

逐次的行動選択のための大脳基底核ループにおける複数の表現系

中原 裕之[†]

銅谷 賢治[‡]

彦坂 興秀^{*}

永野 三郎[§]

[†] 理化学研究所国際フロンティア情報表現研究チーム 〒 351-01 埼玉県和光市広沢 2-1
nakahara@irl.riken.go.jp

[‡] 科学技術振興事業団 川人学習動態脳プロジェクト 〒 619-02 京都府相楽郡精華町 2 - 2
* 順天堂大学医学部第一生理, 〒 113, 文京区本郷 2-1-1

[§] 東京大学総合文化研究科広域科学専攻広域システム科学系, 〒 153, 目黒区駒場 3-8-1

近年、大脳基底核は、ドーパミン神経細胞によって与えられる強化信号を利用して、強化学習を行なっているという仮説が提案されている。大脳基底核-視床-大脳皮質を結ぶ複数のループが知られているが、その機能についてはあまり分かっていない。我々は、この論文で、複数のこれらのループが、複数の表現系を用いて、視覚運動の系列の学習を行なっていることを示す計算論的モデルを提出する。このモデルの中心的な考えは、視覚運動の系列は、空間座標系（例 視覚座標）の方が学習が容易であり、身体座標系（例 関節座標）の方が制御が容易であるという点にある。”2 x 5 課題”と呼ばれる近年の猿での学習実験の行動と神経生理のレベルでの結果を、このモデルは再現できた。

キーワード：大脳基底核、補足運動野、前頭前野、手続き記憶、強化学習、視覚運動学習

Multiple Representations in the Basal Ganglia Loops for Sequential Decision Making

Hiroyuki Nakahara[†]

Kenji Doya[‡]

Okihide Hikosaka^{*}

Saburo Nagano[§]

[†]Lab. For Info. Representation, Frontier Research Program,RIKEN 2-1 Hirosawa, Wako, Saitama 351-01
nakahara@irl.riken.go.jp

[‡] Kawato Dynamic Brain Project, JST, 2-2 Hikoridai, Seika, Soraku, Kyoto 619-02

^{*} Dept. of Physiology, Juntendo Univ., School of Med., 2-1-1 Hongo, Bunkyo, Tokyo 113

[§] Dept. of General Systems Studies, Univ. of Tokyo, 3-8-1 Komaba, Meguro, Tokyo 153

ABSTRACT

The basal ganglia (BG) have been hypothesized to perform reinforcement learning by use of reinforcement signals provided by dopamine neurons. It is well known that there exist multiple BG-thalamocortical loops, but their functions are poorly understood. Here, we propose a computational model of how different BG loops are employed in visuomotor sequence learning using different representations of sequence. The central idea of the model is that a visuomotor sequence is easier to *learn* in spatial representation (e.g. visual coordinates) but is easier to *control* in body-based representation (e.g. joint angle coordinates). The results of simulations of the model replicated both behavioral and neurophysiological findings in recent experimental studies using ”2x5 task”.

keywords: basal ganglia, supplementary motor area, prefrontal cortex, procedural memory, reinforcement learning, visuomotor learning

1 INTRODUCTION

Schultz and his colleagues showed in their experiments [15] that the response tuning of dopamine (DA) neurons in the substantia nigra pars compacta (SNc) in the basal ganglia (BG) shifts from primary reward to conditioned stimuli that predict reward as the conditioning establishes. This fact has led to a hypothesis that a specific form of reinforcement learning (RL), temporal difference (TD) learning [2], occurs in the BG with TD error of reward prediction provided by DA neurons as the reinforcement signal (e.g. [3,6,14]). This hypothesis provides a computational basis for investigating the BG function for sequential motor control, which has been suggested in studies on brain lesions and neuronal recordings.

