

# **Neurocontroller Alternatives for “Fuzzy” Ball-and-Beam Systems With Nonuniform, Nonlinear Friction**

Paul H. Eaton, Danil V. Prokhorov, and Donald C. Wunsch, II  
Applied Computational Intelligence Laboratory  
Department of Electrical Engineering, Box 43102  
Texas Tech University, Lubbock, TX 79409-3102, U.S.A.

## **ABSTRACT**

The ball-and-beam problem is a benchmark for testing control algorithms. In the World Congress On Neural Networks, 1994, Prof. Lotfi Zadeh proposed a twist to the problem, which, he suggested, would require a fuzzy logic controller. This experiment uses a beam, partially covered with a sticky substance, increasing the difficulty of predicting the ball’s motion. We complicated this problem even more by not using any information concerning the ball’s velocity. Although it is common to use the first differences of the ball’s consecutive positions as a measure of velocity and explicit input to the controller, we preferred to exploit recurrent neural networks, inputting only consecutive positions instead. We have used truncated backpropagation through time with the Node-Decoupled Extended Kalman Filter (NDEKF) algorithm to update the weights in the networks. Our best neurocontroller uses a form of approximate dynamic programming called an adaptive critic design. A hierarchy of such designs exists. Our system uses Dual Heuristic Programming (DHP), an upper-level design. To our best knowledge, our results are the first use of DHP to control a physical system. It is also the first system we know of to respond to Zadeh’s challenge. We do not claim this neural network control algorithm is the best approach to this problem, nor do we claim it is better than a fuzzy controller. It is instead a contribution to the scientific dialogue about the boundary between the two overlapping disciplines.

## 1. INTRODUCTION

At the World Congress On Neural Networks, 1994, and later, at the IFAC 13<sup>th</sup> Triennial World Congress in San Francisco, 1996, Prof. Lotfi Zadeh issued a challenge to the neural network community. This challenge was as follows [1]:

*A more complex problem involving dynamic motion planning - a problem which does not lend itself to solution by the methods of classical control - is what might be called the fuzzy ball and beam problem. The main difference between this problem and the standard ball and beam problem is that the beam is assumed to be covered with a strip of fuzzy material, e.g., a strip of thick-pile rug. The fuzziness of the beam's surface is intended to preclude the possibility of setting up the differential equations governing the ball's motion (gliding and rolling) on the beam. This rules out - or at least makes it very difficult to use classical control techniques to derive a way of stabilizing the ball within a prescribed interval. The same obstacle stands in the way of application of neural network and genetic algorithm techniques since a prerequisite to the use of these techniques is the possibility of computer simulation.*

(He provided another example fuzzy material, sticky tape, in his earlier plenary talk at World Congress on Neural Networks, 1994.) Zadeh asserted that the only way to control this system was with fuzzy logic, or at least that other approaches would be difficult. We demonstrate two neurocontrollers to solve this problem. Our approach is to train neurocontrollers on a model of the system off-line, and then deploy them with fixed weights for testing on the actual system on-line. First we verify that this approach is valid using the ordinary beam, i.e., one providing uniform friction for the ball everywhere. Then we show that our approach is applicable to the “fuzzy” beam as well.

Our experience with this problem shows that the presence of a sticky surface does increase the difficulty of control, as asserted in the challenge. A shallow beam angle can allow the ball to become stuck on an undesired location, while a steep angle will cause it to shoot past

the desired point. Setting up differential equations to describe the ball's dynamics is very difficult since these responses are nonlinear and ill-posed. (On a non-sticky beam, these problems do not exist because the ball's movement is smooth and predictable.) Yet, we believe that modeling this problem is not impossible. In fact, the sheer success of our experiments confirms feasibility of relatively simple control solutions attained with *qualitatively* accurate models.

We begin by describing the experimental apparatus, neural network architecture, and training algorithms (Sections 2 and 3). In Section 4, we discuss the creation of a model of the system for off-line simulation and training. In Section 5.1, a conventional neurocontroller is trained, using backpropagation through time, to balance the ball. Adaptive critic approaches to train neurocontrollers are discussed in Section 5.2. We also provide references in that section to papers with full implementation details of these algorithms. In Section 6, we present results from our test on ordinary (non-sticky) and "fuzzy" (sticky) beams for both controllers, shown in computer simulations and actual hardware tests.

## 2. INSTALLATION FOR EXPERIMENTS

In the usual designs for the ball-and-beam, the position sensors are mounted on the beam [7, 20, 21]. For instance, photodiodes mounted along the beam are used as position sensors in [7]. When the ball rolls, it blocks an external light incident on a photodiode thereby signaling its relative position on the beam. We cover the beam with a nontransparent sticky tape, but this obviously would interfere with photosensors as described above. Therefore we built the ball-and-beam setup with displaced fulcrum as shown in Figure 1. The beam moves like a cradle, with a pivot point high above it. A ball position sensor is moved above the beam such that the ball is always easily detectable, whether the beam is covered with a sticky tape or not.

The displaced fulcrum design just described does make the physics of the problem easier than that of the usual ball-and-beam. The version with nonlinear friction (beam partially covered

with sticky tape) thus inherits this advantage. However, our design retains the intrinsic instability of balancing the ball *in an arbitrary off-center location* on the beam as the essential property of the usual ball-and-beam problems. In any case, it is a necessary step, enabling experimentation with the “fuzzy” ball-and-beam, since the nonlinear friction is obtained by covering the surface where photodetectors would be used in other implementations.

The ball position sensor is a 512-element linear Charge-coupled Device (CCD). The light reflecting off the ball and beam is focused onto the CCD using a lens. The ball is painted black and the beam is white. The software reads in the CCD pixels serially, with sampling period of 0.055 seconds, looking for the dark spot and determining the ball’s position on the beam.

A servo motor, using Pulsed Width Modulation (PWM) encoding, moves the beam. The PWM signal is a digital signal with a fixed frequency (the same as the inverse of the sampling period) and a variable duty cycle, with changes in the duty cycle changing the servo motor position.

The ball is 2 cm in diameter. The controller can balance it at any position on the beam, which is 40 cm long. A foam bumper is at each edge of the beam. The CCD reads in the beam from pixel 189 to 412, allowing for 223 pixels to scan the 40 cm beam. This allows an 0.18 cm accuracy limit in determining the ball’s position. Actual accuracy is 0.54 cm after experiments showed an added  $\pm 1$  pixel from system noise. The servo motor can be set to any value from 0 to 255. Due to hardware limitations, only the range between 50 and 140 is used to move the beam. This range produces 90 different possible positions, which set the beam at an angle between 79 degrees and 103 degrees, with 90 degrees corresponding to the beam being level. This allows for the beam to be set to positions in 0.26 degree increments.

The “fuzzy” beam used in this paper was created by laying a sticky tape across half of the beam (sticky surface facing upward). One half of the beam (approximately -1 to 0 in normalized units, and 189 to 290 in the hardware) remained the same as the ordinary beam, while another half (0 to +1 in normalized units, and 290 to 405 in the hardware) was covered with the tape.

### 3. NEURAL NETWORKS USED AND THEIR TRAINING ALGORITHM

Our training method is BackPropagation Through Time truncated with depth  $h$  (BPTT( $h$ )) with Node-Decoupled Extended Kalman Filtering (NDEKF). BPTT( $h$ ) uses  $h$  copies of the network with the same weights but different node activations corresponding to different time steps. We typically used  $h=10$ . The Recurrent Neural Networks (RNN) we use (Figure 2) are Discrete-Time Recurrent Multi-layer Perceptrons [5].

In BPTT( $h$ ), multiple copies of the same network weights are made. This can be pictured as one larger network, with each time step representing a hidden layer passing its output to the next layer. While the copies' weights are usually the same, the node activations are different and time dependent. Thus, unfolding the network provides a natural way to preserve this time dependence. Starting at the beginning of the data, each copy calculates its outputs. At the end, the error is calculated for the last copy, and the weights are updated using the chain rule for ordered derivatives. The copies are then moved forward one time step and the process is repeated [14, 15, 16].<sup>1</sup>

EKF training is a parameter identification technique for an RNN [4], which adapts weights of the network pattern-by-pattern, accumulating training information in approximate error covariance matrices and providing individually adjusted updates for the network's weights.

All weights of the RNN are assembled in a vector  $W$  of length  $M$ . This is split into  $G$  disjoint groups.  $W_i$  is a vector of the  $i$ -th group of weights.  $W_1 \cup W_2 \cup \dots \cup W_G = W$ , and  $\text{size}(W_1) + \text{size}(W_2) + \dots + \text{size}(W_G) = M$ .

The EKF method requires that we compute derivatives of the RNN's  $N$  outputs, rather than output errors, with respect to the weights. These derivatives are obtained using BPTT( $h$ ) [16]. We store them in a set of matrices  $H_i$ , where each  $H_i$  has dimensions  $\text{size}(W_i) \times N$ .

