

Reward modulated Hebbian plasticity as leverage for partially embodied control in compliant robotics

Ken Caluwaerts, Jeroen Burms and Joni Dambre

Journal Name:	Frontiers in Neurorobotics
ISSN:	1662-5218
Article type:	Original Research Article
First received on:	26 Jan 2015
Frontiers website link:	www.frontiersin.org



1

Reward modulated Hebbian plasticity as leverage for partially embodied control in compliant robotics.

K. Caluwaerts^{1,2}, J. Burms¹ and J. Dambre¹

¹ Computing Systems Lab (Reservoir Team), Electronics and Information Systems Department (ELIS), Ghent University, Ghent, Belgium

² Intelligent Robotics Group, NASA Ames Research Center/Oak Ridge Associated Universities, Moffett Field CA, United States

Correspondence*:

Joni Dambre - Joni.Dambre@UGent.be

Neural plasticity for rich and uncertain robotic information streams

2 ABSTRACT

We present a dual contribution to the Hebbian learning field. First, we demonstrate that a simple Hebbian-like learning rule can be used to optimize the connection weights of recurrent neural networks. We then apply this learning rule to compliant Tensegrity robots by optimizing static feedback controllers that directly exploit the dynamics of the robot body. This leads to partially embodied controllers, i.e., hybrid controllers that naturally integrate the computations that are performed by the robot body into a neural network architecture.

The goal of this work is also twofold. First, it solves a crucial missing aspect of our recent work in which we showed that static feedback controllers can be used to induce robust locomotion in compliant robots. Previously, the problem of finding the desired motor signals was solved by using Evolutionary Algorithms. This study directly optimizes the feedback controller, bypassing this problem.

Secondly, this study strengthens our belief that compliant robots should or can be seen as computational units, instead of dumb hardware that needs a complex controller. This link between compliant robotics and neural networks is also the main reason for our search for simple universal learning rules for both neural networks and robotics.

Keywords: distal reward, tensegrity, compliant robotics, Hebbian plasticity, recurrent neural networks, morphological computation

1 INTRODUCTION

Hebbian theory has been around for over half a century **Hebb** (1949), but it still sparks the interest of today's researchers. Small changes to the basic correlation learning rule result in various well known algorithms, such as principal **Oja** (1982); **Sanger** (1989) or independent component **Clopath et al.** (2008); **Hyvriinen and Oja** (2000) extractor networks. The basic rule is biologically plausible as are some of its variations **Mazzoni et al.** (1991); **Loewenstein and Seung** (2006). Whereas all these approaches belong to the general category of unsupervised learning approaches, *reward modulated Hebbian* (RMH) learning is similar to reinforcement learning in that it can be used to tune a neural system to solve a specific task without the need for the desired output signals at the neural level **Hoerzer et al.** (2012);

27 **Fiete and Seung** (2006); **Legenstein et al.** (2010); **Soltoggio and Steil** (2013). When using RMH
28 learning in a robotics context, a reward can, e.g., be computed by comparing the sensory inputs with
29 the desired observations. The use of RMH learning for optimizing robot motor control has several
30 additional advantages. Firstly, the basic learning rule is simple. There is no need for complex mathematical
31 operations and it can therefore be efficiently implemented on various platforms in hardware and software.
32 Secondly, it allows for a distributed implementation: a central unit can be responsible for a global reward,
33 which can then be broadcast to the learning units of local controllers. Finally, RMH learning is an online
34 learning approach. If the reward mechanism remains active, the controller can adapt to changes in the
35 robot morphology or dynamics, e.g., due to wear or damage.

36 Class one tensegrity structures **Skelton and de Oliveira** (2009); **Caluwaerts et al.** (2014) consist of
37 compression members held together by tension members in such a way that compression members are
38 never directly connected. This results in flexible pin-jointed structures that can make efficient use of
39 materials. Tensegrities have been researched from various perspectives, from architecture and art **Snelson**
40 (1965) over mathematics **Connelly and Back** (1998) to biology **Ingber** (1997).

41 In previous work (**Caluwaerts et al.** (2012)), we demonstrated that the motor control of a tensegrity
42 robot can be drastically simplified by using its body as a computational resource. This approach
43 originated from the concepts of *physical reservoir computing* **Verstraeten et al.** (2007) and *morphological
44 computation* **Pfeifer and Bongard** (2007), both of which treat the use of physical systems or *bodies* as
45 a computational resource in so-called embodied computation. In **Caluwaerts et al.** (2012), we mainly
46 focused on approximating motor signals through a single layer linear neural network acting as a feedback
47 controller. The flexibility and lack of joints of our tensegrity robot allowed for simple learning rules, as
48 the risk of failure due to mechanical stress or hard constraints is minimal. As a consequence, the feedback
49 weights were learned by applying online supervised learning rules to approximate the target motor signals,
50 among which a supervised version of RMH learning. However, the target motor signals themselves were
51 generated using genetic algorithms.

52 We now extend this work by incorporating a similar feedback controller into a two-level control
53 hierarchy for end-effector control in a highly compliant class one tensegrity robot. The primary controller,
54 a simple feed forward kinematic controller, generates control signals derived from a very rough static
55 inverse model of the relationship between end effector positions and actuator signals. The secondary
56 embodied controller, consisting only of the robot body and 'neural' linear feedback weights, handles the
57 dynamics, i.e., it tunes the primary control signals to result in smooth and stable trajectories. In this task,
58 only the desired end-effector trajectories are known, not the control signals required to generate them.
59 We use RMH learning to train the feedback weights in the secondary controller, demonstrating that a
60 biologically plausible learning rule can indeed be used to control a compliant system. For this purpose, we
61 first use traditional neural network models to show how RMH learning can be extended to neural systems
62 with partial embodiment. This leads to a version of the learning rule tailored to robotics applications in
63 which the number of observable variables is limited compared to neural networks.

2 METHODS

64 In this section, we introduce the basic learning rule used throughout the paper, we derive a stabilized
65 version of this rule, and we discuss additional changes to the rule to make it more suitable for the targeted
66 application in partially embodied control of a tensegrity robot.

67 Throughout this work, we will employ the term *observations* instead of state or network activity, because
68 the reward can be based on a general performance criterion of the neural network state or the robot
69 behavior.

2.1 HEBBIAN LEARNING IN ANALOG RECURRENT NEURAL NETWORKS

70 Hebbian plasticity is a biologically plausible learning methodology for neural networks. A learning rule is
 71 called Hebbian if it modifies the weights between a set of presynaptic neurons \mathbf{x} and postsynaptic neurons
 72 \mathbf{y} as a function of their joint activity. Although, Hebb did not provide a precise mathematical formulation
 73 of his postulate, a relatively general form can be written as **Hebb** (1949):

$$\Delta \mathbf{W}_{\text{Hebb}} = f(\mathbf{X}, \mathbf{Y}). \quad (1)$$

74 Note that we have used capital \mathbf{X} and \mathbf{Y} ¹ to indicate that the weight updates in the learning rule can
 75 depend on multiple time steps, i.e., the history of the pre- and postsynaptic neuron activations. As we are
 76 using fully connected recurrent neural networks, the sets of presynaptic and postsynaptic neurons are the
 77 same, shifted by a single time step, i.e. $\mathbf{y}[k] = \mathbf{x}[k + 1]$.

78 To apply Hebbian theory in a reinforcement learning setting, we have to introduce the notion of a reward
 79 r into the learning rule. Indeed, reinforcement learning aims at making behavior that optimizes the reward
 80 more likely to happen. However, learning new behaviors necessitates another tool, namely exploration.

81 We use noise injected at the postsynaptic neurons for exploration. Noise causes variations $\delta \mathbf{y}$ of the
 82 postsynaptic neurons, which are essentially drifts from the expected value $\bar{\mathbf{y}}$. We therefore introduce the
 83 exploration noise variable $\mathbf{z} = \delta \mathbf{y} = (\mathbf{y} - \bar{\mathbf{y}})$.

84 If the exploratory noise causes an improvement in behavior, this will result in a higher reward (and vice
 85 versa). A basic learning rule based on this idea is:

$$\Delta \mathbf{W} = r \mathbf{z} \mathbf{x}^T. \quad (2)$$

86 This rule, however, suffers from a number of basic flaws to be able to solve distal reward problems.
 87 First, we need to add some efficient notion of memory of the relationship between exploration noise and
 88 the presynaptic neurons. This can be achieved by estimating the covariance between the exploration and
 89 the states of the presynaptic neurons. Secondly, we note that in its current form, any significant bias of the
 90 reward r will cause unfavorable results. The solution to this is to predict the reward and subtract this from
 91 the obtained reward, resulting in a learning rule of the form:

$$\Delta \mathbf{W} = \alpha(r - \bar{r}) \mathbf{Z}^T \mathbf{X}, \quad (3)$$

92 in which \bar{r} is the predicted reward and where we have added a learning rate parameter α .

93 The predicted reward is sometimes ambiguously referred to as the (short term) average reward.
 94 More precisely, it is the average (or expected) reward when noise is present in the system. As we
 95 will demonstrate, the average reward is typically highly dependent on the noise level of the system.
 96 The learning rule therefore optimizes the expected reward while noise is present in the system (i.e.
 97 $\max \mathbb{E}[r|\mathbf{z}]$), under the assumption that this also optimizes the performance when the exploration noise is
 98 removed (i.e. $\max \mathbb{E}[r|\mathbf{z} = 0]$).

99 In this work, we apply the learning rule on a trial-to-trial basis, by evaluating the sample covariance
 100 $\mathbf{X}^T \mathbf{Z}$ of a number of samples and then applying the learning rule once. It can be seen that - in case the
 101 reward signal can be computed at every time step - one could apply the algorithm at every time step,
 102 similar to the Rare Correlation learning rule of **Soltoggio and Steil** (2013).

2.2 A STABILIZED REWARD MODULATED HEBBIAN RULE

103 Oja's rule **Oja** (1982) and its extension, the Generalized Hebbian Algorithm or Sanger's rule **Sanger**
 104 (1989), provide a single layer neural network implementation to compute principal components. Contrary

¹ We use the notation x to denote a scalar, \mathbf{x} for a vector, and \mathbf{X} for a matrix. In general \mathbf{x}_i is the i th row of \mathbf{X} and x_i is the i th element of \mathbf{x} .

105 to pure Hebbian plasticity, the learning rules are stable, because they force the norm of the weight vectors
 106 to unity. Unlike in the unsupervised learning case, reward modulated rules tend to be stable in practice (i.e.
 107 the trained weights remain bounded). However, it can still be useful to control the norm of the weights
 108 as this can have practical implications. For example, in a robotics application, this would allow to limit
 109 the required feedback gain and thus the required motor power. From a theoretical point of view it is
 110 also instructive to see how the learning rules from the previous sections resemble the now classic rule
 111 discovered by Sanger over 20 years ago. In this section, we provide a similar derivation for the RMH
 112 learning rule we studied in this work.

113 To simplify the notation, we start by defining a number of variables:

$$\mathbf{E} = \mathbf{Z}^T \mathbf{X} \quad (4)$$

$$r' = r - \bar{r}. \quad (5)$$

114 The basic learning rule we studied can now be written in element-wise form as:

$$w_{ij}[n+1] = w_{ij}[n] + \alpha r' e_{ij} \quad (6)$$

115 Oja's rule is a first order approximation to the normalization of the weights at every update step:

$$w_{ij}[n+1] = b_i \frac{w_{ij}[n] + \alpha r' e_{ij}}{\|\mathbf{w}_i[n] + \alpha r' \mathbf{e}_i\|}, \quad (7)$$

116 where \mathbf{b} contains the desired L_2 norms of the weight vectors.

117 We now consider the linearization of this rule for small learning rates α . To further simplify the notation,
 118 we drop the time index and consider a single output dimension:

$$w_j \approx b \frac{w_j + \alpha r e_j}{\|\mathbf{w} + \alpha r \mathbf{e}\|} \Bigg|_{\alpha=0} + \alpha b \left(\frac{\partial}{\partial \alpha} \frac{w_j + \alpha r e_j}{\|\mathbf{w} + \alpha r \mathbf{e}\|} \right) \Bigg|_{\alpha=0} \quad (8)$$

119 A straightforward calculation shows that the part inside the parentheses of the second term can be written
 120 as:

$$\frac{\partial}{\partial \alpha} \frac{w_j + \alpha r e_j}{\|\mathbf{w} + \alpha r \mathbf{e}\|} \Bigg|_{\alpha=0, \|\mathbf{w}\|=b} = \frac{1}{b} r e_j - \left(\frac{r \mathbf{w} \cdot \mathbf{e}}{b^3} \right) w_j, \quad (9)$$

121 assuming $\|\mathbf{w}\| = b$ and $\alpha = 0$.