The BG-thalamocortical circuit is one of the major cortico-subcortical circuits for motor control, with the other being the cerebellocortical circuit. A striking feature of the BG circuit is its separate, closed loop organization. The BG receives projections from almost the entire cerebral cortex but each part of BG projects to specific area of the frontal cortex. At least four BG loops have been identified, including motor, oculomotor and dorsolateral prefrontal loops [1]. Functions of some of these loops have been proposed, but it is very poorly understood how these BG loops work *together* in sequential motor control tasks.

In recent experiments of sequential arm reaching task called “2x5 task” [5], Hikosaka and his colleagues have found that different parts of the BG and the frontal cortex are involved differentially in acquisition and execution of sequential movement [10,12]. Motivated by their results, we propose a computational model of how different BG-cortical loops are involved in different stages of learning of sequential movement. The key idea is that a visuomotor sequence like that in the 2x5 task is easier to *learn* in spatial representation (e.g. visual coordinates) but is easier to *control* in body-based representa-

tion (e.g. joint angle coordinates) [13]. Specifically, we propose that the dorsolateral prefrontal (DLPF) loop learns a sequence in spatial representation and the supplementary motor area loop learns a sequence in body-based representation. Both loops learn concurrently using the reinforcement signal carried by DA neurons, but relative ease of learning and effectiveness in control causes the differential involvement of these two loops. Simulation of the proposed model replicated both behavioral and neurophysiological findings of the 2x5 task experiments.

2 2x5 TASK AND ITS FINDING

Figure 1 shows an example of the sequence of events in a single trial of the 2x5 task. When the animal pressed the home key at the start of a trial, two out of the 16 LED buttons were turned on simultaneously, which is called a ‘set’ of stimulus. The animal had to press the illuminated buttons in a predetermined order, which he had to find out by trial-and-error. If successful, another pair of LEDs, the second set, was illuminated which the monkey had to press again in a predetermined order. A fixed sequence of 5 sets, called a ‘hyperset,’ was presented in a trial. When the animal pressed a wrong button, all LED buttons were illuminated briefly with an unpleasant beep sound, and the trial was aborted without any reward. The animal then had to start over a new trial by pressing the home key. After each successful set, the animal was given a liquid reward.

The same hyperset was used throughout a ‘block’ of experiments until completing a certain number of successful trials (criterion). On each day during a training period, the monkey performed several blocks of trials with different hypersets. Some hypersets were used every day, called ‘learned’ hypersets, and others were randomly generated and used only once, called ‘new’ hypersets. Order of presentation of

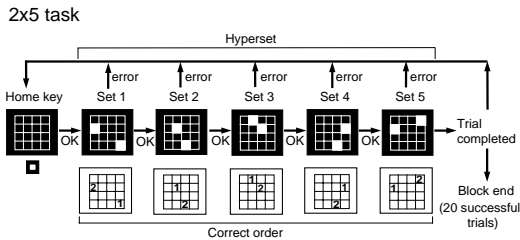


Figure 1: Procedure of 2x5 task with an example of a hyperset . To complete a trial, a monkey has to press 10 buttons (2 buttons x 5 sets) in a correct (predetermined) order.

learned and new hypersets are randomized everyday.

Learning in 2x5 tasks is measured by the decrease in the number of trials to criterion and/or the decrease in the performance time. Hikosaka et al. [5] observed both *short-term* learning and *long-term* learning. Short-term learning is indicated by improved performance during a block of experiment and long-term learning is indicated by improved performance across days (See Figure 2).

Functional differentiation was observed in blockade experiment by muscimol injection between anterior BG (caudate head) and posterior BG (putamen) [10], as well as the presupplementary motor area (pre-SMA) and the supplementary motor area (SMA) [12]. Number of error trials to criterion was significantly increased by blockade of the anterior BG or the pre-SMA for new hypersets but not for learned ones. Number of error trials to criterion was significantly increased by blockade of the posterior BG for learned hypersets but not for new ones. Blockade of the SMA affected both learned and new ones, but only mildly. In brief, these lines of experimental evidence suggest that the anterior BG and the pre-SMA is more involved in early acquisition process of sequences, and that the posterior BG and, possibly, the SMA is more involved in maintenance and retrieval of acquired sequential memory.

