Published in Psychological Research (2009) Vol. 73, No. 4, pp. 527-44

Brain mechanisms for predictive control by switching internal models: implications for higher-order cognitive functions

Hiroshi Imamizu^{1, 2} and Mitsuo Kawato²

5

15

20

¹ Biological Information and Communications Technology Group, National Institute of Information and Communications Technology, 2-2-2, Hikaridai, Keihanna Science City, Kyoto 6190288, Japan.

² Computational Neuroscience Laboratories, Advanced Telecommunications Research Institutes International, 2-2-2, Hikaridai, Keihanna Science City, Kyoto 6190288, Japan.

Note: Supplemental movie is temporally located on our web server

(http://www.cns.atr.jp/~imamizu/multi_functions.mpg) only for review purposes

Corresponding author: Hiroshi Imamizu

National Institute of Information and Communications Technology 2-2-2 Hikaridai, Keihanna Science City, Kyoto 619-0288, Japan Telephone: +81 774 95 1220, Fax: +81 774 95 1236 e-mail: imamizu@gmail.com

Abstract

5

10

15

Humans can guide their actions toward the realization of their intentions. Flexible, rapid and precise realization of intentions and goals relies on the brain learning to control its actions on external objects and to predict the consequences of this control. Neural mechanisms that mimic the input-output properties of our own body and other objects can be used to support prediction and control, and such mechanisms are called internal models. We first summarize functional neuroimaging, behavioral and computational studies of the brain mechanisms related to acquisition, modular organization, and the predictive switching of internal models mainly for tool use. These mechanisms support predictive control and flexible switching of intentional actions. We then review recent studies demonstrating that internal models are crucial for the execution of not only immediate actions but also higher-order cognitive functions, including optimization of behaviors toward long-term goals, social interactions based on prediction of others' actions and mental states, and language processing. These studies suggest that a concept of internal models can consistently explain the neural mechanisms and computational principles needed for

fundamental sensorimotor functions as well as higher-order cognitive functions.

Introduction

10

Two opposing approaches have been proposed to understanding higher-order cognitive functions, such as tool use, social interaction, and language, that are generally specific to humans. One approach hypothesizes that these functions are unique faculties of humans, which should be investigated independently of studies on the cognitive functions of related faculties in other species, such as non-human primates. This approach is often taken in studies on language (e.g. Hauser, Chomsky, & Fitch, 2002). In contrast, the other approach hypothesizes continuity between the cognitive functions of humans and those of other species from an evolutionary point of view, where the "human" faculties are based on computational principles in common with those of other species. For example, humans and non-human primates share common computational principles in fundamental sensorimotor control (e.g. reaching and grasping). The latter approach attempts to explore how these principles can also form the basis for higher-order cognitive functions.

The current special issue of *Psychological Research* focuses on goal-directedness, mirror systems and internal models, each of which are closely related to common computational principles between basic sensorimotor functions and higher-order cognitive functions, as well as between humans and non-human primates. In particular, mirror systems and internal models have been considered key concepts for elucidating enigmas in neural mechanisms that support the human abilities of social interactions and long-term planning of behaviors based on predictions. This article reviews our studies on internal models while explaining the relationships between internal models and goal-directedness, or mirror systems, and discusses how our studies on internal models could be extended to understanding higher-order cognitive functions.

Internal models are promising concepts for explaining neural mechanisms and computational principles supporting the flexible abilities of prediction and learning in cognitive functions. Abilities to learn relationships between actions and resultant changes in states of external objects are particularly important for planning of goal-directed behaviors. Such abilities are largely dependent on neural mechanisms that can model or simulate the relationships between an action and its consequences before the action's execution. For example, skilled manipulation of a computer mouse requires the ability to predict how the mouse should be moved in order to move a cursor to a particular position on the screen (predictive control: Fig. 1A) and how the cursor will move on the screen if the mouse is moved in a particular direction (prediction of feedback: Fig. 1B). Neural mechanisms that mimic the input-output properties of controlled objects can support the predictive control and prediction of sensorimotor feedback, and these mechanisms are called internal models (Kawato, 1999; Kawato, Furukawa, & Suzuki, 1987; Wolpert, Ghahramani, & Jordan, 1995). Although the concept of internal models was developed in motor neuroscience, many studies have suggested that it can be extended to explain the fundamental computational principles of higher-order cognitive functions, such as goal-directed

10

15

20

In the first part of this article, we review studies on internal models in the context of sensorimotor learning and use of tools. These studies revealed brain mechanisms related to acquisition, modular organization, and switching of internal models. Next, we summarize studies suggesting that internal models contribute to the cognitive functions discussed above. We believe that it is important to understand how the acquisition, modularity and switching of internal models contribute to the cognitive functions and that such an understanding can lead to the construction of a computational framework, which can

behaviors, mirror systems, social interactions, communication, and languages.

-4-

consistently explain the neural basis for fundamental as well as higher-order cognitive functions such as sensorimotor control, tool use, social interaction, and language.

Cerebellar activity related to an internal model

5

10

15

20

To investigate the acquisition process of internal models in the human brain, we measured brain activity using functional magnetic resonance imaging (fMRI) when human subjects learned to use a novel tool (Imamizu, et al., 2000). Based on previous neurophysiological and computational studies (e.g. Ito, 1984; Kawato, et al., 1987; Kitazawa, Kimura, & Yin, 1998; Shidara, Kawano, Gomi, & Kawato, 1993), we focused on the cerebellum and conducted detailed analysis of changes in cerebellar activity during learning. Subjects manipulated a computer mouse in a magnetic resonance (MR) scanner so that the corresponding cursor followed a randomly moving target on a screen (tracking task). In test periods, the cursor appeared in a position rotated 120° around the center of the screen to necessitate subject learning (novel mouse; Fig. 2), while in baseline periods it was not rotated (normal mouse). Each subject's performance was measured by tracking errors, i.e., the distance between the cursor and the target.

The errors in the test periods significantly decreased as the number of sessions increased (Fig. 3A), suggesting that learning progressed. When we investigated cerebellar activity that significantly and positively correlated with tracking error, we identified a strong correlation ($r^2 = 0.82$) between activity and error in the large part of the lateral cerebellum (white regions in Fig. 3C), suggesting that most of the activity in the cerebellum reflects the error. However, in a further experiment, we found that activity in some parts of the cerebellum was not explained by the error. In this experiment, we increased target velocity in baseline periods so that the errors there were equalized to the error in the test period (Fig. 3B). Then, we subtracted the activity in the baseline period from that in the test period. We could still find significant activity in the hatched regions in Fig. 3C, suggesting that the activity in theses regions cannot be explained solely by error.

We investigated how activity in the white regions and that in the hatched regions in Fig. 3C changed during training sessions. Activity in the white regions drastically decreased, as shown in the middle panel of Fig. 3D, while activity in the hatched regions did not markedly decrease as shown in the left panel of Fig. 3D. This suggests that activity in the hatched region includes activity that cannot be explained by the error. By subtracting the middle curve from the left curve, we found that the activity unrelated to the error increased at the beginning and remained high during the training sessions. This activity was thought to reflect the acquired internal model representing the input-output property of the novel mouse.

10

Although change in activity with learning was not investigated, an fMRI study indicated that the lateral cerebellum contributes to an internal model of a complex dynamics (Milner, Franklin, Imamizu, & Kawato, 2007). Subjects manipulated an object with a complex dynamics (balancing an inverted pendulum created by attaching weights to a flexible ruler) in a complex condition, and they manipulated an object with a simple dynamics (squeezing a soft foam ball) in a simple condition. Muscle activation was

precisely matched between these conditions. Consequently, a significant difference in activity between the complex and simple conditions was found only in the lateral

-6-

cerebellum among regions where a significant increase in activity was found by comparing the complex (or the simple) condition with a rest condition.

In the above studies, we focused on changes in cerebellar activity based on previous neurophysiological and computational studies. Many studies have reported changes in whole-brain activity with sensorimotor learning when a force field alters limb dynamics (Shadmehr & Holcomb, 1997) or when a screen controlled by a computer program kinematically alters visual feedback of hand position (Krakauer, et al., 2004) or joystick position (Graydon, Friston, Thomas, Brooks, & Menon, 2005). These studies identified different cerebral regions related to learning as a consequence of the differences in experimental methods (e.g. adaptation to a force field or altered visual feedback; different effectors such as the arm, a computer mouse and a joystick; tracking a continuously moving target or aiming at a static target). However, their results are consistent with ours in that significant activity is found in the lateral cerebellum after learning. These results suggest that the cerebellum is one of the regions where internal models representing input-output properties of controlled objects are most likely acquired.

Forward and inverse internal models

5

10

15

20

It is thought that the central nervous system (CNS) uses two forms of internal models. Inverse models transform intended actions or goals into the motor commands to reach those goals (Kawato, et al., 1987; Fig. 1A). Forward models transform efference copies of motor command into the resultant trajectory or sensorimotor feedback (Kawato, et al., 1987; Miall, Weir, Wolpert, & Stein, 1993; Wolpert, et al., 1995; Fig. 1B). The above imaging studies

-7-

investigating neural correlates of internal models do not take these two forms into account. Because both forward and inverse models are thought to be necessary for rapid and smooth movements, the above brain activity probably reflected activities of both models. Neurophysiological studies have shown data indicating that Purkinje cells in the cerebellar cortex contribute to inverse models of motor systems (Gomi, et al., 1998; Shidara, et al., 1993). Many functional neuroimaging studies have shown data circumstantially indicating that the cerebellum contributes to forward models as described below.