122 The complete learning rule can therefore be written as:

$$w_j = b \frac{w_j}{b} + \alpha b \left(\frac{1}{b} r e_j - \left(\frac{r \mathbf{w} \cdot \mathbf{e}}{b^3} \right) w_j \right) \quad (10)$$

$$= w_j + \alpha r \left(e_j - \frac{w_j}{b^2} \mathbf{w} \cdot \mathbf{e} \right). \quad (11)$$

123 The matrix form of the rule for multiple outputs is given by:

$$\Delta \mathbf{W} = \alpha r' (\mathbf{E} - \text{diag}_v(\mathbf{b})^{-2} \text{diag}_v(\text{diag}_m(\mathbf{E} \mathbf{W}^T)) \mathbf{W}), \quad (12)$$

124 where the diag_v operator transforms a vector into a diagonal matrix, and the diag_m operator transforms
 125 a matrix into a vector, containing its diagonal elements. The time complexity to this rule is of the order
 126 $\mathcal{O}(mn)$, with m the number of outputs and n the number of inputs. The stabilized learning rule therefore
 127 takes the same form as Sanger's rule.

2.3 DECORRELATED LEARNING RULE FOR ROBOTICS EXPERIMENTS

128 In partially embodied control, the dynamics of the robot body are used directly as a computational
 129 resource. In our RMH learning setup, this is equivalent to replacing part of the neural network by the
 130 robot body, which receives inputs from the remaining neurons and the observations of which are fed
 131 back into those neurons. The training procedure is now heavily constrained, as it can only adapt synaptic
 132 weights of the remaining neurons, whereas the part of the network that is replaced by the body remains
 133 unchanged.

134 Although this situation is similar to RMH learning in neural networks, it differs in the fact that most of
 135 the physical state of the robot remains hidden to the observer and the number of observable signals that
 136 can be fed into the trainable neural network is relatively small. Furthermore, the dynamics of the observed
 137 variables tend to be highly correlated. For example, stiffening the structure typically causes an increase in
 138 all sensor values.

139 As RMH learning relies on the varying influence of exploration noise on the observed variables.
 140 A common influence (increase or decrease) of the noise on all observed variables therefore reduces
 141 the effectiveness of the learning rule. A simple approach to overcome this issue is to decorrelate the
 142 observations. In (Caluwaerts et al., 2012, Appendix I), we showed that a decorrelation layer that uses
 143 Sanger's rule offers a biologically plausible solution for this. In this work, we take a more pragmatic
 144 approach and decorrelate on a trial-by-trial basis.

145 The resulting learning rule is:

$$\Delta \mathbf{W}^T = \alpha H(r - \bar{r})(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Z}, \quad (13)$$

146 where $H(\cdot)$ represents the Heaviside step function. In this variant of the learning rule, weight updates only
 147 occur when the observed reward is better than expected. This is not strictly necessary, but we found that
 148 it slightly improved our results.

149 This learning rule is similar to ridge regression, the difference being that the algorithm will try to
 150 reproduce the noise \mathbf{Z} (instead of a desired output) proportionally to the relative performance with the
 151 injected noise w.r.t. the expected reward. We used a large regularization parameter ($\lambda = 1$) which results
 152 in only a minor decorrelation of the sensor variables, yet is enough to allow for efficient learning. A high
 153 regularization parameter allows for simple covariance estimators, in case a more biologically plausible
 154 version of the rule is desired.

3 EXPERIMENTS

3.1 NEURAL NETWORK EXPERIMENTAL SETUP

155 Before presenting our results in a robotics context, we first study RMH plasticity in discrete time recurrent
 156 neural networks with hyperbolic tangent activation functions. They receive input, which we denote \mathbf{U} and
 157 a readout function provides observations of the network state. Additionally, exploration noise \mathbf{Z} is injected
 158 into the network. The network update equation is given by:

$$\mathbf{x}[k+1] = \tanh(\mathbf{W}\mathbf{x}[k] + \mathbf{W}_{\text{in}}\mathbf{u}[k+1] + \mathbf{z}[k+1]). \quad (14)$$

159 For each of our neural network experiments, the network is initialized according to the reservoir
 160 computing approach Verstraeten et al. (2007). This implies that we initialized the weights randomly
 161 (i.i.d. standard normally distributed samples) and then tune the network dynamics in a useful regime.
 162 This is achieved by rescaling the weight matrix \mathbf{W} such that its spectral radius² is such that the learning

² The spectral radius is the largest amongst the absolute values of the eigenvalues of \mathbf{W} .

163 converges. For the experiments we will describe, we obtained good performance for initial spectral radii in
 164 [0.95, 1.2], i.e., stable or almost stable networks. This differs from related approaches such as **Legenstein**
 165 **et al.** (2010) and **Hoerzer et al.** (2012), where initially chaotic networks are used. The input weight
 166 matrix \mathbf{W}_{in} was sparsely initialized (20% non-zero elements) with i.i.d. normally distributed values with
 167 standard deviation 0.05. All networks contained 100 neurons.

168 Fig. 1 shows our learning setup for neural networks. The neural network to be trained is the central
 169 element. A reward is provided after a trial based on the network observations. A trial is defined as a
 170 number of time steps which corresponds to a period of time in which the network tries to perform a task of
 171 interest. In parallel to the network, a reward prediction system estimates the expected reward based on the
 172 network inputs. The RMH learning rule finally combines information from the network state, exploration
 173 noise, reward and estimated expected reward to compute an update $\Delta \mathbf{W}$ of the network weights.

3.2 NEURAL NETWORK EXPERIMENTS

174 The networks used for the 3 tasks described below only differ in the way observed outputs are generated
 175 and in the subset of weights that can be modified by the learning rule. In the first two networks, two
 176 neurons are selected as output-generating neurons, and the output is computed as the sum of their states.
 177 The third network has three output-generating neurons and has an output equal to the product of these
 178 neurons' states.

179 In all networks, the weights to and from the observed neurons are fixed and recurrent. This prevents
 180 the learning algorithm from generating solutions in which the observation neurons become a pure output
 181 layer, which does not influence the state of the rest of the network. In the networks for tasks 1 and 3,
 182 all other internal weights are modifiable, i.e., there are 98^2 and 97^2 trainable connections in tasks 1 and
 183 3, respectively. In task two, the training is further restricted by fixing the input weights for half of the
 184 remaining neurons, i.e., there are only 49 ($48 + 49$) adjustable weights. This means that about half of
 185 the network is a random recurrent neural network. In a neural or neurorobotics context, we can see this
 186 as a rudimentary model for a trainable network interacting with an untrained dynamical system, such as
 187 another brain area or a physical body, e.g., the partially embodied control of a robot arm with a neural
 188 network.

189 Note that having a learning rule that handles delayed rewards is crucial, as the input does not directly,
 190 nor immediately define the observations.

191 In what follows, we describe the three neural network tasks in more detail. We first consider problems
 192 with discrete input spaces. More precisely, we solve the 2 bit delayed XOR problem and a 3 bit delayed
 193 classification task. Our third example has a continuous input space and a more complex readout function.

194 3.2.1 Task 1: proof-of-principle

195 *Inputs* The input signal for this task represents a single bit stream. A zero bit is coded as the negative half
 196 period of a sine wave, a one bit as the positive half (see the top row of Fig. 7). For each trial, we randomly
 197 select one of the 4 possible 2-bit input sequences.

198 *Desired output* The neural network has to compute the so-called 2-bit delayed XOR task, i.e., the
 199 exclusive OR function applied to the last two bits of its input stream, represented as binary values. More
 200 concretely, the output of the network should be as close to plus one or minus one during the last half of
 201 the second bit.

202 The XOR task is a common test or benchmark, because the inputs are not linearly separable. A linear
 203 network cannot obtain optimal performance for all inputs simultaneously. Therefore, this task is a simple
 204 test to verify if the learning rule can exploit the non-linear effects of the network. The task also requires

205 the network to remember a specific part of the input, while ignoring inputs that occurred more than one
 206 bit length in the past.

207 *Reward function* Throughout this manuscript, we use different reward functions. The main reasons for
 208 this is that some reward functions are more appropriate for a specific task and to show that the learning
 209 rule does not depend on a specific reward function (e.g. square loss). For the results presented for the 2-bit
 210 XOR task, we used minus the mean squared hinge loss, as the hinge loss is a more appropriate reward
 211 function for a binary classification task:

$$r_n^{2\text{bit}} = \frac{-1}{5} \sum_{k=15}^{19} \max(0, 1 - t[k] (x_0^o[20n+k] + x_1^o[20n+k]))^2, \quad (15)$$

212 where n indicates the number of the current trial, x^o are the neurons that generate the observations and
 213 $t[k]$ are the desired observations. However, we have also successfully applied the mean squared error
 214 reward function for this task.

215 *Prediction of the expected reward* Estimating the expected reward is trivial in the case of a modest number
 216 of different inputs. For the results presented here, we averaged the last 50 rewards per input sequence.

217 3.2.2 Task 2: partially embodied computation

218 *Inputs* For this task, the input is the same as for the previous task. For each trial, we now randomly select
 219 one of 8 possible 3-bit input sequences.

220 *Desired outputs* The network now has to act as a 3-bit digital-to-analog decoder, i.e., it has to produce
 221 one of 8 equidistant analog values in the range $[-1, 1]$, corresponding to the decimal interpretation of the
 222 last three encoded bits it received. Similar to the previous task, the desired value has to be present on the
 223 output during the second half of the third bit.

224 This task is more complex and nonlinear than the previous one and it requires more memory as well.
 225 The target values are indicated by red crosses in the bottom row of Fig. 7.

226 *Reward function* For the 3 bit task, the reward value $r^{3\text{bit}}$ is defined as minus the mean squared error of
 227 the observations during the last five time steps:

$$r_n^{3\text{bit}} = \frac{-1}{5} \sum_{k=25}^{29} (t[k] - x_0^o[30n+k] - x_1^o[30n+k])^2. \quad (16)$$

228 *Prediction of the expected reward* The rewards were estimated in the same way as in the previous task.

229 3.2.3 Task 3: nonlinear observation function

The task at hand is to reproduce part of the input in
 230 reverse after a delay (see Figure 3).

231 *Inputs* The network receives two input signals. The first input signal (see example in Figure 3) consists of
 232 sequences of 12 time steps per trial. The first 5 of these steps form a linear segment between two values,
 233 which are sampled with uniform probability from $[0, 1]$ for each trial. The final value is held constant for
 234 two more time steps. The final 5 time steps again form a linear segment that is obtained by connecting the
 235 last value from the first 5 steps with a third random sample from $[0, 1]$.

236 Because the trials are fed into the system one by one, it is not clear to the network when a trial starts.
 237 The second input, a binary signal, is used to inform the network when it has to generate the desired output.

238 *Desired outputs* The network must learn to recall the input during the first 5 steps and reproduce them in
 239 reverse order during the last 5 steps of the trial. The system must ignore the remaining 7 input samples.

240 *Reward function* The reward function used here is minus the mean absolute error of the observations:

$$r_n^{\text{cont}} = -\frac{1}{5} \sum_{k=0}^4 |u[12n+k] - \prod_{j=0}^2 x_j^0[12n+11-k]|. \quad (17)$$

241 *Neural network structure* The complete network structure is shown in Fig. 2(right). The observations are
 242 computed by multiplying the states of three internal neurons, which have fixed (i.e., untrained) incoming
 243 and outgoing connections.

244 *Prediction of the expected reward* The expected reward \bar{r} estimates the performance of the system given
 245 the noise level σ . Furthermore, the reward is input dependent, therefore \bar{r} estimates the following quantity:

$$\bar{r} = \mathbb{E}[r|\mathbf{u}, \sigma]. \quad (18)$$

246 Various algorithms can be used to estimate this quantity. We employed the well known Recursive Least
 247 Squares algorithm **Kailath et al.** (2000) to learn a simple online estimator of this quantity, which we
 248 applied to the input sequences \mathbf{u} and the network state \mathbf{x} at the end of a trial.

3.3 THE TENSEGRITY ROBOT

249 The general setup of our tensegrity robot control problem is shown in Figure 4. It is similar to the neural
 250 network setup represented in Fig. 1, but the recurrent neural network has been replaced with the tensegrity
 251 robot.

252 The tensegrity structure used for our experiments has 4 struts and is shown in Figure 5. It is based
 253 on the standard 3 strut tensegrity prism to which a shorter rod has been added that acts as a compliant
 254 end-effector. The bottom 3 nodes of the original prism have been fixed through ball-joints. The resulting
 255 structure has $17 k = 20 \text{ N m}^{-1}$ springs, 14 of which are actuated (the length of the 3 bottom springs is
 256 fixed). The controller time step was 50 ms and gravity was not modeled.