The following equations form the basis of the EKF training procedure:

---

<sup>1</sup> Our notation BPTT( $h$ ) corresponds to BPTT( $h-1$ ) of [15] and [24].

$$A(k) = \left( \frac{I}{\eta(k)} + \sum_{i=1}^G H_i^T(k) P_i(k) H_i(k) \right)^{-1}, \quad (1)$$

$$K_i(k) = P_i(k) H_i(k) A(k), \quad (2)$$

$$W_i(k+1) = W_i(k) + K_i(k) e(k), \quad (3)$$

$$P_i(k+1) = P_i(k) - K_i(k) H_i^T(k) P_i(k) + Q_i(k), \quad (4)$$

where  $\eta(k)$  is a scalar learning rate,  $K_i(k)$  is the Kalman gain matrix for the  $i$ -th group of weights,  $e(k) = d(k) - y(k)$  is the  $N \times 1$  error vector,  $d(k)$  is a vector of the desired outputs, and  $y(k)$  is the actual output vector of the RNN;  $e^T(k)e(k)/2$  forms the squared error,  $P_i(k)$  is the  $\text{size}(W_i) \times \text{size}(W_i)$  approximate error covariance matrix which models correlation between each pair within the  $i$ -th group of weights, and  $Q_i(k)$  is a positive diagonal matrix that helps avoid numerical divergence of the procedure and prevents getting stuck in a local minimum [4].

Grouping of weights can be done in a variety of ways. We employ grouping by node, i.e. weights belonging to the same neuron are grouped together. Thus, we ignore frequently unnecessary correlation between weights belonging to different neurons. This results in a significant reduction of computational complexity and storage requirements since the size of each error covariance matrix  $P_i$  can then be made much smaller than  $M^2$ , the size in the case when  $G = 1$ .

The matrices  $P_i(0)$  are initialized as diagonal matrices with large diagonal elements (values between 100 and 10,000). User-specified values of  $\eta(k)$  are usually increased from 0.01 to 1 whereas diagonal components of  $Q_i(k)$  are decreased from 0.01 to  $10^{-6}$  as training progresses.

The set of matrices  $H_i$  in (1), (2), (4) is obtained by truncated backpropagation through time with depth 10. This means that we do not use more than 10 copies of the network to accumulate appropriate derivatives in the matrices  $H_i$ , hence the notation BPTT(10). Training usually lasts for around one hundred passes, where one pass corresponds to a complete processing of the whole training set.

## 4. SYSTEM SIMULATION

Our approach to balancing the ball on the beam belongs to the category of indirect adaptive control [2]. Problems of this category are characterized by the presence of a system model, which simulates behavior of the actual system. The model naturally has to be accurate enough in predicting the system behavior over a sufficiently long time horizon. The length of the time horizon depends on the training mode. If training is performed on-line, the model is typically used for short-term predictions of the system outputs. For example, in the case of one-step-ahead predictions the model gets outputs of the actual system from the previous time step as its inputs. The model outputs are predictions of the current system outputs. Such model-system configuration is known under different names including series-parallel model [3] and teacher forcing [15]. We model the ball-and-beam system via a recurrent neural network henceforth referred to as an ID (identification) network. After this network is trained, its weights are fixed, and it *replaces* the system in subsequent off-line training of the neurocontrollers. We emphasize that the complete replacement of the actual system is performed. Our configuration is therefore an extreme case of so-called parallel model [3]. Unlike the series-parallel model, it features an ID network inputting its own predictions for the previous time step instead of the system outputs. The parallel model is used when the ID network is to simulate long-term behavior of the system, which is precisely what we need. The parallel method clearly puts more demands on the neural network than the series-parallel method. This is justified by the convenience of off-line training, which precludes the latter.

Regardless of whether the parallel or series-parallel model is used, the main function of the model for training controllers is to provide sensitivity signals, i.e. gradients of outputs with respect to inputs. These gradients are used to adapt adjustable parameters of the controller. It is convenient to use the backpropagation algorithm to obtain the gradients, in which case one recovers backpropagation through the forward model as proposed in [8].

Our ID network uses two inputs, the ball position,  $B(t)$ , and the PWM,  $PWM(t)$ , and computes the ball position  $B(t+1)$ . It is a one-hidden-layer RNN. The hidden layer has two neurons which read in the inputs and the recurrent values from the previous time, and pass their outputs to one output neuron.

About forty data files were created for training the ID network using various techniques such as a human randomly changing the angle, a human trying to balance the ball, the computer randomly changing the beam angles, and the computer changing the beam angle in specific patterns. The best models were created using a combination of the last two. In the beginning of gathering the data file, the computer was moving the beam from the largest angle in one direction, followed by the largest angle in the other direction. The beam then oscillated back and forth, slowly decreasing the angle until the beam was horizontal. This portion of the data file was followed by the computer randomly choosing a beam angle and staying in this position for a certain amount of time. An example of a typical data set recorded in this way is shown in Figure 3.

The parallel model shown in Figure 4 feeds the output from each step as the input for the next time step. Only the first few ball positions are read from a data file of consecutive ball positions and corresponding PWMs providing the initial ball position, after which the network uses its previous estimate of the ball position as the current input. The parallel method is trained in 70 to 100 passes using BPTT(10) with NDEKF as described in Section 3. A typical comparison of the actual  $B(t)$  and the corresponding predictions of the ID network on the training data is shown in Figure 5. An example test performance of the ID network is illustrated in Figure 6. While discrepancies between the predicted and the actual  $B(t)$  appear to be large, they turn out to be small enough to enable acceptable off-line training of our neurocontrollers, to be discussed in Sections 5-6. More details on the system simulation using series-parallel and parallel models are given in [17, 18].

## 5. NEUROCONTROLLERS

Once we have the ID network, we can proceed with designing and training of the system's controller. Section 5.1 discusses the conventional neurocontroller, i.e. one which is trained using BPTT with NDEKF in the closed loop with the ID network. In Section 5.2, we discuss an adaptive critic-based neurocontroller. To distinguish between the two neurocontrollers we used the designator "conventional" for the former. We will later show that our best performance for this problem was obtained with the adaptive critic known as DHP. We provide some motivation for this design and references to complete implementation details.

### 5.1. CONVENTIONAL NEUROCONTROLLER

The controller is a one-hidden-layer RNN, with two nodes in the hidden layer and one node in the output layer. The neurocontroller uses the two inputs, the current ball position ( $B$ ), and the desired ball position ( $B_{desired}$ ), to produce an output which is the PWM signal used by the ID network.

The desired ball position is a value provided in a data file. The current ball position comes from the output of the ID network from the previous time step. The neurocontroller is shown with its connections to the ID network in Figure 7. The output of the neurocontroller ( $PWM$ ) is a sigmoid node, where  $PWM = -1$  sets the beam to the highest angle possible in one direction, and  $PWM = +1$  moves to the highest angle in the other direction. The controller is trained using BPTT(10) with an error term corresponding to the difference between the current ball position and the desired ball position at each time step. This error is backpropagated through the ID network to the neurocontroller, providing  $\partial E(t)/\partial PWM(t)$ . The weights of the controller are then updated using NDEKF and the appropriate derivatives with respect to the controller weights.

### 5.2. ADAPTIVE CRITIC-BASED NEUROCONTROLLER

A family of approaches united under the common name adaptive critic designs (ACDs) have recently emerged as a synthesis of reinforcement learning, dynamic programming, and backpropagation [10, 19, 22, 23]. This section discusses the most salient issues of ACDs with respect to this problem. A typical ACD consists of the critic network, the ID network, and the controller network, also called the action network in ACD literature. The critic is trained to estimate the cost-to-go function  $J(t)$  of the Bellman equation of classical dynamic programming [13]. Similarly to the off-line training algorithm Dyna [9], in our case the critic is added to the ID network and the conventional controller of Section 5.1 to obtain an adaptive critic-based neurocontroller. The addition of the critic turned out to be crucial in improving the balancing performance, as shown in Section 6.

The critic is a one-hidden-layer RNN with three hidden neurons and one output neuron. The critic inputs are the actual ball position  $B(t)$  and the corresponding desired ball position  $B_{desired}(t)$ . Two ACDs, Heuristic Dynamic Programming (HDP) and Dual Heuristic Programming (DHP), were implemented. HDP outputs the cost-to-go function  $J(t)$ , while DHP outputs the derivative of  $J(t)$ . The reader is referred to [11, 12, 19] for the detailed descriptions of these ACDs. As we shall see, DHP works well for this problem, whereas HDP does not.

We nevertheless begin by introducing the architecture of HDP, because DHP is best explained by its contrast with HDP, and the shortcomings of HDP for this problem serve to motivate the use of DHP. The HDP critic produces the  $J(t)$  function described above as the output, calculated using:

$$Error(t) = J(t) - \gamma J(t+1) - U(t), \quad (5)$$

where

$$U(t) = \frac{1}{2}(B(t) - B_{desired}(t))^2, \quad (6)$$

The error is calculated by finding the output error,  $U(t)$ , and the  $J(t)$  and  $J(t+1)$  terms for the current and next time step. By using this cost-to-go function, the effect of the current neurocontroller output on future positions can be seen. For instance, if the beam is set to a very

high angle, the ball could roll to the desired position faster than from a small angle, but it will be more difficult to stop the ball when it arrives. Learning to incorporate this tradeoff is why reinforcement learning allows good control decisions in spite of the inability to create an accurate model. This advantageous property of ACDs was also noted in [22].