Miall and colleagues investigated brain activity related to eye-hand coodrdination using a tracking task (Miall, Reckess, & Imamizu, 2001). Subjects followed a moving target with their eyes while simultaneously moving a joystick to control the cursor. The 10 temporal offset between targets for eye and hand motions caused parametric variation of the degree of eye-hand coordination. The behavioral data indicated that manual tracking performance was optimal when the target for eye motion anticipated the target for hand motion by 38 ms. Synchronous movements of two effectors with such a small offset cannot be achieved simply by reaction to reafferents or visual input. This suggests that a forward 15 model predicts the movement outcome based on a motor command and that the predicted outcome is sent to the oculomotor system for programming or modifying the manual movements. The fMRI data found a parametric increase in activity of the lateral cerebellum and the oculomotor vermis as eye-hand coordination increased, suggesting a contribution of the cerebellum to prediction of the movement outcome.

20

Behavioral studies on grip force-load force coupling have found convincing evidence that the CNS makes use of forward models in sensory motor control. When an object is held in a precision grip (e.g. a grasp between the tips of the thumb and forefinger) and

moved by voluntary movements (e.g. arm movements), the grip force perpendicular to the contact surface changes in phase with the load force induced by the movements (Johansson & Westling, 1988). The coupling between the two forces prevents the object from slipping while using minimal grip force. This grip force modulation is anticipatory in the sense that changes in the grip force occur at the same time as, or even prior to, changes in the load force. Based on theoretical analysis of behavioral data (Flanagan & Wing, 1997) suggesting that output signals from a forward model of arm movements are used for control of grip force, an fMRI experiment examined brain activity related to coordination of grip force and load force (Kawato, et al., 2003). The results indicated that parts of the anterior lobule in the cerebellum contribute to the coordination.

10

15

Other studies have indicated that the cerebellar forward models contribute to prediction of sensorimotor feedback in various situations, such as cancellation of tactile sensation during self-tickling (Blakemore, Frith, & Wolpert, 2001; Blakemore, Wolpert, & Frith, 1998) and state-dependent control of arm and finger movements (Diedrichsen, Criscimagna-Hemminger, & Shadmehr, 2007). It has been suggested that the cerebellum contributes to a prediction of change in the state of external objects that is not caused by its own motor commands (O'Reilly, Mesulam, & Nobre, 2008). Regarding neurophysiological studies, Miall and colleagues proposed that simple spike activity of Purkinje cells represents prediction of sensory feedback and is corrected by complex spike activity

²⁰ representing a discrepancy between the prediction and actual feedback (Miall, et al., 1993). They found, as supporting evidence, that the interval between an increase in simple spike activity and the resulting complex spike activity is about 150 ms, which is equivalent to visuomotor feedback delay and necessary for synchronizing the prediction and feedback (Miall, Keating, Malkmus, & Thach, 1998). Recently, it has been suggested that simple

-9-

spike discharge of Purkinje cells has several characteristics of a forward internal model of the arm (Ebner & Pasalar, 2008).

These functional imaging and neurophysiological studies suggest that the cerebellum is related to both forward and inverse internal models, but it is unknown how these two forms of internal models are organized in the cerebellum.

Modular organization of internal models

5

10

Humans interact with myriad objects and environments that often change in a discrete manner. If the CNS maintains only a small number of global internal models, relearning is needed whenever manipulated objects and environments change. However, if the CNS maintains a large number of internal models or modules for different objects and environments, less relearning is needed and thus learning interference is avoided. Moreover, initial learning of objects and environments may be facilitated by a combination of stored modules.

Many lines of behavioral studies have shown the multiplicity and modularity of internal models. For example, it has been demonstrated that humans can independently learn dynamic properties of their own arms altered by weights and kinematic properties altered by the rotation of visual feedback of their hand position (Krakauer, Ghilardi, & Ghez, 1999). This result suggests that some types of internal models are independently acquired and do not interfere with each other. Ghahramani and colleagues indicated that the CNS can appropriately combine output signals from stored internal models for different

-10-

sensorimotor mappings (Ghahramani & Wolpert, 1997). Flanagan and colleagues made subjects learn a kinematic transformation (visuomotor rotation), a dynamic transformation (force field), and a combination of these transformations (Flanagan, et al., 1999). When the subjects learned the combined transformation, reaching errors were smaller if the subject first learned the separate kinematic and dynamic transformations. These results suggest the ability of subjects to combine internal models as needed, depending on the situation.

In functional imaging studies, we investigated cerebellar activity after subjects learned to use a velocity-control mouse in which cursor velocity was proportional to the mouse position. Here, we examined the difference in activated regions between when subjects used the velocity-control mouse and when they used the rotated mouse (see above). By subtraction of activity when subjects manipulated the normal mouse (baseline condition) from activity when subjects used the rotated or velocity-control mouse (test conditions), we derived a map specific to each type of mouse. Figure 4 shows three-dimensional displays of the maps. Similar regions in the lateral cerebellum were activated, but the rotated-mouse activations (yellow) tend to be located more anteriorly and laterally than the velocity-control mouse activations (blue). The different tools evoked activities in distinct locations with small overlap (2.1% of the total activated volume), demonstrating the modularity and multiplicity of internal models for tools.

20

10

15

Higuchi and colleagues measured cerebellar activity when subjects used sixteen
common tools (scissors, a hammer, chopsticks and so on) and when they mentally imagined
using the tools without actual hand movements (Higuchi, Imamizu, & Kawato, 2007).
Figure 5 shows t-value-weighted centroids of activation when subjects actually used
individual tools (Fig. 5A) or when they imagined the tools' use (Fig. 5B) in comparison to

rest condition. Activities during the actual use tend to be located in the anterior lobule of the cerebellum. In contrast, activities during the imaginary use tend to be located more laterally, in the posterior lobule, than those during the actual use. We measured the distance of the centroids from the fourth ventricle. Because the fourth ventricle is the most anterior and medial part of the cerebellum, the longer the distance is, the more posteriorly and laterally the centroid is located (Fig. 5C). As shown in Figure 5E, the mean distance across the tools during the imaginary use was significantly longer than that during actual use (t(28))= 2.66, P < 0.05), suggesting that activities during the imaginary use were located posteriorly and laterally. Figure 5E shows lines connecting the centroids during the actual use (rectangles) with those during the imaginary use (circles). Regarding the tools that evoked activities in the posterior lobule during the imaginary use (thick circles in Fig 5B), the lines are often orthogonal to the primary fissure between the anterior and the posterior lobules (thick lines in Fig. 5E). This suggests that the order of the centroids along the primary fissure for the tools is almost identical between the anterior and posterior lobules. Activities in the anterior lobule are probably evoked by activities of limb muscles and sensory feedbacks (Grodd, Hulsmann, Lotze, Wildgruber, & Erb, 2001), while activities in the posterior lobule may reflect internal models for use of the tools. This result suggests that internal models contributing to skillful use of common tools are modularly organized, that is, different parts of the lateral cerebellum contribute to the use of different tools.

10

15

Our functional imaging study suggests the ability of the CNS to combine output signals from internal models (Imamizu, Higuchi, Toda, & Kawato, 2007). Subjects sufficiently learned to use 60° and 160° rotated mice, in each of which the cursor appeared in a position rotated 60° or 160° around the center of the screen. Then we investigated brain activity when subjects learned to use a 110° rotated mouse (an intermediate angle between

-12-

 60° and 160°). In the early and late stages of learning the 110° mouse, we measured cerebellar activity specific to the 60° , 110° or 160° mouse according to the same method we used when measuring activity specific to the rotated mouse or the velocity-control mouse (see above). In the early stage of learning, activated volumes for 60° , 110° and 160° were 7.5 cm³, 5.2 cm³ and 11.3 cm³, respectively (across-subjects mean based on individual activity maps at P < 0.001 uncorrected for multiple comparisons), suggesting that the volume was the smallest for 110° . In contrast, in the late stage, the volumes were 3.3 cm³, 10.7 cm³ and 9.3 cm³. Although these differences in volumes among conditions did not reach statistically significant levels, we observed that the volume for 110° became the largest while the volumes for 60° and 160° decreased. Possible explanations of these changes in activity are as follows. In the early stage, an internal model for 110° was not acquired yet, and the CNS combined output signals from internal models for 60° and 160° to cope with the novel 110° mouse. However, in the late stage, an internal model for 110° had been acquired, and the necessity of internal models for 60° and 160° decreased. The CNS may be able to combine acquired internal models according to the degree of

acquisition of a new internal model.

Studies reviewed in this section suggest that the CNS maintains multiple internal models for different objects and environments in a modular fashion and that it can combine output signals from the stored internal models depending on the situation.