257 Instead of observing the state of the neurons, we measure the spring forces:

$$x_i = f_i \quad (19)$$

$$= \max(k(l_i - l_i^0), 0). \quad (20)$$

258 Actuators changing the equilibrium lengths of the springs replace the postsynaptic neurons and motor
 259 babbling takes on the role of the exploration noise:

$$l^0 = l^{\text{init}} + \mathbf{W}\mathbf{x} + \mathbf{u} + \mathbf{z}. \quad (21)$$

260 All tensegrity experiments were performed in our tensegrity simulator which is based on an Euler-
 261 Lagrange formulation of the tensegrity dynamics (**Skelton and de Oliveira**, 2009, Chapter 5). For more
 262 details on the simulation setup, we refer to **Caluwaerts et al.** (2012).

3.4 HIERARCHICAL END-EFFECTOR CONTROL FOR THE TENSEGRITY ROBOT

263 The task we consider is writing characters with the top node of the rod suspended in the tensegrity
 264 structure. More precisely, the node has to trace letters in a horizontal (XY) plane. The characters

265 were taken from UCI Character Trajectories data set **Bache and Lichman** (2013), integrated and then
266 subsampled and rescaled.

267 The robot is controlled by combining a feed forward kinematic controller and a learned static linear
268 feedback controller. The kinematic controller provides the input signals \mathbf{u} in Equation 21. We sampled 100
269 random spring lengths to create a set of configurations for the kinematic controller. To write a character,
270 the kinematic controller selects a combination of spring lengths that move the end-effector as close as
271 possible to the desired position when the structure is in equilibrium.

272 The reward function used for the next experiments tries to bring the end-effector close to the desired
273 trajectory:

$$r^{\text{traj}} = \frac{-1}{s} \sum_{k=0}^{s-1} \max(\|\mathbf{n}[k] - \mathbf{c}[k]\| - 0.01, 0), \quad (22)$$

274 where s is the number of steps of the current character, $\mathbf{c}[k]$ the vector containing the target position at
275 time k (relative to the beginning of the trial) and $\mathbf{n}[k]$ the position in the XY plane of the end-effector at
276 time k . This reward function will cause the learning rule to stop improving a feedback controller w.r.t. a
277 point on the trajectory in case the end-effector is within 1cm of the target position.

3.5 ROBUSTNESS AGAINST FAILURES

278 To demonstrate the robustness of the controllers as well as a more practical application of the learning
279 rule, we simulated various actuator failures. In this case, we used a more realistic feed forward controller.
280 Again starting from a simple kinematic controller, we now optimized the inputs \mathbf{u} in Equation (21) at
281 each time step using a basic exploration method. More precisely, a small amount of noise z is injected at
282 each time step and the change in expected reward is observed. If an improvement of the expected reward
283 is observed after a trial, we reproduce the noise in the feed forward controller $\mathbf{u} = \mathbf{u}_{\text{kin}} + \mathbf{u}_{\text{expl}}$ by using
284 $\Delta \mathbf{U}_{\text{expl}} = \mathbf{Z}$ when the expected reward (of the trial) improved and $\Delta \mathbf{U}_{\text{expl}} = 0$ when it did not. In the
285 previous equations, \mathbf{u}_{kin} refers to the original kinematic controller discussed in the previous section.

286 In this setup, we now simulate actuator failures by resetting one or more actuators to the original
287 kinematic controller instead of the optimized ones, i.e. $\mathbf{u} = \mathbf{u}_{\text{kin}}$.

4 RESULTS

4.1 NEURAL NETWORK EXPERIMENTS

288 We first demonstrate on an academic benchmark task (task 1) that the proposed learning rule can be
289 successfully applied to neural networks. The left panel of Figure 6 shows the evolution of the reward and
290 the spectral radius while the system learns the 2 bit delayed XOR task. The noise level was $\sigma = 0.01$
291 and the learning rate was $\alpha = 0.05$. We observe that the network indeed learns to solve the task almost
292 perfectly, as the reward converges to its maximal value of 0.

293 In this experiment, we set the initial spectral radius to 0.95 to show how the spectral radius evolves as
294 the average reward increases. Every run of the algorithm for different initial random weights resulted in a
295 final spectral radius of approximately 1.2, which indicates that the learning rule indeed tunes the memory
296 of network. The right panel of Fig. 6 shows a detailed view of the reward and expected reward for a single
297 input sequence.

298 For the particular run shown in this Figure, we see a drop in performance after about 125 000 trials.
299 This is probably a significant change in network behavior (a bifurcation) due to a minor weight change, a

300 typical phenomenon in non-linear recurrent neural networks. However, the algorithm continues learning
 301 successfully afterwards.

302 Next, we address task 2, which is similar to task 1 (the classification of sequentially offered bit patterns),
 303 but more difficult, as the sequences now consist of 3 bits. However, in this case the learning rule is only
 304 allowed to modify half of the internal weights of the recurrent neural network. The network structure can
 305 be considered as a model for partially embodied computation, i.e., we replace part of the original trainable
 306 network by a fixed one, which acts as a dummy for a physical body.

307 The results presented here show that, again, the learning rule is successful. Indeed, Figure 7 shows
 308 the results of overlaying 50 random orderings of the input sequences. Note that the classification result
 309 must only be available at the output during 5 time steps at the end of each trial (indicated by red crosses
 310 in the figure). As trials follow each other continuously, the variability of the state trajectories is due to
 311 the different initial states. Clearly, the network has learned the correct time window, as the classification
 312 result is available at the right time and only depends on the two previous bits and not on older inputs (the
 313 random inputs were fed into the system sequentially). We can identify a number of separate trajectories
 314 that keep track of the possible outcomes. It can be seen that within a trial, each additional bit that is offered
 315 at the input reduces the number of possible states of the system by half. Interestingly, the two neurons
 316 that generate the observations have different state trajectories, because the learning rule only quantifies
 317 the performance based on the observations, without directly enforcing a specific behavior of the neurons
 318 responsible for the observations.

319 Whereas the network was trained without any noise on the input signals, the resulting behaviour is
 320 robust against such noise. In Figure 8, we show the behavior of the network when input noise is present.
 321 This plot was generated by first applying k-means clustering on the trajectories and then estimating the
 322 variance of each centroid. Shown are the various centroids and the standard deviation of each. We see
 323 that the network is robust against high amounts of noise on the input data (σ up to 0.5), as the original
 324 trajectories are maintained.

325 Finally, we move to the last and most elaborate task. Figure 9 visualises the rewards during testing for
 326 various state noise levels, using 100 000 random input trials per noise level. The noise was added to the
 327 internal neurons of the network, but not to the 3 neurons which generate the observations. The network
 328 shown was first trained for 1 000 000 time steps with an exploration noise level $\sigma = 0.05$ and a learning
 329 rate of 0.005. Each graph in the top panel shows the average rewards of the trained networks, across the
 330 whole spectrum of possible input sequences for a given level of input noise (increasing from left to right).
 331 In the top left graph, we see that without noise, the average reward remains close to its optimal value
 332 of 0 for most input patterns, although some regions of the input pattern space seem to be slightly more
 333 difficult. This demonstrates that the learning rule also works, although less perfectly than in the previous
 334 cases, when the relation between the internal states in the network and the way they are translated into
 335 actions and rewards is highly nonlinear and when the input patterns do not fall into discrete categories.

336 The bottom panel shows the reward distribution, averaged across the different input trials, for each noise
 337 level. It becomes clear that \bar{r} is indeed the expected reward under the influence of noise for the given input
 338 sequence. The expected reward shifts to the left as a function of the amount of noise in the network. The
 339 learning rule evaluates the performance with respect to the expected reward for the current noise level.

4.2 TENSEGRITY EXPERIMENTS

340 The left panel of Figure 10 shows how the tensegrity robot performs when drawing characters between
 341 20 and 68 time steps long (1s to 3.4s). A different set of feedback weights was learned for each character,
 342 therefore it is easy to predict the expected reward. To clarify, we estimated the expected reward by
 343 averaging the rewards obtained during the previous 30 trials. As can be seen from the top row, the initial
 344 performance of the system with only the kinematic controller is very low, whereas the combination of
 345 both controllers, using our RMH learning rule, performs considerably better.

346 The plot on the right of Fig. 10 shows the learning curves for each character, indicating that for most
347 characters good results were obtained after 1000 to 1500 trials, which would be equivalent to less than
348 one hour real robot time for most characters. It is possible to accelerate learning by further tuning the
349 learning parameters. We used a conservative level of exploration noise ($\sigma = 5$ mm) and a learning rate
350 $\alpha = 1$, which consistently resulted in stable feedback controllers. The learning rule did not achieve the
351 same final reward for all characters (e.g. the “m”). This is due to physical limitations enforced on the
352 motor commands.

353 Finally, we simulate actuator failures. The results of these experiments are presented in Fig. 11. As
354 could be expected, the performance immediately drops significantly after each failure. By applying the
355 RMH learning rule to the feedback controller, the system is able to recover from the various failures. To
356 investigate the stability of the learning rule, each experiment was performed 30 times, with similar results.

5 DISCUSSION

357 Hebbian theory is a well-established approach to explain synaptic plasticity between neurons. Over the
358 years, many variations of the basic learning rule have been developed. Each of these had a specific
359 application, ranging from unsupervised feature extraction to reinforcement learning. In this work, we
360 presented Hebbian like learning rules in which the synaptic plasticity is based on the correlation of the
361 presynaptic neurons and an exploratory noise signal. The plasticity was modulated by a distal reward
362 signal, resulting in a learning method that maximizes the expected reward of a trial.

363 While similar rules have already been presented, we focused on distal reward learning in recurrent neural
364 networks and compliant robots. The rationale for this is our belief that both can be seen as computational
365 resources and can therefore benefit from similar learning techniques. We had already shown this before,
366 but in this work we also presented a method to efficiently train the substrates.

367 Our work builds upon Legenstein's **Legenstein et al.** (2010), who considered simulated motor control
368 tasks in combination with an instantaneous reward signal in an initially chaotic neural network. One
369 significant difference with respect to our experimental setup is that Legenstein estimated the exploration
370 noise as well as the expected reward. This allows for uncontrolled or unknown noise sources to be used,
371 which adds to the biological plausibility of the method learning **Faisal et al.** (2008). Covariance and noise
372 based rules have a strong biological foundation **Loewenstein and Seung** (2006); **Loewenstein** (2008);
373 **Soltani and Wang** (2006). For example, it is well known that neural networks in biology have intrinsic
374 noise sources **Faisal et al.** (2008), which could be used for learning **Maass** (2014). While this type of noise
375 can sometimes be measured by external means (e.g. voltage clamps), a plasticity rule within the biological
376 substrate cannot generally observe the noise signals, hence the importance of the noise estimator in
377 Legenstein's rule. In this work, we did not consider such a noise estimator, as significant uncontrolled
378 noise sources seem unlikely in robotics. Finally, remark that Legenstein's learning rule extends various
379 earlier techniques with similar mathematical formulations **Fiete and Seung** (2006); **Loewenstein** (2008);
380 **Loewenstein and Seung** (2006).

381 Our study of RMH learning in neural networks shows that our version of the learning rule, using
382 distal rewards, works for a wide range of conditions (trainability and observation functions) and for very
383 different tasks. In addition, we have interpreted the network configuration of task 2, in which an entire
384 subnetwork was not trainable, as a model for partially embodied computation. An interesting question in
385 this case is whether learning in the trainable part of the neural network was fundamentally necessary in
386 order to recognize the input patterns, or, in contrast, whether the necessary computations were already
387 present in the network dynamics. In this case, as in a traditional reservoir computing setup, the learning
388 only needed to provide a suitable mapping from the internal dynamics to the observation function.

389 The presented results show that after training, the system state evolves along a fixed number of highly
390 robust trajectories. This phenomenon is not commonly observed in reservoirs without trained feedback
391 weights, indicating that training half of the network at least provides feedback loops in addition to a

392 suitable observation function. However, whether an actual trainable neural network has added value on
393 top of these functionalities is not clear from the current experimental results. A more detailed analysis of
394 the learning outcomes in assemblies of fixed and trainable substrates, as a model for partially embodied
395 computation, is the subject of ongoing work.

