



2015). In contrast, DA receptor activation in the striatum, one of the main target structures of SN and ventral tegmental area projections, has been linked to flexible updating of relevant information in working memory (e.g., D'Ardenne et al., 2012; Cools, Sheridan, Jacobs, & D'Esposito, 2007; Bilder, Volavka, Lachman, & Grace, 2004; Frank, Loughry, & O'Reilly, 2001). This has been described as a selective gating process (Chatham & Badre, 2015; Badre, 2012), in which input gating of the striatum ensures that behaviorally relevant information enters cortical working memory. The first questions of this study therefore is whether striatal activity relates exclusively to flexible updating driven by relevant changes or rather reflects the processing of unpredicted sensory information in general.

Beyond this moment-by-moment perspective, the volatility, that is, the rate of change in the environment, plays a role in regulating the interplay between flexible and stable states. Temporary extended increases in prediction errors indicate volatile environments and lead to a fast adjustment of internal predictions and behavior (Jiang, Beck, Heller, & Egner, 2015; Chumbley et al., 2014; Schiffer, Ahlheim, Wurm, & Schubotz, 2012; Friston, Daunizeau, & Kiebel, 2009; Behrens, Woolrich, Walton, & Rushworth, 2007). Recent studies provide evidence that longer timescale tonic DA modulates the extent to which prior action outcome biases phasic DA release and hence future action selection (Yu, FitzGerald, & Friston, 2013; Humphries, Khamassi, & Gurney, 2012; Beeler, Daw, Frazier, & Zhuang, 2010). This idea has been formalized in the dual-state theory of working memory, according to which representations in pFC are regulated by so-called attractor networks (Durstewitz & Seamans, 2008; Durstewitz, Seamans, & Sejnowski, 2000). These can assume either high- or low-energy barriers, corresponding to a shielding or a destabilization of current working memory representations, respectively. The adjustment of each of the two states relies on tonic DA release in the dopaminergic midbrain, driven by the interplay of previous event predictability and the history of behavioral outcomes. This leads to the hypothesis that the failure to flexibly adapt in volatile environments would lead to increased DA levels thereby promoting transition into a flexible cognitive state. Conversely, the failure to ignore drifts and to maintain stable responding in stable environments should be accompanied by decreased DA activity, thereby promoting transition into a more stable state. Whereas the first question of this study thus pertains to the immediate response to different types of prediction errors (i.e., flexible switching or stable maintenance), the second question is whether these responses are modulated by cognitive states pertaining to recent performance history resulting from different levels of volatility.

We employed fMRI while participants performed a task that required monitoring of a digit sequence for structure-violating items. Switches between predictable sequences had to be indicated via button press (cognitive flexibility),

whereas sequence omissions (drifts) had to be ignored (cognitive stability). Switch and drift probabilities varied across the experiment. Implementing switches and drifts in the same design allowed us to use correlational analysis on the rate of correctly detected switches and ignored drifts to show that cognitive flexibility and stability are functionally independent (Cools & D'Esposito, 2011). This way, we could further test our hypothesis that model updating and stabilization as functionally independent processes would recruit different cortical regions: Model update and retrieval from episodic memory was expected to be reflected in activity of medial prefrontal areas and the hippocampus (Schlichting & Preston, 2015; Preston & Eichenbaum, 2013). In contrast, stabilization of prediction was hypothesized to be accompanied by premotor and lateral pFC activation (Cohen, Braver, & Brown, 2002; Miller & Cohen, 2001), whereby either dorsal or ventral activation would reflect the way a model content is stored in working memory, that is, spatially or verbally, respectively (Rottschy et al., 2012). With regard to the role of the striatum in immediate response selection, we hypothesized caudate activity in response to both switches and drifts (Schiffer et al., 2012). Crucially, we expected the degree of activation increase to be correlated with the ability to discriminate between both events, showing that the striatum is related to selecting correct responses toward different types of prediction errors (Chatham & Badre, 2015; Badre, 2012). Finally, to test the idea that the dopaminergic midbrain is involved in the transition between cognitive states (Humphries et al., 2012; Durstewitz & Seamans, 2008), we measured adaptation of flexible and stable states as determined by previous performance within time windows that were individually fitted by levels of volatility predicting participants' behavioral data.

## METHODS

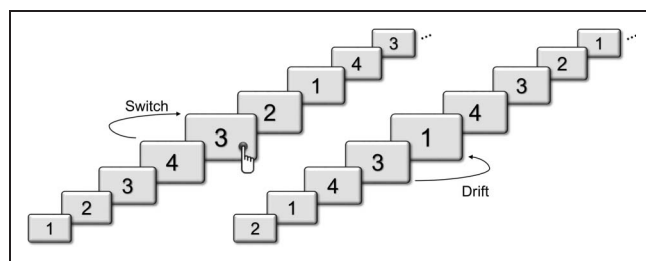
### Participants

Twenty-two right-handed healthy participants (16 women;  $23.26 \pm 2.19$  years old, range = 19–27 years) with normal or corrected-to-normal vision participated in the study. None of them reported a history of medical, neurological, or psychiatric disorders or substance abuse. The study protocol was conducted in accordance with ethical standards of the Declaration of Helsinki and approved by the local ethics committee of the University of Münster. Each participant submitted a signed informed consent notification and received reimbursement or course credits for their participation afterwards. For further assessment, participants were given the Barratt Impulsiveness Scale (BIS-11; Patton, Stanford, & Barratt, 1995). Two participants were excluded (because of pressing the wrong response button and an incidental finding of brain abnormality). Thus, a total of 20 participants (15 women;  $23.57 \pm 2.60$  years old, range = 19–27) were included in the analysis.

## Task

Participants were presented with two different digit sequences, which allowed them to generate an internal model predicting forthcoming input (ascending model: 1–2–3–4, descending model: 4–3–2–1). Digits continuously followed one another and were presented on the screen for 900 msec, separated by an ISI of 100 msec (Figure 1). Sequences repeated constantly to enable the participants to predict the regular sequence. Occasionally, switches between the models, that is, directional changes, occurred at a random position within the current sequence. In addition, single digits were omitted sometimes at variable positions without experiencing a temporal gap (drifts, hereafter). Derived from signal detection theory measures, the participant's task was to indicate a switch from one model to the other by button press (switch detection) but to ignore the sequential omissions (drift rejection). They had to respond as fast and accurately as possible and received individual performance feedback in three breaks of 14 sec every four blocks.

