell implementing stochastic error descent is given in Figure 8. The incremental update  $-\eta\pi_i\hat{\mathcal{E}}$  to be performed in (5) is first decomposed in amplitude and sign components. This allows for a hybrid digital-analog implementation of the learning cell, in which amplitudes of certain operands are processed in analog format, and their polarities implemented in logic. Since  $|\pi_i| \equiv 1$ , the amplitude  $\eta|\hat{\mathcal{E}}|$  is conveniently communicated as a global signal to all cells, in the form of two gate voltages  $V_{bn}$  and  $V_{bp}$ . The (inverted) polarity  $\overline{\text{POL}}$  is obtained as the (inverted) exclusive-or combination of the perturbation  $\pi_i$  and the polarity of  $\hat{\mathcal{E}}$ . The decomposition of sign and amplitude ensures proper convergence of the learning increments in the presence of mismatches and offsets in the physical implementation of the learning cell. This is because the polarities of the increments are more accurately implemented through logic-controlled circuitry, which are independent of analog mismatches in the implementation.

The perturbation  $\pi_i$  is applied to  $p_i$  in three phases (6) by capacitive coupling onto the storage node  $C$ . The binary state of the local perturbation  $\pi_i$  selects one of two global perturbation signals to couple onto  $C$ . The perturbation signals ( $V_\sigma^+$  and its complement  $V_\sigma^-$ ) globally control the three phases  $\phi^0$ ,  $\phi^+$  and  $\phi^-$  of (6), and set the perturbation amplitude  $\sigma$ . The simple configuration using a one-bit multiplexer is possible because each perturbation component can only take one of two values  $\pm\sigma$ .



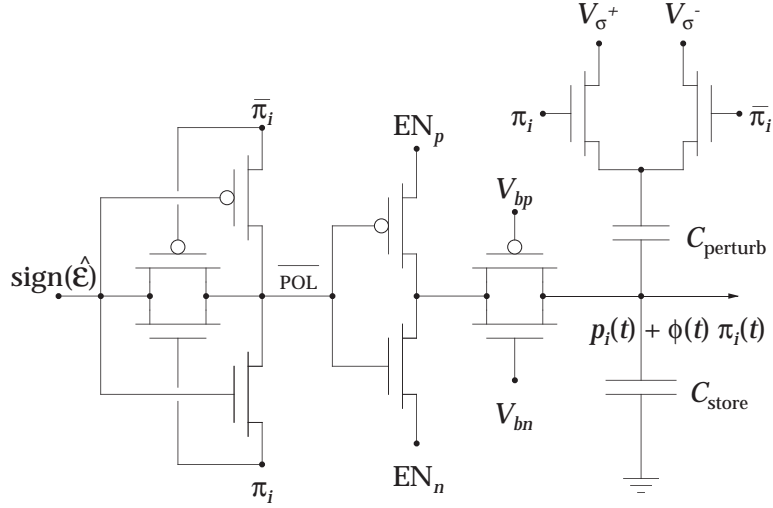


Figure 8: Circuit schematic of a learning cell implementing stochastic error descent, using the charge pump adaptive element.

## 4.2 Learning Systems

### 4.2.1 Continuous-time trajectory learning in an analog VLSI recurrent neural network

On-chip learning of continuous-time recurrent dynamics has been demonstrated in an analog VLSI neural network, using stochastic error descent [31],[32]. We briefly summarize the architecture, operation and results here. The chip contains an integrated network of six fully interconnected continuous-time neurons

$$\tau \frac{d}{dt} x_i = -x_i + \sum_{j=1}^6 W_{ij} \sigma(x_j - \theta_j) + y_i \quad , \quad (22)$$

with  $x_i(t)$  the neuron states representing the outputs of the network,  $y_i(t)$  the external inputs to the network, and  $\sigma(\cdot)$  a sigmoidal activation function. The 36 connection strengths  $W_{ij}$  and 6 thresholds  $\theta_j$  constitute the free parameters to be learned, and the time constant  $\tau$  is kept fixed and identical for all neurons.