396 Given the promising results we obtained in our simulations and the limited number of assumptions
397 we had to make to obtain successful results, this work paves the way towards more complex control
398 hierarchies for robot motor control, in which each level refines the output of the previous one. In this
399 context, too, our results raise some interesting questions, for instance, about the exact role of the, very
400 poorly performing, kinematic controller in our experiments. In fact, the main goal of the kinematic
401 controller is not to have the optimal performance, but rather to inject energy into the structure. In our
402 previous work, we showed an example with instantaneous reward function in which we first trained
403 a feedback controller with known target signals using recursive least squares, and then proceeded to
404 learn additional feedback signals using a reward modulated Hebbian rule (Caluwaerts et al., 2012,
405 Section 5.1.3). The reason why an additional energy source is employed in both cases, is that it is hard to
406 consistently learn pure feedback controllers with simple Hebbian like learning rules. A small change in a
407 feedback weight can cause the system dynamics to fade out, which often results in instability. Therefore,
408 an easy and efficient solution is to use an additional input that consistently pumps energy into the system.
409 In principle, this can be accomplished using a feedback controller as in our previous work, a simple
410 feed forward controller as we use here, or another controller such as a Central Pattern Generator. More
411 extensive research is needed to determine how much of the workload can be offloaded to the lowest level,
412 partially embodied feedback controller and how this scales to more complex control tasks, e.g. involving
413 more complex robot bodies.

414 In summary, our main conclusion is that reward modulated Hebbian plasticity with distal rewards
415 provides a simple, yet effective tool for bridging learning in recurrent neural networks and the exploitation
416 of the own dynamics of compliant robots. This strengthens our belief that both the body and the neural
417 network can be used as computational tools and that they should be combined in a self-organising way
418 into partially embodied hierarchical controllers.

DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

419 The authors declare that the research was conducted in the absence of any commercial or financial
420 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

421 K. Caluwaerts performed the experiments. All authors contributed to the interpretation and structuring of
422 the results and to writing the paper.

ACKNOWLEDGEMENT

423 We would like to thank Andrea Soltoggio for useful discussions in the preparation of this letter.

424 *Funding:* This research was funded by a Ph.D. fellowship of the Research Foundation - Flanders (FWO)
425 and the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreements No
426 248311 (AMARSi), and No 604102 (Human Brain Project).

REFERENCES

- 427 Bache, K. and Lichman, M. (2013), UCI machine learning repository
- 428 Caluwaerts, K., Despraz, J., Iscen, A., Sabelhaus, A., Bruce, J., Schrauwen, B., et al. (2014), Design and
429 Control of Compliant Tensegrity Robots through Simulation and Hardware Validation, *Journal of the
430 Royal Society Interface*, 98, 11
- 431 Caluwaerts, K., D'Haene, M., Verstraeten, D., and Schrauwen, B. (2012), Locomotion without a brain:
432 Physical reservoir computing in tensegrity structures, *Artificial Life*, 19, 1, 35–66, doi:10.1162/ARTL_
433 a_00080
- 434 Clopath, C., Longtin, A., and Gerstner, W. (2008), An online hebbian learning rule that
435 performs independent component analysis, *BMC Neuroscience*, 9, Suppl 1, O13, doi:10.1186/
436 1471-2202-9-S1-O13
- 437 Connelly, R. and Back, A. (1998), Mathematics and tensegrity, *American Scientist*, 86, 2, 142–151,
438 doi:10.1511/1998.2.142
- 439 Faisal, A. A., Selen, L. P., and Wolpert, D. M. (2008), Noise in the nervous system, *Nature Reviews
440 Neuroscience*, 9, 4, 292–303
- 441 Fiete, I. R. and Seung, H. S. (2006), Gradient learning in spiking neural networks by dynamic perturbation
442 of conductances., *Physical Review Letters*, 97, 4, 5
- 443 Hebb, D. O. (1949), The organization of behavior (Wiley)
- 444 Hoerzer, G. M., Legenstein, R., and Maass, W. (2012), Emergence of complex computational structures
445 from chaotic neural networks through reward-modulated hebbian learning, *Cerebral Cortex*, doi:10.
446 1093/cercor/bhs348
- 447 Hyvriinen, A. and Oja, E. (2000), Independent component analysis: algorithms and applications, *Neural
448 Networks*, 13, 45, 411 – 430, doi:[http://dx.doi.org/10.1016/S0893-6080\(00\)00026-5](http://dx.doi.org/10.1016/S0893-6080(00)00026-5)
- 449 Ingber, D. E. (1997), Tensegrity: The architectural basis of cellular mechanotransduction., *Annual Review
450 of Physiology*, 59, 1, 575–599
- 451 Kailath, T., Sayed, A. H., and Hassibi, B. (2000), Linear estimation, number 3 in Information and System
452 Sciences (Prentice Hall)
- 453 Legenstein, R., Chase, S. M., Schwartz, A. B., and Maass, W. (2010), A reward-modulated Hebbian
454 learning rule can explain experimentally observed network reorganization in a brain control task.,
455 *Journal of Neuroscience*, 30, 25, 8400–8410
- 456 Loewenstein, Y. (2008), Robustness of learning that is based on covariance-driven synaptic plasticity,
457 *PLoS Computational Biology*, 4, 3, 10
- 458 Loewenstein, Y. and Seung, H. S. (2006), Operant matching is a generic outcome of synaptic plasticity
459 based on the covariance between reward and neural activity, *Proceedings of the National Academy of
460 Sciences*, 103, 41, 15224–15229, doi:10.1073/pnas.0505220103
- 461 Maass, W. (2014), Noise as a Resource for Computation and Learning in Networks of Spiking Neurons,
462 *PROCEEDINGS OF THE IEEE*, 102, 5, SI, 860–880, doi:{10.1109/JPROC.2014.2310593}
- 463 Mazzoni, P., Andersen, R. A., and Jordan, M. I. (1991), A More Biologically Plausible Learning Rule for
464 Neural Networks, *Proceedings of the National Academy of Sciences of the United States of America*,
465 88, 10, 4433–4437
- 466 Oja, E. (1982), A simplified neuron model as a principal component analyzer., *Journal of Mathematical
467 Biology*, 15, 3, 267–273
- 468 Pfeifer, R. and Bongard, J. (2007), How the body shapes the way we think: A new view of intelligence
469 (MIT Press)
- 470 Sanger, T. (1989), Optimal unsupervised learning in a single-layer linear feedforward neural network,
471 *Neural Networks*, 2, 6, 459–473, doi:10.1016/0893-6080(89)90044-0
- 472 Skelton, R. E. and de Oliveira, M. C. (2009), Tensegrity systems (Springer)
- 473 Snellen, K. (1965), Continuous tension, discontinuous compression structures
- 474 Soltani, A. and Wang, X.-J. (2006), A biophysically based neural model of matching law behavior:
475 melioration by stochastic synapses., *Journal of Neuroscience*, 26, 14, 3731–3744
- 476 Soltoggio, A. and Steil, J. (2013), Solving the distal reward problem with rare correlations, *Neural
477 Computation*, 25, 978, doi:10.1162/NECO_a_00419

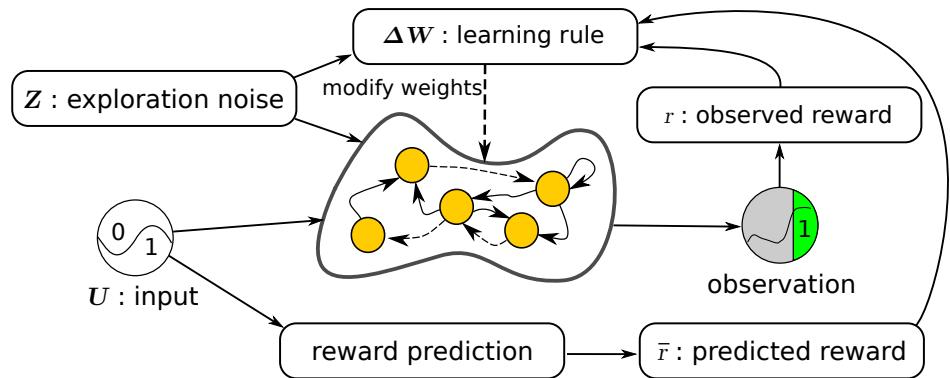


Figure 1. Overview of the learning setup for recurrent neural networks. The initially random recurrent neural network receives the inputs U and the exploration noise Z . The state of the postsynaptic neurons is computed by applying the hyperbolic tangent function to the sum of the inputs, the noise and the weighted sum of the presynaptic neurons. Observations are made of the state of the network and after every trial (fixed number of time steps), a reward is computed, based on the observations made during the last trial. In parallel, a simple reward prediction network predicts the expected reward for the given input. The learning rule then updates the weights between the presynaptic and postsynaptic neurons, by using the reward, the expected reward, the exploration noise and the states of the presynaptic neurons.

478 Verstraeten, D., Schrauwen, B., D'Haene, M., and Stroobandt, D. (2007), 2007 special issue: An
479 experimental unification of reservoir computing methods, *Neural Networks*, 20, 3, 391–403

FIGURES

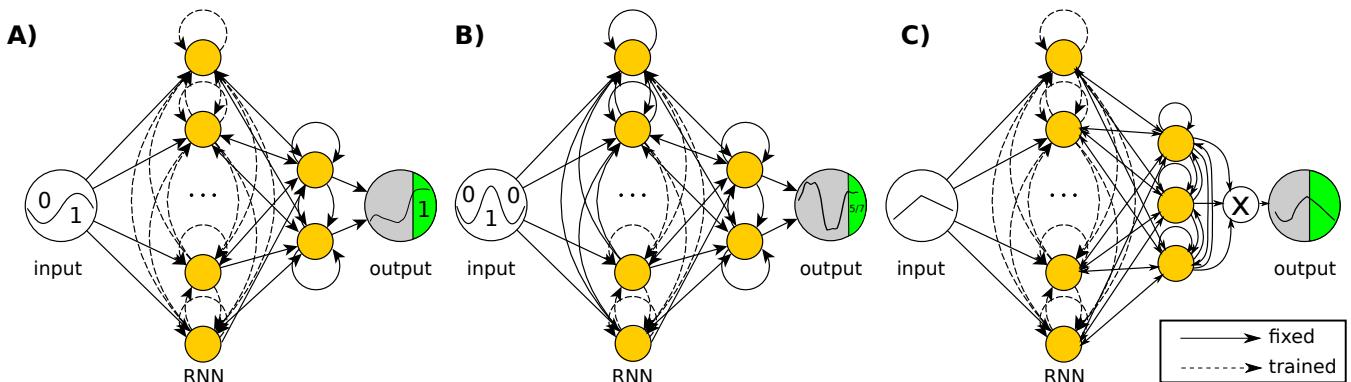


Figure 2. Network structures for the Recurrent Neural Network tasks. The networks are simulated in discrete time with hyperbolic tangent neurons (yellow nodes). Full lines are fixed connections, while dashed lines are trained. The reward is evaluated at the output neuron over the green period of time (one reward per trial).

(left) Task 1 (2-bit delayed XOR): the output equals the sum of two neurons (yellow nodes at the right), but only the internal connections can be modified by the algorithm. The RNN consisted of 100 neurons and the bit period was 10 time steps. (center) Task 2 (3 bit delayed classification): The output equals the sum of two neurons as in the XOR task, but now only half of the internal weights can be modified. The RNN consisted of 100 neurons and the bit period was 10 time steps. (right) continuous input task: the network has to reproduce the reversed input from the first 5 time steps of each trial. A trial consists of 12 time steps in total, of which the last 7 have to be ignored by the network.

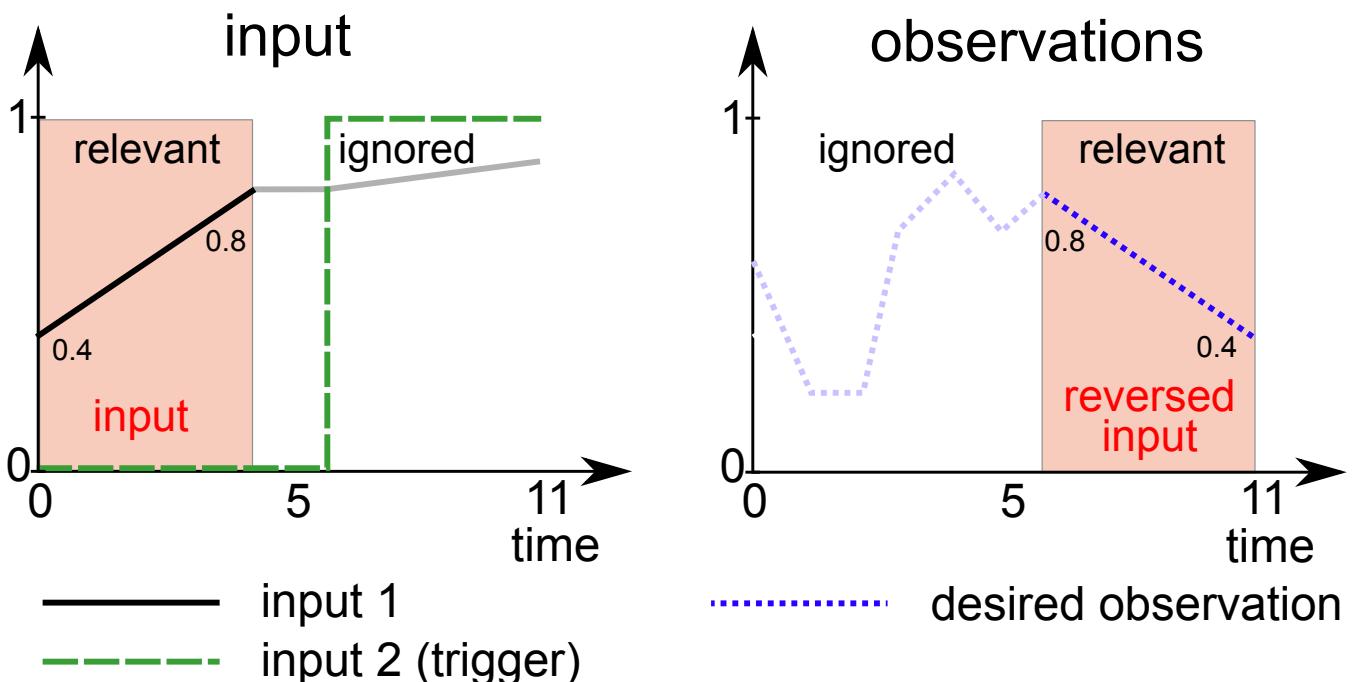


Figure 3. Overview of task 3: the network has to reproduce part of the first input (dashed black line) in reverse at the end of the trial (dashed blue line). More precisely, the first 5 steps of the first input (dashed line) are to be reproduced in reverse at the end of the trial (12 time steps total). The input space consists of all straight line originating and ending in $[0, 1]$. A second input indicates when the network has to start producing the desired observations.