The task consisted of 12 blocks with an average number of 160 trials ( $SD = 6.82$ ) in a full-factorial 2 (Probability: high vs. low)  $\times$  2 (Event: switch vs. drift) design. This means that blocks either had a high or low probability of switches, paired with a high or low probability of drifts. Transitions between block types resulting from factor combination were balanced across the entire session. Event probabilities were individually staircased before the scanner session and probabilities for switches and drifts always adapted to the same extent. Here, participants performed 10 blocks with 80 trials each, starting with an event frequency of 17.5%. For a block performance higher (lower) than 75%, event frequency increased (decreased) with a rate of 2.5% in the subsequent block. The maximum reached frequency of events across the entire staircase session ( $M = 24.5\%$ ,  $SD = 1.7$ ) served as maximum event frequency in unmixed blocks of the main experiment, in which switches and drifts occurred with the same frequency. Minimum event frequency was set to



**Figure 1.** Schematic diagram of the task. Stimuli of a simple four-digit sequence continuously followed each other with a frequency of 1 Hz. Participants had to indicate a change in the direction of a sequence (switch), as displayed in the upper row, via button press. At the same time, they had to ignore the omission of a single digit (drift), as displayed in the lower row. The probability of occurrence of each respective type of sequence violation changed from block to block.

approximately one third of the respective individual maximum frequency. In mixed (i.e., high-switch and low-drift or vice versa) blocks, the difference between maximum frequency for both event types and minimum frequency for one event type served as maximum event frequency, whereas minimum frequency remained equal. In this way, difficulty level in terms of overall probability of events was kept constant across the experiment (with the exception of unmixed low-frequency blocks).

Stimulus presentation per block was pseudorandomized by using the stochastic universal sampling method (Baker, 1987). This method ensured a balanced distribution of event types across the block so that the observed event frequencies were in line with the expected frequencies.

## fMRI Data Acquisition

Whole-brain imaging data were collected on a 3-T Siemens Magnetom Prisma MR tomograph (Siemens, Erlangen, Germany) using a 20-channel head coil. To minimize head motion, the head was tightly fixated with cushions. Functional images were acquired using a gradient T2\*-weighted single-shot EPI sequence sensitive to BOLD contrast ( $64 \times 64$  data acquisition matrix, 192 mm field of view,  $90^\circ$  flip angle, repetition time = 2000 msec, echo time = 30 msec). Each volume consisted of 30 adjacent axial slices with a slice thickness of 4 mm and a gap of 1 mm, resulting in a voxel size of  $3 \times 3 \times 5$  mm. Images were acquired in ascending order along the AC–PC plane to provide a whole-brain coverage. Structural data were acquired for each participant using a standard Siemens 3-D T1-weighted MPRAGE sequence for detailed reconstruction of anatomy with isotropic voxels ( $1 \times 1 \times 1$  mm) in a 256-mm field of view ( $256 \times 256$  matrix, 192 slices, repetition time = 2130, echo time = 2.28). Stimuli were projected on a screen that was positioned behind the participant's head. They were presented in the center of the field of vision by a video projector, and participants viewed the screen by a  $45^\circ$  mirror, which was fixated on the top of the head coil and adjusted for each participant to provide a good view of the entire screen.

## Behavioral Data Analysis

Performance on the task was defined by hits (correct detection of switches), correct rejections of drifts, and correspondingly, switch misses and false alarms at drifts. Discrimination index ( $P_r$ ; probability of recognition of switches and drifts, i.e.,  $(\text{hits} + 0.5/\text{number of switches} + 1) - (\text{false alarms} + 0.5/\text{number of drifts} + 1)$ ) and bias index ( $B_r$ ; response probability in an uncertain state, i.e.,  $(\text{false alarms} + 0.5/\text{number of drifts} + 1)/(1 - P_r)$ ) were calculated (Snodgrass & Corwin, 1988). Hit and correct rejection rate and RTs at hits and false alarms were compared by Student's paired  $t$  tests. To assess the relationship between hits and correct rejections, we calculated Pearson's correlation coefficient. To show that there was no relationship between the two measures, we additionally calculated

the attenuation adjusted correlation coefficient, which provides an estimate of the strength of a correlation assuming no measurement error. For this purpose, the correlation of the two variables ( $r_{xy}$ ) was divided by the square root of their multiplied reliabilities ( $r_{xx}$  and  $r_{yy}$ ; see Spearman, 1904). Here, we calculated split-half reliability, that is, Pearson's correlation coefficients between scores of the two halves of the test for both switch detection and drift rejection. If not stated otherwise, significance tests were performed at  $\alpha = .05$ , two-sided.

## fMRI Data Analysis

### fMRI Data Preprocessing

Brain image preprocessing and basic statistical analyses were conducted using LIPSIA software package, version 3.0 (Lohmann et al., 2001). As a first step, spikes in time series were corrected by interpolating them with adjacent time points. To correct for temporal offsets between the slices acquired in one scan, a cubic-spline interpolation was used. Additionally, individual functional magnetic resonance (EPI) images were motion-corrected with the first time-step as reference and six degrees of freedom (three rotational, three translational). Then, the average across all time points of this corrected data was used as reference scan for a second pass of motion correction. Motion correction estimates were inspected visually. A rigid linear registration with six degrees of freedom (three rotational, three translational) was performed to align the functional data slices with a 3-D stereotactic coordinate reference system. Rotational and translational parameters were acquired by coregistration of the first EPI magnetic resonance time step to the individual 3-D MPRAGE reference set. Anatomical datasets were normalized to the ICBM/MNI space by linear scaling. The resulting parameters were then used to transform all functional slices employing a trilinear interpolation. Resulting data had a spatial resolution of  $3 \times 3 \times 3$  mm ( $27$  mm<sup>3</sup>). Normalized functional images were spatially smoothed with a Gaussian kernel of 6 mm FWHM. A temporal high-pass filter of 1/128 HZ was applied to the data to remove low-frequency noise such as scanner drift. To prevent effects of physiological noise in the midbrain (e.g., pulsation artifacts), the component-based noise correction method (CompCor) was applied to the epoch-related analysis (see below) to reduce the temporal standard deviation of the BOLD signal (Behzadi, Restom, Liao, & Liu, 2007). During this application, the first few principal components of regions with high temporal variance are obtained and factored out via linear regression.