The term  $\gamma$  in (5) is called the discount factor and is set to a value between 0 and 1. When  $\gamma$  is close to 0, the critic network will be trained with a very small predictability horizon (if  $\gamma=0$ , this horizon is just the current time step error  $U(t)$ ), while having values of  $\gamma$  close to 1 is equivalent to extending prediction horizon of the critic up to infinity. Once the critic network is trained, its weights are fixed and the critic is used to train the action network. This is done by using the derivative of the cost-to-go function as the error term

$$0 - \frac{\partial J(t)}{\partial B(t)}. \quad (7)$$

Simultaneously requiring smooth estimates and a minimal final cost is proven to minimize cost over time [6]. Backpropagation through (now fixed) weights of the critic and model networks yields an error signal for the action network. The training of the action network is performed until its error reaches a minimum, and the weights of the action network are fixed and used to retrain the critic. This process is continually repeated until the neurocontroller performance is no longer improved over the previous training pass. The training process follows guidelines of the general training procedure for ACDs proposed in [19], and the NDEKF algorithm is used for updating both critic and action weights.

For this problem, HDP, one of the simplest ACDs, produced unsatisfactory results. This was expected, because, as we shall explain below, DHP computes derivatives that are essential to training the action network that makes control decisions. To illustrate the limitations of HDP with respect to this issue, we trained the HDP critic using  $\gamma=0$  in (5), i.e.

$$Error(t) = J(t) - U(t). \quad (8)$$

The trained critic showed it could produce results with the error (8) close to zero implying that

$$J(t) \approx U(t) \quad (9)$$

However, the problem occurred when the derivative term  $\partial J(t)/\partial B(t)$  was examined. This derivative term was very inaccurate when compared to  $\partial U(t)/\partial B(t)$ . But an inaccurate  $\partial J(t)/\partial B(t)$  causes improper training of the action network due to (7). The choice of  $\gamma=0$  turns the task of approximating the cost-to-go function  $J(t)$  into the much simpler task of recreating the one-step cost function  $U(t)$ . It is therefore hard to expect a more accurate  $\partial J(t)/\partial B(t)$  for  $\gamma>0$ . This inaccuracy of  $\partial J(t)/\partial B(t)$  in the most benign case demonstrated that even if a network can produce a close estimate of a desired function, it is still possible for the network to produce a poor approximation of the derivatives of this function, unless it is trained to estimate them explicitly.

The explicit approximation of  $\partial J(t)/\partial B(t)$  is precisely the goal of the more sophisticated ACD called Dual Heuristic Programming (DHP). DHP uses a critic whose output tries to minimize

$$Error(t) = \frac{\partial J(t)}{\partial B(t)} - \gamma \frac{\partial J(t+1)}{\partial B(t)} - \frac{\partial U(t)}{\partial B(t)}. \quad (10)$$

The advantage of this ACD is that the critic produces the  $\partial J(t)/\partial B(t)$  term as its output, to be directly used to adapt the action network (see (7)). The disadvantage lies in the necessity to perform extra computations for the  $\partial J(t+1)/\partial B(t)$  term, but this is a minor technical difficulty. The DHP arrangement of all three networks for two consecutive time steps is shown in Figure 8. DHP is an algorithm of moderate complexity, and all our experience with ACDs to date recommends it as a good trade-off between complexity and power. Reference [22] contains implementation details pertaining to DHP, while references [19, 23] give an analysis of the whole family of ACDs, including an algorithm which includes HDP and DHP as special cases. While ACDs are sufficiently complex to preclude even an overview section here, these papers, especially [23], provide the complete discussion of relevant issues.

## 6. RESULTS

## 6.1. ORDINARY (NON-STICKY) BALL-AND BEAM

### 6.1.1. Conventional neurocontroller

The conventional neurocontroller was trained using the ID network simulating the ball-and-beam system with the ordinary beam, first around a single position and then around multiple positions. The controller trained relatively fast on the single position, and depending on the distance between the balancing positions, could also train quickly on multiple balancing points. The network was repeatedly trained on many different positions using random initial weights with the best neurocontroller being chosen and applied to test simulations. An example test simulation for this type of controller is shown in Figure 9. Once it was shown that the neurocontroller could balance the ball in the off-line simulations, the weights were saved and applied to the actual hardware, with the test performance as shown in Figure 10.

### 6.1.2. Critic-based neurocontroller

Neurocontroller training within the DHP framework began with random weights in both the neurocontroller and critic networks. The critic network was trained first, using  $\gamma=0$  in (10) (see Section 5.2). This training made the first critic act very similar to the explicit error function used to train the conventional neurocontroller. The critic quickly trained to almost zero error within the first 10 passes. Once the error reached a minimum, the critic weights were saved and the critic was used to train the action network. The results of this were successful, but not as good as the conventional neurocontroller. This was because  $\partial J(t)/\partial B(t)$  is an approximation of  $\partial U(t)/\partial B(t)$ , and there is a slight difference between the two. When the minimum was reached, the first neurocontroller's weights were saved, and the new critic was trained using the first neurocontroller, this time with  $\gamma=0.6$ . This process was repeated, with gradually increasing  $\gamma$  (final value 0.9), until the neurocontroller performance was no longer improving. Figure 11 shows the results of the test simulation for the critic-based neurocontroller and Figure 12 shows the hardware test results.

### 6.1.3. Comparison of the two neurocontrollers response characteristics

We observe that the conventional neurocontroller used smaller changes in the PWM signal to move the ball. This makes the controller slower when making a large change in *PWM*, causing overshooting to be observed in many hardware tests. The conventional neurocontroller appears to have a greater difficulty in correcting this than the critic-based neurocontroller. The critic-based neurocontroller used larger PWM signals, usually causing overshoot. However, it could compensate for this quickly resulting in a shorter settling time compared to that of the conventional neurocontroller.

As seen in Figures 9-12, performances of both controllers are comparable, with small discrepancies between the results of off-line simulations (Figures 9 and 11) and those of runs on the actual system (Figures 10 and 12) being observed.

Table 1 quantifies performance of both neurocontrollers in terms of settling time, rise time, and the overshoot. All numbers represent averages over 20 test trials of balancing the ball.

Table 1. Comparative results of the neurocontrollers for the ordinary ball-and-beam.

Neurocontroller	Settling time, sec	Rise time, sec	Overshoot, %
Conventional	5.5	1.4	30
Adaptive critic-based	4.5	1.1	35

Thus, we have verified that our approach of training neurocontrollers off-line on the model (ID network) of the actual system is valid when applied to balancing the ball on the ordinary beam. In the next section, we demonstrate that the approach remains applicable to balancing the ball on the “fuzzy” beam as well.

## 6.2. “FUZZY” (STICKY) BALL-AND-BEAM

### 6.2.1. Conventional neurocontroller

Similar with section 6.1.1, the conventional neurocontroller was trained using the ID network. But this time the ID network was trained in advance to simulate dynamics of the system with the “fuzzy” beam. A typical test simulation of the best neurocontroller is shown in Figure 13. A typical performance of the same controller when applied to the actual hardware is illustrated in Figure 14.

### 6.2.2. Critic-based neurocontroller

Similar with Section 6.1.2, after the critic was trained first (using  $\gamma=0$ ), critic weights were fixed and used to train the neurocontroller. On the first cycle the controller reached a reasonable minimum of the RMS error term quickly, after two or three passes through a 300-point data file. The critic was then retrained using  $\gamma=0.6$ , and used for further training of the controller. This retraining was repeated for several cycles, usually three or four, with the final  $\gamma=0.9$ . By this process, the neurocontroller eventually functioned at a performance level that could not be improved by further training. An example test simulation is shown in Figure 15. This neurocontroller was then applied to the actual system, with the typical results shown in Figure 16.

### 6.2.3. Comparison of the two neurocontrollers’ response characteristics

Comparing Figures 13-14 and 15-16, we conclude that the critic-based neurocontroller outperforms the conventional neurocontroller both in the off-line simulation and hardware tests, but in both cases there are differences between the performance in the simulations and corresponding actual behavior. These differences are more salient than those observed in the case of the ordinary ball-and-beam (see Section 6.1.3). In particular, the ball is seen to be stuck in the time interval from 77 to 172 in Figure 14 whereas it should have oscillated according to Figure 13. Alternatively, the ball keeps oscillating in the time interval from 229 to the end in

Figure 14 whereas it should have settled down according to Figure 13. Such discrepancies are not very surprising. After all, the sticky tape is intentionally used to complicate accurate long-term predictability of the ball movement.

Table 2 quantifies performance of both neurocontrollers in terms of settling time, rise time, and the overshoot. All numbers represent averages over 20 test trials of balancing the ball. Unlike the results in Table 1, the settling time of the conventional controller is not reported since in most of the trials the ball either got stuck far from the desired balancing point or never stopped oscillating. We note that the steady-state error of the conventional controller, however, is not particularly large. (The critic-based controller has some steady-state error too.) We therefore do not claim that a conventional neurocontrol solution better than ours can not be designed -- in fact, we suspect that it can be done.

Table 2. Comparative results of the neurocontrollers for the “fuzzy” ball-and-beam

Neurocontroller	Settling time, sec	Rise time, sec	Overshoot, %
Conventional	Did not settle	1.3	37
Adaptive critic-based	6.5	1.5	31

### 6.3. BALANCING IN A GIVEN TIME INTERVAL

It is interesting to change the problem so as to permit the user to specify *when* it is desirable to have the ball balanced at a particular location on the beam. To implement this, a simple modification of the proposed neurocontrollers turns out to be sufficient.