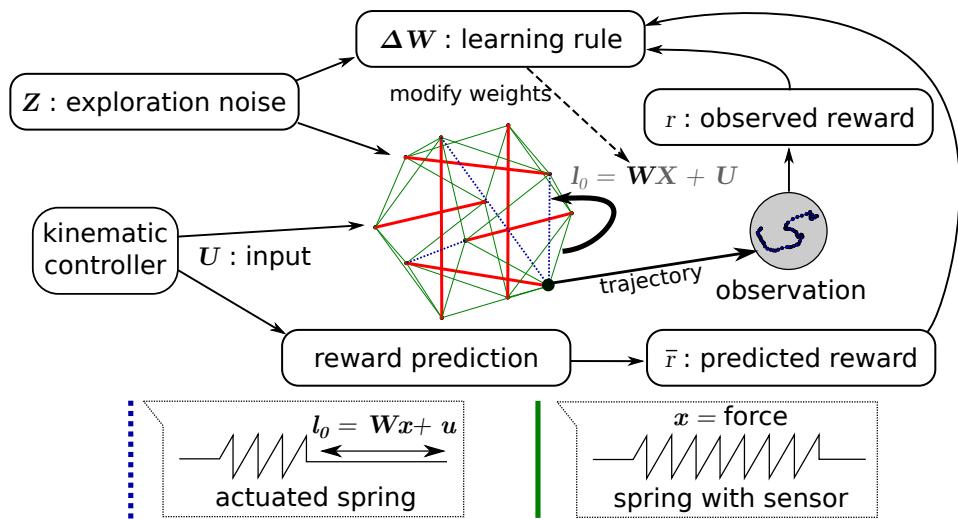


Figure 4. Overview of the way the learning rule is applied to compliant Tensegrity structures. The setup is similar to the recurrent neural network setup of Fig. 1. The neural network has been replaced by the combination of the compliant robot body and the neural linear feedback weights. It now receives input from the kinematic controller. Force sensors on the springs act as presynaptic neurons for the trained weights and the actuator signals correspond to the postsynaptic neurons. The learning rule adapts the feedback weights from the force sensors to the motor signals. The observations used for reward computation are based on the trajectories of an end-effector.

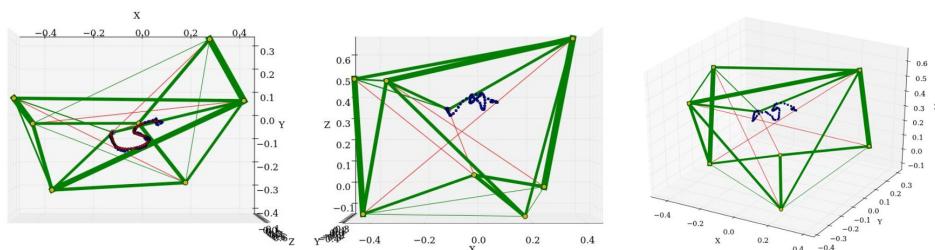


Figure 5. Tensegrity structure used for the experiments. The top node of the center rod is used as an end-effector to draw in the XY plane. In this example, the robot draws an "S" as can be seen on the left. The center and right figures show other perspectives to demonstrate that the reward does not depend on the vertical position.

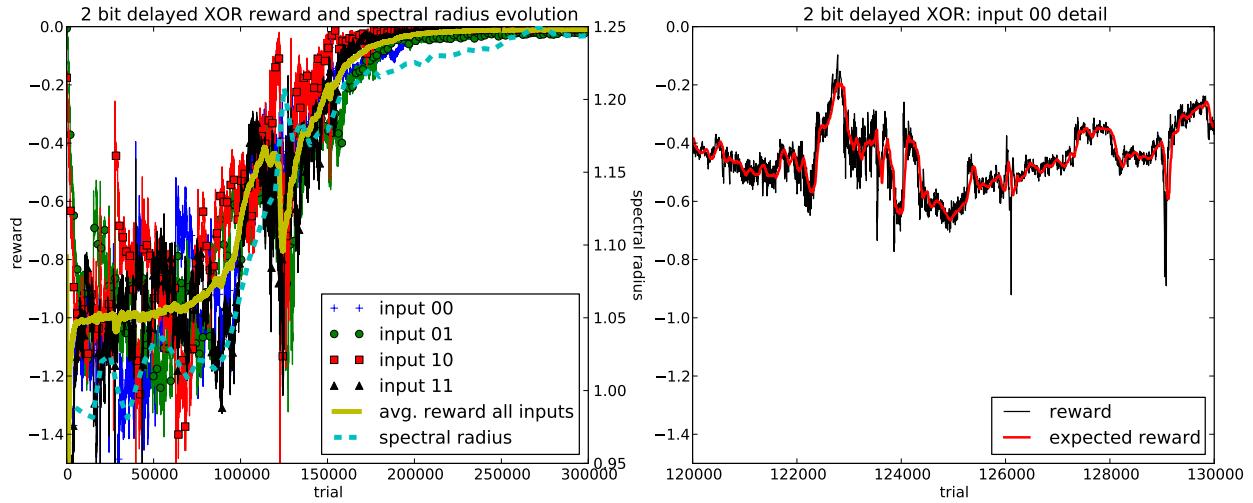


Figure 6.(Left) Evolution of the reward and the spectral radius for the 2 bit delayed XOR task and (Right) detail of the reward and expected reward for a single input sequence.

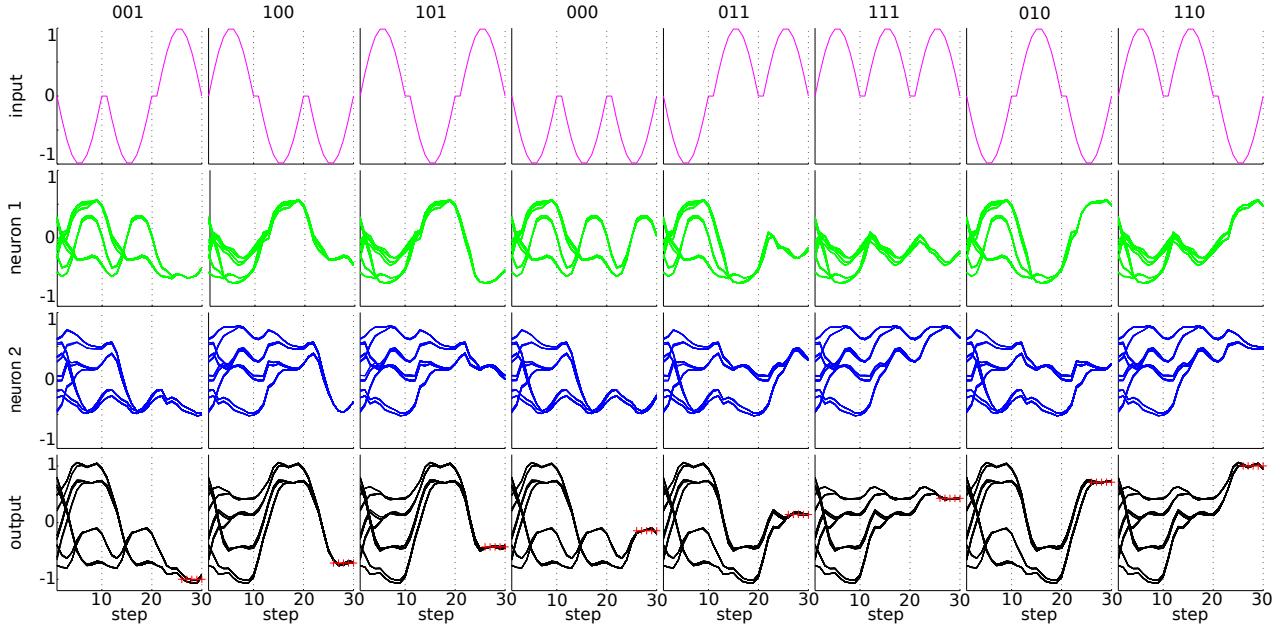


Figure 7. State trajectories for the 3 bit classification task. The top row shows the 8 input sequences. The next two rows show the state trajectories of the 2 neurons which are summed to compute the observations (see Fig. 2(center)) for the given input sequence of the top row. The bottom row shows the observations (the sum of the two middle rows). The desired observation at the end of the trial is indicated by red crosses. It can be seen that each additional bit reduces the number of possible states by half. Interestingly, the two neurons that generate the observations have different state trajectories, because the learning rule only quantifies the performance based on the observations, without directly enforcing a specific behavior of the neurons responsible for the observations. The plots were generated by overlaying 50 random orderings of the input sequences.

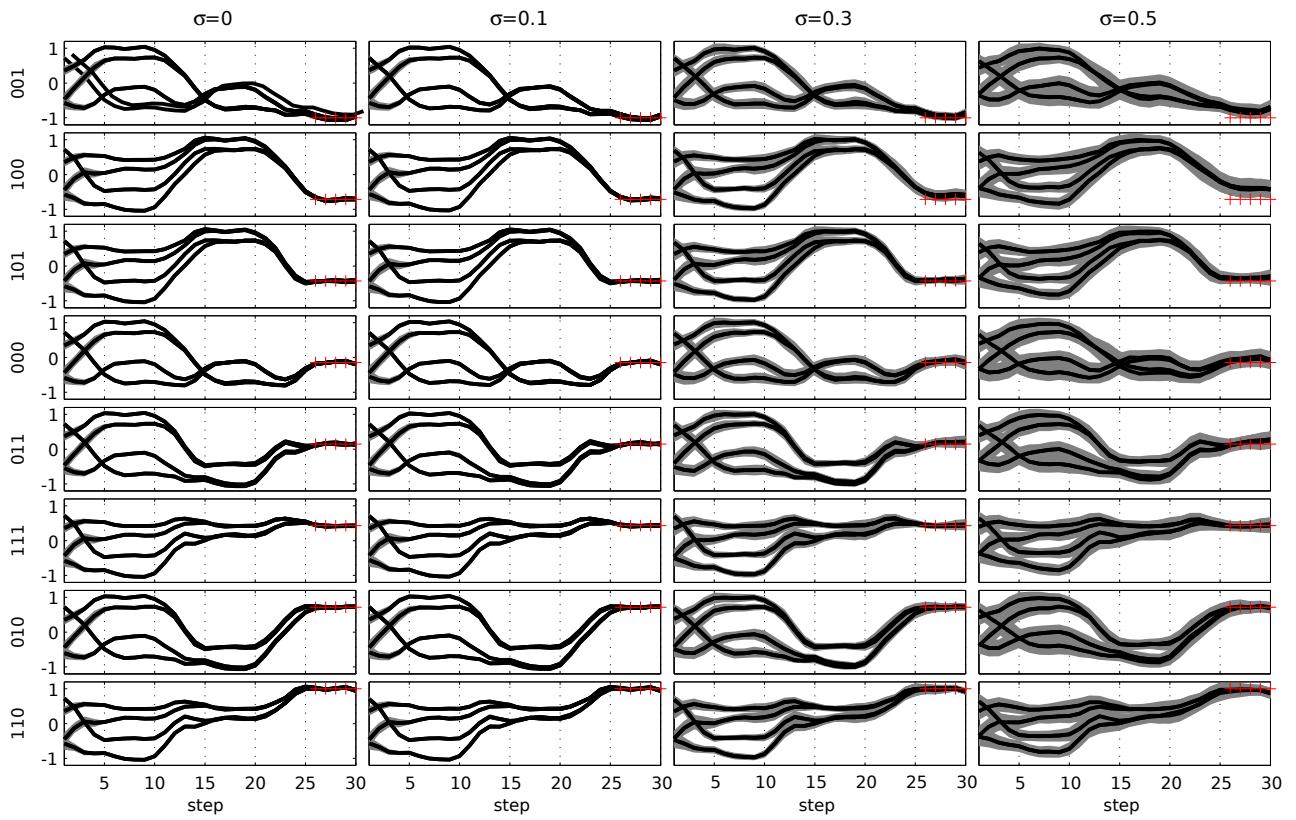


Figure 8. Evaluation of a network trained for the 3 bit classification task under the influence of input noise. As exploration noise drives the learning rule, the trained networks are generally robust against both input and state noise. The noise amplitude increased from left to right. The 8 possible inputs are shown from top to bottom. The gray area indicates one standard deviation around the observations for each of the different state trajectories (thick lines). The red crosses indicate the target observations at the end of the trial.