### Design Specification

Statistical analysis was based on a least squares estimation using the general linear model (GLM) for serially auto-correlated observations (Friston et al., 1995; Worsley &

Friston, 1995). Event- and epoch-related analyses were conducted in separated GLMs: The event-related analysis focused on BOLD signal changes during single trials to assess prediction error processing, whereas the epoch-related analysis included entire periods where either unstable or inflexible states were adjusted. Single trials and epochs were modeled as delta and box-car functions, respectively, and convolved with a canonical hemodynamic response function. In both models, the subject-specific six rigid body transformations obtained from residual motion correction were included as further covariates of no interest.

*Flexible and stable responses to prediction errors.* To analyze common and differential neural signatures of responses to both types of prediction errors, that is, switches and drifts, we calculated first-level regression models containing the specific events with an amplitude of one, that is, standard digits (STD), switches (SW), drifts (DR), and breaks of 14 sec. Only correct responses were analyzed, which comprised the full duration of the presented trial (1 sec). Because of the high event density, only events at a distance of at least two trials to the next modeled event were included. Furthermore, the GLM contained a separate regressor subsuming all button presses, that is, hits and false alarms. This controlled for motor response activity during switch versus drift processing inasmuch as only the first event type required a motor response whereas the second did not. This way, the switch effect, which would otherwise have concealed the drift effect, could be leveled. To provide further verification for the existence of the two networks, we calculated an additional GLM containing the specific error types, that is, switch misses and false alarms at drifts, which were contrasted with hits and correct rejections, respectively (for a detailed analysis, see Supplementary Material; <https://figshare.com/s/e0bfd1e93bf80e57f6f0>). Contrast images, that is, beta value estimates of the raw score differences between specified conditions, were generated for each participant.

*Thresholds for flexible and stable state transitions driven by recent performance history.* We expected midbrain dopaminergic activity to correspond to the experience-driven transitions to flexible and stable states. We therefore conducted an epoch-related analysis to assess neural activity emerging from performance history in response to different switch and drift probabilities in the recent past within an individually fitted time window. To estimate the participant-specific length of the time window, across which switch and drift probability were accumulated, logistic regressions were conducted for each subject. The regression model estimated the degree to which window length-based Shannon's surprise  $I(x_i)$  for switches and drifts (Shannon, 1948) could predict (variance in) response accuracy for both types of events (hits and correct rejections).

Shannon's surprise can be used as a measure of anticipation success because it is based on the frequency of an



event  $p(x_i)$  normalized by the sum of all event types over a defined window of recent events (see Equation 1).

Calculation of event probability:

$$p(x_i) = \frac{n(x_i) + 1}{\sum x_i + 1} \quad (1)$$

The surprise  $I(x_i)$  of each event given by the negative logarithm of this probability quantifies the amount of information provided by the current stimulus dependent on the history of previous stimuli (see Equation 2).

Calculation of Shannon's surprise:

$$I(x_i) = -\ln p(x_i) \quad (2)$$

We conducted these regressions for sliding windows with a minimum length of 20 trials and a trial-wise increase up to the mean predefined block length of 160 trials. Subsequently, the respective window length that provided the minimum deviance, that is, the difference between the log-likelihood of the fitted model and the maximum possible likelihood, was chosen as subject-specific epoch duration in the fMRI analysis. A paired  $t$  test between these individually obtained deviance values and deviances of models determining surprise as a function of mean block length of 160 trials was calculated. This was done to validate the significance of fit improvement by the individual sliding window length.

On the basis of the derived window length, we then calculated the difference of error type per window (number of misses – number of false alarms) to reflect the bias of prior performance. High values correspond to rare responding, resulting in a high miss as well as a high correct rejection rate. In contrast, low values reflect a bias toward frequent responses resulting in a higher hit, but also higher false alarm rate. The GLM underlying the fMRI analysis included the parametric effect of this bias on neural activity within the respective subsequent period as the main regressor of interest. We included mean surprise at switches and at drifts per window as nuisance regressors in the GLM. This ensured that the effects of the previous error bias on activity within a current period could not be attributed to the probability of critical events within this period.

### Group Analysis

To obtain group statistics, the resulting contrast images of all participants were entered into a second-level random-effects analysis using a one-sample  $t$  test across participants to test for significant deviation from zero. To assess differences between switch and drift processing, hits at switches were contrasted with drift rejections (SW > DR) and vice versa (DR > SW) in the GLM, which controlled for motor responses. Furthermore, we calculated the block-wise parametric effects for the error bias (BIAS). We corrected for multiple comparisons across all voxels using the threshold-free cluster enhancement (TFCE)

method (Smith & Nichols, 2009). The significance level for whole-brain activations was set to  $p < .05$  TFCE-corrected. Default TFCE parameters  $H = 2$  and  $E = 0.5$  were used.

### ROI Analysis

To test for a specific role of the caudate nucleus in global prediction error processing, we assessed overlapping effects of switch- and drift-elicited activity using small volume correction (SVC) on the additive conjunction analysis [(SW > STD)  $\cap$  (DR > STD)] at  $p < .05$  TFCE-SVC-corrected. Anatomical masks of left and right caudate were defined based on the Automated Anatomical Labeling (AAL) atlas (Tzourio-Mazoyer et al., 2002). Beta values of these ROIs, that is, left and right caudate, were extracted and correlated with  $P_r$  index. Alpha level was Bonferroni-corrected.

To investigate effects of error bias on dopaminergic midbrain activity, left and right SN ROIs were used for SVC. These ROIs were derived from the probabilistic atlas of the BG (ATAG; Keuken et al., 2014).