20

5

10

15

Neural mechanism for selection and switching of internal models

In this section, we review behavioral and imaging studies investigating neural mechanisms that select or switch internal models according to changes in the environment and controlled objects.

Behavioral studies have shown that humans can switch internal models based on contextual information. For example, an auditory tone cue can induce context-dependent adaptation to prismatic displacement in the opposite directions (Kravitz & Yaffe, 1972). It had long been thought that simultaneous adaptations to opposing force fields are impossible (Brashers-Krug, Shadmehr, & Bizzi, 1996; Gandolfo, Mussa-Ivaldi, & Bizzi, 1996; Karniel & Mussa-Ivaldi, 2002). However, it was recently demonstrated that cognitive cues such as color and shape, and random and frequent presentation of the force fields, contribute to simultaneous learning and predictive switching of internal models for the opposing fields (Osu, Hirai, Yoshioka, & Kawato, 2004).

Using a continuous tracking task in which subjects used a computer mouse, we investigated brain activity related to switching of internal models (Imamizu, Kuroda,

Yoshioka, & Kawato, 2004). Subjects sufficiently learned to use three types of computer mouse with different input-output properties (rotated, velocity-control and normal mice) before their brain activities were scanned. During the tracking task in the MR scanner, the input-output property changed at random timing (from rotated to velocity-control, from normal to rotated, and so on). We investigated activity that increased immediately after the change and found that activities in the dorsolateral prefrontal cortex (DLPF; area 46), the insula, the anterior parts of the intra-parietal regions, and the lateral cerebellum are related to switching of internal models.

-14-

Our close examination of activation time courses revealed that there exist two types of temporal profiles in activity change depending on the brain region. One type of profile transiently increased immediately after the switch, but the levels of sustained background activity 20 sec after switching were almost the same as those before switching, suggesting that the dominant component of this profile is transient response (Fig. 6A). This transient response is probably related to the switching of internal models corresponding to the change of mouse type. This type of profile was found in area 46 and the insula. In the second type of profile, we could observe not only a transient increase of activity but also a change in the sustained activity level. In Figure 6B, the level of activation was low before the switch when the subjects used the normal mouse (open circles). It transiently increased immediately after switching and then remained high as long as the subjects used the rotated mouse (filled circles). Thus, this type of profile consisted of both transient response and sustained response.

10

Our further analysis found that the rotated and the velocity-control mice evoked sustained activity in distinct regions of the lateral cerebellum, suggesting that the activity is related to internal models. This type of profile was mainly observed in the cerebellum and the anterior part of the intra-parietal regions. We quantitatively investigated magnitudes of the transient and sustained responses in time courses of individual regions using a linear regression analysis and then calculated the ratio of the magnitude of the sustained response to that of the transient response (Fig. 6C). The results confirmed that transient response is dominant in the frontal regions (area 46 and the insula), while both responses are contained in activity in the parietal regions and the cerebellum. We also investigated the spatial overlap between the transient response related to the switching and the sustained response related to the internal models and found a significant overlap in the cerebellum, suggesting

-15-

that internal models contribute to the switching. As we discuss below in relation to computational models, this result suggests that internal models play an important role in switching mechanisms in the parietal and cerebellar regions.

Predictive switching of internal models

10

15

20

Empirically, two types of information are crucial for the switching of internal models: contextual information, such as color or shape of the objects that can be perceived before movement execution, and information on the difference between actual and predicted sensorimotor feedbacks calculated during or after execution. For example, when we lift a transparent bottle, the CNS can switch between internal models for light and heavy objects in a predictive fashion, since we know whether the bottle is empty or full beforehand. However, when lifting up a milk carton, we cannot estimate the weight, and the CNS relies on the error between actual and predicted sensorimotor feedbacks (prediction error). It is probably important for anticipatory adjustment of behavior that a mechanism for predictive switching of internal models can work before movement execution independently from a postdictive mechanism based on prediction error.

We conducted a behavioral experiment to investigate whether the predictive mechanism is functionally independent from the postdictive mechanism (Imamizu, Sugimoto, et al., 2007). Subjects learned to move their index fingers to targets while visual feedback of the finger movements was rotated clockwise (CW) or counterclockwise (CCW) by 40° around the initial position. When subjects adapted to alternating blocks of opposing rotations, we investigated the effects on the subjects' performances due to contextual information (a verbal instruction) on the forthcoming direction of rotation. We measured the effect of such contextual information on the predictive mechanism by measuring the performance error at the beginning of each block and that on the postdictive mechanism by measuring the speed of gradual decrease of the error within blocks. Consequently, the contextual information selectively improved predictive switching performance but did not affect postdictive switching performance based on prediction errors, suggesting the existence of functionally independent mechanisms. Based on the results of our behavioral study, we planned an fMRI experiment to examine whether these two mechanisms are based on separate neural substrates.

The experimental design of our previous fMRI study (Imamizu, et al., 2004) did not allow us to distinguish the activity related to predictive switching from that related to postdictive switching for the following reasons. While subjects tracked a target continuously moving at high speed on a screen, the mouse type was changed and, simultaneously, cognitive cues (change of cursor color and letters indicating the mouse type) were presented. In this way subjects simultaneously obtained cognitive cues for predictive switching and sensorimotor feedback for postdictive switching; consequently, cue-related activity temporally overlapped feedback-related activity.

In our new experiment, discrete pointing movements and event-related fMRI were used to separate activity related to the presentation of the cognitive cue from that related to sensorimotor feedback (Imamizu & Kawato, 2008). The task for subjects followed that in our behavioral study, and subjects sufficiently learned the 40° CW and 40° CCW visuomotor rotations before scanning of brain activity. During the fMRI experiment, the direction of rotation changed in a block-random fashion. A cue was presented at the

20

-17-

beginning of each trial and before movement initiation. The color of the cue corresponded to the direction of rotation of the feedback in an instructed condition, and thus predictive switching was possible. However, the color did not correspond to the direction in the non-instructed condition, and thus subjects relied on prediction errors calculated from sensorimotor feedback for switching in a non-instructed condition. Switching-related activity was identified as activity that transiently increased after the direction of rotation was changed. The switching-related activity in cue periods in the instructed condition, when a predictive switch is possible, was observed in the superior parietal lobule (SPL). However, the switching-related activity in feedback periods in the non-instructed condition, when prediction error is crucial for the postdictive switch, was observed in the inferior parietal lobule (IPL) and prefrontal cortex (PFC). These results clearly demonstrate regional differences in neural substrates between the predictive and postdictive mechanisms.

10

The above study suggests that the SPL contributes to predictive switching when the CNS has to select internal models. By contrast, Bursztyn and colleagues investigated brain activity related to predictive loading of an internal model when only one type of skill or an internal model was required throughout their experiment (Bursztyn, Ganesh, Imamizu, Kawato, & Flanagan, 2006). In their experiment, subjects learned to compensate a novel dynamics applied to their wrist movement. After learning, brain activity was measured during the interval between the cue and the initiation of movement. Their analysis revealed activity in supplementary motor areas (SMA), the primary motor (M1) regions, the dorsal premotor (PMd) regions, and the cerebellum. These results suggest that regions directly related to motor control are involved in internal-model recruitment in preparation for movement execution when selection of internal models is not needed.

-18-

Computational models for task switching

5

10

15

20

A mixture-of-experts architecture (Fig. 7A) was previously proposed for a computational model for task switching, including switching of internal models (Ghahramani & Wolpert, 1997; Jacobs, Jordan, Nowlan, & Hinton, 1991). In this architecture, expert modules (i.e., internal models) are trained so as to split the input data into subparts in which particular experts are specialized. For example, an expert module is specialized for the input-output property of each tool. Depending on the context, a gating module weights the contribution of the output of each expert module to the final output. A computational model for simultaneous learning and switching of internal models (MOSAIC model: Modular Selection and Identification for Control model) has recently been proposed (Haruno, Wolpert, & Kawato, 2001; Wolpert & Kawato, 1998). This model can explain the above results of behavioral and imaging studies in a consistent manner. The MOSAIC model (Fig. 7B) has two features that are largely different from the mixture-of-experts architecture.

First, in a mixture-of-experts architecture, the switching function is centralized in the gating module and segregated from the internal models. By contrast, in the MOSAIC model, internal models themselves play crucial roles in switching as follows. Multiple pairs of forward internal models (predictors: "F" in Fig. 7B) and inverse internal models (controllers: "I" in the figure) are tightly coupled as functional units in the MOSAIC model. For example, when we use a new tool, forward models of various types of similar tools simultaneously predict sensory feedback from an efference copy of motor commands. The prediction of each forward model is then compared with actual feedback. The smaller the error, the more likely it is that the forward model was an appropriate predictor in the current context. The inverse model paired with the appropriate predictor is considered an appropriate controller. Accordingly, the selection mechanism depends on the internal models, and forward models must be active when switching internal models. Therefore, the MOSAIC predicts that the switching activity spatially and temporally overlaps the internal model activity. Our fMRI study (Imamizu, et al., 2004) indicated that activity in the anterior parts of the intra-parietal regions and the lateral cerebellum contains both transient response related to switching and sustained response related to internal models. This result suggests that the MOSAIC model can well explain the switching mechanisms in these regions. Especially in the cerebellum, the transient response was observed in regions related to an internal model for the rotated mouse and one for the velocity-control mouse. This suggests that the transient response reflects activity of forward internal models for both types of mice, simultaneously predicting sensory feedback, and that the sustained response reflects activity of the selected internal models.