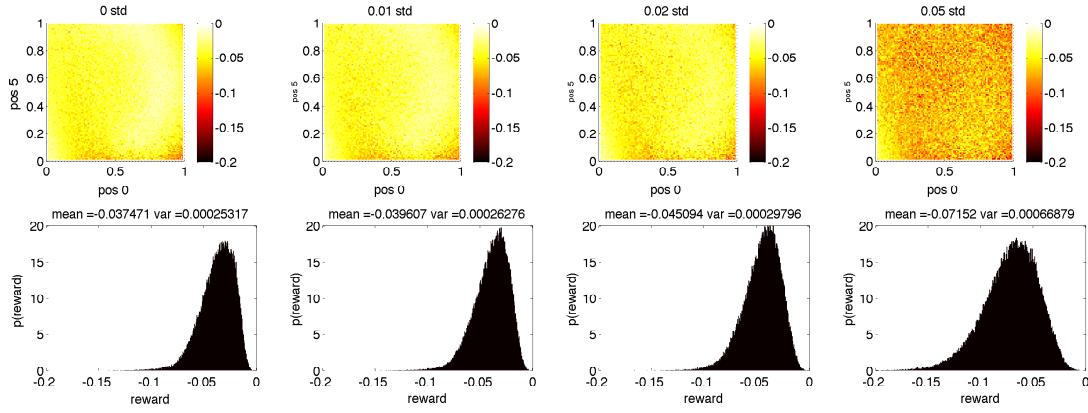


Figure 9. Reward distribution for Task 3 during testing with different state noise levels: (Top) average reward (negative mean absolute error) for the whole range of input combinations. The horizontal and vertical axes of each plot indicate the initial and final values of the linear segments that need to be reproduced, in reverse order, at the output (Bottom) sample distribution of the rewards for all inputs. The amount of injected state noise increases from left to right. We can see the average reward level shift to the left (lower rewards).

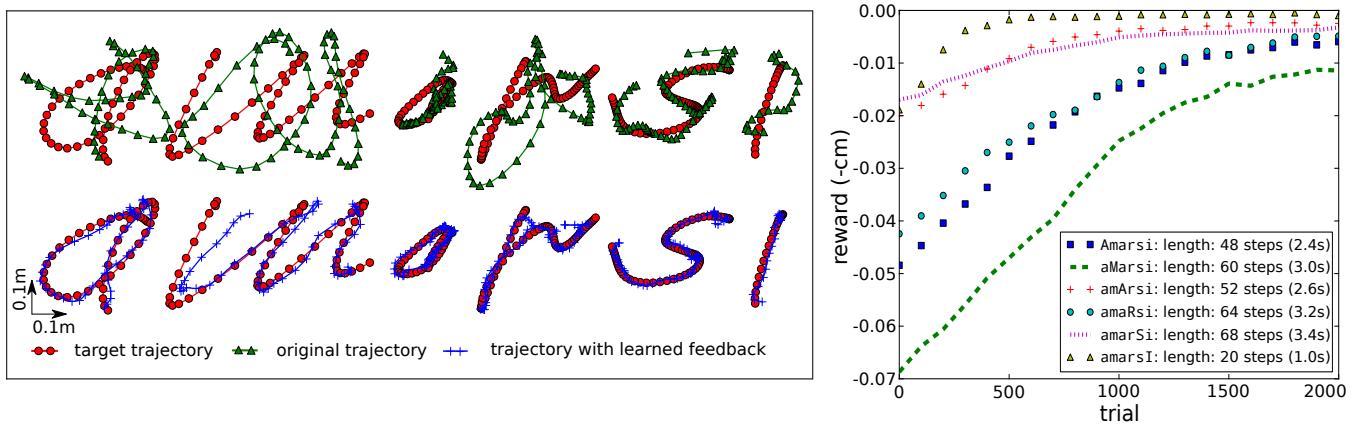


Figure 10. Writing characters with a tensegrity end-effector. Top left: characters drawn with only the kinematic feed forward controller active. Bottom left: characters drawn with the kinematic feed forward controller and the learned feedback controller active. Right: learning curves for the different characters. The legend indicates the length of a trial.

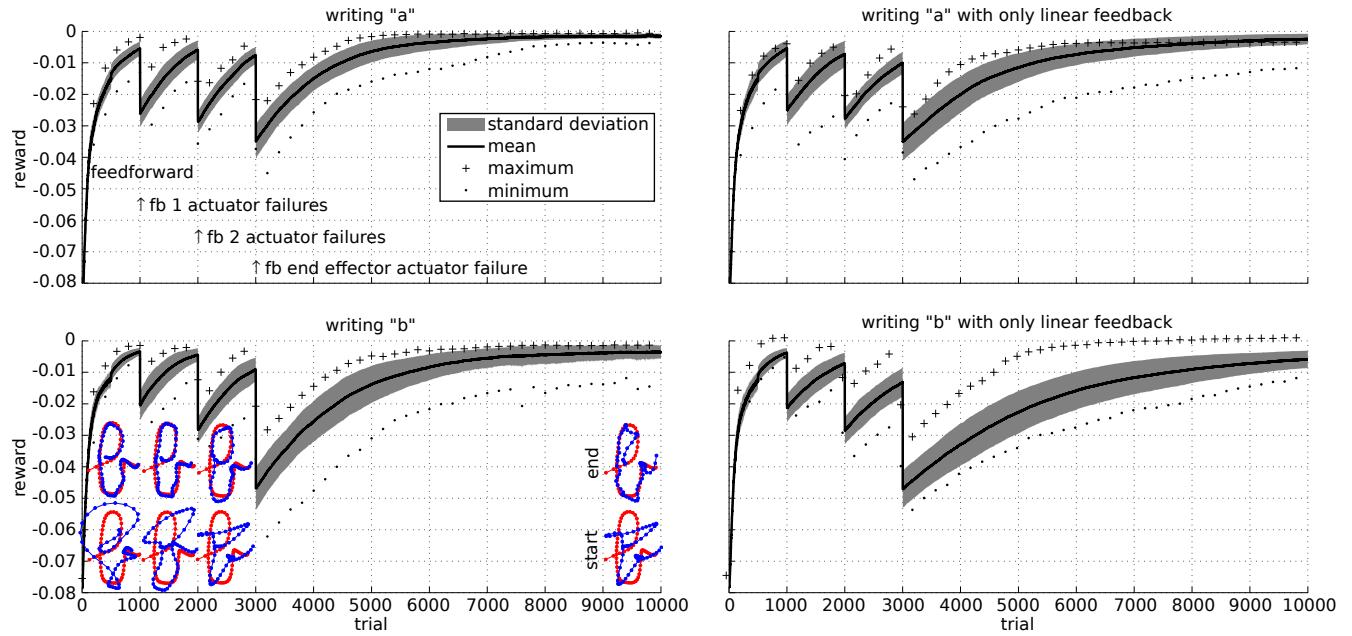


Figure 11. Robustness of the learning rule for the writing task. Initially a feed forward controller is optimized using gradient ascent. We then simulate a failure by making a single actuator follow its initial trajectory from trial 1000 onwards. At the same time, the learning rule starts learning a set of feedback weights to compensate for the actuator failure. Similarly, we simulate a 2 actuator failure after trial 2000. At time 3000, we simulate a failure of an actuator directly attached to the end-effector. The top row shows the results for writing an "a" character, while the bottom row shows the results for a "b". The left column show the results when the feedback includes the spring forces and the square of the spring forces, while the right column only includes the spring forces.

