



Universiteit
Leiden
The Netherlands

Control of complex actions in humans and robots

Kleijn, R.E. de

Citation

Kleijn, R. E. de. (2017, November 23). *Control of complex actions in humans and robots*. Retrieved from <https://hdl.handle.net/1887/57382>

Version: Not Applicable (or Unknown)

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/57382>

Note: To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/57382> holds various files of this Leiden University dissertation

Author: Kleijn, Roy de

Title: Control of complex actions in humans and robots

Date: 2017-11-23

CHAPTER 6

Optimized behavior in a robot model of sequential action

SEQUENTIAL ACTION is one of the cornerstones of human everyday action. Most of our everyday activities, such as coffee making or driving a car, can be regarded as complex but sequential actions. How humans perform these sequential actions has been the subject of study for at least a century. As described in Chapter 2, sequential action can be represented on a symbolic (what will my next action be?) level, as well as a subsymbolic, sensorimotor (what motor parameters should I use?) level [173]. Interaction effects between the two levels of representation have been observed, and integration between the two is necessary to produce smooth sequential action. Due to their embeddedness (i.e. an implementation in a physical environment), (virtual, humanoid) robots are suitable subjects for developing and investigating models of behavior in which interaction with the environment is important (see [9] for an extensive overview). Robot paradigms have been successfully used to investigate psychological phenomena that require such embeddedness like hand–eye coordination [84], object handling [67], and imitation learning

This chapter is an adaptation of the article *de Kleijn, R., Kachergis, G., & Hommel, B. (in preparation). Optimized behavior in a robot model of sequential action.*

[138]. Used in the proper way, they hold promise to investigate the relation between symbolic planning of actions and the subsymbolic execution of these actions.

6.1.1 Optimization of motor control

The specific motor parameters used in the execution of motor commands is influenced by several effects and constraints. A good example is the *end-state comfort effect* [30], in which the grasp location of an object is a function of the expected end state of the arm. In other words, the arm end state is optimized. Other optimization is seen in the form of contextual lip rounding [35], where the lips are rounded in preparation for pronouncing the /u/ sound well in advance, and bending of mouse trajectories when sequentially reaching for stimuli with a mouse cursor by predicting its future location [75].

Other authors have investigated such predictive movements using similar measures. In earlier work, Dale et al. [34] used a paradigm similar to the one used in Chapters 4 and 5, with different levels of sequence complexity¹. As sequence complexity decreased, participants were found to make larger predictive movements (i.e. movements toward the next stimulus) and be more likely to have explicit sequence knowledge. Participants not making predictive movements were observed to move their mouse cursor to the center of the screen, equidistant from all stimuli. The authors mention that “even participants with low pattern awareness engaged in this form of behavior” (p. 204), but our findings described in Chapter 5 show that it is specifically this group without explicit sequence awareness that engages in this type of behavior. Duran and Dale [38] agree with this finding, and report that this centering strategy is likely employed to compensate for lack of sequence knowledge, making it impossible to accurately predict the next target. In those circumstances, moving the mouse cursor to a position equidistant to all alternatives would be an effective strategy.

¹More specifically, a measure of grammatical regularity was used inverse to the first-order entropy of the sequence, as used in [70].

6.1.2 The current study

In the current study, we directly manipulated prediction quality in a sequential reaching task with a virtual robot hand controlled by an artificial neural network. The task was similar in nature to the task described by Dale et al. [34] and Kachergis et al. [77]: reaching for targets that appeared or changed color in a repeating sequence.

In any modeling problem using artificial neural networks, the connection weights between the artificial neurons (or units) have to be optimized. In other words, the goal is to find those connection weights that cause the artificial agent to produce the behavior that most closely approaches the required behavior as measured by a fitness or cost function determined by the researcher. One of the most popular methods for determining suitable connection weights is known as *backpropagation* [132], in which the network is presented with an input vector, after which the output produced by the network is compared to the desired output, and network weights are then updated according to their error value, starting with the output units and working back through the network.