10

Second, the MOSAIC model has two architectures, each of which is specialized for
 predictive switching based on contextual information or postdictive switching based on
 error of prediction derived from sensorimotor feedback. This is consistent with our
 behavioral (Imamizu, Sugimoto, et al., 2007) and fMRI (Imamizu & Kawato, 2008) studies
 indicating the two independent switching mechanisms.

²⁰ Cerebellar activity in reinforcement-learning tasks: contributions of internal models to goal-directed behaviors

-20-

In the above sections, we reviewed behavioral and neuroimaging studies investigating internal models for control of peripheral objects (e.g. tools or objects in the hand) toward immediate goals in time (e.g. moving a cursor to a target). However, humans often have to guide their behaviors toward distal goals in time such as maximizing a reward that will be obtained in a long-term future under complicated stochastic environments. Learning based on reward has been investigated in a framework of reinforcement learning models (Sutton & Barto, 1998). Neurophysiological (W. Schultz, Apicella, & Ljungberg, 1993) and neuroimaging studies have shown that the basal ganglia and prefrontal regions play a key role in such types of learning. However, some studies have shown involvement of the lateral cerebellum as well as the basal ganglia in tasks designed for investigation of reinforcement learning (Doya, Okada, Ueda, Okamoto, & Yamawaki, 2001; Haruno, et al., 2004). Figure 8 shows examples of cerebellar regions activated in reinforcement-learning tasks (see also a supplemental movie:

10

http://www.cns.atr.jp/~imamizu/multi_functions.mpg). Red regions were activated when subjects conducted a stochastic decision task maximizing monetary rewards, in which subjects had to learn behaviors involving different task difficulties that were controlled by probability (Haruno, et al., 2004). Blue regions were activated when subjects planned their behaviors predicting a log-term reward in a Markov decision problem (Doya, et al., 2001).

The above activations of the lateral cerebellum suggest that internal models are needed for goal-directed behaviors in complex environments. Reinforcement-learning algorithms can be effective for optimizing a chain of actions in small-scale stochastic environments. However, many studies indicated limitations of the model-free approach adopted by plain reinforcement-learning algorithms and suggested the necessity of complementary use of model-based approaches. Doya has suggested that the cerebellum is specialized for supervised learning (model-based approach), which is guided by the error signal, while the basal ganglia are specialized for reinforcement learning (model-free approach), which is guided by the reward signal, and that each neural mechanism plays complementary roles in motor control and cognitive functions (Doya, 1999, 2000).

Anatomical connectivity between the basal ganglia and the cerebellum (Hoshi, Tremblay, Feger, Carras, & Strick, 2005) may support interplay between the cerebellar internal model and reinforcement-learning mechanisms in the basal ganglia. In the same line of thought, Daw and colleagues proposed a computational model consisting of two parallel reinforcement-learning modules in the brain: a model-free module associated with the dorsolateral striatum in the basal ganglia and a model-based module associated with the prefrontal cortex (Daw, Niv, & Dayan, 2005).

Kawato and Samejima theoretically pointed out the inefficiency of a plain reinforcement-learning algorithm when applied to practical problems including multiple degrees of freedoms, nonlinearity, and large delays (Kawato & Samejima, 2007). Such problems are often encountered in optimization of most goal-directed behaviors based on learning associations between motor commands and resultant trajectory (or sensorimotor feedback) and the associations between actions and resultant rewards in practical environments. They suggested that internal models contribute to dividing a complex task into simple subtasks, each of which is learned by separate reinforcement-learning modules. In extending reinforcement-learning tasks, it is fair to state that humans need good 'models' that can predict long-term changes in environments when they efficiently plan and select

behaviors toward distal goals in complex environments. Thus, internal models are thought

to be important for goal-directed behaviors in general.

15

20

-22-

The above results indicate that cerebellar internal models contribute to reinforcement-learning tasks based on long-term reward. To illustrate regional differences between activities related to the reinforcement-learning tasks and those related to sensorimotor control, several types of activities reviewed in the earlier part of this article were superimposed onto Figure 8. Green regions were activated when subjects manipulated an object with a complex dynamics (Milner, et al., 2007). Cyan regions are related to coordination of grip-force and load-force (Kawato, et al., 2003). Yellow regions indicate activity related to use of various common tools (Higuchi, et al., 2007; Fig. 5). Here, we averaged activities when subjects imagined use of the different tools. As the figure shows, activities related to the sensorimotor control tend to be located in superior and medial parts, while those related to the reinforcement-learning tasks tend to be located in inferior and lateral parts. We found activity reflecting an internal model of a novel tool (a 120° rotated mouse; Imamizu et al., 2000) in magenta regions. These results suggest that activity related to relatively higher cognitive functions (i.e., maximizing a long-term reward and use of a nevel tool) exist in inferior and lateral parts.

novel tool) exist in inferior and lateral parts.

10

15

20

Contributions of internal models to mirror system, social interactions, communication, and language

Many neurons in the PMv (F5) of macaque monkeys show activity in correlation with the grasp type being executed. A subpopulation of these neurons, the mirror neurons, responds to observation of goal-directed movements performed by another monkey or an experimenter (Gallese, Fadiga, Fogassi, & Rizzolatti, 1996). The mirror neurons and mirror

-23-

systems (not individual neurons) have been thought to support the abilities of communication and social interaction.

Experimental and theoretical studies suggest that internal models support information processing in mirror systems. Anatomical connectivity has been found between the ventral premotor (PMv) regions and the cerebellar output nucleus (the dentate nucleus) (Middleton 5 & Strick, 1997). Corresponding to this connection, our fMRI study suggests that functional connectivity between the lateral cerebellum and the PMv regions increased after acquisition of internal models (Tamada, Miyauchi, Imamizu, Yoshioka, & Kawato, 1999). This study used the novel mouse and confirmed similar change in activity to our previous study in the lateral cerebellum (Imamizu, et al., 2000). Moreover, the study investigated change in 10 activity in cerebral regions and found a significant effect of learning on the activated volume (decrease or increase of volume) in frontal and occipital regions. The authors found that activity in the left lateral cerebellum increased after learning in comparison to the right cerebellum. In their analysis, they adopted a hypothesis that increase in activity in the right cerebral region should be observed in comparison to the left homologous region if the 15 region has functional connectivity with the lateral cerebellum. They found that the activity in the right PMv regions increased in comparison to the left homologous regions after learning. Miall suggested that inverse models in the cerebellum and projections from parietal regions to the PMv via the cerebellum contribute to converting observation of another's action into one's own motor control signals (Miall, 2003). Oztop and colleagues 20 conducted computer simulations of tasks that are closely related to mirror systems (Oztop, Kawato, & Arbib, 2006; Oztop, Wolpert, & Kawato, 2005). In their simulations, an

observer estimated the goal of the reaching movements or the intention of the agent

performing grasping movements. The results theoretically indicated that internal models for sensorimotor control are effective in inferring the goals or mental states of others.

Estimation of mental states of others is essential for communication and social interactions. "Theory of mind," the ability to conceive the intentions and beliefs of others, has become another important key concept for understanding the mechanisms involved in 5 the estimation (Baron-Cohen, 1997; Frith & Frith, 1999). Many functional imaging studies have suggested involvement of regions near the superior temporal sulcus (STS) and other prefrontal regions in theory of mind (e.g. Decety, Jackson, Sommerville, Chaminade, & Meltzoff, 2004; Tankersley, Stowe, & Huettel, 2007). Our study found that the STS regions are activated during observation of geometrical shapes whose movements appear intentional or goal-directed (J. Schultz, Imamizu, Kawato, & Frith, 2004). Recently, Haruno and Kawato indicated that the strength of activation in the STS regions reflects individuals' competence to construct internal models of others' mental states (Haruno & Kawato, 2005). In their experiment, subjects were categorized into two groups according to strategies adopted by subjects for maximizing monetary rewards in a social interaction game (the "prisoner's dilemma" game). Subjects in a group tried to learn the association between one's own action and reward independently from the strategy of the other agent. Their strategies were well explained by a plain (model-free) reinforcement algorithm. By contrast, the strategies of the other group could not be explained by such an algorithm, and behavioral data suggested that the subjects exploited the agent's strategies to predict the agent's behavior in response to the subjects' own action (forward internal model). Their imaging data indicated that activity in STS regions in the latter group was significantly

stronger than that in the former group.

10

15

20

-25-

In our study investigating switching mechanisms for internal models (Imamizu & Kawato, 2008), activity in the STS regions increased as accuracy of subjects' performances increased after alteration of environments, and thus we concluded that the STS is one of the regions that receive output signals from the acquired internal models. Although the above experiment using the social interaction game did not investigate the learning process, we speculate that the STS region plays an important role in predicting the agent's behavior at the initial stage of learning but that the cerebellum acquires internal models of the agent after repeated games with the same agent. However, to fully understand the roles of the STS and the cerebellum in social interactions, it would be necessary to reconstruct or decode what information is presented in these regions and how the reconstructed information changes with learning.