Figure 1.JPG

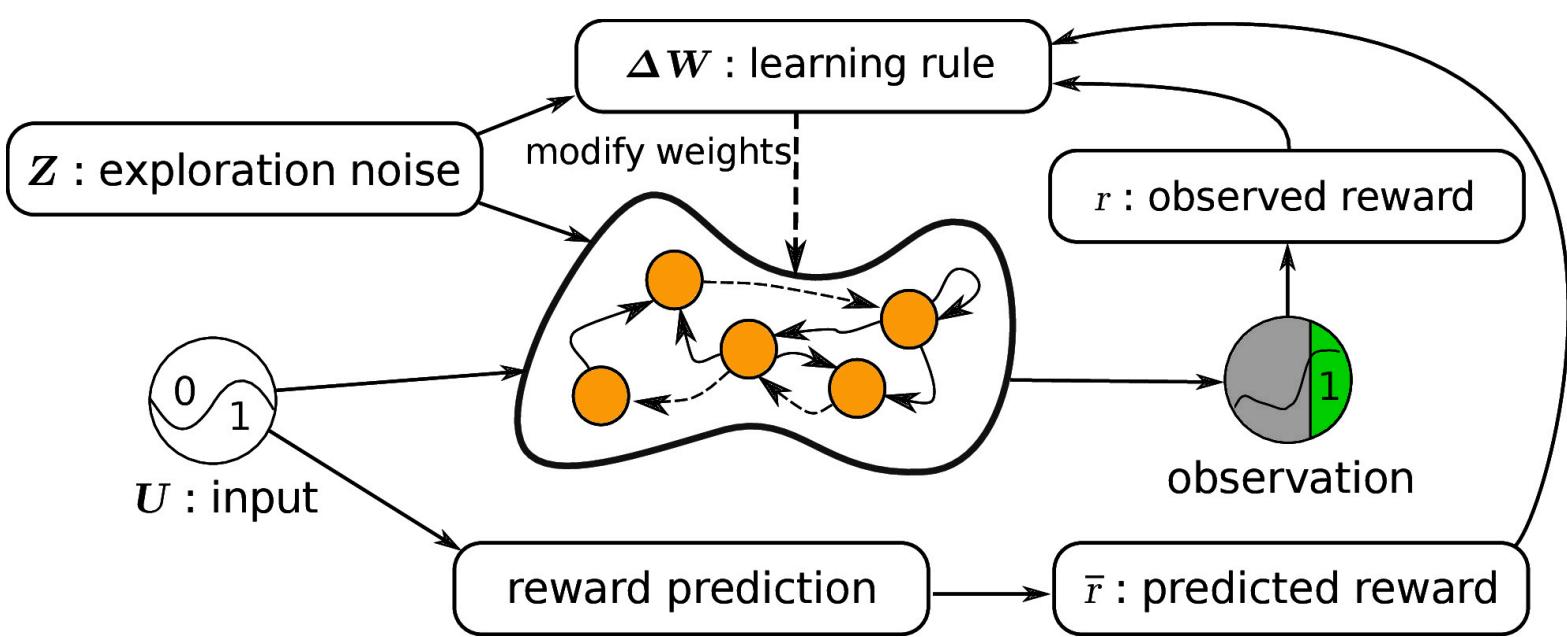


Figure 2.JPG

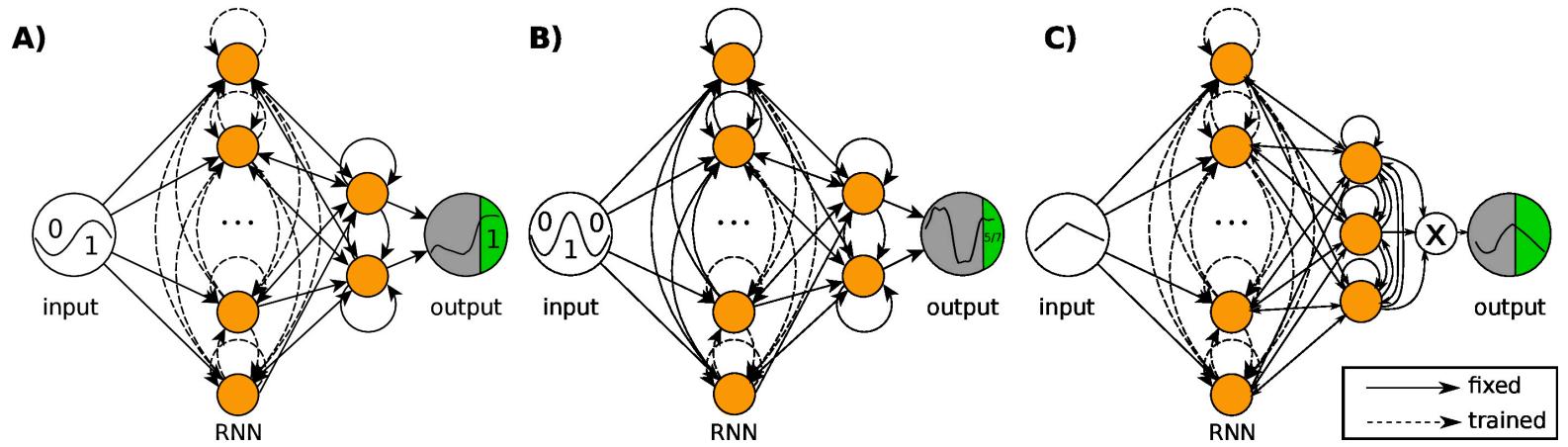


Figure 3.JPG

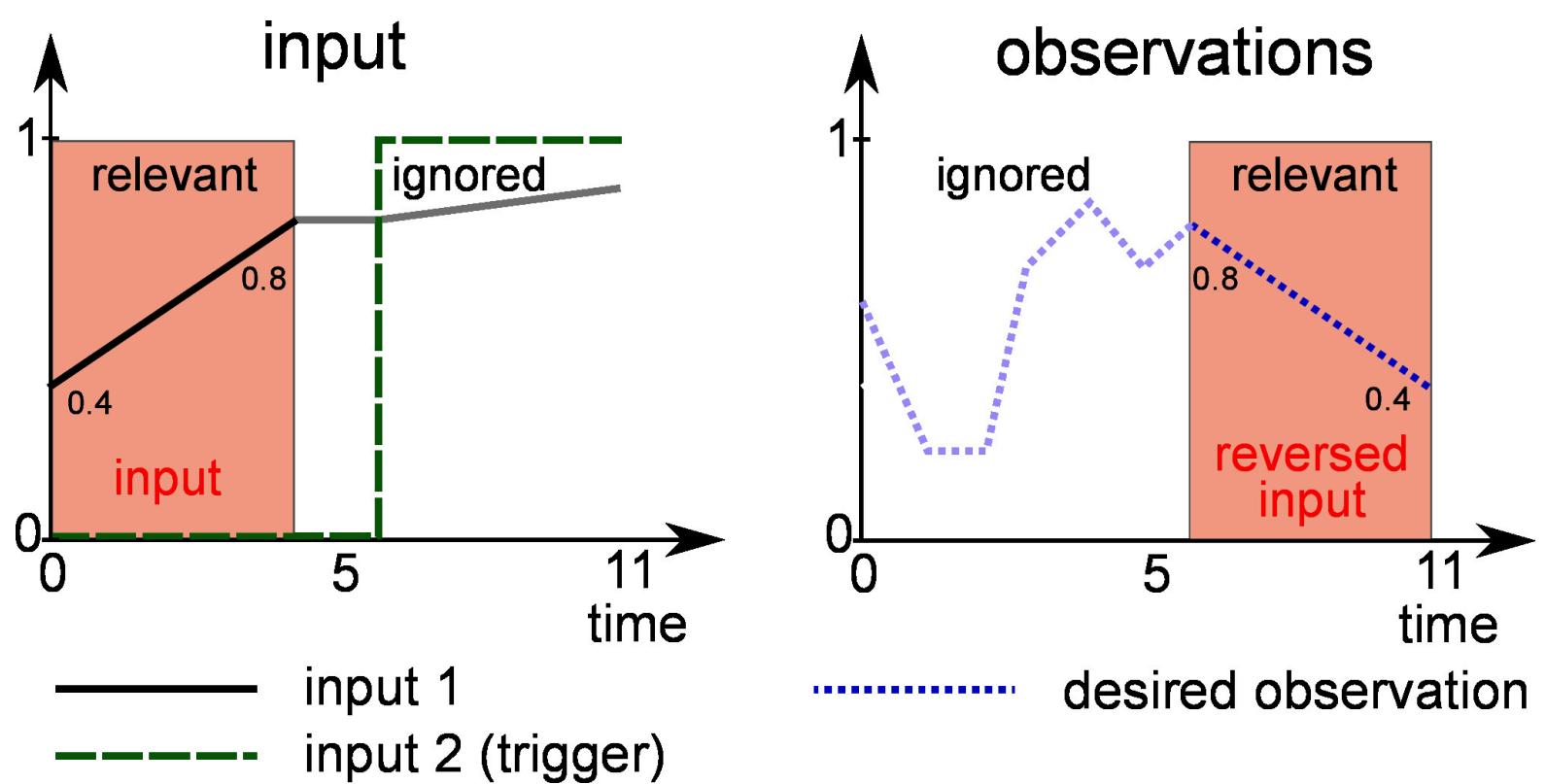


Figure 4.JPG

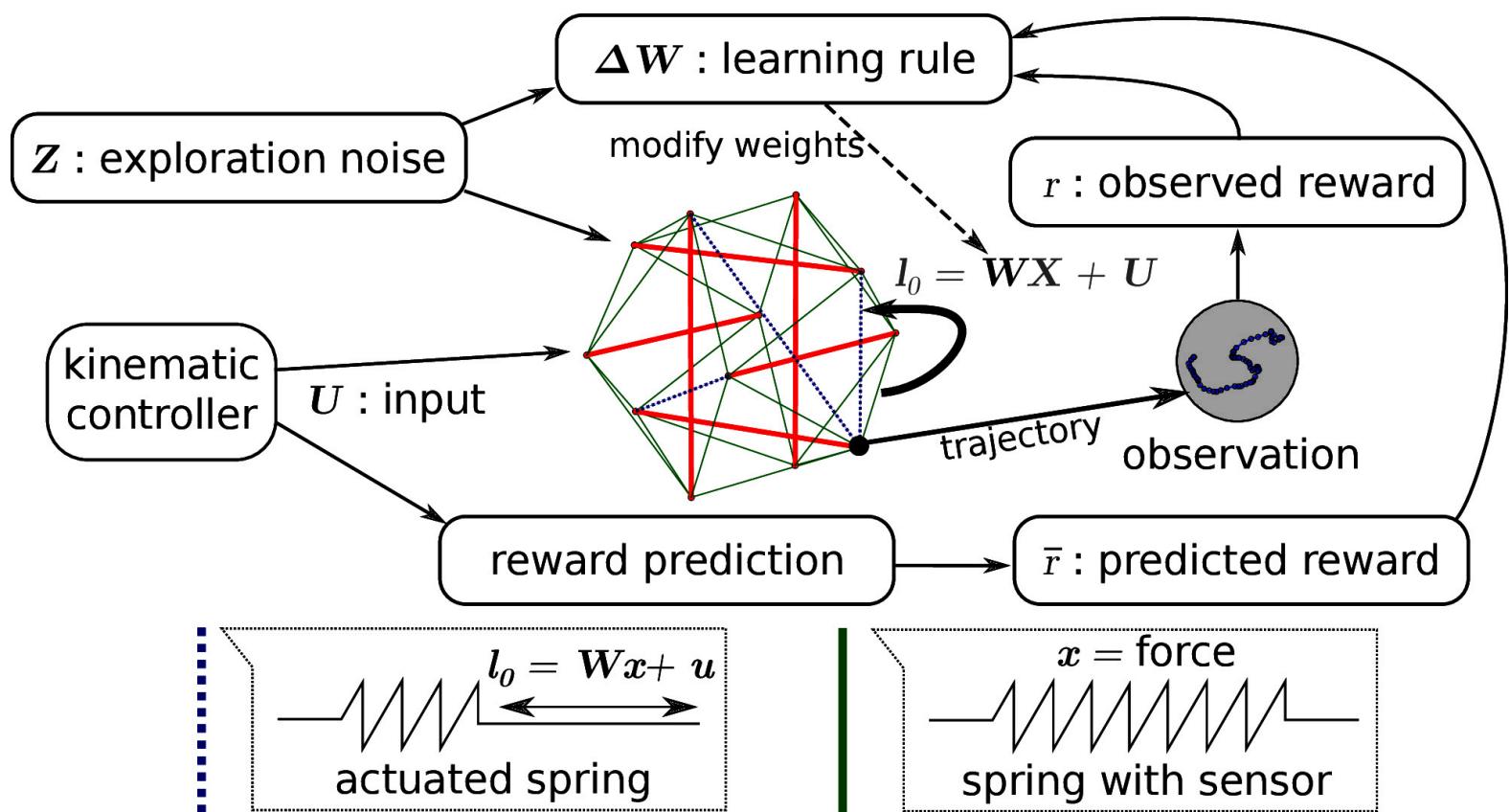


Figure 5.JPG

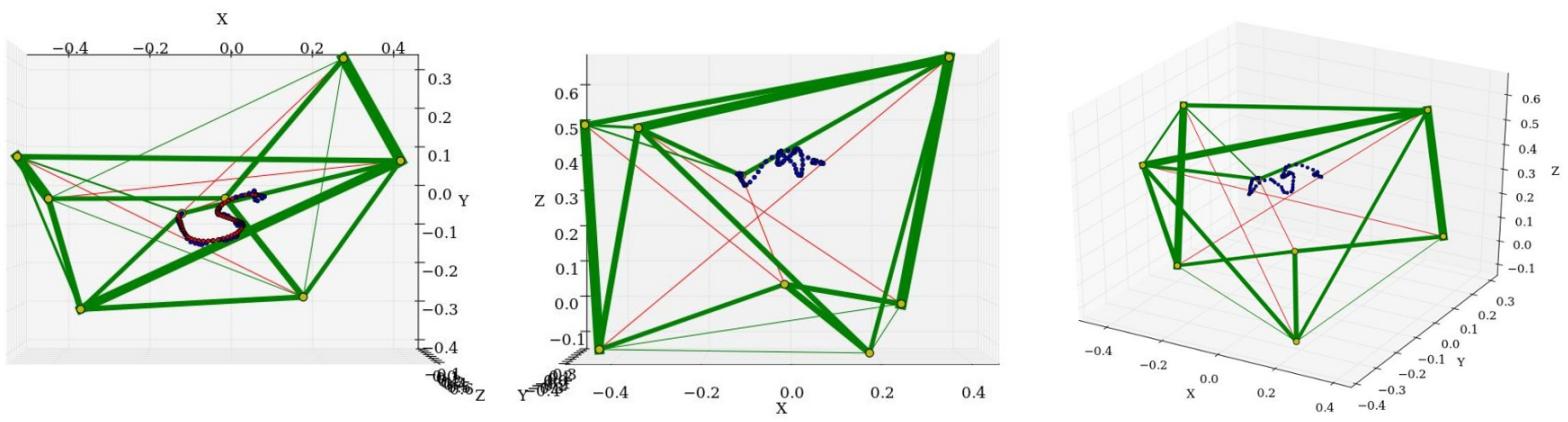


Figure 6.JPG