Let us consider a structure described by the following equations

$$T(t) = T(t-1) + u(t)\Delta, \quad (11)$$

$$I(t) = \begin{cases} 0, & \text{if } T(t) < \Theta \\ 1, & \text{if } T(t) \geq \Theta \end{cases}, \quad (12)$$

where

$$\Theta = t_b - t_s. \quad (13)$$

Here  $u(t)$  is the unit step function,  $\Delta$  is an input weight equal to the sampling interval (0.055 sec, as in Section 2),  $t_b$  is a user-specified time (in seconds) marking the end of the balancing interval, and  $t_s$  is the average settling time (see Tables 1-2). In the equation (11) one can recognize a linear recurrent node, with its only external connection  $u(t)$ . In addition, the equations (12) and (13) describe a threshold node. Both of these nodes are connected in series, with the recurrent node feeding the threshold node. At each time step  $t$  all states and inputs of both the ID network and the controller are multiplied by the output  $I(t)$  of the threshold node. As a result, no balancing happens as long as  $I(t)$  stays at zero. The duration of the interval when  $I(t)$  remains zero (in time steps) is determined by the user specifying  $t_b$  ( $\Delta$  and  $t_s$  are the predefined system parameters). The balancing interval begins at  $t=0$  with initialization of the state  $T$  to zero. Due to using the average settling time  $t_s$ , the balancing may only approximately be completed by the end of the time interval from 0 to  $t_b$ . In addition, if  $t_b < t_s$  the balancing begins instantly (i.e., in the very first time step) but, in this case it would be impossible to meet such an invalid user specification.

The equations (11)-(13) are interpreted as a special neural network. This network is added to either of the neurocontrollers discussed in Sections 6.1 and 6.2 to make them immediately applicable for balancing the ball within the time specification.

## 7. CONCLUSION

In this paper we have demonstrated that the control of the “fuzzy” ball-and-beam system, as described here, can be done without fuzzy logic, using only neural networks. It is, to our knowledge, the first reported experimental study to respond to Prof. Zadeh's challenge. Furthermore, to the best of our knowledge, this is the first accomplished implementation of DHP to control a physical system. These results showcase the capabilities of DHP in particular and neural network controllers in general, although we do not claim that the results presented here

represent the best possible approach to this particular problem. We do believe that reinforcement learning approaches offer a valuable alternative when simulation of the system to be controlled can yield only qualitatively accurate predictions, which was the real crux of this challenge.

#### ACKNOWLEDGEMENTS

The authors gratefully acknowledge support from National Science Foundation (Neuroengineering Program Grant #ESC-9413120 and Research Equipment Grant Program), Ford Research Laboratories, and the Texas Tech Center for Applied Research. We wish to note that our response to this challenge is intended to be understood in the same cordial spirit in which the challenge was given, and appreciate the thoughts of Prof. Lotfi Zadeh in creating this worthy challenge. We also appreciate support of Colby Hoffmann and John Casall in conducting experiments, and useful suggestions of the anonymous reviewers.

## REFERENCES

- [1] Zadeh, L., "Fuzzy Control: Issues, Contentions, And Perspectives," In *Proc. of IFAC 13<sup>th</sup> Triennial World Congress*, San Fransisco, USA, 1996. pp. 35-38.
- [2] K. S. Narendra, "Adaptive control using neural networks," in W.T. Miller, R. Sutton, and P. Werbos, Eds., *Neural networks for Control*. Cambridge, MA: MIT Press, 1990. pp. 115-143.
- [3] Narendra, K. S., & Parthasarathy, P. (1990). "Identification and Control of Dynamical Systems Using Neural Networks." *IEEE Trans. Neural Networks*, **1**(1), 4-27.
- [4] Puskorius, G., and Feldkamp, L. "Neurocontrol of Nonlinear Dynamical Systems with Kalman Filter-Trained Recurrent Networks," *IEEE Trans. Neural Networks*, March 1994.
- [5] Hush, D. R., and Horne B. "Progress in Supervised Neural Networks - What's New since Lippman?," *IEEE Signal Processing Mag.*, Jan. 1993. pp. 8-34.
- [6] Sutton, R. S., "Learning to Predict by the Methods of Temporal Differences," *Machine Learning*, 3:9-44, Kluwer Academic Publishers, Boston, 1988. pp. 10-43
- [7] Godfrey N., H. Li, Y. Ji, W. Marcy. "Real Time Fuzzy Logic Controller For Balancing A Beam-And-Ball System," in H. Li and M. Gupta, Eds., *Fuzzy Logic and Intelligent Systems*. Kluwer Academic Publishers, Norwell, Ma, 1995. pp. 157-148.
- [8] Jordan, M. I. (1988). Supervised Learning and Systems with Excess Degrees of Freedom. COINS Technical Report 88-27. Department of Computer Science and Information, University of Massachusetts, Amherst.
- [9] Sutton, R. S. "First Results with Dyna, an Integrated Architecture for Learning, Planning and Reacting" in W.T. Miller, R. Sutton, and P. Werbos, Eds., *Neural networks for Control*. Cambridge, MA: MIT Press, 1990. pp. 179-189.
- [10] Werbos, P. J. "A Menu of Designs for Reinforcement Learning Over Time," in W.T. Miller, R. Sutton, and P. Werbos, Eds., *Neural networks for Control*. Cambridge, MA: MIT Press, 1990. pp. 67-95.
- [11] Werbos, P. J. (1992). Neurocontrol and Supervised Learning: An Overview and Evaluation. In D. A. White & D. A. Sofge (Eds.), *Handbook of Intelligent Control: Neural, Fuzzy and Adaptive Approaches* (pp. 65-89). New York, NY: Van Nostrand Reinhold.
- [12] Werbos, P. J. (1992). Approximate Dynamic Programming for Real-Time Control and Neural Modeling. In D. A. White & D. A. Sofge (Eds.), *Handbook of Intelligent Control: Neural, Fuzzy and Adaptive Approaches* (pp. 493-525). New York, NY: Van Nostrand Reinhold.
- [13] Bellman, R. E. *Dynamic Programming*, Princeton University Press, Princeton, NJ, 1957. pp. 1-54.
- [14] Werbos, P. J. "Backpropagation Through Time: What it Does and How to Do It," *Proc. of the IEEE*, vol. 78, no. 10, 1990, pp. 1550-1560.
- [15] Williams, R., and D. Zipser, "Gradient-Based Learning Algorithms for Recurrent Networks and Their Computational Complexity," Ch. 13 in *Backpropagation: Theory, Architecture, and Applications* (Chauvin and Rumelhart, Eds.), LEA, 1995. pp. 311-340.

- [16] Puskorius, G., & Feldkamp, L. (1994). Truncated Backpropagation Through Time and Kalman Filter Training for Neurocontrol. In *Proc. of the IEEE International Conference on Neural Networks (ICNN)* (pp. 2488-2493). Orlando, FL.
- [17] Paul Eaton, *Control of a Ball-and-Beam System with Neural Networks*, M.S. thesis, Department of Electrical Engineering, Texas Tech University, Lubbock, TX, 1996.
- [18] Eaton, P., Prokhorov, D., & Wunsch, D. (1996). "Neurocontrollers for Ball-and-Beam Systems," in C. Dagli et. al. (Eds.), *Intelligent Engineering Systems Through Artificial Neural Networks (Proc. Conf. Artificial Neural Networks in Engineering)*, V. 6 (pp. 551-557). New York: ASME Press, 1996.
- [19] Prokhorov, D., & Wunsch, D. (1997). "Adaptive Critic Designs." *IEEE Trans. Neural Networks*, **8**(5), 997-1007.
- [20] Jiang, Y., McCorkell, C., Zmod, R., "Application of Neural Networks for Real Time Control of a Ball-Beam System," *proceedings of the IEEE International Conference on Neural Networks, Western Australia*. December 1995. Vol. 5. pp. 2397-2402.
- [21] Ng, K., Trivide, M., "Neural Integrated Fuzzy Controller and Real-Time Implementation of A Ball Balancing Beam," *Proceedings of the 1996 IEEE International Conference on Robotics and Automation*. Minneapolis, Minnesota, April 1996. Vol. 2. pp. 1590-1595.
- [22] Prokhorov Danil, Santiago Roberto and Wunsch, Donald C., "Adaptive Critic Designs: A Case Study for Neurocontrol," *Neural Networks*, Vol. 8, No. 9, pp. 1367-1372, December 1995.
- [23] Danil V. Prokhorov, *Adaptive Critic Designs and Their Applications*, Ph.D. dissertation, Department of Electrical Engineering, Texas Tech University, Lubbock, 1997; also available by request from the author ([dprokhor@ford.com](mailto:dprokhor@ford.com)).
- [24] Feldkamp, L., & Puskorius, G., "A Signal Processing Framework Based on Dynamic Neural Networks with Application to Problems in Adaptation, Filtering, and Classification," *Proc. of the IEEE*, vol. 86, no. 11, 1998, pp. 2259-2277.