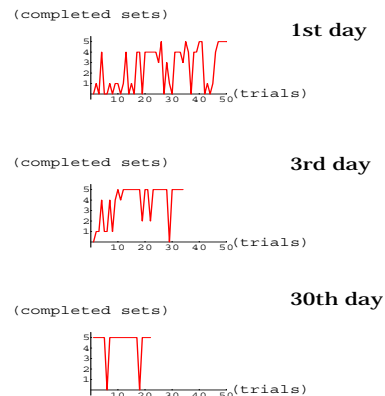


Figure 2: Experimental result of learning a hyperset across days. The change in the number of completed sets (ordinate) across trials (abscissa) is compared among the 1st day (top), the 3rd day (middle) and the 30th day (bottom). Taken from Hikosaka et al. [5]

3 HYPOTHESIS ON MULTIPLE REPRESENTATIONS IN THE BASAL GANGLIA LOOPS

In visuomotor task such as visually-guided reaching, problems in inverse kinematics and inverse dynamics must be solved to reach a target based on visual information [7]. We propose that it is easier to learn a visuomotor sequence in visual coordinates particularly when the sequence is learned by trial-and-error, whereas it is faster and easier to execute the sequence in body-based coordinates once it is acquired. We further propose that it is advantageous to have multiple learning processes with different speeds concurrently. A quick acquisition of a sequence in visual coordinates would help the animals solve new problems that they encounter in everyday life; a slow acquisition of a sequence in body coordinates would help the animals store the memory of the frequently-used sequence robustly and retrieve it quickly. Thus, it is plausi-

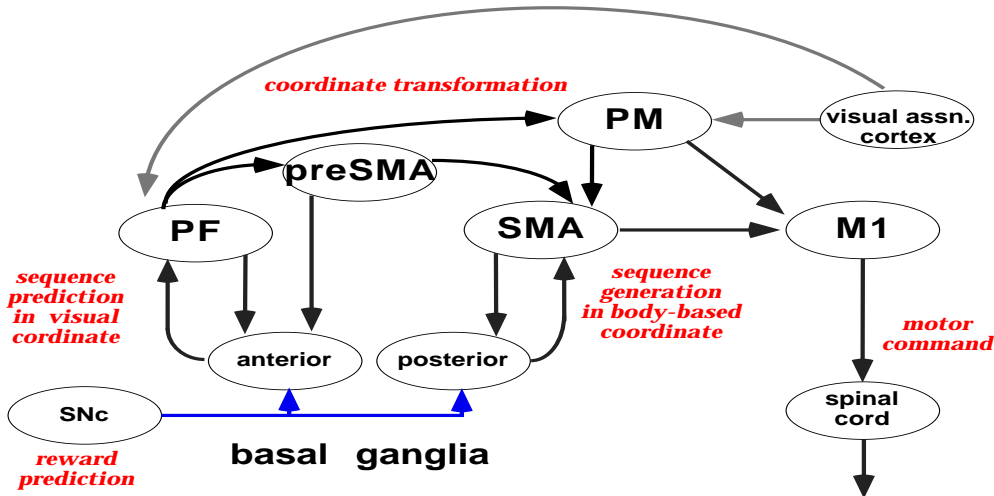


Figure 3: Hypothesized scheme of the basal ganglia-thalamocortical loops for sequential motor control.

ble to consider from a computational viewpoint that a sequence in visual coordinates is more suitable for early acquisition stage, whereas a sequence in body-based coordinates is more suitable for execution with robust maintenance and quick retrieval of its memory.

The dorsolateral prefrontal loop: the projection from DLPF to the BG is primarily to anterior striatum, including head of the caudate nuclei (CD) and and the rostral putamen. The DLPF is well known to be involved in visuospatial memory and is considered to play a role in control and planning of sequential movements [4]. Posterior parietal cortex that is connected with the DLPF also projects to the dorsolateral head of the CD [1]. Based on these facts, it is likely that the DLPF loop, including the DLPF and the anterior striatum, engages computation in visual coordinates (Figure 3).