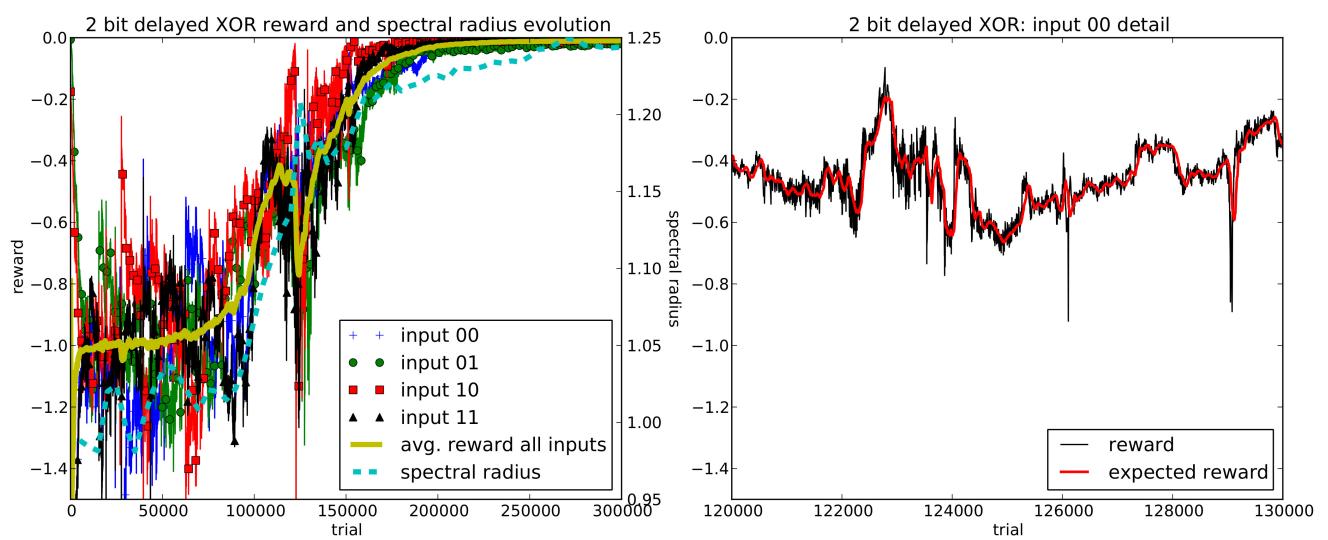


Figure 7.JPG

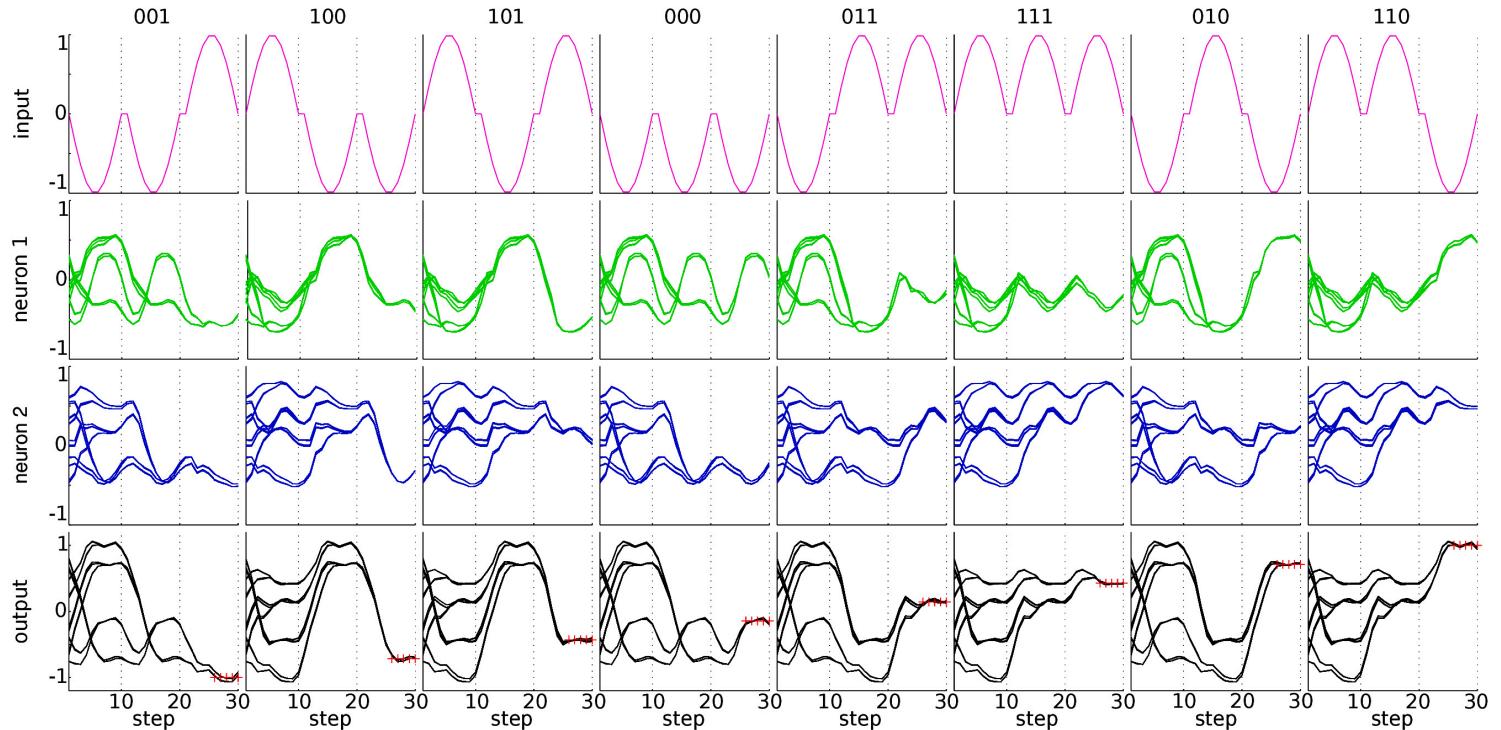


Figure 8.JPG

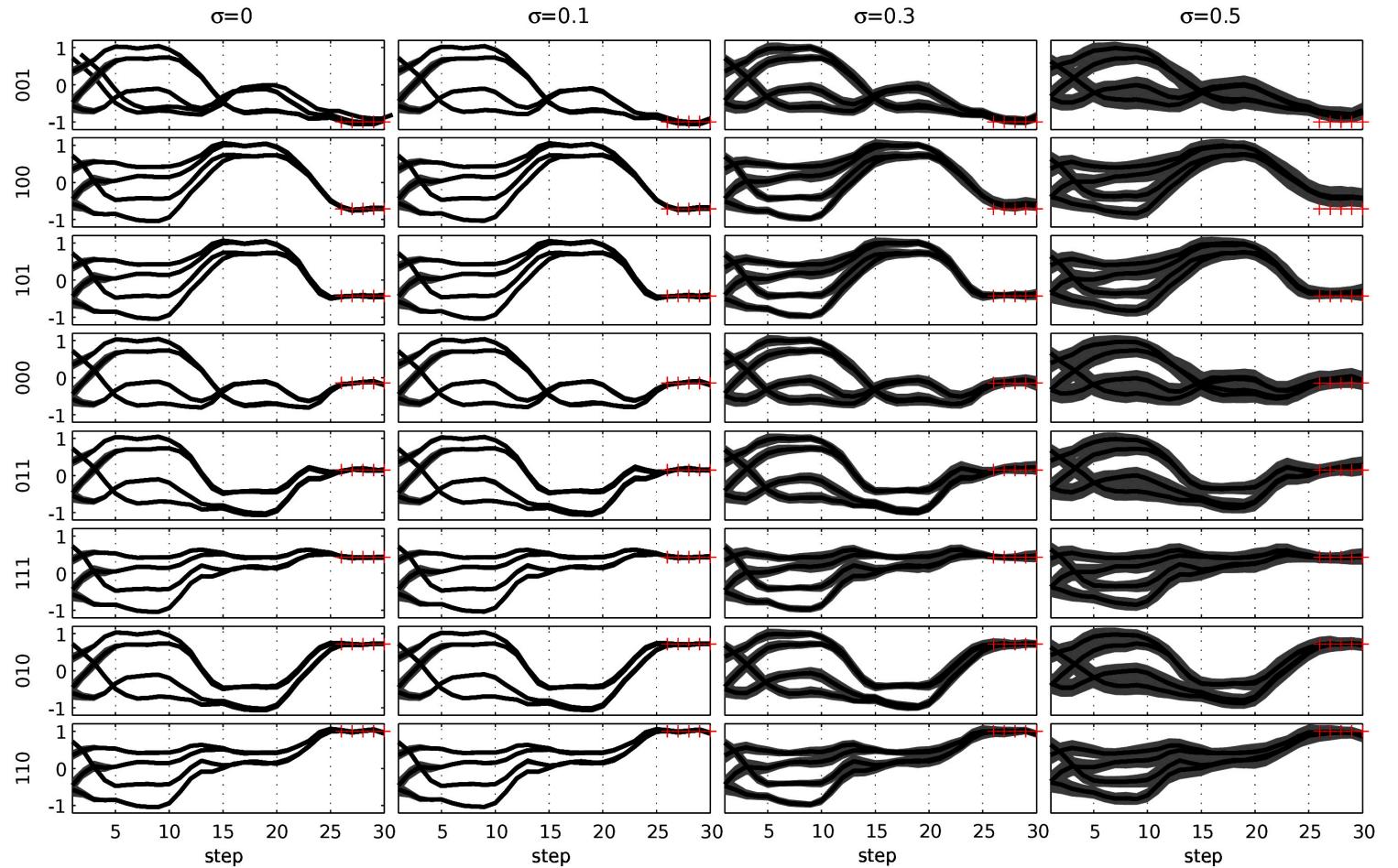


Figure 9.JPG

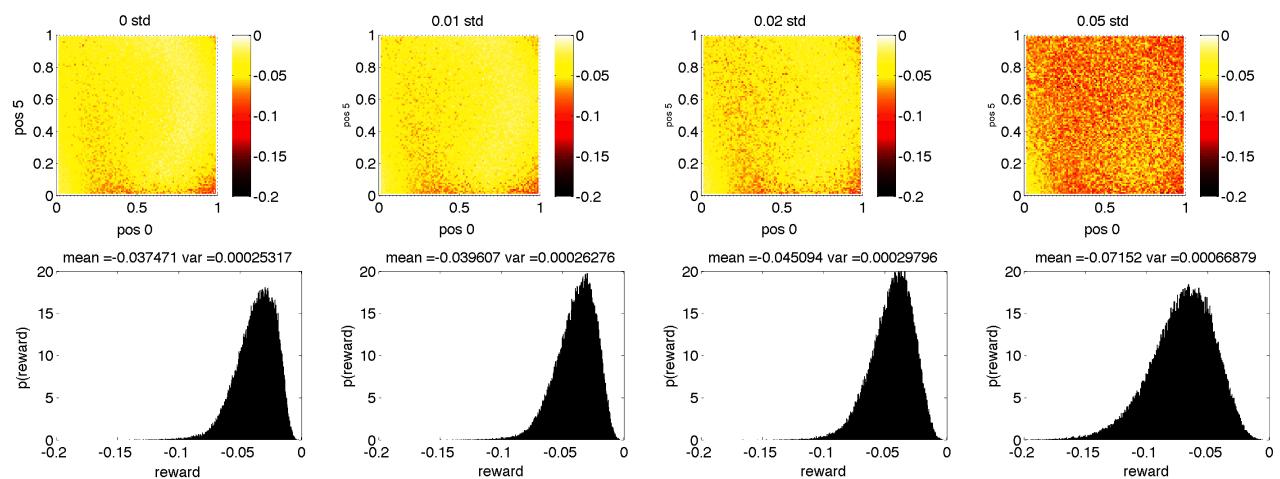


Figure 10.JPG

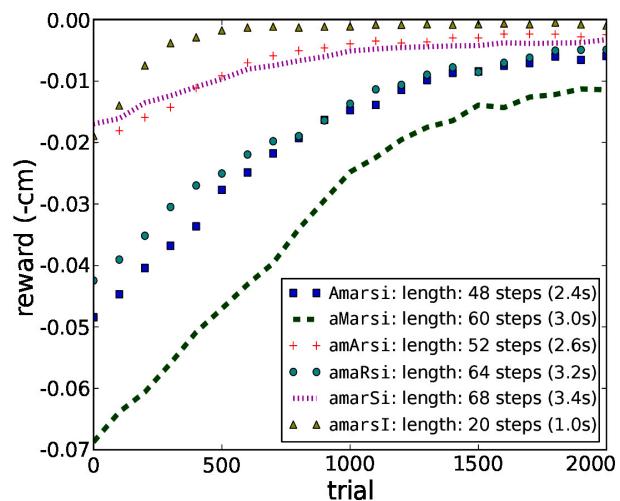
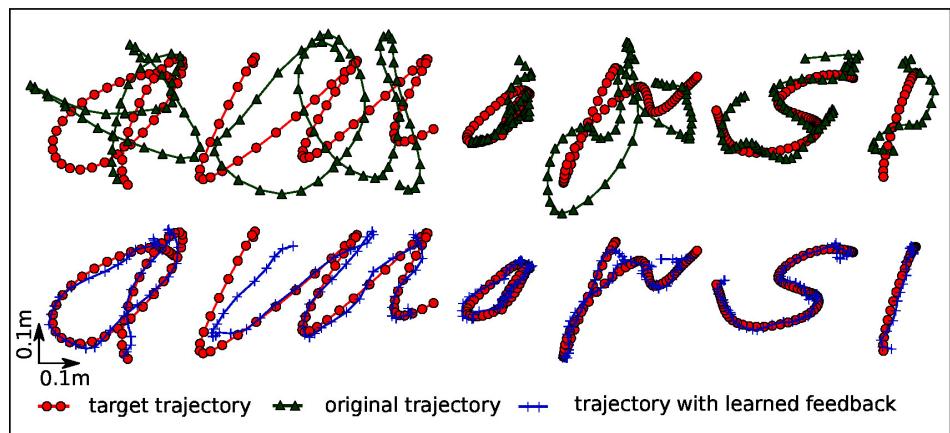


Figure 11.JPG