## RESULTS

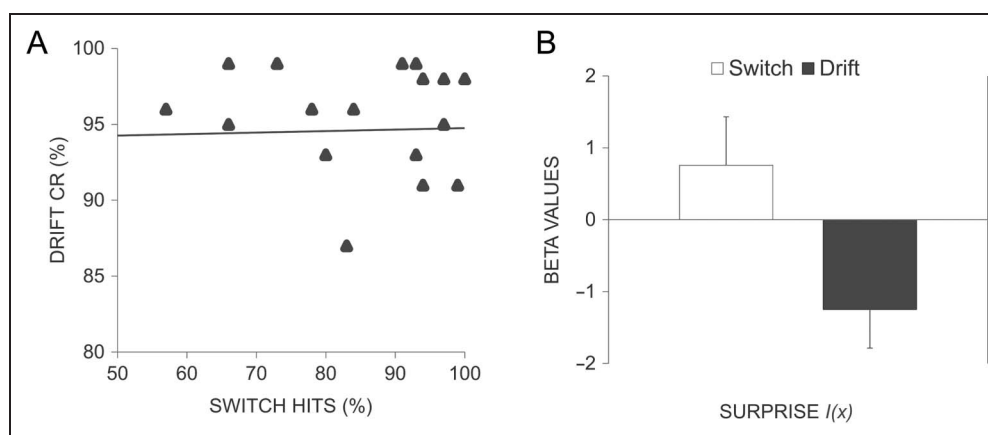
### Behavioral Results

To assess participants' performance and response tendency on the task, we used signal detection theory measures mean discrimination index ( $P_r$ ) and bias index ( $B_r$ ; Snodgrass & Corwin, 1988). Mean  $P_r$  index was  $M = 0.78$  ( $SD = 0.14$ ), and mean  $B_r$  index was  $M = 0.35$  ( $SD = 0.27$ ). This bias toward conservative response thresholds across the whole group also translated into a significant difference between hits at switches  $M = 84.31\%$  ( $SD = 12.60$ ) and correct rejections at drifts ( $M = 94.26\%$ ,  $SD = 5.73$ ),  $t(19) = 3.243$ ,  $p = .004$ . Mean RT at hits ( $M = 953$  msec,  $SD = 177$ ) did not differ significantly from mean RT at false alarms ( $M = 890$  msec,  $SD = 186$ ),  $t(19) = 1.569$ ,  $p = .133$ . As we expected the  $B_r$  value to reflect trait impulsivity of the participants, we conducted a correlational analysis on  $B_r$  score with BIS-11 total score that, however, did not reveal a significant association ( $r = .15$ ,  $p = .36$ ).

To test our hypothesis that cognitive flexibility and cognitive stability would be functionally independent, we performed a correlation analysis between the rate of hits (detected switches) and the rate of correct rejections (ignored drifts). This analysis revealed no significant relationship between the two variables ( $r = .025$ ,  $a(r) = 0.16$ ,  $p = .918$ ; Figure 2A), despite attenuation correction.

Individual drift and switch surprise rates were modeled on the basis of an individually fitted window of recent events for each subject. Mean trial number of the sliding window was  $M = 40.76$  ( $SD = 36.81$ ). The individually optimized logistic regressions fitted the data significantly better than those that were based on the mean block length of 160 trials ( $t(19) = 7.31$ ,  $p < .001$ ). Although participants showed an anticipation effect of drift processing,

**Figure 2.** (A) Scatter plot of nonsignificant correlation between performance on switches and drifts, measured as percent correct responses toward switches (hits) and percentage of correctly withheld responses toward drifts (correct rejections, CR). (B) Beta values signifying the relationship between error rate on switches and drifts and switch and drift surprise  $I(x)$ . Switch surprise was positively correlated with correct responses; a lower drift surprise, (i.e., anticipation of drifts) was accompanied by a lower error rate.



participants did not adapt to environments with a high amount of switches: High surprise of drifts was related to a high chance of making an error, whereas high switch surprise was accompanied by lower error chance (Figure 2B).

### Imaging Results

To assess differential cortical processing of switches and drifts, we contrasted BOLD signal changes between successfully detected switches and successfully rejected drifts and vice versa. We expected to find medial pFC (MPFC) and hippocampus in response to switches as a sign of updating predictions from long-term memory. In contrast, we hypothesized drift-specific activation in lateral prefrontal and premotor regions reflecting heightened demands on model stabilization.

#### Differential Cortical Networks for Flexible and Stable Responses

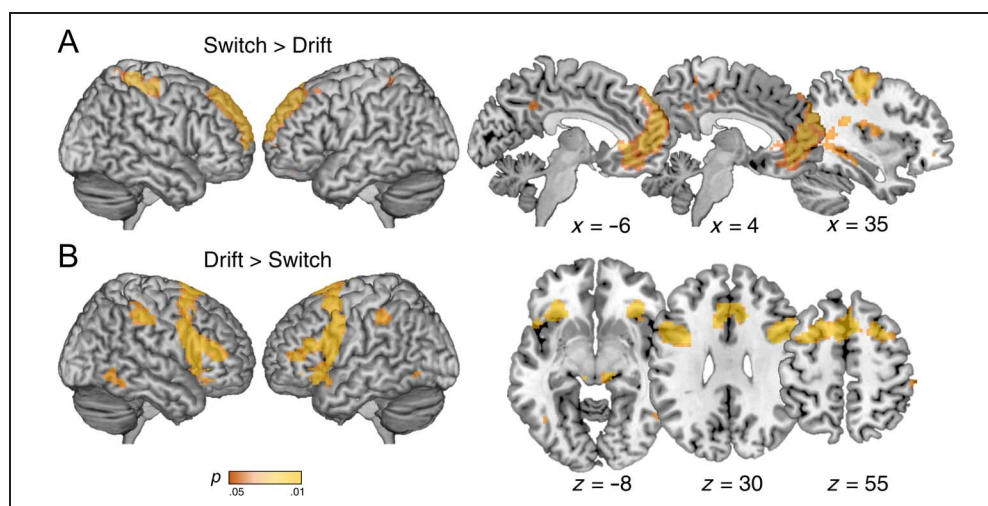
The switch-specific contrast estimate (SW > DR) revealed prefrontal activation of bilateral MPFC, including

BA 9 and BA 10, extending into dorsal ACC, as well as of the right hippocampus (Figure 3A). The reverse contrast, that is, successful rejection of drifts (DR > SW), yielded BOLD signal changes in bilateral SMA and premotor cortex, portions of the triangular and opercular parts of the inferior frontal gyrus (IFG), and anterior insula (Figure 3B; Table 1). Results of an alternative GLM excluding motor responses as regressor of no interest and of the GLM, in which error types were contrasted with correct responses, are included in the Supplementary Material.

#### Striatal Activity toward Prediction Errors Associated with Task Performance

To assess BOLD signal changes during the processing of prediction errors, an additive conjunction of the contrasts switch versus standard digit and drift versus standard digit (SW > STD)  $\cap$  (DR > STD) was computed. This contrast revealed the expected caudate nucleus activation at a threshold of  $p < .05$ , TFCE-corrected (R:  $x = 8, y = 18, z = -2$ ; L:  $x = -10, y = 18, z = -2$ ).