The network is trained with target output signals  $x_1^T(t)$  and  $x_2^T(t)$  for two neuron outputs,  $i = 1, 2$ . The other four neurons are hidden to the output, and the internal dynamics of these hidden neuron state variables play an important part in optimizing the output. Learning consists of minimizing the time-averaged error (7) with respect to the parameters  $W_{ij}$  and  $\theta_j$ , using stochastic error descent. For a consistent evaluation of the stochastic gradient, the perturbed function measurements  $\mathcal{E}(\mathbf{p} \pm \boldsymbol{\pi})$  are performed on a time scale significantly (60 times) larger than the period of the target signals.

All local learning functions, including the generation of pseudo-random perturbations and the stochastic learning update, are embedded with the synaptic functions (22) in a scalable

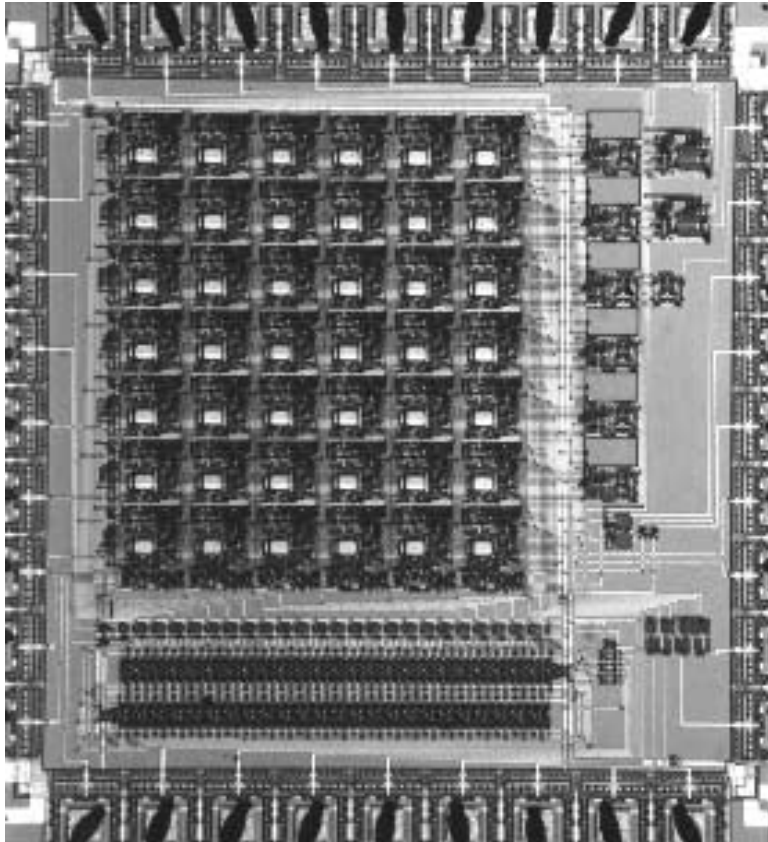


Figure 9: *Micrograph of an analog VLSI recurrent neural network chip that learns continuous-time internal dynamics using stochastic error-descent. Center:  $6 \times 7$  array of weight and threshold parameter cells with integrated learning and storage functions. Bottom: random binary array generator providing the parallel parameter perturbations.*

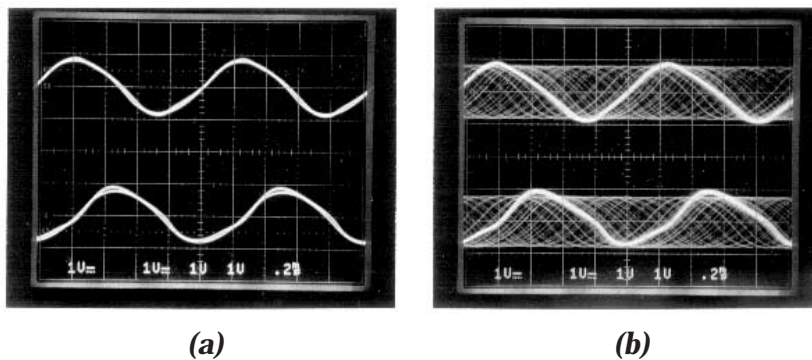


Figure 10: *Oscillograms of the network outputs and target signals after learning, (a) under weak residual teacher forcing, and (b) with teacher forcing removed. Top traces:  $x_1(t)$  and  $x_1^T(t)$ . Bottom traces:  $x_2(t)$  and  $x_2^T(t)$ .*

2-D array of parameter cells  $W_{ij}$  and  $\theta_j$ . The circuitry implementing the learning functions is essentially that of Figure 8. The dynamic refresh scheme described above is incorporated locally in the parameter cell for long-term storage of the parameters. A micrograph of the chip is shown in Figure 9. Power dissipation is 1.2 mW from a 5 V supply for a 1 kHz signal being trained.