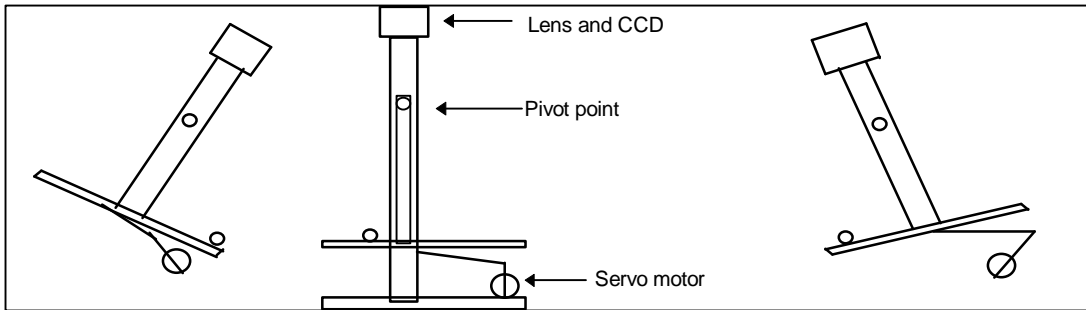


Figure 1 - Hardware setup. The beam moves like a cradle, with a pivot point high above the beam.

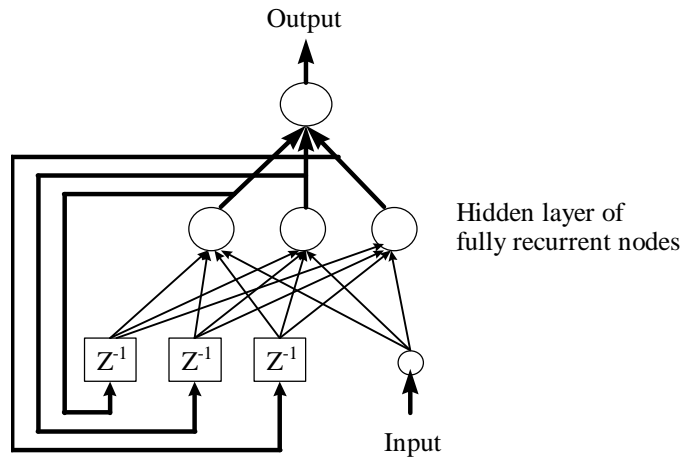
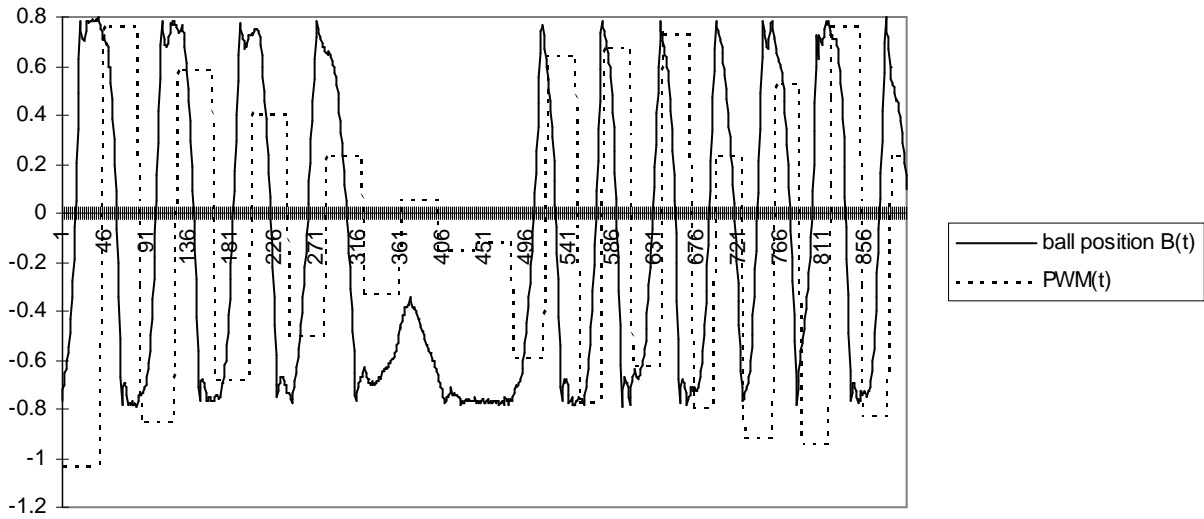
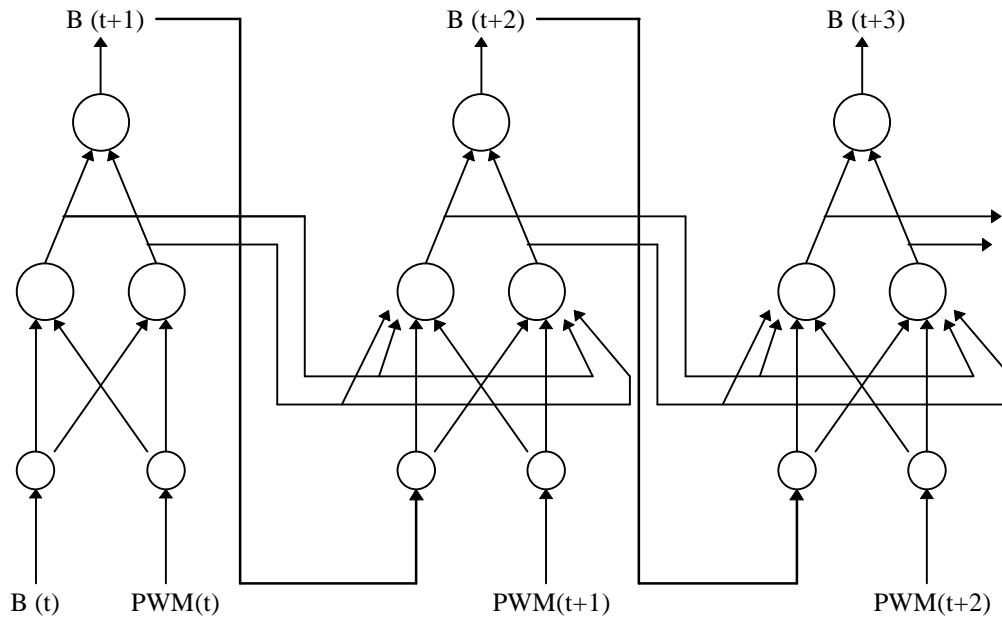


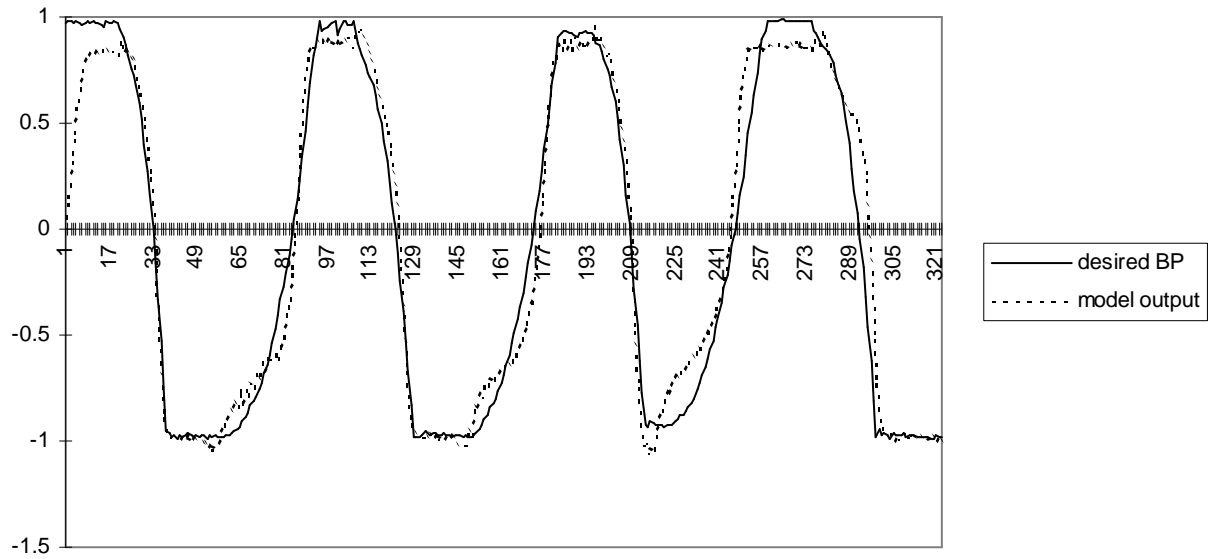
Figure 2 - Recurrent network architecture. The output of each neuron in the hidden layer from the previous time step is used as input for the current step.