**Figure 3.** fMRI main effects at  $p < .05$ , whole-brain TFCE-corrected. (A) There was statistically significantly increased activation during switches in medial prefrontal and cingulate areas and in the right hippocampus. (B) Drifts elicited significant activation of BA 6 extending into IFG and anterior insula.



**Table 1.** fMRI Activations

Region	Side	BA	Cluster Size	MNI Coordinates			$p^a$
				$x$	$y$	$z$	
<i>Switch &gt; Drift</i>							
Superior frontal gyrus	L	9	432	-21	57	33	$10^{-4}$
Anterior cingulate cortex	L	10	67986	-6	48	0	$10^{-3}$
Postcentral gyrus <sup>b</sup>	R	3	65745	21	-36	63	$10^{-3}$
Cuneus	R	18	3024	-6	-93	27	$10^{-3}$
Inferior temporal gyrus	R	20	135	54	-27	-24	$10^{-2}$
<i>Drift &gt; Switch</i>							
Frontal inferior gyrus	R	48	1377	39	12	23	$10^{-4}$
Supplementary motor area	R/L	32	79002	0	15	48	$10^{-3}$
Supramarginal gyrus	R	40	999	54	-33	42	$10^{-3}$
Premotor cortex	L	6	81	-15	3	75	$10^{-2}$
Inferior parietal lobule	L	21	540	-48	-36	42	$10^{-2}$
Inferior temporal gyrus	R	20	189	51	-48	-15	$10^{-2}$
Angular gyrus	R	7	621	30	-57	45	$10^{-2}$

R = right; L = left;  $x$ ,  $y$ ,  $z$  = MNI coordinates of peak voxel activation.

<sup>a</sup>TFCE-corrected for multiple comparison.

<sup>b</sup>Extending into bilateral posterior cingulate cortex and right hippocampus.

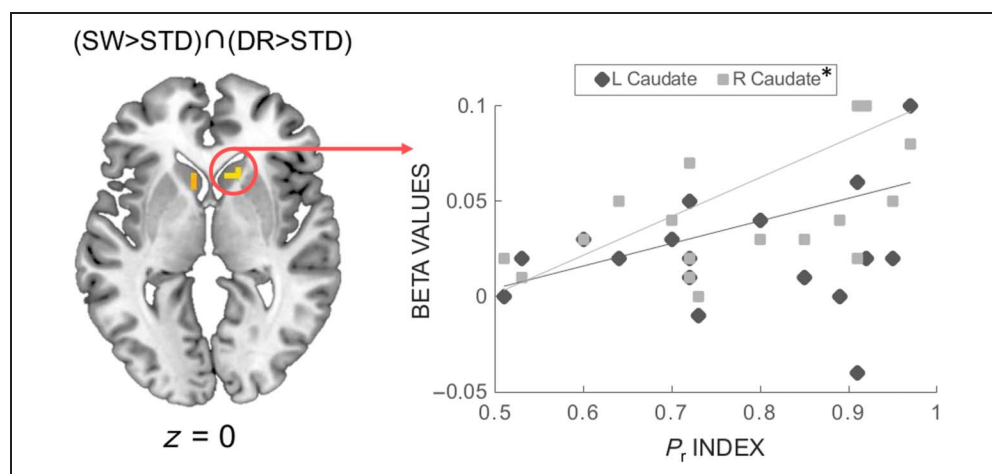
Parameter estimates of right caudate activity significantly correlated with  $P_r$  ( $r = .536$ ,  $p = .007$ , one-tailed) but not with  $B_r$  index ( $r = .275$ ,  $p = .121$ , one-tailed) at an adjusted alpha level of .025. Left caudate activity was not related to either of the two measures ( $P_r$ :  $r = .33$ ,  $p = .075$ , one-tailed;  $B_r$ :  $r = .24$ ,  $p = .153$ , one-tailed; Figure 4). This correlation of striatal activation with  $P_r$  value suggests that prediction error processing in the striatum is associated with the ability to make the correct re-

sponse toward surprising stimuli. This finding stands in contrast to the proposal that the striatum is limited to signaling any deviation from predictions per se.

#### *SN Activity Reflects Adaptation of Flexible and Stable State Transitions*

The epoch-based analysis revealed that activity in the left SN ( $x = -12$ ,  $y = -21$ ,  $z = -18$ ) significantly correlated

**Figure 4.** Results of caudate ROI analysis, reported at  $p < .05$ , TFCE-SVE-corrected. In the caudate ROIs, there was higher activation for both switches and drifts compared with standard digits. Activation during these critical events scaled with participants' ability to discriminate between events as indicated by behavioral  $P_r$  index.







the ability to shield the internal model from temporary violation, whereas flexibility comprises the ability to update the internal model because of permanent changes. If flexibility and stability were two extremes on the same dimension, we would expect that being inflexible leads to missing relevant events while helping to ignore irrelevant changes. At the same time, the ability to detect relevant events would lead to responses toward irrelevant prediction errors. In the present task, participants' performance could be driven by a set motor response threshold as shown in a negative correlation between switch detection and drift rejection. However, we did not find a corresponding relationship between flexible and stable responses, thus providing evidence that cognitive flexibility and stability rely on functionally independent processes. Because of the limited variability in our measure of cognitive stability that might have prevented us from finding a significant correlation between hits at switches and correct rejections of drifts, the finding of a missing relationship needs to be confirmed in future studies using a similar task.

### **Different Cortical Substrates of Cognitive Flexibility and Stability**

Our hypothesis that cognitive flexibility and stability rely on independent processes is substantiated by the differential processing of switches and drifts on the cortical level. As hypothesized, we found a network comprising medial prefrontal regions, that is, BA 9 and BA 10, and the right hippocampal formation in response to correctly detected model switches. This suggests that this network plays a role in the top-down adaptation of working memory representations by retrieval of an alternative predictive model from long-term memory. This finding significantly extends previous studies where mesial BA 9 was found in response to predictable events compared with destabilized predictions (Kühn & Schubotz, 2012). More rostral portions of the MPFC have been implicated in guiding hippocampal encoding and retrieval when new information is integrated into existing knowledge (Schlichting & Preston, 2015; Preston & Eichenbaum, 2013; van Kesteren, Ruiters, Fernández, & Henson, 2012). Our findings suggest that this network provides top-down predictions, which initiate an immediate update of currently valid internal models.