The results of training the chip with a periodic analog target signal representing a quadrature-phase oscillator are illustrated in Figure 10. Learning is achieved in roughly 1500 training cycles of 60 msec each, using “teacher forcing” during training for synchronization between network and target dynamics, and by careful but unbiased choice of initial conditions for the weight parameters to avoid local minima. These conditions are less critical in more general applications of nonlinear system identification where the network during training is presented input signals to be associated with the target output signals.

#### 4.2.2 Reinforcement learning in a VLSI neural classifier for nonlinear noise-shaping delta-sigma modulation

A VLSI classifier consisting of 64 locally tuned, hard-thresholding neurons was trained using reinforcement learning to produce stable noise-shaping modulation of orders one and two [36]. While this system does not implement the stochastic version of reinforcement learning studied above, it presents a particularly simple VLSI implementation and serves to demonstrate some of the properties also expected of more advanced implementations that incorporate stochastic learning with continuous neurons.

Similar to the “boxes-system” used in [18], the classifier implements a look-up table from a binary address-encoded representation of the state space spanned by  $u(t)$  and  $x_i(t)$ . In particular,  $y(t) = y_{\chi(t)}$  and  $q(t) = q_{\chi(t)}$  where  $\chi(t)$  is the index of the address determined by hard-limiting thresholding operations on the components  $u(t)$  and  $x_i(t)$ . Each neuron cell, identified by address  $k$ , locally stores the two parameters  $y_k$  and  $q_k$  in analog format, and

updates them according to the external reinforcement signal defined in (20).

In its simplest form, the implemented reinforcement learning performs updates in the eligible  $y_k$  parameters opposite to their thresholded output values, each time failure occurs ( $r(t) = -1$ ). Hysteresis is included in the dynamics of the  $y_k$  updates to ensure some degree of stability under persistent negative reinforcement during training, even without use of the adaptive critic  $q_k$ . Although this simple form of reinforcement learning with non-adaptive hysteretic critic is not meant to be adequate for more complex tasks, it has proven sufficient to train the VLSI neural classifier to produce noise-shaping modulation of orders 1 and 2.

The integrated system contains a cascade of 6 integrators, an 11-bit address state encoder, and an address-encoded classifier with 64 reinforcement learning neurons on a  $2.2 \text{ mm} \times 2.2 \text{ mm}$  chip in  $2 \mu\text{m}$  CMOS technology. A record of a learning experiment reinforcing first-order noise-shaping modulation in the first integrator, using 2-bit address encoding  $\chi(t)$  of the polarities of  $u(t)$  and  $x_1(t)$ , is shown in Figure 11. As in the simulations above, the input sequence  $u(t)$  during training is uniformly random with half full-scale maximum amplitude (1 V pp), and the integrator variables  $x_i(t)$  as well as the eligibilities  $e_k(t)$  are reset to zero after every occurrence of failure,  $r(t) = -1$ . The dynamics of the state variables and parameters recorded in Figure 11 shows convergence after roughly 150 input presentations. The time step in the experiments was  $T = 2.5 \text{ msec}$ , limited by the bandwidth of the instrumentation equipment in the recording. Notice that the learned pattern of  $y_k$  at convergence conforms to that of a standard first-order delta-sigma modulator [37], which it should in this rather simple case. Learning succeeded at various values of the learning constants  $\delta$  and  $\alpha$ , affecting mainly the rate of convergence.

Tests for higher-order noise-shaping modulation on the same learning system only succeeded for order  $n = 2$ , using a total of 8 parameters  $y_k$ . For higher orders of noise-shaping, a continuous neuron representation and learning are required, as the above simulations of the stochastic reinforcement system indicate.