The motor loop: in motor loop, most of the projections to the BG originate from the cortical motor areas, among which the SMA is of particular interest, and principally terminate in the posterior striatum (the bulk of putamen) [1]. The SMA is well connected with other motor cortical areas and has been known long as involved in sequential movements, particularly internally-generated complex ones [17]. Thus,

the motor loop, including the SMA and the posterior striatum, may engage computation in body-based coordinates (Figure 3).

It is noteworthy that not the SMA but the pre-SMA is connected with the DLPF [17]. The pre-SMA has the projections to the anterior striatum. Hence, the pre-SMA is heavily interacted with the DLPF loop. It is also experimentally shown that the pre-SMA is generally more activated in the period after receiving sensory inputs and before starting movements [9, 17].

We hypothesize that the DLPF loop, perhaps together with the pre-SMA, learns the sequence using visual coordinates. The DLPF loop, hence, is more critical in early acquisition stage of sequences. In contrast, the motor loop learns the sequence using body-based coordinates. The motor loop, hence, is more involved in execution of well-learned sequences. A unique feature of our model is that both loops concurrently learn sequences based on reinforcement signal provided by DA neurons (Figure 3).

4 SIMULATION ON 2x5 TASK

In order to test the behavior of a model based on our hypothesis, we built a neural network

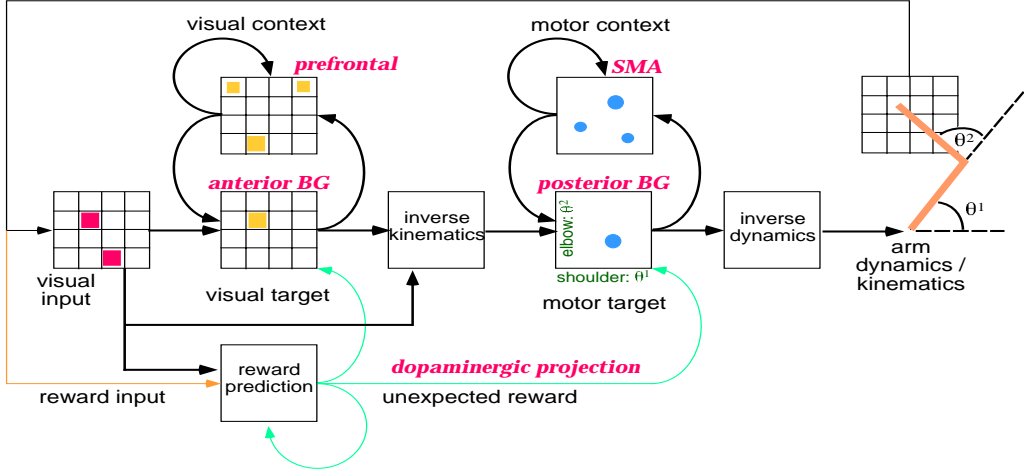


Figure 4: Diagram of the proposed model in context of the 2x5 task

of 2x5 task learning. The specific questions we asked were whether the model would have the short-term and long-term learning of performance and whether the model could replicate the results of blockade experiments for the DLPF loop.

4.1 Network Architecture

Figure 4 illustrates the overall structure of the network model, which consisted of the visual network (DLPF-BG loop), the motor network (SMA-BG loop), the critic network (DA system), and inverse kinematics modules.