We further hypothesized drift-specific activation of lateral pFC, which would reflect model stabilization (D'Esposito, 2007; Bilder et al., 2004; Cohen et al., 2002; Miller & Cohen, 2001). Either dorsal or ventral prefrontal portions were expected to be activated, depending on whether model content was stored visuospatially, that is, by circular sequence representation, or phonologically (see Rottschy et al., 2012, for a meta-analysis). Furthermore, previous studies have shown that the prediction of sequences primarily activates premotor and not necessarily prefrontal regions (Schubotz, 2007; Schubotz

& von Cramon, 2003). Even increasing the demand on mnemonic representation of the sequence by occluding some of its stimuli did not recruit prefrontal areas but rather further boosted premotor activity (Schönberger, Hagelweide, Pelzer, Fink, & Schubotz, 2015). Thus, we expected the premotor network to be activated in response to irrelevant prediction errors. In line with our hypotheses, we observed that mental sequence completion and stabilization in case of drifts activated a network comprising SMA, premotor cortex, and portions of the IFG extending along the operculum into anterior insula. Because the common task strategy was subvocalization of the digit sequence, activation of IFG presumably reflects heightened verbal working memory load because of sequential interruptions (Fegen, Buchsbaum, & D'Esposito 2015; Shergill et al., 2002). Our results thus further substantiate that a network comprising not only lateral prefrontal but also premotor regions plays a basic role in the stabilization of predictive internal models that are stored in working memory.

### **Recognition of Prediction Error Types in the Striatum**

Although the behavioral consequences of different types of prediction errors requiring updating or shielding of predictions are processed in different neural networks, our study further provides evidence that all types of surprising events are captured by the same (pre-) attention control system, presumably gated through the striatum. A number of recent studies suggest that the striatum is responsible for flexible updating and adaptation of cortical representations (e.g., Stelzel, Fiebach, Cools, Tafazoli, & D'Esposito, 2013; Cools et al., 2007; Dreisbach & Goschke, 2004). Moreover, there is evidence for a basic stimulus selection function of the striatum when faced with unpredicted sensory input (e.g., den Ouden, Danizeau, Roiser, Friston, & Stephan, 2010; Corlett et al., 2004; O'Doherty et al., 2004; Seymour et al., 2004).

In this study, we observed that caudate activity is related to selecting the correct response following prediction errors rather than being driven by a global response probability bias. This suggests that specific behavioral implications of both event types have been or are recognized at this processing stage. This novel finding extends the proposed input gating function to control cognitive and motor representations in the pFC (Chatham & Badre, 2015; Badre, 2012; Cools, 2011). Our data thus integrate the above-mentioned findings, suggesting that flexible stimulus processing in the striatum might comprise a selection process, which can entail adaptation, stabilization, or building of internal models.

### **Adaptation of Flexible and Stable States in the SN**

Our results support the hypothesis that activity in the SN is associated with the transition between states of

cognitive flexibility and cognitive stability, driven by participants' performance in response to recent environmental demands. This finding delivers evidence in favor of recent computational models (Durstewitz & Seamans, 2008), which propose a dopaminergic modulation of cognitive states: A D1 receptor dominated state is associated with active maintenance, whereas DA action on D2 receptors promotes flexibility of the system. Furthermore, it has been proposed that tonic DA release balances transitions between different states by modulating the degree to which prior learning biases action selection (Humphries et al. 2012). The adjustment of each of the two states might thus be realized by presynaptic effects of tonic DA release: DA released in this manner acts as an inhibitory feedback signal and changes responsivity of the DA system in such a way as phasic responses in the dopaminergic midbrain become curtailed (O'Reilly & Frank, 2006; Schmitz, Benoit-Marand, Gonon, & Sulzer, 2003; Grace, 1991). Therefore, in this study, we associate the model of prefrontal function proposed by Durstewitz et al. (2000) with activation in the midbrain as a source of dopaminergic activity in pFC. We are the first to show that the disposition toward (too) stable or (too) flexible states differs not only between individuals but that these states also vary intraindividually dependent on the recent performance history in response to environmental demands.

## Conclusion

Taken together, our data provide evidence that flexible and stable responses in short-term prediction error processing correspond to functionally independent functions but share a common neural substrate in the striatum, responsible for stimulus discrimination and corresponding response decisions. MPFC is associated with model update, whereas lateral pFC stabilizes working memory when faced with distraction. Furthermore, the adaptation of different cognitive states, resulting from performance history, is modulated by SN activation, emphasizing the likely role of tonic DA in setting a threshold for corresponding state transitions. Future studies can build on these findings, especially as they might shed new light on DA-related diseases, for example, Parkinson disease, in which a deficient interplay of flexibility and stability may contribute to the phenotype.

## Acknowledgments

We thank Monika Mertens for her great help in conducting this investigation. Advice given by Irina Kaltwasser and Daniel Kluger has been a big help in improving the manuscript. We also thank Gabriele Lohmann and Johannes Stelzer for their great support in analyzing our data. We thank the German Research Foundation (Clinical Research Group KFO219 "Basal-Ganglia-Cortex-Loops: Mechanisms of Pathological Interactions and Therapeutic Modulation", SCHU 1439/5-2) for financially supporting the project.

Reprint requests should be sent to Ima Trempler, Fliednerstr. 21, 48149 Münster, Germany, or via e-mail: ima.trempler@uni-muenster.de.