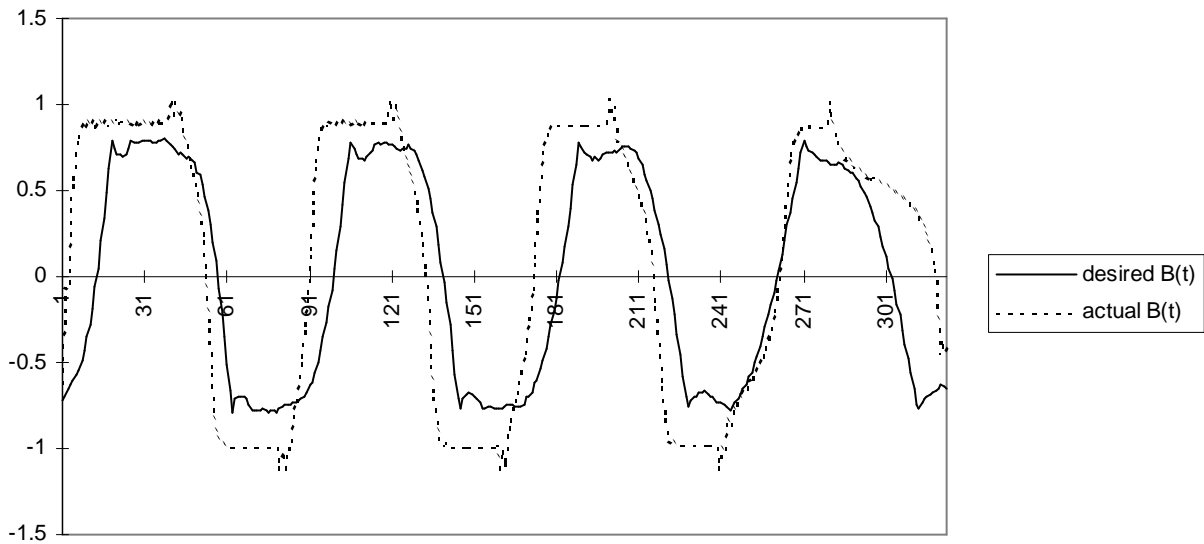
**Figure 3 - Example training data set. The first 350 steps are a predetermined pattern executed moving the beam through its entire range of motion. The next 650 steps are randomly generated PWM values and the corresponding ball positions.**



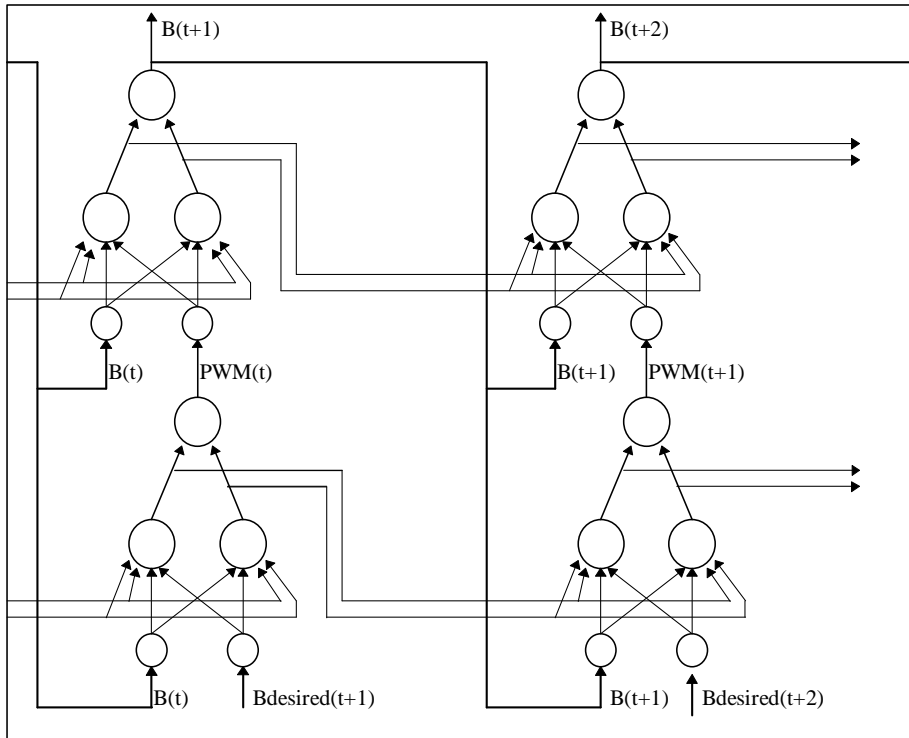
**Figure 4 - Architecture for the parallel model. The ID network uses its previous prediction of the ball position and the PWM signal from a data file as its inputs. This creates a reliance on the accuracy of previous predictions in the current position calculation.**



**Figure 5 – Training results of the ID network trained as the parallel model.**



**Figure 6 - Test results for the ID network trained as the parallel model. This typical result demonstrates that the ID network indeed captures essential long-term dynamics of the actual ball-and-beam system.**



**Figure 7 - Controller network connected to the ID network. Both networks contain two hidden neurons and one output neuron. The controller uses the current  $B(t)$  from the previous time step and the desired position  $B_{desired}(t)$  from a data file, calculating the needed  $PWM(t)$  value. The ID network receives the current  $B(t)$  and the controllers output  $PWM(t)$  value to determine the next ball position  $B(t+1)$ .**

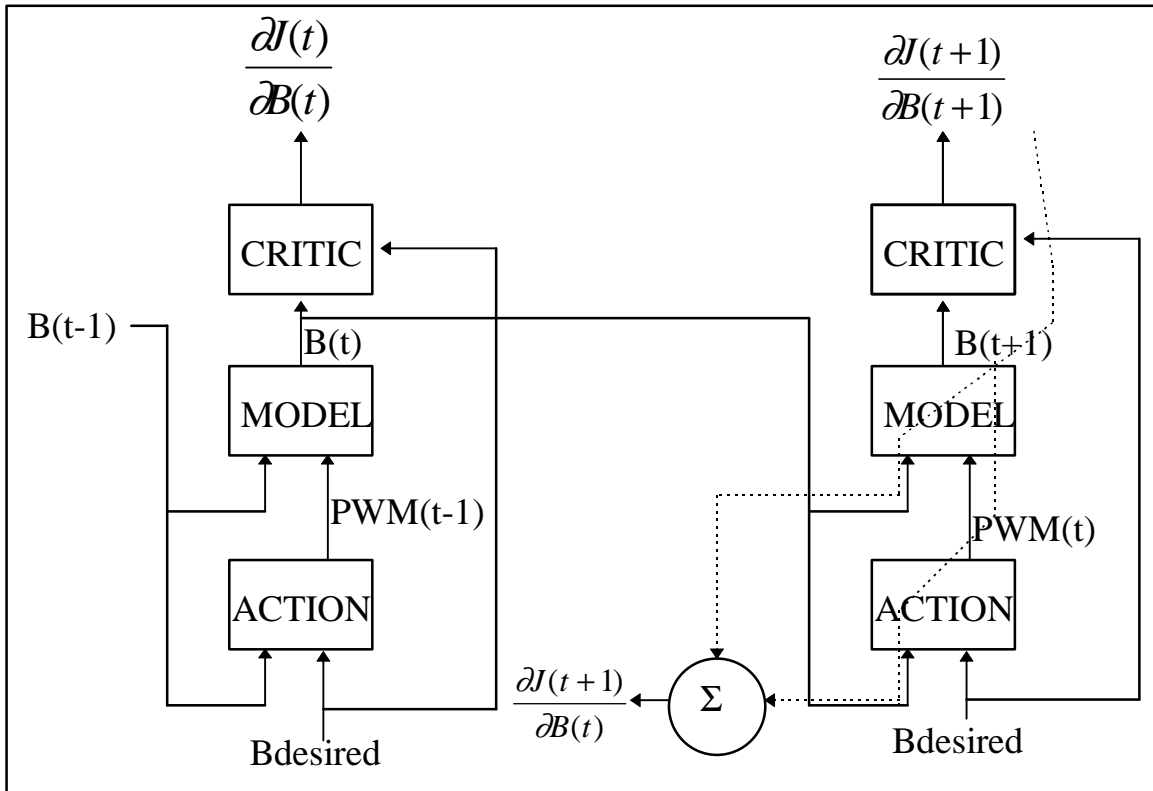
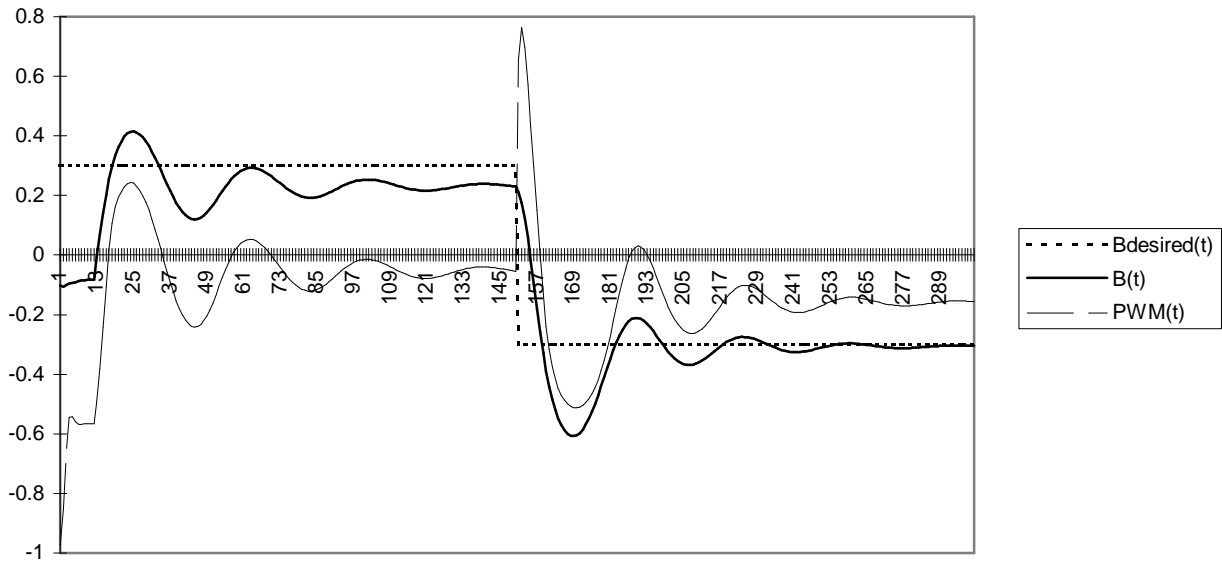
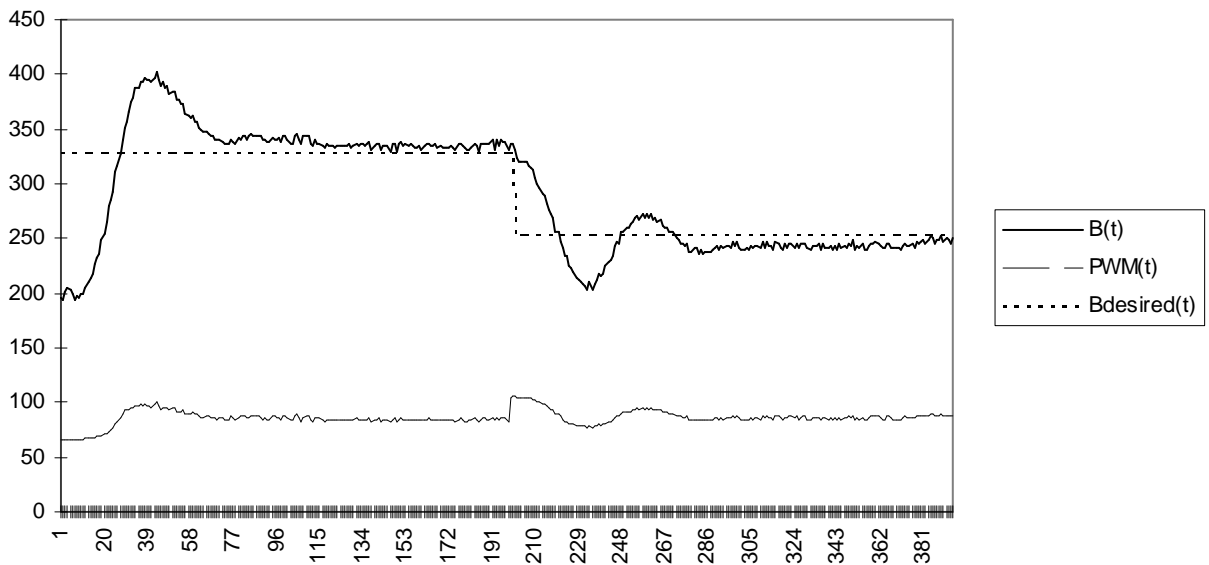