### 4.3 Structural Properties

The general structure of neuromorphic information processing systems has some properties differentiating them from some more conventional human-engineered computing machinery, which are typically optimized for general-purpose digital programming. Some of the desirable properties for neuromorphic architectures are: fault-tolerance and robustness to noise through a redundant distributed representation, robustness to changes in operating conditions through on-line adaptation, real-time bandwidth through massive parallelism, and modularity as a result of locality in space and time. We illustrate these properties in the two architectures for supervised and reinforcement learning in Figures 1 and 2. Since both architectures are similar, a distinction between them will not explicitly be made.

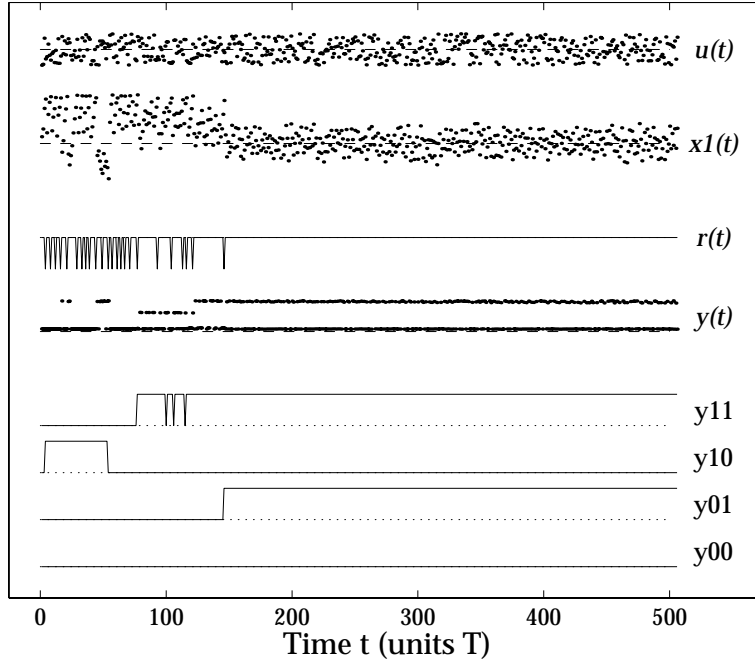


Figure 11: First-order modulator experiments: recorded dynamics of state variables and parameters during on-chip learning.

#### 4.3.1 Fault-tolerance through statistical averaging in a distributed representation:

Direct implementation of gradient descent, based on an explicit model of the network, is prone to errors due to unaccounted discrepancies in the network model and mismatches in the physical implementation of the gradient. This is due to the localized representation in the computation of the gradient as calculated from the model, in which any discrepancy in one part may drastically affect the final result. For this and other reasons, it is unlikely that biology performs explicit gradient calculation on complex systems such as recurrent neural networks with continuous-time dynamics. Stochastic error descent avoids errors of various kinds by physically probing the gradient onto the system rather than deriving it. Using simultaneous and uncorrelated parallel perturbations of the weights, the effect of a single error on the outcome is thus significantly reduced, by virtue of the statistical nature of the computation. However, critical in the accuracy of the implemented learning system is the precise derivation and faithful distribution of the global learning signals  $\hat{\mathcal{E}}(t)$  and  $\hat{r}(t)$ . Strictly speaking, it is essential only to guarantee the correct *polarity* and not the exact amplitude of the global learning signals, as implemented in Figure 8.

### 4.3.2 Robustness to changes in the environment through on-line adaptation:

This property is inherent to the on-line incremental nature of the studied supervised and reinforcement learning algorithms, which track structural changes in  $\hat{\mathcal{E}}(t)$  or  $\hat{r}(t)$  on a characteristic time scale determined by learning rate constants such as  $\eta$  or  $\alpha$  and  $\beta$ . Learning rates can be reduced as convergence is approached, as in the popular notion in cognitive neuroscience that neural plasticity decreases with age and experience.