Evolutionary algorithms such as *neuroevolution* (e.g. [5]) can find suitable network weights not by directly calculating an error measure for each input-output pair presented to the network, but by quantifying the performance of agents controlled by the network. In its most simple form, the method of neuroevolution generates a large number of agents with randomly initialized networks and quantifies how well they perform on the required task during a fixed period of time. Next, the best performing agents are allowed to “reproduce”, and are copied to the following generation in a slightly modified way (e.g. by adding random noise to the connection weights). In subsequent generations, this procedure is repeated until some predefined fitness criterion is reached. Neuroevolution is considered an efficient approach to solving reinforcement learning problems. Past studies have shown neuroevolution to be faster and more efficient than reinforcement learning methods such as Q-learning (see Chapter 4) on several tasks, including robot arm control [99, 100, 150].

Evolutionary algorithms have been used to simulate a wide range of psychological phenomena, ranging from reciprocity [4] to selective attention [114] and category learning [101].

6.2 Method

6.2.1 Task design

The task used for the virtual robots was analogous to the task described by Kachergis et al. [77]. It was designed as an environment of size 50×50 represented in continuous space (i.e. as floating-point values). Over the course of one run of 500 discrete time steps, target stimuli appeared sequentially in one of the four corners of the environment (distance 10 from the environment border), following a simple repeating 1-2-3-4 sequence. In one condition, networks were provided with accurate information about the next stimulus. In a second condition, the information was not predictive of the next stimulus. In a third condition, no information about the next stimulus was provided to the network. The exact implementation is described below under *Network design*.

A virtual robot arm was to touch the target (come within a square of size 6×6 centered on the target) as quickly as possible. After touching a target, no targets were visible for 20 time steps as an inter-stimulus interval (ISI), after which the next target would appear. Every run (one network-controlled virtual robot arm performing the task for 500 time steps), the starting location was initialized to the center of the screen. During each run, the amount of targets touched and the total distance moved was calculated. Also, to encourage fast movement, a reward with decaying value was associated with each target. Rewards were initialized to value 100, decreasing by 1 with each time step. After completion of the run, network fitness was calculated by

$$\text{fitness} = \text{touched stimuli} + \text{total reward} - (.0001 \times \text{distance moved})$$

An agent with perfect prediction capability (i.e. immediately touching the stimulus that just appeared by already being in its location) would

therefore be able to reach a theoretical maximum fitness score of 2525.

6.2.2 Network design

The virtual robot arm was controlled by a two-layer feedforward neural network with four sensory neurons, two prediction neurons, eight internal (hidden) neurons, and two motor neurons (see Figure 6.1). All sensory and prediction neurons were normalized in the range $[0.0, 1.0]$, with Gaussian noise sampled from $N(0, .05)$ added to the input. The two motor neurons were truncated to the range $[-2.0, 2.0]$, and allowed for movement in the two-dimensional plane. For simplicity we did not model the kinematics of an articulated effector.

The input to the two prediction neurons was constant (i.e. also present during the ISI) and represented either (1) the correct location of the next stimulus, (2) the location of one of the four stimuli, randomly chosen, or (3) a constant input of $[0.0, 0.0]$. So although in the second condition the prediction neurons were provided with the location of a stimulus, this location was not informative of the actual location of the next stimulus. These conditions will be referred to as accurate prediction, random prediction, and no prediction, respectively.

The output O_j of a hidden or motor neuron j was determined by the sigmoid activation function

$$O_j = \frac{1}{1 + \exp(-\sum_{i=1}^N w_{ij} O_i - b_j)} \quad (6.1)$$

in which N represents the number of input neurons i , O_i their output, w_{ij} the connection weight from i to j , and b_j the bias. Of the four sensory neurons, two were used for sensing the target, and two for sensing the location of the agent.