Both the visual and motor networks had the same structure and had an output layer and a context layer, each corresponding to the BG and the cortex. First, the state y_i of the output layer was updated by the softmax function of the weighted sum of the input u_j and the context x_j with a tuning parameter for sharpness of the softmax, denoted by ζ .

$$s_i(t) = \sum_j w_{ij}^I u_j(t) + w_{ij}^C x_j(t) \quad (1)$$

$$y_i(t) = \frac{\exp[\zeta s_i(t)]}{\sum_k \exp[\zeta s_k(t)]} \quad (2)$$

Then, one of the output units was stochastically set as $z_i(t) = 1$ with a probability

$$\text{Prob}(z_i(t) = 1) = y_i(t). \quad (3)$$

The context layer was updated, with a time constant, τ , by the following equation.

$$x_i(t+1) = x_i(t) + \frac{1}{\tau}(z_i(t) - x_i(t)) \quad (4)$$

The weight matrices w^I and w^C were updated by a Hebbian rule, weighted by the reinforcement signal $\hat{r}(t)$, as defined below;

$$w_{ij}^I(t+1) = w_{ij}^I(t) + \eta^I \hat{r}(t) z_i(t) u_j(t) \quad (5)$$

$$w_{ij}^C(t+1) = w_{ij}^C(t) + \eta^C \hat{r}(t) z_i(t) x_j(t) \quad (6)$$

Initially, the input weights were set by an identity matrix so that the default behavior is to press one of the lit buttons.

In the visual network, the state was encoded by the 16 units corresponding to the positions of 16 buttons in the Cartesian space. The input u^V to the visual network was a 16 dimensional vector of 0 and 1 encoding whether the corresponding buttons were lit. Its output z^V represented the button to be pressed, or the movement target in the visual space.

In the motor network, state was encoded by a population vector of joint angles. Each unit had a Gaussian softmax activation function

$$a_j(\theta) = \exp\left[-\frac{1}{2} \sum_{i=1}^2 \left(\frac{\theta^i - \theta_j^i}{\delta^i}\right)^2\right], \quad (7)$$

$$b_j(\theta) = \frac{a_j(\theta)}{\sum_k a_k(\theta)} \quad (8)$$

where θ^1 and θ^2 denote shoulder and elbow angles and θ_j^i denotes the preferred joint angle for the j -th unit.

An inverse kinematic model, which was derived analytically from the geometry of the arm, was used to transform the visual representation into corresponding motor representation. In the normal operation, the input to the motor network was a sum of two activation vectors, one calculated from the current visual input and the other from the output of the visual network.

The critic network provides the estimated value function given the current state, or the current sensory input, denoted by u^R . Its encoding is the same as the encoding of the state in the visual network. The estimated value of a state, $P(t)$, is defined by

$$P(t) = \sum_j w_j^R u_j^R + b^R \quad (9)$$

where w_j^R and b^R are the weight matrix and bias for the critic. Using this estimated value function, the reinforcement signal, or TD error, $\hat{r}(t)$, is computed by

$$\hat{r}(t) = r(t) + \gamma P(t+1) - P(t) \quad (10)$$

The critic weight matrix, w^R , and bias, b^R , is updated by use of TD error as defined below:

$$w_j^R(t+1) = w_j^R(t) + \eta^R \hat{r}(t) u_j^R \quad (11)$$

$$b^R(t+1) = b^R(t) + \eta^b \hat{r}(t) \quad (12)$$

Thus, note that we used TD(0) learning in simulation we report below for sake of simplicity and that the overall architecture of the network above is based on actor-critic scheme [2].

4.2 Simulation Setup

We let the model learn 2 learned hypersets and 1 new hyperset per a simulated day, approximately keeping ratio between learned and new hypersets in experiment. Simulation is run for 10 days to train the model networks. Learning process during this period is examined to

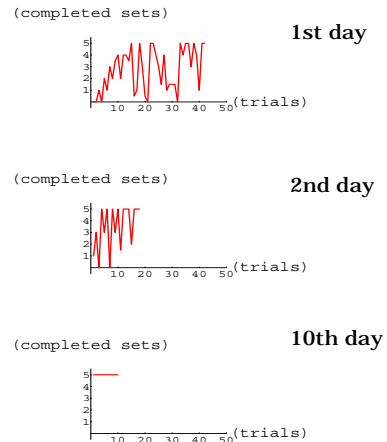


Figure 5: Simulation result of learning a hyperset across days. The change in the number of completed sets (ordinate) across trials (abscissa) is compared among the 1st day (top), the 2nd day (middle) and the 10th day (bottom).

test whether the model exhibits different learning levels. Then, using the trained network parameters, we tested the performance of the network for learned and new hypersets in case of the blockade of the DLPF loop, or the visual network. In simulation, the blockade of the visual network was realized by inhibiting input from the visual network to the motor network.