### 4.3.3 Real-time bandwidth through parallelism:

All learning operations are performed in parallel, with exception of the three-phase perturbation scheme (6) or (18) to obtain the differential index  $\hat{\mathcal{E}}$  under sequential activation of complementary perturbations  $\boldsymbol{\pi}$  and  $-\boldsymbol{\pi}$ . We note that the synchronous three-phase scheme is not essential and could be replaced by an asynchronous perturbation scheme as in [16] and [42]. While this probably resembles biology more closely, the synchronous gradient estimate (4) using complementary perturbations is computationally more efficient as it cancels error terms up to second order in the perturbation strength  $\sigma$  [24]. In the asynchronous scheme, one could envision the role of random noise naturally present in biological systems as a source of perturbations, although it is not clear how noise sources can be effectively isolated to produce the correlation measures necessary for gradient estimation.

### 4.3.4 Modular architecture with local connectivity:

The learning operations are local in the sense that a need for excessive global interconnects between distant cells is avoided. The global signals are few in number and common for all cells, which implies that no signal interconnects are required *between* cells across the learning architecture, but all global signals are communicated uniformly *across* cells instead. This allows to embed the learning cells directly into the network (or adaptive critic) architecture, where they interface physically with the synapses they adapt, as in biological systems. The common global signals include the differential index  $\hat{\mathcal{E}}(t)$  and reinforcement signal  $\hat{r}(t)$ , besides common bias levels and timing signals.  $\hat{\mathcal{E}}(t)$  and  $\hat{r}(t)$  are obtained by any global mechanism that quantifies the “fitness” of the network response in terms of teacher target values or discrete rewards (punishments). Physiological experiments support evidence of local (hebbian [5]) and sparsely globally interconnected (reinforcement [6]) mechanisms of learning and adaptation in biological neural systems [3],[4].

## 5 Conclusion

Neuromorphic analog VLSI architectures for a general class of learning tasks have been presented, along with key components in their analog VLSI circuit implementation. The architectures make use of distributed stochastic techniques for robust estimation of gradient information, accurately probing the effect of parameter changes on the performance of the

network. Two architectures have been presented: one implementing stochastic error-descent for supervised learning, and the other implementing a stochastic variant on a generalized form of reinforcement learning. The two architectures are similar in structure, and both are suitable for scalable and robust analog VLSI implementation.

While both learning architectures can operate on (and be integrated in) arbitrary systems of which the characteristics and structure does not need to be known, the reinforcement learning architecture additionally supports a more general form of learning, using an arbitrary, externally supplied, reward or punishment signal. This allows the development of more powerful, generally applicable devices for “black-box” sensor-motor control which make no prior assumptions on the structure of the network and the specifics of the desired network response.

We presented results that demonstrate the effectiveness of perturbative stochastic gradient estimation for supervised learning and reinforcement learning, applied to nonlinear system identification and adaptive oversampled data conversion. A recurrent neural network was trained to generate internal dynamics producing a target periodic orbit at the outputs. A neural classifier controlling the second-order noise-shaping modulator was trained for optimal performance with no more evaluative feedback than a discrete failure signal indicating whenever any of the modulation integrators saturate. The critical part in the VLSI implementation of adaptive systems of this type is the precision of the *polarity*, rather than the amplitude, of the implemented weight parameter updates. A binary controlled charge-pump provides voltage increments and decrements of precise polarity spanning four orders of magnitude in amplitude, with charge transfers down to a few hundred electrons.

## References

- [1] C.A. Mead, “Neuromorphic Electronic Systems,” *Proceedings of the IEEE*, vol. **78** (10), pp 1629-1639, 1990.
- [2] C.A. Mead, *Analog VLSI and Neural Systems*, Reading, MA: Addison-Wesley, 1989.
- [3] G.M. Shepherd, *The Synaptic Organization of the Brain*, 3rd ed., New York, NY: Oxford Univ. Press, 1992.
- [4] P.S. Churchland and T.J. Sejnowski, *The Computational Brain*, Cambridge, MA: MIT Press, 1990.
- [5] S.R. Kelso and T.H. Brown, “Differential Conditioning of Associative Synaptic Enhancement in Hippocampal Brain Slices,” *Science*, vol. **232**, pp 85-87, 1986.
- [6] R.D. Hawkins, T.W. Abrams, T.J. Carew, and E.R. Kandell, “A Cellular Mechanism of Classical Conditioning in *Aplysia*: Activity-Dependent Amplification of Presynaptic Facilitation,” *Science*, iss. **219**, pp 400-405, 1983.