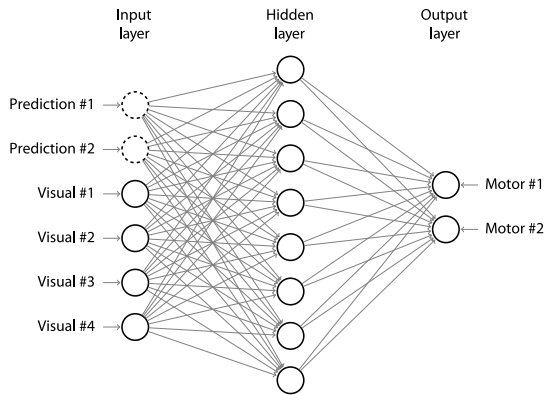


Figure 6.1 | Two-layer feedforward network architecture used. Six input units (two prediction units and four sensory units), eight hidden units, and two output units controlled the virtual robot arm.

6.2.3 Evolution of the network

Network weights were optimized using a neuroevolution algorithm using a direct encoding scheme (i.e. there was a one-to-one mapping of genotype to phenotype) similar to Nolfi et al. [109]. Although direct encoding schemes have been criticized for being biologically implausible [108], and having difficulties with scalability², direct encoding provided a good trade-off between simplicity and performance for the relatively simple networks used in this study. The initial population consisted of 100 networks with weights uniformly random $\in [-2.0, 2.0]$. For each subsequent generation, the twenty networks with the highest fitness value were allowed to reproduce by generating four copies each, with Gaussian noise sampled from $N(0, .3)$ added to the network weights. In addition, each of the twenty best networks was kept unmodified and added to the next generation, keeping the population size a constant 100. In pseu-

²The search space in direct encoding schemes increases exponentially with network size.

decode, the evolutionary algorithm was:

Algorithm 1: High-level description of neuroevolution algorithm

```

initialize 100 networks with random weights;
for 1000 generations do
  foreach network do
    | evaluate fitness;
  end
  sort networks by fitness;
  for 20 best networks do
    | generate 4 mutated copies;
    | generate 1 identical copy;
  end
end

```

All simulations were run 30 times per condition, so a total of 90 simulations were run.

6.3 Results

Maximum fitness of the networks differed between conditions, $F(2, 87) = 9.29$, $p < .001$, $\eta_G^2 = .176$. Post-hoc pairwise t -tests showed that networks with accurate predictions fed into the prediction neurons developed a higher maximum fitness ($M = 1868$) than networks with no prediction ($M = 1262$), $t(58) = 2.76$, $p = .008$, $d = .72$, and than networks with random prediction ($M = 947$), $t(58) = 4.12$, $p < .001$, $d = 1.08$. These differences remained significant after Holm–Bonferroni correction.

Figure 6.2 shows the evolution of fitness over time. Although the networks with no prediction evolved somewhat faster than networks with accurate prediction, maximum fitness leveled off after ~ 250 generations. For the networks with accurate prediction the network weights evolved slower, but surpassed the fitness of the non-predicting networks after 320 generations and continued to increase. Networks with random prediction evolved slower overall, and attained lowest maximum fitness.

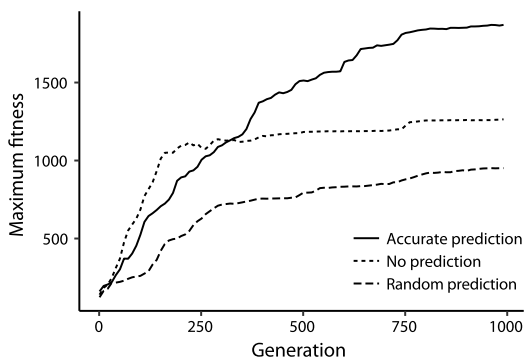


Figure 6.2 | Networks with accurate prediction attained higher maximum fitness than networks with no prediction or random prediction. These networks evolved to make efficient use of the information from the prediction neurons. Displayed are means over 30 simulations per condition.

Centering behavior differed between conditions, $F(2, 86) = 8.09, p < .001, \eta_G^2 = .158$. Post-hoc pairwise t -tests showed that networks with accurate prediction spent a smaller proportion of ITI time in the center 10×10 units ($M = .195$) than both networks with no prediction ($M = .415$), $t(57) = 4.64, p < .001, d = 1.23$, and networks with random prediction ($M = .340$), $t(58) = 2.96, p = .004, d = .778$. These differences remained significant after Holm–Bonferroni correction. The networks with no prediction and random prediction did not differ significantly, $p = .277$. Results are shown in Figure 6.3.