4.3 Results

4.3.1 Different learning levels

Hikosaka et al. [5] indicated *short-term* learning by improved performance during a block and *long-term* learning by improved performance across days. It is clear in Figure 6 (right) that the model makes more errors in the first half of total trials than in the second half, in particular for the 1st and 2nd days. This indicates short-term learning (See also Figure 5). Figure 6 (left) shows that the model improved its performance for learned hypersets across days particularly for first few days, indicating long-term learning as observed in the experiments (See also

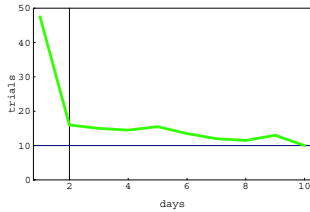
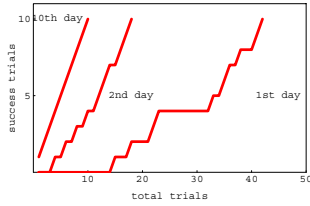


Figure 6: Performance of the proposed model: (left) example of learning a learned hyperset across days by the model. The number of successful trials is plotted against the total number of trials for the 1st, 2nd, and 10th days. (right) the averaged performance for learned hypersets across 10 days. The mean number of trials to criterion is plotted across days.

Figure 5) [5].

4.3.2 Blockade of the anterior basal ganglia and the presupplementary motor area

It is observed that the performance for the learned hyperset was not deteriorated by the blockade of the visual network (Table 1), similarly to the results of the blockade experiments in the 2x5 task [10, 12]. In the present network architecture, it was observed in training period that the network could not learn some types of new hypersets at all and the same phenomena was sometimes observed in the blockade condition as well, which will be discussed in the next section.

Table 1: Results on the blockade of visual network: number of trials to criterion for learned hypersets

	NORMAL	BLOCKADE
Learned	11.4(± 1.07)	11.7(± 0.95)

5 DISCUSSION

We have hypothesized that the BG contribute to sequential motor control based on RL with multiple representations. The DLPF loop uses visual coordinates and contributes to learning sequences in early stage. The motor loop uses body-based coordinates and provides robustness with the sequential memory once acquired, even though the motor loop learns them slowly. The results of the model based on this hypothesis resembled some of experimental results in the 2x5 task: different levels of learning at behavior level and functional differentiation at neurophysiological level.

As noted above, it was observed in the present simulation that acquired sequential memory severely interfered learning for some types of new hypersets in training period as well as in blockade condition. This is probably because, in the present architecture, the output from the visual network is forwarded only as an input to the motor network so that the output from the visual network cannot directly influence on hand movements. We suggest that the present architecture should be extended to include explicit mechanism to let the output from the visual network directly influence on producing a final output and, for this purpose, to weight, or gate, outputs from both of the visual and motor networks. The recent study on neural activities in the pre-SMA suggests that such mechanism may exist and that the pre-SMA may be a part of it [11, 16]. This point is currently investigated. In addition, the functions of the cerebellum was not explicitly addressed in the present study, however, there are experimental results on the

cerebellum in the 2x5 task, indicating that the dentate nucleus is involved in maintenance of sequential memory [8]. Integrating the functions of the cerebellum into the proposed hypothesis is a fruitful future research.

Acknowledgments

Supported by JSPS and SPRF of RIKEN to HN, Uehara Memorial Foundation and JSPS Research for the Future Program to OH. HN acknowledges many insights from S. Miyachi, K. Miyashita, S. Sakai, and X. Lu in development of this study.