- [7] P.R. Montague, P. Dayan, C. Person and T.J. Sejnowski, "Bee Foraging in Uncertain Environments Using Predictive Hebbian Learning," *Nature*, vol. **377** (6551), pp. 725-728, 1996.
- [8] C.A. Mead and M. Ismail, Eds., *Analog VLSI Implementation of Neural Systems*, Norwell, MA: Kluwer, 1989.
- [9] C. Diorio, P. Hassler, B. Minch and C.A. Mead, "A Single-Transistor Silicon Synapse," to appear in *IEEE Transactions on Electron Devices*.
- [10] C.A. Mead, "Adaptive Retina," in *Analog VLSI Implementation of Neural Systems*, C. Mead and M. Ismail, Eds., Norwell, MA: Kluwer Academic Pub., 1989, pp 239-246.
- [11] G. Cauwenberghs, C.F. Neugebauer, and A. Yariv, "Analysis and Verification of an Analog VLSI Outer-Product Incremental Learning System," *IEEE Transactions on Neural Networks*, vol. **3** (3), pp 488-497, 1992.
- [12] Y. Horio, and S. Nakamura, "Analog Memories for VLSI Neurocomputing," in *Artificial Neural Networks: Paradigms, Applications, and Hardware Implementations*, C. Lau and E. Sanchez-Sinencio, Eds., IEEE Press, 1992, pp 344-363.
- [13] E. Vittoz, H. Oguey, M.A. Maher, O. Nys, E. Dijkstra, and M. Chevroulet, "Analog Storage of Adjustable Synaptic Weights," in *VLSI Design of Neural Networks*, Norwell MA: Kluwer Academic, pp 47-63, 1991.
- [14] G. Cauwenberghs, and A. Yariv, "Fault-Tolerant Dynamic Multi-Level Storage in Analog VLSI," *IEEE Transactions on Circuits and Systems II*, vol. **41** (12), pp 827-829, 1994.
- [15] G. Cauwenberghs, "Analog VLSI Long-Term Dynamic Storage," in *Proceedings of the International Symposium on Circuits and Systems*, Atlanta, GA, 1996.
- [16] A. Dembo and T. Kailath, "Model-Free Distributed Learning," *IEEE Transactions on Neural Networks*, vol. **1** (1), pp 58-70, 1990.
- [17] S. Grossberg, "A Neural Model of Attention, Reinforcement, and Discrimination Learning," *International Review of Neurobiology*, vol. **18**, pp 263-327, 1975.
- [18] A.G. Barto, R.S. Sutton, and C.W. Anderson, "Neuronlike Adaptive Elements That Can Solve Difficult Learning Control Problems," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. **13** (5), pp 834-846, 1983.
- [19] R.S. Sutton, "Learning to Predict by the Methods of Temporal Differences," *Machine Learning*, vol. **3**, pp 9-44, 1988.
- [20] C. Watkins and P. Dayan, "Q-Learning," *Machine Learning*, vol. **8**, pp 279-292, 1992.