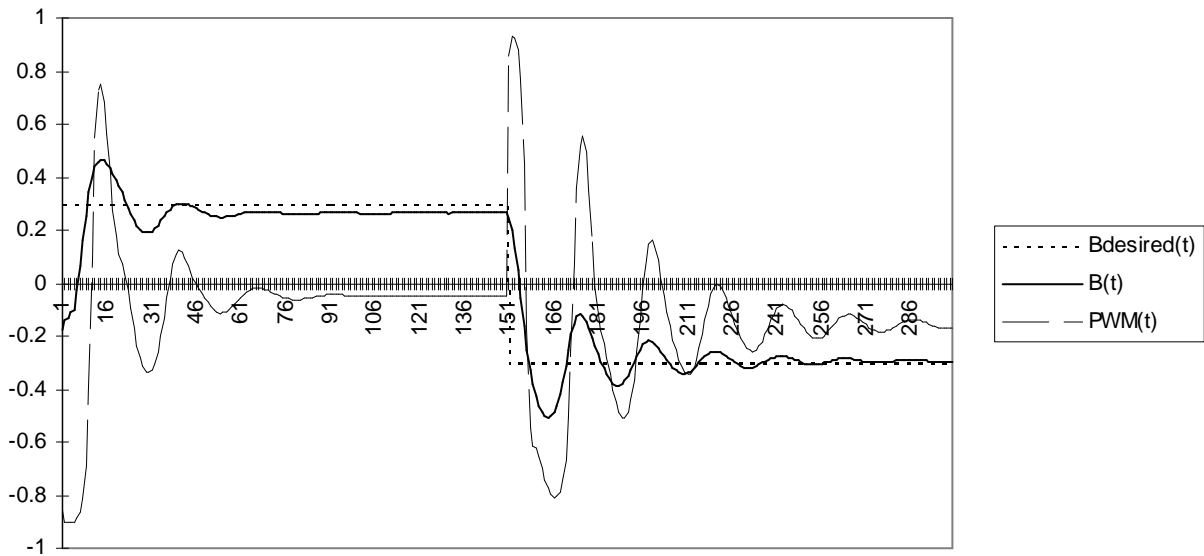
Figure 8 - Diagram of the connections between networks for DHP. This method uses a critic trained to produce a derivative of the cost-to-go function. The derivative of the cost-to-go function is backpropagated through the networks as shown by the dashed lines. (The left branch captures the dependency upon  $B(t)$ , and is thus the source of error signal in the previous time step. The right branch captures the dependency upon the control signal  $PWM(t)$ . This enables training of the Action net as shown.) The model (ID) network and the action (controller) network are the same as those in Figure 7.



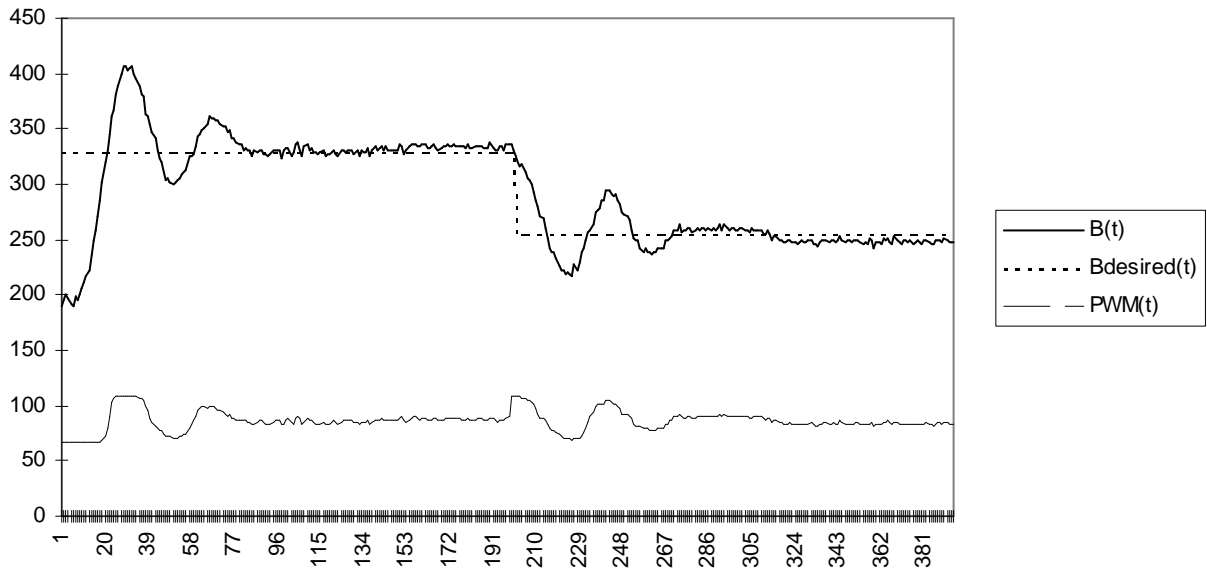
**Figure 9 - Computer simulation of conventional controller test performance on ordinary beam. The results show the conventional neurocontroller approach is capable of moving the ball to different desired positions in a relatively short period of time.**



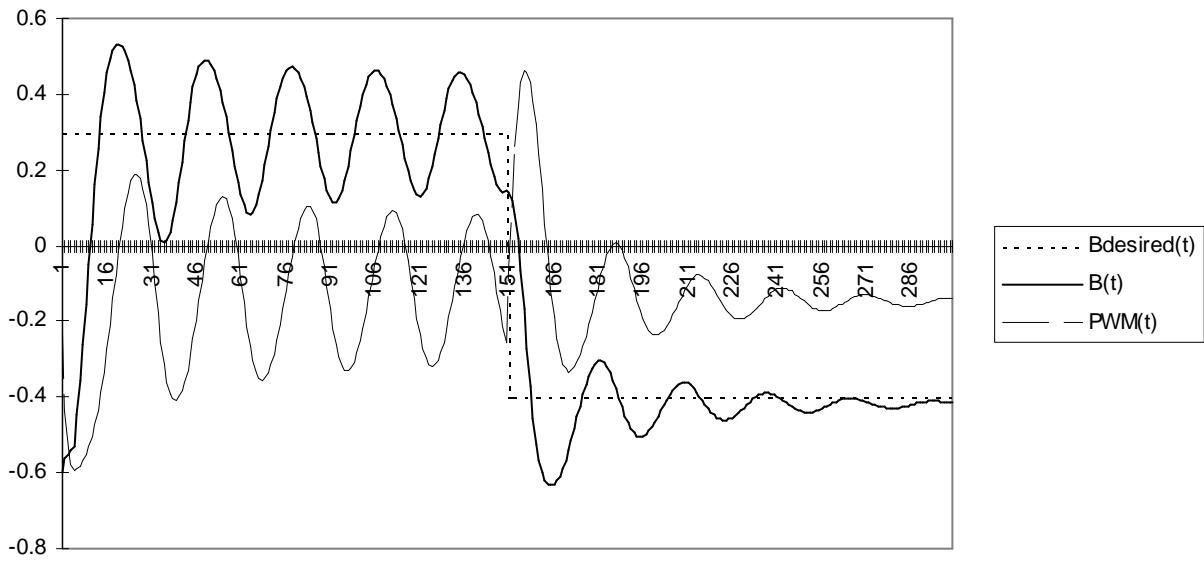
**Figure 10 - Hardware results of conventional controller test performance on ordinary beam. The results show the conventional neurocontroller is able to move and balance the ball around different positions on the beam. These are the same desired balancing points (dotted) as in Figure 9 but the overall duration is about 100 time steps longer. This difference between the off-line simulations and actual hardware runs is present in all the figures below.**