10

15

How do internal models contribute to imitations and theory of minds? An essential problem in imitations is to infer covert information in others' brains (such as motor commands and intentions) from observation of their action. An effective way for this inference is to utilize one's own internal inverse models that translate action to motor commands or forward models (Oztop, et al., 2006; Oztop et al., 2005). Similarly, one's own internal models that translate actions and communication signals to intention and belief can be utilized for inference of others' intentions and beliefs that we cannot directly observe. Learning internal models of other people using one's own internal models may largely rely on the similarity of musculoskeletal systems and brains across people. Therefore, it becomes more difficult to learn internal models of people with different social and cultural

on the similarity of musculoskeletal systems and brains across people. Therefore, it becomes more difficult to learn internal models of people with different social and cultural backgrounds compared to those of people with a common background. Wolpert, Doya and Kawato pointed out similar computational difficulties, such as the large amount of noise, nonlinear properties, high dimensionality, and delayed feedback encountered in sensorimotor control and social communication, including imitations and theory of minds, although these difficulties are more severe in social communication than in sensorimotor control (Wolpert, Doya, & Kawato, 2003). It has been suggested that "mental simulations" using forward and inverse models (Oztop, et al., 2005) and hierarchical organization of internal models (Wolpert, Doya, & Kawato, 2003) can increase the inference accuracy of intentions and beliefs despite these difficulties.

The contribution of the cerebellum to language has been suggested by activation in the lateral cerebellar cortex during a verbal response selection task (saying an appropriate verb for a visually presented noun) (Raichle, et al., 1994). Our fMRI study (see above) revealed functional connectivity between the lateral cerebellum and PMv regions, parts of 10 which are known as Broca's area (Tamada, et al., 1999). Recently, Higuchi and colleagues found an overlap of brain activity for language and tool use in Broca's area (Higuchi, Imamizu, Chaminade, & Kawato, 2004). Their tool-use task required subjects to perform hierarchical manipulation of objects and tools, e.g. moving an object while holding it with chopsticks. The overlap was found in the dorsal parts of area 44 (a part of Broca's area). It 15 has been suggested that area 44 is involved in the syntactic aspects of language (Sakai, 2005) and specifically complex hierarchical processing (e.g. understanding of embedded sentences) (Friederici, Bahlmann, Heim, Schubotz, & Anwander, 2006). The location of this overlap suggests that language and tool use may share computational principles for processing hierarchical structures common to these two distinct abilities. In combination with a study indicating involvement of the PMv regions in monkeys during tool use (Obayashi, et al., 2001), this study suggests that neural processes for computation of hierarchical structures exist in primates and evolved secondarily to support human grammatical ability.

20

-27-

Beyond syntactic aspects, internal models are thought to contribute to the semantic aspects of words related to actions and manipulation. It is known that semantic memory is represented by distributed brain networks of sensory and motor regions. Recent functional brain imaging studies have intensively investigated semantic memory of tools (for review, see Martin & Chao, 2001) and found that the medial fusiform gyrus stores the form of tools, the left posterior middle gyrus represents the visual motion related to tool use, and the PMv regions represent the tool-use-associated action. Input-output properties are important semantic aspects of tools, and thus internal models representing these properties are key parts of the distributed network used for the semantic memory of tools. Our studies on the existence of internal models in the cerebellum (e.g. Higuchi, et al., 2007; Imamizu, Kuroda, Miyauchi, Yoshioka, & Kawato, 2003; Imamizu et al., 2000) suggest that the cerebellum also contributes to the semantic representation of words related to actions and manipulation, such as words for tools.

10

Hurley proposed a "Shared Circuit Model" in which cognitive functions such as
¹⁵ mirror systems, imitation, mental simulation of social interactions, and mind reading
(theory of minds) use internal simulation loops for sensorimotor control and additional
systems that inhibit motor outputs and generate virtual sensory inputs during the simulation
(Hurley, 2008). In our understanding, the internal simulation loops correspond to
combinations of forward and inverse internal models. Therefore, the studies reviewed in
²⁰ this section are consistent with the Shared Circuit Model in that internal models contribute
to mirror system, social interactions, communication, and language processing using the

same computational principles involved in sensorimotor control.

-28-

Discussion

10

15

20

Humans acquire internal models of the environment and external objects for effective realization of goal-directed behaviors. Neural substrates of internal models had been investigated by neurophysiological studies on other animals. However, recent advances in non-invasive functional neuroimaging methods such as PET and fMRI have enabled us to investigate how internal models are acquired and organized in the human brain. This review article first presented neuroimaging studies indicating how internal models are acquired in the brain network, including the cerebellum. Environments and objects with which humans interact often change in a discrete manner. Behavioral and imaging studies have indicated that the CNS acquires multiple internal models in a modular fashion and flexibly copes with such discrete changes by reducing interference and combining acquired internal models. A switching mechanism of internal models is also important for flexible adaptation under rapid and frequent environmental changes. Our studies suggested that neural mechanisms in the parietal regions (the SPL and IPL) and prefrontal regions contribute to the selection of appropriate internal models.

We then presented studies indicating the contribution of internal models to higher-order cognitive functions. Many studies have suggested that internal models are involved in optimization of goal-directed behaviors such as maximizing long-term rewards in collaboration with neural mechanisms for reinforcement learning. Our analysis of functional connectivity between the lateral cerebellum and the PMv suggests the contribution of internal models to mirror systems and faculties of language. Theoretical and simulation studies supported such a contribution to the mirror systems. Our recent imaging study demonstrated that regions probably receiving output signals from internal models for tool use are closely related to neural mechanisms for language processing and speech production (Broca's area). Finally, theoretical and empirical studies have suggested that internal models are involved in the theory of mind during social interactions by predicting others' behaviors in response to one's own behaviors.

Figure 9 shows schematic diagrams of functional pathways between the cerebral regions and the cerebellum based on the principal studies in this review article. In our study investigating predictive and postdictive mechanisms for switching of internal models (Imamizu & Kawato, 2008), we conducted analysis of functional connectivity using a method called dynamical causal modeling (Friston, Harrison, & Penny, 2003).

Consequently, we identified a significant increase in the influence of the SPL on the lateral 10 cerebellum during predictive switching based on contextual information and an increased influence of the IPL on the lateral cerebellum during postdictive switching based on error in the prediction of sensorimotor feedback (Fig. 9A). We hypothesized that the increased influence of the cerebellum on the IPL corresponded to the prediction of sensorimotor feedback, which is computed by forward models and necessary for calculation of prediction 15 error. Although our connectivity analysis could not find a statistically significant increase, an anatomical study on monkeys indicated that a region in the IPL (area 7b in monkeys) is the target of output from the cerebellum (Clower, West, Lynch, & Strick, 2001). These diagrams can be mapped onto the MOSAIC model (Fig. 7B), as indicated by circled numbers in the figures. 20

Furthermore, regarding postdictive switching, we found increased influence of the IPL on the SPL. This increase in influence suggests that the error of prediction for sensorimotor feedback was used as contextual information in the next trial because it is important

information on changes in the environment. This information flow may be analogous to those underlying behavioral adjustment after conflict or error in cognitive control tasks such as the Stroop color-naming task. Kerns and colleagues (Kerns, et al., 2004) found that an increase in activity in the anterior cingulate cortex (ACC) in an error trial leads to an increase in activity in the PFC (areas 8 and 9) in the subsequent trial, and they suggested that the ACC monitors the conflict and that the PFC produces behavioral adjustments based on detection of the conflict. It can be postulated that the IPL is involved in the monitoring of error and that the SPL contributes to subsequent behavioral adjustment by predictive switching of internal models. We also found an increase in bidirectional influences between the IPL and the DLPFC around area 46, suggesting that the DLPFC contributes to behavioral adjustment through interaction with the IPL in the switching of internal models.

10

20

In addition to the above analysis of functional connectivity, we found that activity increased in the lateral occipito-temporal cortices (LOTC), the supplementary motor area (SMA), the dorsal premotor (PMd) region, and the primary motor cortex (M1) as subjects' performances improved after alteration of the environment (direction of visuomotor 15 rotation). According to previous studies described below, these cerebral regions are closely related to internal models, and they are assumed to receive output signals from internal models (Fig. 9B). The LOTC is related to biological-motion perception (Bonda, Petrides, Ostry, & Evans, 1996), imitation (Iacoboni, et al., 2001), trajectory learning (Maquet, Schwartz, Passingham, & Frith, 2003) and smooth pursuit eye movements (Schmid, Rees, Frith, & Barnes, 2001). Using fMRI and computational modeling, Kawawaki et al. (Kawawaki, Shibata, Goda, Doya, & Kawato, 2006) indicated the contribution of the LOTC to prediction of target motion during visual pursuit. Output signals from forward

internal models have been suggested to play an important role in prediction and observation

-31-

of movements of objects and other persons (Blakemore & Decety, 2001; Frith, Blakemore, & Wolpert, 2000). Consistent with these studies, Haruno and Kawato found that the STS region, which is adjacent to the LOTC, is related to internal forward models of others' behaviors during human-human interaction (Haruno & Kawato, 2005).