- [21] P.J. Werbos, “A Menu of Designs for Reinforcement Learning Over Time,” in *Neural Networks for Control*, W.T. Miller, R.S. Sutton and P.J. Werbos, Eds., Cambridge, MA: MIT Press, 1990, pp 67-95.
- [22] G. Cauwenberghs, “A Fast Stochastic Error-Descent Algorithm for Supervised Learning and Optimization,” in *Advances in Neural Information Processing Systems*, San Mateo, CA: Morgan Kaufman, vol. **5**, pp 244-251, 1993.
- [23] P. Werbos, *Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences*. Ph.D. dissertation, 1974. Reprinted in P. Werbos, *The Roots of Backpropagation*. New York: Wiley, 1993.
- [24] H.J. Kushner, and D.S. Clark, *Stochastic Approximation Methods for Constrained and Unconstrained Systems*, New York, NY: Springer-Verlag, 1978.
- [25] H. Robins and S. Monro, “A Stochastic Approximation Method,” *Annals of Mathematical Statistics*, vol. **22**, pp 400-407, 1951.
- [26] J.C. Spall, “A Stochastic Approximation Technique for Generating Maximum Likelihood Parameter Estimates,” Proceedings of the 1987 American Control Conference, Minneapolis, MN, 1987.
- [27] M.A. Styblinski and T.-S. Tang, “Experiments in Nonconvex Optimization: Stochastic Approximation with Function Smoothing and Simulated Annealing,” *Neural Networks*, vol. **3** (4), pp 467-483, 1990.
- [28] M. Jabri and B. Flower, “Weight Perturbation: An Optimal Architecture and Learning Technique for Analog VLSI Feedforward and Recurrent Multilayered Networks,” *IEEE Transactions on Neural Networks*, vol. **3** (1), pp 154-157, 1992.
- [29] J. Alspector, R. Meir, B. Yuhua, and A. Jayakumar, “A Parallel Gradient Descent Method for Learning in Analog VLSI Neural Networks,” in *Advances in Neural Information Processing Systems*, San Mateo, CA: Morgan Kaufman, vol. **5**, pp 836-844, 1993.
- [30] B. Flower and M. Jabri, “Summed Weight Neuron Perturbation: An  $\mathcal{O}(n)$  Improvement over Weight Perturbation,” in *Advances in Neural Information Processing Systems*, San Mateo, CA: Morgan Kaufman, vol. **5**, pp 212-219, 1993.
- [31] G. Cauwenberghs, “A Learning Analog Neural Network Chip with Continuous-Recurrent Dynamics”, in *Advances in Neural Information Processing Systems*, San Mateo, CA: Morgan Kaufman, vol. **6**, pp 858-865, 1994.
- [32] G. Cauwenberghs, “An Analog VLSI Recurrent Neural Network Learning a Continuous-Time Trajectory,” *IEEE Transactions on Neural Networks*, vol. **7** (2), March 1996.
- [33] F. Pineda, “Mean-Field Theory for Batched-TD( $\lambda$ ),” submitted to *Neural Computation*, 1996.

- [34] S. Grossberg and D.S. Levine, "Neural Dynamics of Attentionally Modulated Pavlovian Conditioning: Blocking, Inter-Stimulus Interval, and Secondary Reinforcement," *Applied Optics*, vol. **26**, pp 5015-5030, 1987.
- [35] E. Niebur and C. Koch, "A Model for the Neuronal Implementation of Selective Visual Attention Based on Temporal Correlation Among Neurons," *Journal of Computational Neuroscience*, vol. **1**, pp 141-158, 1994.
- [36] G. Cauwenberghs, "Reinforcement Learning in a Nonlinear Noise Shaping Oversampled A/D Converter," to appear in *Proc. Int. Symp. Circuits and Systems*, Hong Kong, June 1997.
- [37] J.C. Candy and G.C. Temes, "Oversampled Methods for A/D and D/A Conversion," in *Oversampled Delta-Sigma Data Converters*, IEEE Press, pp 1-29, 1992.
- [38] E. Vittoz and J. Fellrath, "CMOS Analog Integrated Circuits Based on Weak Inversion Operation," *IEEE Journal on Solid-State Circuits*, vol. **12** (3), pp 224-231, 1977.
- [39] A.G. Andreou, K.A. Boahen, P.O. Pouliquen, A. Pavasovic, R.E. Jenkins, and K. Strohhahn, "Current-Mode Subthreshold MOS Circuits for Analog VLSI Neural Systems," *IEEE Transactions on Neural Networks*, vol. **2** (2), pp 205-213, 1991.
- [40] A.L. Hodgkin and A.F. Huxley, "Current Carried by Sodium and Potassium Ions Through the Membrane of the Giant Axon of *Loligo*," *Journal of Physiology*, vol. **116**, pp. 449-472, 1952.
- [41] G. Cauwenberghs, "A Micropower CMOS Algorithmic A/D/A Converter," *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, vol. **42** (11), pp 913-919, 1995.
- [42] D. Kirk, D. Kerns, K. Fleischer, and A. Barr, "Analog VLSI Implementation of Gradient Descent," in *Advances in Neural Information Processing Systems*, San Mateo, CA: Morgan Kaufman, vol. **5**, pp 789-796, 1993.