## REFERENCES

- Armbruster, D. J., Ueltzhöffer, K., Basten, U., & Fiebach, C. J. (2012). Prefrontal cortical mechanisms underlying individual differences in cognitive flexibility and stability. *Journal of Cognitive Neuroscience*, *24*, 2385–2399.
- Baddeley, A. D. (1986). *Working memory*. Oxford, UK: Oxford University Press.
- Badre, D. (2012). Opening the gate to working memory. *Proceedings of the National Academy of Sciences, U.S.A.*, *109*, 19878–19879.
- Baker, J. E. (1987). Reducing bias and inefficiency in the selection algorithm. In Grefenstette J. (Ed.), *Proceedings of the Second International Conference on Genetic Algorithms and their Application* (pp. 14–21). Hillsdale, NJ: Erlbaum.
- Beeler, J. A., Daw, N., Frazier, C. R., & Zhuang, X. (2010). Tonic dopamine modulates exploitation of reward learning. *Frontiers in Behavioral Neuroscience*, *4*, 170.
- Behrens, T. E., Woolrich, M. W., Walton, M. E., & Rushworth, M. F. (2007). Learning the value of information in an uncertain world. *Nature Neuroscience*, *10*, 1214–1221.
- Behzadi, Y., Restom, K., Liu, J., & Liu, T. T. (2007). A component based noise correction method (CompCor) for BOLD and perfusion based fMRI. *Neuroimage*, *37*, 90–101.
- Bilder, R. M., Volavka, J., Lachman, H. M., & Grace, A. A. (2004). The catechol-O-methyltransferase polymorphism: Relations to the tonic-phasic dopamine hypothesis and neuropsychiatric phenotypes. *Neuropsychopharmacology*, *29*, 1943–1961.
- Chatham, C. H., & Badre, D. (2015). Multiple gates on working memory. *Current Opinion in Behavioral Sciences*, *1*, 23–31.
- Chumbley, J. R., Burke, C. J., Stephan, K. E., Friston, K. J., Tobler, P. N., & Fehr, E. (2014). Surprise beyond prediction error. *Human Brain Mapping*, *35*, 4805–4814.
- Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, *36*, 181–204.
- Cohen, J. D., Braver, T. S., & Brown, J. W. (2002). Computational perspectives on dopamine function in prefrontal cortex. *Current Opinion in Neurobiology*, *12*, 223–229.
- Cools, R. (2011). Dopaminergic control of the striatum for high-level cognition. *Current Opinion in Neurobiology*, *21*, 402–407.
- Cools, R., & D'Esposito, M. (2011). Inverted-U-shaped dopamine actions on human working memory and cognitive control. *Biological Psychiatry*, *69*, e113–e125.
- Cools, R., Sheridan, M., Jacobs, E., & D'Esposito, M. (2007). Impulsive personality predicts dopamine-dependent changes in frontostriatal activity during component processes of working memory. *Journal of Neuroscience*, *27*, 5506–5514.
- Corlett, P. R., Aitken, M. R., Dickinson, A., Shanks, D. R., Honey, G. D., Honey, R. A., et al. (2004). Prediction error during retrospective revaluation of causal associations in humans: fMRI evidence in favor of an associative model of learning. *Neuron*, *44*, 877–888.
- D'Ardenne, K., Eshel, N., Luka, J., Lenartowicz, A., Nystrom, L. E., & Cohen, J. D. (2012). Role of prefrontal cortex and the midbrain dopamine system in working memory updating. *Proceedings of the National Academy of Sciences, U.S.A.*, *109*, 19900–19909.
- den Ouden, H. E., Danizeau, J., Roiser, J., Friston, K. J., & Stephan, K. E. (2010). Striatal prediction error modulates cortical coupling. *Journal of Neuroscience*, *30*, 3210–3219.