The SMA, PMd and M1 are involved in motor control and likely receive output signals from internal inverse models. This is consistent with a study finding activity in these regions related to preparatory loading of information stored in internal models for compensation of a novel dynamics (Bursztyn, et al., 2006). In addition to these regions, output signals from inverse models are probably sent to the PMv regions (and Broca's area). This was suggested by our study finding an increase in functional connectivity after acquisition of internal models (Tamada, et al., 1999).

10

15

20

The studies we reviewed in the earlier sections mainly investigated internal models for rapid and smooth control of our bodies and tools to realize relatively immediate goals. However, some characteristics of internal models revealed by these studies are postulated to play key roles in supporting higher-order cognitive functions.

Modular organization of internal models is essential for effective organization of behavior in complex environments. If internal models were modularly organized, many novel situations that we encounter could be dealt with as combinations of previously experienced contexts. By modulating the contribution of the output signals from individual internal models to the final output signal, an enormous repertoire of behaviors could be generated (MOSAIC; Wolpert & Kawato, 1998). Our fMRI study has demonstrated fundamental neural mechanisms supporting such an ability in a relatively simple task, that is, use of three types of computer mouse with different input-output properties (Imamizu, Higuchi, et al., 2007). A modular decomposition strategy is effective for tackling a complex task by dividing it into simple subtasks. It has been suggested that internal models can contribute not only to learning subtasks but also to dividing a complex task into simple subtasks, each of which can be learned by model-free or model-based

reinforcement-learning modules (Kawato & Samejima, 2007). Furthermore, realization of distal goals in time under complex environments often needs multiple steps of actions that should be organized in a hierarchical fashion. Increasing the accuracy of hierarchical plans of actions requires precise internal models for individual actions and hierarchical organization of these internal models. Modularity of internal models is essential for such organization of action plans. Because realization of distal goals in time often needs step-by-step actions and a long time to accomplish them, environments often change during this process. For flexible reorganization of action plans depending on changes in the environment, it is important that internal models be modularly and hierarchically organized

¹⁵ Modularity and hierarchy are also thought to be essential for language processing.

20

and that they can be flexibly switched depending on available contextual information.

Bidirectional and recursive information processing is also important for higher-order cognitive functions. As we reviewed, many studies have suggested the existence of forward and inverse internal models in the CNS. Using these two forms of internal models, closed-loop circuits can be constructed in the CNS without relying on feedback loops in the external world (Hurley, 2008). These internal circuits support "mental simulations" of interactions between one's own actions and the resultant changes in the environment, and they can increase accuracy based on recursive computation in planning and selection of behaviors toward distal goals. It has been suggested that internal circuits including forward

-33-

and inverse models are essential for the inference of others' mental states in computer simulations related to mirror systems (Oztop, et al., 2005).

Functional connectivity between the cerebellum and various cerebral regions, as illustrated in Figure 9, indicates that the cerebellar internal models contribute to not only motor control but also various cognitive functions. In particular, STS, LOTC and PMv have been suggested to be involved in prediction of movements of external objects and actions of others, making inferences about intentions and goals of others, and language processing. Neurophysiological and anatomical studies have shown functional connections between the lateral cerebellum and both prefrontal and parietal regions (e.g. Clower, et al., 2001; Middleton & Strick, 2001; Sasaki, Oka, Kawaguchi, Jinnai, & Yasuda, 1977). However, previous studies, including ours, are mainly based on temporal correlations in activities between the regions or anatomical connectivity revealed by virus-based tracers, and thus little is known about the exact information exchanged between the cerebellum and the

cerebral cortices. We can make inferences about types of information based on our knowledge of the functions of particular cerebral regions; however, it would be necessary 15 to reconstruct or decode what information is presented in these cerebellar and the cerebral regions in the human brain to exactly understand the roles of the cerebellum and internal models in higher-order cognitive functions.

20

5

10

As opposed to the theoretical studies and computer simulations reviewed above, a small number of experimental studies have directly investigated the contributions of internal models to cognitive functions. Here, we propose several possible experimental and robotic studies. First, sensorimotor tasks could be used to investigate the hierarchical organization of internal models, where subjects would learn to hierarchically combine

several types of tools or sensorimotor transformations. This work would extend a study by Higuchi and colleagues (Higuchi, et al., 2007), and the results could be compared to those for activity related the hierarchical aspects of language, such as understanding of embedded sentences. To study social interactions, two fMRI-compatible manipulandums and fMRI scanners could be used to scan the brain activities of two subjects while they play interactive force-exerting games. Here, it would be possible to investigate how activity changes when the subject must learn different properties of the opponent, i.e., these properties change from motor dynamics such as force levels to higher-order cognitive properties such as strategies and personalities. Such a study would help us to understand the continuity or discontinuity of internal models between sensorimotor control and cognitive functions. Regarding robotic experiments, some robots have already been made for

http://www.irc.atr.jp/productRobovie/robovie-r2-e.html). Using these robots as a starting point, we could build new robots that possess the internal models of several types of people and as well as the ability to autonomously refine these internal models based on the

interacting with people (e.g. our institute's Robovie:

- and as well as the ability to autonomously refine these internal models based on the
 feedback obtained from people who have actually interacted with the robots in experiments.
 This would allow us to examine how their abilities, flexibility, and impressions they give of
 their intelligence improve in comparison to previous robots that react to people simply
 based on a database such as a lookup table of questions and answers.
- A series of our fMRI studies was motivated by the need to investigate sensorimotor learning mechanisms under novel environments. However, our results revealed some essential characteristics of internal models that can be generalized to understanding higher-order cognitive functions such as optimization of behaviors toward long-term goals, social interactions based on prediction of others' actions and mental states, and language

processing.

Acknowledgement

5

We thank Toshinori Yoshioka for developing the software used for the three-dimensional display of multiple brain activities (Figs. 4 and 7, and supplemental movie). This MATLAB(R) based software is freely available at: http://www.cns.atr.jp/multi_color

References

5

- Baron-Cohen, S. (1997). *Mindblindness: An essay on autism and theory of mind (learning, development and conceptual change)*. Cambridge: MIT Press.
- Blakemore, S. J., & Decety, J. (2001). From the perception of action to the understanding of intention. *Nat Rev Neurosci, 2*(8), 561-567.
- Blakemore, S. J., Frith, C. D., & Wolpert, D. M. (2001). The cerebellum is involved in predicting the sensory consequences of action. *Neuroreport*, *12*(9), 1879-1884.
- Blakemore, S. J., Wolpert, D. M., & Frith, C. D. (1998). Central cancellation of self-produced tickle sensation. *Nat Neurosci, 1*(7), 635-640.
- ¹⁰ Bonda, E., Petrides, M., Ostry, D., & Evans, A. (1996). Specific involvement of human parietal systems and the amygdala in the perception of biological motion. *J Neurosci, 16*(11), 3737-3744.
 - Brashers-Krug, T., Shadmehr, R., & Bizzi, E. (1996). Consolidation in human motor memory. *Nature*, 382(6588), 252-255.
- ¹⁵ Bursztyn, L. L., Ganesh, G., Imamizu, H., Kawato, M., & Flanagan, J. R. (2006). Neural correlates of internal-model loading. *Curr Biol*, *16*(24), 2440-2445.
 - Clower, D. M., West, R. A., Lynch, J. C., & Strick, P. L. (2001). The inferior parietal lobule is the target of output from the superior colliculus, hippocampus, and cerebellum. *J Neurosci*, 21(16), 6283-6291.
- Daw, N. D., Niv, Y., & Dayan, P. (2005). Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nat Neurosci,* 8(12), 1704-1711.
 - Decety, J., Jackson, P. L., Sommerville, J. A., Chaminade, T., & Meltzoff, A. N. (2004). The neural bases of cooperation and competition: an fMRI investigation. *Neuroimage*, 23(2), 744-751.
 - Diedrichsen, J., Criscimagna-Hemminger, S. E., & Shadmehr, R. (2007). Dissociating timing and coordination as functions of the cerebellum. *J Neurosci*, 27(23), 6291-6301.

- Doya, K. (1999). What are the computations of the cerebellum, the basal ganglia and the cerebral cortex? *Neural Netw, 12*(7-8), 961-974.
- Doya, K. (2000). Complementary roles of basal ganglia and cerebellum in learning and motor control. *Curr Opin Neurobiol*, *10*(6), 732-739.
- Doya, K., Okada, G., Ueda, K., Okamoto, Y., & Yamawaki, S. (2001). Pediction of shortand long-term reward: A functional MRI study with a Markov decision problem.
 Paper presented at the Annual Meeting Society for Neuroscience.
 - Ebner, T. J., & Pasalar, S. (2008). Cerebellum predicts the future motor state. *Cerebellum*, 7(4), 583-588.
- ¹⁰ Flanagan, J. R., Nakano, E., Imamizu, H., Osu, R., Yoshioka, T., & Kawato, M. (1999). Composition and decomposition of internal models in motor learning under altered kinematic and dynamic environments. *J Neurosci, 19*(20), RC34.
 - Flanagan, J. R., & Wing, A. M. (1997). The role of internal models in motion planning and control: evidence from grip force adjustments during movements of hand-held loads. *J Neurosci*, 17(4), 1519-1528.