References

- [1] G. E. Alexander, M. D. Crutcher, and M. R. DeLong. Basal ganglia-thalamocortical circuits: Parallel substrates for motor, oculomotor, “prefrontal” and “limbic” functions. In H. Uylings, C. Van Eden, J. De Bruin, M. Corner, and M. Feenstra, editors, *Progress in Brain Research*, volume 85, chapter 6, pages 119–146. Elsevier Science Publishers B.V., 1990.
- [2] A. Barto, R. Sutton, and C. Anderson. Neuron-like adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics*, 13:834–846, 1983.
- [3] K. Doya. Efficient nonlinear control with actor-tutor architecture. In D. Touretzky, M. Mozer, and M. Hasselmo, editors, *Advances in Neural Information Processing Systems*, volume 9, 1997.
- [4] P. S. Goldman-Rakic. Circuitry of primate prefrontal cortex and regulation of behavior by representational memory. In F. Plum and V. Mountcastle, editors, *Handbook of Physiology - The Nervous System V*, volume 5, chapter 9, pages 373–417. 1987.
- [5] O. Hikosaka, M. Kato, S. Miyachi, and K. Miyashita. Learning of sequential movements in the monkey: Process of learning and retention of memory. *Journal of Neurophysiology*, 74(4):1652–1661, 1995.
- [6] J. C. Houk, J. L. Adams, and A. G. Barto. A model of how the basal ganglia generate and use neural signals that predict reinforcement. In *Models of Information Processing in the Basal Ganglia*, pages 249–270. MIT Press, Cambridge, Massachusetts, 1995.
- [7] M. Kawato. Computational schemes and neural network models for formation and control of multijoint arm trajectory. In *Neural Networks for Control*. MIT Press, Cambridge, MA, 1990.
- [8] X. Lu, O. Hikosaka, and S. Miyachi. Role of monkey cerebellar dentate nucleus in procedural memory. In *Proceedings of 26th Annual Meeting, Society for Neuroscience.*, page 546.13, 1996.
- [9] Y. Matsuzaka, H. Aizawa, and J. Tanji. A motor area rostral to the supplementary motor area (presupplementary motor area) in the monkey: neuronal activity during a learned motor task. *Journal of Neurophysiology*, 68(3):653–662, 1992.
- [10] S. Miyachi, O. Hikosaka, K. Miyashita, Z. Karadi, and M. Kato. Different roles of monkey striatum in learning of sequential movements. *Experimental Brain Research*, in press, 1996.
- [11] K. Miyashita and O. Hikosaka. Activity of pre-SMA neuron for the performance of sequential movement in monkeys. In *Proceedings of 27th Annual Meeting, Society for Neuroscience.*, (submitted) 1997.
- [12] K. Miyashita, K. Sakai, and O. Hikosaka. Effects of SMA and pre-SMA inactivation on learning of sequential movements in monkey. In *Proceedings of 26th Annual Meeting, Society for Neuroscience.*, page 731.3, 1996.
- [13] H. Nakahara. *Sequential Decision Making in Biological Systems: The Role of Nonlinear Dynamical Phenomena in Working Memory and Reinforcement Learning in Long-Term Memory*. PhD thesis, Univ. of Tokyo, 1997.
- [14] W. Schultz, P. Dayan, and R. Montague. A neural substrate of prediction and reward. *Science*, 275:1593–1599, 1997.
- [15] W. Schultz, R. Romo, T. Ljungberg, J. Mirenowicz, J. R. Hollerman, and A. Dickinson. Rewarded-related signals carried by dopamine neurons. In *Models of Information Processing in the Basal Ganglia*, pages 231–248. MIT Press, Cambridge, Massachusetts, 1995.
- [16] K. Shima, H. Mushiake, N. Saito, and J. Tanji. Role for cells in the presupplementary motor area in updating motor plans. *Proceedings of the National Academy of Sciences of the United States of America*, 93(16):8694–8, Aug 6 1996.
- [17] J. Tanji. The supplementary motor area in the cerebral cortex. *Neuroscience Research*, 19(3):251–68, May 1994.