- D'Esposito, M. (2007). From cognitive to neural models of working memory. *Philosophical Transactions of the Royal Society of London, Series B, Biological Sciences*, *362*, 761–772.
- D'Esposito, M., & Postle, B. R. (2015). The cognitive neuroscience of working memory. *Annual Review of Psychology*, *66*, 115–142.
- Draheim, C., Hicks, K. L., & Engle, R. W. (2016). Combining reaction time and accuracy: The relationship between working memory capacity and task switching as a case example. *Perspectives on Psychological Science*, *11*, 133–155.
- Dreisbach, G., & Goschke, T. (2004). How positive affect modulates cognitive control: Reduced perseveration at the cost of increased distractibility. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *30*, 343–353.
- Durstewitz, D., & Seamans, J. K. (2008). The dual-state theory of prefrontal cortex dopamine function with relevance to catechol-o-methyltransferase genotypes and schizophrenia. *Biological Psychiatry*, *64*, 739–749.
- Durstewitz, D., Seamans, J. K., & Sejnowski, T. J. (2000). Neurocomputational models of working memory. *Nature Neuroscience*, *3*, 1184–1191.
- Fegen, D., Buchsbaum, B. R., & D'Esposito, M. (2015). The effect of rehearsal rate and memory load on verbal working memory. *NeuroImage*, *105*, 120–131.
- Frank, M. J., Loughry, B., & O'Reilly, R. C. (2001). Interactions between the frontal cortex and basal ganglia in working memory: A computational model. *Cognitive, Affective, and Behavioral Neuroscience*, *1*, 137–160.
- Friston, K. J., Daunizeau, J., & Kiebel, S. J. (2009). Reinforcement learning or active inference? *PLoS One*, *4*, e6421.
- Friston, K. J., Holmes, A. P., Worsley, K. J., Poline, J.-P., Frith, C. D., & Frackowiak, R. S. J. (1995). Statistical parametric maps in functional imaging: A general linear approach. *Human Brain Mapping*, *2*, 189–210.
- Goschke, T., & Dreisbach, G. (2008). Conflict-triggered goal-shielding attenuates background-monitoring for prospective memory cues. *Psychological Science*, *19*, 25–32.
- Grace, A. A. (1991). Phasic versus tonic dopamine release and the modulation of dopamine system responsivity: A hypothesis for the etiology of schizophrenia. *Neuroscience*, *41*, 1–24.
- Hedden, T., & Gabrieli, J. D. (2015). Shared and selective neural correlates of inhibition, facilitation, and shifting processes during executive control. *NeuroImage*, *51*, 421–431.
- Humphries, M. D., Khamassi, M., & Gurney, K. (2012). Dopaminergic control of the exploration-exploitation trade-off via the basal ganglia. *Frontiers in Neuroscience*, *6*, 9.
- Jiang, J., Beck, J., Heller, K., & Egner, T. (2015). An insula-frontostriatal network mediates flexible cognitive control by adaptively predicting changing control demands. *Nature Communications*, *6*, 8165.
- Keuken, M. C., Bazin, P. L., Crown, L., Hootsmans, J., Laufer, A., Müller-Axt, C., et al. (2014). Quantifying inter-individual anatomical variability in the subcortex using 7 T structural MRI. *NeuroImage*, *94*, 40–46.
- Kühn, A. B., & Schubotz, R. I. (2012). Temporally remote destabilization of prediction after rare breaches of expectancy. *Human Brain Mapping*, *33*, 1812–1820.
- Lavie, N. (2005). Distracted and confused?: Selective attention under load. *Trends in Cognitive Sciences*, *9*, 75–82.
- Lohmann, G., Müller, K., Bosch, V., Mentzel, H., Hessler, S., Chen, L., et al. (2001). LIPSI—A new software system for the evaluation of functional magnetic resonance images of the human brain. *Computerized Medical Imaging and Graphics*, *25*, 449–457.
- Miller, E. K., & Cohen, J. D. (2001). An integrative theory of prefrontal cortex function. *Annual Review of Neuroscience*, *24*, 167–202.
- Miyake, A., Friedman, N. P., Emerson, M. J., Witzki, A. H., Howerter, A., & Wager, T. D. (2000). The unity and diversity of executive functions and their contributions to complex “frontal lobe” tasks: A latent variable analysis. *Cognitive Psychology*, *41*, 49–100.
- Müller, J., Dreisbach, G., Goschke, T., Hensch, T., Lesch, K. P., & Brocke, B. (2007). Dopamine and cognitive control: The prospect of monetary gains influences the balance between flexibility and stability in a set-shifting paradigm. *European Journal of Neuroscience*, *26*, 3661–3668.
- O'Doherty, J., Dayan, P., Schultz, J., Deichmann, R., Friston, K., & Dolan, R. J. (2004). Dissociable roles of ventral and dorsal striatum in instrumental conditioning. *Science*, *304*, 452–454.
- O'Reilly, R. C., & Frank, M. J. (2006). Making working memory work: A computational model of learning in the prefrontal cortex and basal ganglia. *Neural Computation*, *18*, 283–328.
- Patton, J. H., Stanford, M. S., & Barratt, E. S. (1995). Factor structure of the Barratt Impulsiveness Scale. *Journal of Clinical Psychology*, *51*, 768–774.
- Preston, A. R., & Eichenbaum, H. (2013). Interplay of hippocampus and prefrontal cortex in memory. *Current Biology*, *23*, R764–R773.
- Redgrave, P., & Gurney, K. (2006). The short-latency dopamine signal: A role in discovering novel actions? *Nature Reviews Neuroscience*, *7*, 967–975.
- Rottschy, C., Langner, R., Dogan, I., Reetz, K., Laird, A. R., Schulz, J. B., et al. (2012). Modelling neural correlates of working memory: A coordinate-based meta-analysis. *NeuroImage*, *60*, 830–846.
- Schiffer, A. M., Ahlheim, C., Wurm, M. F., & Schubotz, R. I. (2012). Surprised at all the entropy: Hippocampal, caudate and midbrain contributions to learning from prediction errors. *PLoS One*, *7*, e36445.
- Schiffer, A. M., Waszak, F., & Yeung, N. (2015). The role of prediction and outcomes in adaptive cognitive control. *Journal of Physiology Paris*, *109*, 38–52.
- Schlichting, M. L., & Preston, A. R. (2015). Hippocampal-medial prefrontal circuit supports memory updating during learning and post-encoding rest. *Neurobiology of Learning and Memory*, *134*, 91–106.
- Schmitz, Y., Benoit-Marand, M., Gonon, F., & Sulzer, D. (2003). Presynaptic regulation of dopaminergic neurotransmission. *Journal of Neurochemistry*, *87*, 273–289.
- Schönberger, A. R., Hagelweide, K., Pelzer, E. A., Fink, G. R., & Schubotz, R. I. (2015). Motor loop dysfunction causes impaired cognitive sequencing in patients suffering from Parkinson's disease. *Neuropsychologia*, *77*, 409–420.
- Schubotz, R. I. (2007). Prediction of external events with our motor system: Towards a new framework. *Trends in Cognitive Sciences*, *11*, 211–218.
- Schubotz, R. I., & von Cramon, D. Y. (2003). Functional-anatomical concepts of human premotor cortex: Evidence from fMRI and PET studies. *NeuroImage*, *20*, 120–131.
- Schultz, W., & Dickinson, A. (2000). Neuronal coding of prediction errors. *Annual Review of Neuroscience*, *23*, 473–500.
- Seymour, B., O'Doherty, J. P., Dayan, P., Koltzenburg, M., Jones, A. K., Dolan, R. J., et al. (2004). Temporal difference models describe higher-order learning in humans. *Nature*, *429*, 664–667.
- Shannon, C. E. (1948). A mathematical theory of communication. *The Bell System Technical Journal*, *27*, 379–423.
- Shergill, S. S., Brammer, M. J., Fukuda, R., Bullmore, E., Amaro, E., Jr., Murray, R. M., et al. (2002). Modulation of activity in

- temporal cortex during generation of inner speech. *Human Brain Mapping*, *16*, 219–227.
- Smith, S. M., & Nichols, T. E. (2009). Threshold-free cluster enhancement: Addressing problems of smoothing, threshold dependence and localisation in cluster inference. *Neuroimage*, *44*, 83–98.
- Snodgrass, J. G., & Corwin, J. (1988). Pragmatics of measuring recognition memory: Applications to dementia and amnesia. *Journal of Experimental Psychology: General*, *117*, 34–50.
- Spearman, C. (1904). The proof and measurement of association between two things. *American Journal of Psychology*, *15*, 72–101.
- Stelzel, C., Fiebach, C. J., Cools, R., Tafazoli, S., & D'Esposito, M. (2013). Dissociable fronto-striatal effects of dopamine D2 receptor stimulation on cognitive versus motor flexibility. *Cortex*, *49*, 2799–2811.
- Tzourio-Mazoyer, N., Landeau, B., Papathanassiou, D., Crivello, F., Etard, O., Delcroix, N., et al. (2002). Automated anatomical labeling of activations in SPM using a macroscopic anatomical parcellation of the MNI MRI single-subject brain. *Neuroimage*, *15*, 273–289.
- van Kesteren, M. T., Ruitter, D. J., Fernández, G., & Henson, R. N. (2012). How schema and novelty augment memory formation. *Trends in Neurosciences*, *35*, 211–219.
- Vandierendonck, A., Liefvooghe, B., & Verbruggen, F. (2010). Task switching: Interplay of reconfiguration and interference control. *Psychological Bulletin*, *136*, 601–626.
- Worsley, K. J., & Friston, K. J. (1995). Analysis of fMRI time-series revisited—Again. *Neuroimage*, *2*, 173–181.
- Yu, Y., FitzGerald, T. H., & Friston, K. J. (2013). Working memory and anticipatory set modulate midbrain and putamen activity. *Journal of Neuroscience*, *33*, 14040–14047.