20

- Friederici, A. D., Bahlmann, J., Heim, S., Schubotz, R. I., & Anwander, A. (2006). The brain differentiates human and non-human grammars: functional localization and structural connectivity. *Proc Natl Acad Sci U S A*, 103(7), 2458-2463.
- Friston, K. J., Harrison, L., & Penny, W. (2003). Dynamic causal modelling. *Neuroimage*, 19, 1273-1302.
- Frith, C. D., Blakemore, S. J., & Wolpert, D. M. (2000). Abnormalities in the awareness and control of action. *Philos Trans R Soc Lond B Biol Sci*, 355(1404), 1771-1788.
- Frith, C. D., & Frith, U. (1999). Interacting minds--a biological basis. *Science*, 286(5445), 1692-1695.
- ²⁵ Gallese, V., Fadiga, L., Fogassi, L., & Rizzolatti, G. (1996). Action recognition in the premotor cortex. *Brain, 119 (Pt 2),* 593-609.
 - Gandolfo, F., Mussa-Ivaldi, F. A., & Bizzi, E. (1996). Motor learning by field approximation. *Proc Natl Acad Sci U S A*, *93*(9), 3843-3846.
 - Ghahramani, Z., & Wolpert, D. M. (1997). Modular decomposition in visuomotor learning. *Nature, 386*(6623), 392-395.

- Gomi, H., Shidara, M., Takemura, A., Inoue, Y., Kawano, K., & Kawato, M. (1998). Temporal firing patterns of purkinje cells in the cerebellar ventral paraflocculus during ocular following responses in monkeys I. Simple spikes [In Process Citation]. J Neurophysiol, 80(2), 818-831.
- Graydon, F. X., Friston, K. J., Thomas, C. G., Brooks, V. B., & Menon, R. S. (2005). Learning-related fMRI activation associated with a rotational visuo-motor transformation. *Brain Res Cogn Brain Res*, 22, 373-383.
 - Grodd, W., Hulsmann, E., Lotze, M., Wildgruber, D., & Erb, M. (2001). Sensorimotor mapping of the human cerebellum: fMRI evidence of somatotopic organization. *Hum Brain Mapp*, *13*(2), 55-73.

20

- Haruno, M., & Kawato, M. (2005). Two groups of subjects with different learning competence in a prisoner's dilemma task exhibit differential activations in the superior temporal sulcus. Paper presented at the Annual Meeting Society for Neuroscience
- Haruno, M., Kuroda, T., Doya, K., Toyama, K., Kimura, M., Samejima, K., et al. (2004). A neural correlate of reward-based behavioral learning in caudate nucleus: a functional magnetic resonance imaging study of a stochastic decision task. J Neurosci, 24(7), 1660-1665.
 - Haruno, M., Wolpert, D. M., & Kawato, M. (2001). Mosaic model for sensorimotor learning and control. *Neural Comput*, *13*(10), 2201-2220.
 - Hauser, M. D., Chomsky, N., & Fitch, W. T. (2002). The faculty of language: what is it, who has it, and how did it evolve? *Science*, *298*(5598), 1569-1579.
 - Higuchi, S., Imamizu, H., Chaminade, T., & Kawato, M. (2004). *Broca's area during tool use and linguistic processing*. Paper presented at the Annual Meeting Society for Neuroscience
 - Higuchi, S., Imamizu, H., & Kawato, M. (2007). Cerebellar activity evoked by common tool-use execution and imagery tasks: an fMRI study. *Cortex*, *43*(3), 350-358.
 - Hoshi, E., Tremblay, L., Feger, J., Carras, P. L., & Strick, P. L. (2005). The cerebellum communicates with the basal ganglia. *Nat Neurosci*, 8(11), 1491-1493.

- Hurley, S. (2008). The shared circuits model (SCM): how control, mirroring, and simulation can enable imitation, deliberation, and mindreading. *Behav Brain Sci*, 31(1), 1-22; discussion 22-58.
- Iacoboni, M., Koski, L. M., Brass, M., Bekkering, H., Woods, R. P., Dubeau, M. C., et al. (2001). Reafferent copies of imitated actions in the right superior temporal cortex. *Proc Natl Acad Sci U S A*, 98(24), 13995-13999.

- Imamizu, H., Higuchi, S., Toda, A., & Kawato, M. (2007). Reorganization of brain activity for multiple internal models after short but intensive training. *Cortex*, *43*(3), 338-349.
- ¹⁰ Imamizu, H., & Kawato, M. (2008). Neural correlates of predictive and postdictive switching mechanisms for internal models. *J Neurosci*, *28*(42), 10751-10765.
 - Imamizu, H., Kuroda, T., Miyauchi, S., Yoshioka, T., & Kawato, M. (2003). Modular organization of internal models of tools in the human cerebellum. *Proc Natl Acad Sci U S A*, 100(9), 5461-5466.
- ¹⁵ Imamizu, H., Kuroda, T., Yoshioka, T., & Kawato, M. (2004). Functional magnetic resonance imaging examination of two modular architectures for switching multiple internal models. *J Neurosci*, *24*(5), 1173-1181.
 - Imamizu, H., Miyauchi, S., Tamada, T., Sasaki, Y., Takino, R., Putz, B., et al. (2000). Human cerebellar activity reflecting an acquired internal model of a new tool. *Nature*, 403(6766), 192-195.
 - Imamizu, H., Sugimoto, N., Osu, R., Tsutsui, K., Sugiyama, K., Wada, Y., et al. (2007). Explicit contextual information selectively contributes to predictive switching of internal models. *Exp Brain Res*, 181(3), 395-408.
 - Ito, M. (1984). The cerebellum and neural motor control. New York: Raven Press.
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., & Hinton, G. E. (1991). Adaptive mixture of local experts. *Neural Comput*, *3*, 79-87.
 - Johansson, R. S., & Westling, G. (1988). Coordinated isometric muscle commands adequately and erroneously programmed for the weight during lifting task with precision grip. *Exp Brain Res*, 71(1), 59-71.

- Karniel, A., & Mussa-Ivaldi, F. A. (2002). Does the motor control system use multiple models and context switching to cope with a variable environment? *Exp Brain Res*, 143(4), 520-524.
- Kawato, M. (1999). Internal models for motor control and trajectory planning. *Curr Opin Neurobiol, 9*(6), 718-727.
- Kawato, M., Furukawa, K., & Suzuki, R. (1987). A hierarchical neural-network model for control and learning of voluntary movement. *Biol Cybern*, *57*(3), 169-185.
- Kawato, M., Kuroda, T., Imamizu, H., Nakano, E., Miyauchi, S., & Yoshioka, T. (2003). Internal forward models in the cerebellum: fMRI study on grip force and load force coupling. *Prog Brain Res*, 142, 171-188.

15

20

- Kawato, M., & Samejima, K. (2007). Efficient reinforcement learning: computational theories, neuroscience and robotics. *Curr Opin Neurobiol, 17*(2), 205-212.
- Kawawaki, D., Shibata, T., Goda, N., Doya, K., & Kawato, M. (2006). Anterior and superior lateral occipito-temporal cortex responsible for target motion prediction during overt and covert visual pursuit. *Neurosci Res*, *54*(2), 112-123.
- Kerns, J. G., Cohen, J. D., MacDonald, A. W., 3rd, Cho, R. Y., Stenger, V. A., & Carter, C. S. (2004). Anterior cingulate conflict monitoring and adjustments in control. *Science*, 303(5660), 1023-1026.
- Kitazawa, S., Kimura, T., & Yin, P. B. (1998). Cerebellar complex spikes encode both destinations and errors in arm movements. *Nature*, *392*(6675), 494-497.
- Krakauer, J. W., Ghilardi, M. F., & Ghez, C. (1999). Independent learning of internal models for kinematic and dynamic control of reaching. *Nat Neurosci*, 2(11), 1026-1031.
- Krakauer, J. W., Ghilardi, M. F., Mentis, M., Barnes, A., Veytsman, M., Eidelberg, D., et al. (2004). Differential cortical and subcortical activations in learning rotations and gains for reaching: a PET study. *J Neurophysiol*, *91*(2), 924-933.
- Kravitz, J. H., & Yaffe, F. L. (1972). Conditionned adaptation to prismatic displacement with a tone as the conditioal stimulus. *Percept Psychophys*, *12*(3), 305-308.

- Maquet, P., Schwartz, S., Passingham, R., & Frith, C. (2003). Sleep-related consolidation of a visuomotor skill: brain mechanisms as assessed by functional magnetic resonance imaging. *J Neurosci*, 23(4), 1432-1440.
- Martin, A., & Chao, L. L. (2001). Semantic memory and the brain: structure and processes. *Curr Opin Neurobiol*, 11(2), 194-201.
- Miall, R. C. (2003). Connecting mirror neurons and forward models. *Neuroreport, 14*(17), 2135-2137.
- Miall, R. C., Keating, J. G., Malkmus, M., & Thach, W. T. (1998). Simple spike activity predicts occurrence of complex spikes in cerebellar Purkinje cells. *Nat Neurosci, 1*(1), 13-15.