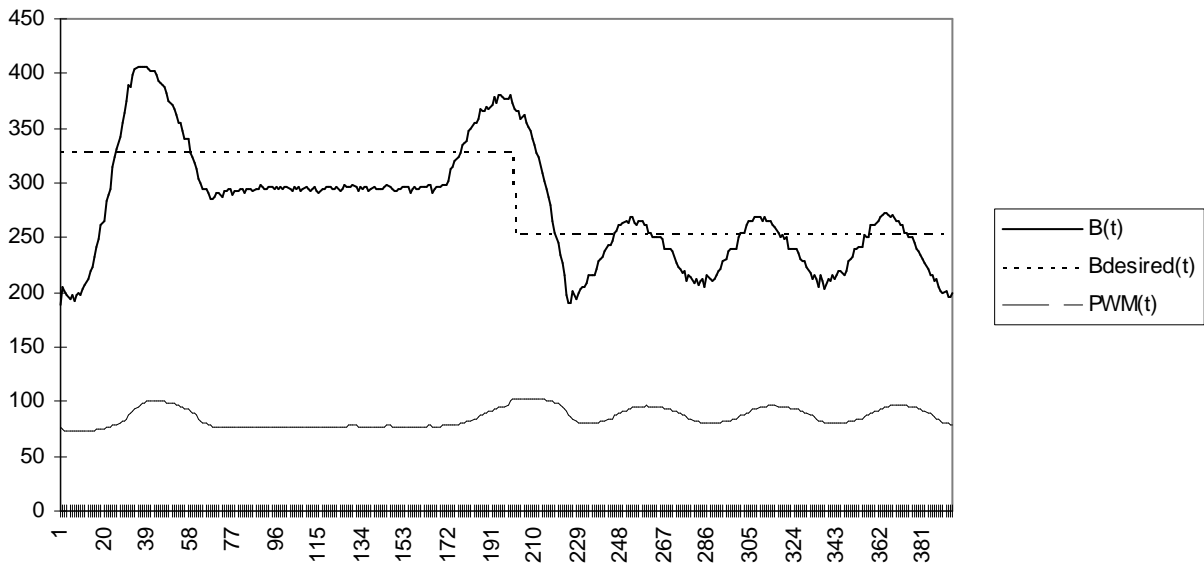
**Figure 11 - Computer simulation of critic-based controller test performance on ordinary beam. The critic based neurocontroller generates larger PWM(t) values compared to the conventional approach when moving the ball. These larger values cause “ringing” effects on the motion of the ball, seen in the second balancing position.**



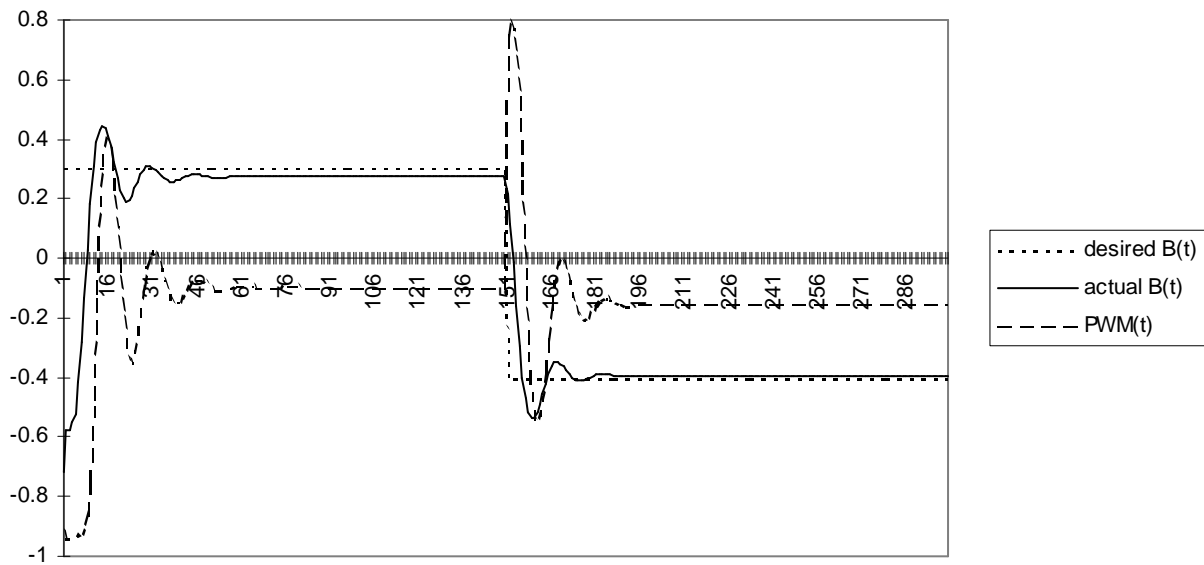
**Figure 12 - Hardware results of critic-based controller test performance on ordinary beam. The hardware results demonstrate the critic based training approach is a viable solution to training neurocontrollers. The larger PWM(t) values generated “ringing” on hardware tests, but reach the final positions in times similar to those of the conventional neurocontroller.**



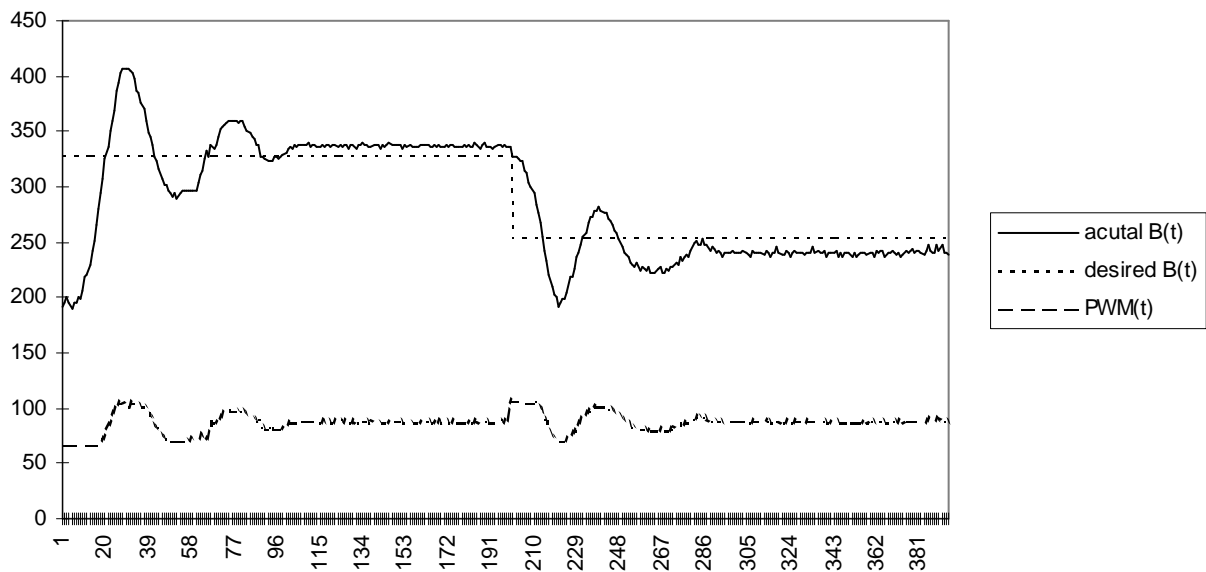
**Figure 13 - Computer simulation of conventional controller test performance on “fuzzy” beam. The conventional neurocontroller training approach was attempted multiple times without producing an adequate controller capable of balancing the ball on the “fuzzy” beam.**



**Figure 14 - Hardware results of conventional controller test performance on “fuzzy” beam. A typical conventional neurocontroller trained on the “fuzzy” model was applied to the hardware, unable to balance the ball.**



**Figure 15 - Computer simulation of critic-based controller test performance on “fuzzy” beam. After several training cycles(training the critic, training the controller, retraining the critic etc), the controller was able to balance the ball on the “fuzzy” surface**



**Figure 16 - Hardware results of critic-based controller test performance on “fuzzy” beam. The critic-based controller was able to balance the ball on the “fuzzy” surface.**