Movement across the environment is displayed in Figure 6.4. The networks with random prediction (Figure 6.4b) learned that the information provided was not informative, and reached their maximum fitness by returning to the center of the environment after touching each stimulus, whereas networks with accurate prediction (Figure 6.4a) moved toward the next target, waiting for it to appear.

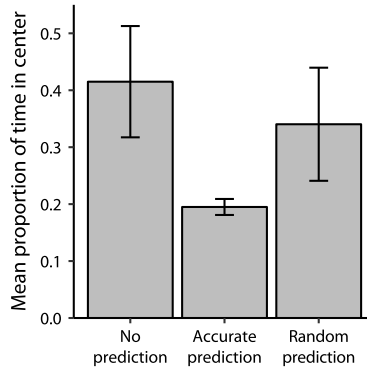
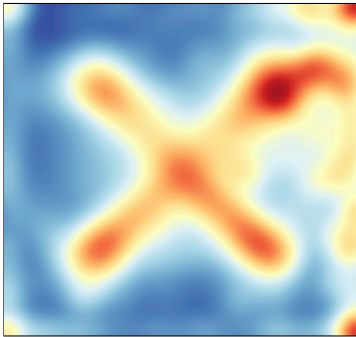
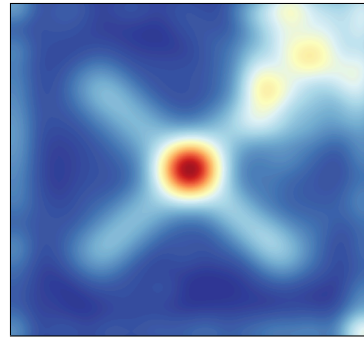


Figure 6.3 | The mean proportion of ITI time spent in the center of the screen for all three conditions. Networks with accurate prediction spent less time in the center. Error bars indicate 95% CI.



(a) In the condition with accurate prediction, position density is clustered around the stimuli, indicating active movement toward stimuli.



(b) With random prediction, the networks evolve to produce centering behavior. Most time is spent in a position equidistant to all targets.

Figure 6.4 | Density heat map showing the relative amount of time spent across locations in the accurate prediction and random prediction conditions, ranging from blue (little time spent) to red (most time spent).

6.4 Discussion

In this study, we investigated the behavior found in earlier work by Duran and Dale [38], Dale et al. [34], and the work described in Chapters 4 and 5. These studies describe a centering behavior in which participants moved their mouse to the center of the screen under some circumstances. In Chapter 5, we describe how this seems to be related to the quality of the action plan, or the capability to predict the next stimulus. This also makes sense on a theoretical level, as a centered position, equidistant to all possible stimuli is optimal under maximum uncertainty.

In the current study we evolved artificial neural networks that controlled a robotic arm, with a task analogous to the one used in Chapters 4 and 5. In one condition, an accurate prediction of the next stimulus was provided to the network as part of the input. In the second condition, the input given was randomly determined, and unrelated to the next stimulus. In a third condition, input to the prediction neurons was kept constant at zero. Under the last two conditions, centering behavior developed, with the networks that were provided random input and networks that were given no input developing the same centering strategy as human participants in Chapters 4 and 5 that had not developed explicit sequence knowledge. In summary, we showed that centering behavior evolved in a robotic arm controlled by an artificial neural network as a function of prediction quality, analogous to the findings described in Chapters 4 and 5.

Future research could shed light on the differences between the random prediction condition and the no prediction condition. From our results, it seems that performance was worse under the random prediction condition (although not significantly so), and developed more slowly. Apparently, the networks had trouble ignoring the dynamic, but uninformative input. In comparative studies with human participants, it would be interesting to distinguish between participants who *know* that they are unaware of the sequence (no prediction), and participants who are

actively, but unsuccessfully, trying to predict the sequence (random, or at least partly incorrect prediction).