- Miall, R. C., Reckess, G. Z., & Imamizu, H. (2001). The cerebellum coordinates eye and hand tracking movements. *Nat Neurosci, 4*(6), 638-644.
- Miall, R. C., Weir, D. J., Wolpert, D. M., & Stein, J. F. (1993). Is the cerebellum a Smith predictor? *Journal of Motor Behavior*, 25, 203-216.
- ¹⁵ Middleton, F. A., & Strick, P. L. (1997). Dentate output channels: Motor and cognitive components. In C. I. de Zeeuw, P. Strata & J. Voogd (Eds.), *The cerebellum: From structure to control* (pp. 553-566): Elsevier Science BV.
 - Middleton, F. A., & Strick, P. L. (2001). Cerebellar projections to the prefrontal cortex of the primate. *J Neurosci*, *21*(2), 700-712.
- Milner, T. E., Franklin, D. W., Imamizu, H., & Kawato, M. (2007). Central control of grasp: manipulation of objects with complex and simple dynamics. *Neuroimage*, 36(2), 388-395.
 - O'Reilly, J. X., Mesulam, M. M., & Nobre, A. C. (2008). The cerebellum predicts the timing of perceptual events. *J Neurosci, 28*(9), 2252-2260.
- ²⁵ Obayashi, S., Suhara, T., Kawabe, K., Okauchi, T., Maeda, J., Akine, Y., et al. (2001). Functional brain mapping of monkey tool use. *Neuroimage*, *14*(4), 853-861.
 - Osu, R., Hirai, S., Yoshioka, T., & Kawato, M. (2004). Random presentation enables subjects to adapt to two opposing forces on the hand. *Nat Neurosci*, 7(2), 111-112.
 - Oztop, E., Kawato, M., & Arbib, M. (2006). Mirror neurons and imitation: a computationally guided review. *Neural Netw*, *19*(3), 254-271.

- Oztop, E., Wolpert, D., & Kawato, M. (2005). Mental state inference using visual control parameters. *Brain Res Cogn Brain Res, 22*(2), 129-151.
- Raichle, M. E., Fiez, J. A., Videen, T. O., MacLeod, A. M., Pardo, J. V., Fox, P. T., et al. (1994). Practice-related changes in human brain functional anatomy during nonmotor learning. *Cereb Cortex*, 4(1), 8-26.
- Sakai, K. L. (2005). Language acquisition and brain development. *Science*, *310*(5749), 815-819.
- Sasaki, K., Oka, H., Kawaguchi, S., Jinnai, K., & Yasuda, T. (1977). Mossy fibre and climbing fibre responses produced in the cerebellar cortex by stimulation of the cerebral cortex in monkeys. *Exp Brain Res, 29*(3-4), 419-428.

- Schmid, A., Rees, G., Frith, C., & Barnes, G. (2001). An fMRI study of anticipation and learning of smooth pursuit eye movements in humans. *Neuroreport, 12*(7), 1409-1414.
- Schultz, J., Imamizu, H., Kawato, M., & Frith, C. D. (2004). Activation of the human superior temporal gyrus during observation of goal attribution by intentional objects. *J Cogn Neurosci, 16*(10), 1695-1705.
 - Schultz, W., Apicella, P., & Ljungberg, T. (1993). Responses of monkey dopamine neurons to reward and conditioned stimuli during successive steps of learning a delayed response task. *J Neurosci*, 13(3), 900-913.
- Shadmehr, R., & Holcomb, H. H. (1997). Neural correlates of motor memory consolidation. *Science*, 277(5327), 821-825.
 - Shidara, M., Kawano, K., Gomi, H., & Kawato, M. (1993). Inverse-dynamics model eye movement control by Purkinje cells in the cerebellum. *Nature*, *365*(6441), 50-52.
 - Sutton, R. S., & Barto, A. G. (1998). Reinforcement learning. Cambridge, MA: MIT Press.
- Tamada, T., Miyauchi, S., Imamizu, H., Yoshioka, T., & Kawato, M. (1999). Cerebro-cerebellar functional connectivity revealed by the laterality index in tool-use learning. *NeuroReport*, 10(2), 325-331.
 - Tankersley, D., Stowe, C. J., & Huettel, S. A. (2007). Altruism is associated with an increased neural response to agency. *Nat Neurosci*, *10*(2), 150-151.

- Wolpert, D. M., Doya, K., & Kawato, M. (2003). A unifying computational framework for motor control and social interaction. *Philos Trans R Soc Lond B Biol Sci*, 358(1431), 593-602.
- Wolpert, D. M., Ghahramani, Z., & Jordan, M. I. (1995). An internal model for sensorimotor integration. *Science*, *269*(5232), 1880-1882.

Wolpert, D. M., & Kawato, M. (1998). Multiple paired forward and inverse models for motor control. *Neural Networks*, 11, 1317-1329.

Figure legends

Figure 1

⁵ Predictive control of a computer mouse (**A**) and prediction of consequence of control (**B**).

Figure 2

Relationship between direction of mouse movement (black arrows) and cursor movement (white arrows) when a cursor position is rotated 120° around the center of a screen (120° rotated mouse).

¹⁰ Figure 3

15

(A) Tracking error (mean ± SD) averaged across subjects as a function of number of training sessions. (B) Tracking error (mean + SD) in an experiment where a target velocity in baseline periods was increased to equalize the errors in baseline period to the errors in test period. (C) Cerebellar regions where activity is related to error signals (white regions enclosed by solid line) and regions where activity contains components unrelated to error (hatched regions). (D) Left panel shows activity change in hatched regions of Fig. 3C. Middle panel shows activity change in white regions of Fig. 3C. Right panel shows subtraction of the activity change in the middle panel from that in the left panel. Each curve indicates the exponential function fitted to the circles.

Figure 4

Cerebellar regions related to manipulation of the novel mice shown in various views. Yellow and blue colors indicate regions where activation was more highly correlated with the manipulation of a rotated and velocity mouse, respectively, than manipulation of a normal mouse. L: left, R: right

Figure 5

10

15

Distribution of t-value-weighted centroid of activation coordinates when subjects actually used common tools (**A**) or when they imagined use of the tools (**B**). Thick circles in Fig. 5B indicate centroids in the posterior lobule. Subjects used sixteen tools but one of the tools (saw) could not evoke significant activation (P < 0.001, uncorrected for multiple comparisons in random effect analysis). Thus, the number of centroids is fifteen. (**C**) Transverse anatomical image of the human brain at the cerebellum. Thick outlines in Figs. 5A, 5B and 5E indicate the region of the right lateral cerebellar hemisphere as shown in Fig. 5C. (**D**) Mean distance of the centroid from the fourth ventricle across tools (+SD) for actual use or imaginary use. (**E**) Lines connecting the centroids during actual use (rectangles) with those during the imagery (circles) for tools that evoked activities in the posterior lobule during the imagery. Thick black lines indicate tools that evoked activities in the posterior lobule during imaginary use, while thin gray lines indicate tools that did not evoke activities in the posterior lobule.

²⁰ Figure 6

(A) Activation time course in area 46 when mouse-type changed from the normal to the rotated mouse. (B) Activation time course in the cerebellum. (C) Schematic representation of a ratio of sustained component to that of transient component in various brain regions.

Figure 7

⁵ Computational models for switching of internal models. (A) Mixture-of-experts model having a single switching mechanism (a gating module). (B) MOSAIC model having separate switching mechanisms for predictive switching based on contextual information and postdictive switching based on the prediction error of sensorimotor feedback. Circled numbers indicate correspondence between information flows in the model and neural pathways in Figure 9.

Figure 8

15

20

Cerebellar regions activated in various kinds of tasks. (**A**) Activated regions shown in superior–posterior–lateral view. (**B**) Activated regions projected onto the sagittal (left), the coronal (right), and the transverse (bottom) planes. A gray object indicates outline of the cerebellum from the same view as Fig. 8A. **Red** regions were activated when subjects conducted a stochastic decision task maximizing monetary rewards. **Blue** regions were activated when subjects predicted a log-term reward. **Green** regions were activated when subjects manipulated an object with complex dynamics. **Cyan** regions were related to coordination of grip-force and load-force. **Yellow** regions indicate activity related to use of various common tools. **Magenta** regions were related to an internal model of a novel tool. Schematic diagrams of functional pathways between the cerebral regions and the cerebellum based on representative studies in this review article. (**A**) Pathways related to predictive or postdictive switching of internal models based on our functional connectivity analysis (Imamizu & Kawato, 2008). DLPFC: dorsolateral prefrontal cortex, IPL: inferior parietal lobule, SPL: superior parietal lobule. (**B**) Output pathways from cerebellar forward and inverse internal models. SMA: supplementary motor area, PMd: dorsal premotor region, PMv: ventral premotor region, BA: Broca's area, M1: primary motor cortex, STS: superior temporal sulcus, LOTC: lateral occipito-temporal cortices. Circled numbers indicate correspondences of the pathways to information flows in the MOSAIC model (Fig. 7B).

5



Fig. 1



Fig. 2







Fig. 4







Fig. 6



Fig. 7



Fig. 8

