

Mnemonic Introspection in Macaques Is Dependent on Superior Dorsolateral Prefrontal Cortex But Not Orbitofrontal Cortex

Sze Chai Kwok (郭思齊),^{1,2,3} Yudian Cai (蔡禹甸),¹ and Mark J. Buckley⁴

¹Shanghai Key Laboratory of Brain Functional Genomics, Key Laboratory of Brain Functional Genomics Ministry of Education, School of Psychology and Cognitive Science, East China Normal University, Shanghai 200062, China, ²Shanghai Key Laboratory of Magnetic Resonance, East China Normal University, Shanghai 200062, China, ³NYU-ECNU Institute of Brain and Cognitive Science at NYU Shanghai, Shanghai 200062, China, and ⁴Department of Experimental Psychology, University of Oxford, Oxford OX1 3SR, United Kingdom

The human PFC has been associated more with meta-perceptual as opposed to meta-memory decisions from correlational neuroimaging investigations. Recently, metacognitive abilities have also been shown to be causally dependent upon anterior and dorsal PFC in nonhuman primate lesion studies. Two studies, using postdecision wagering paradigms and reversible inactivation, challenged this meta-perceptual versus meta-memory notion and showed that dorsal and anterior prefrontal areas are associated with metamemory for experienced objects and awareness of ignorance, respectively. Causal investigations are important but scarce; nothing is known, for example, about the causal contributions of prefrontal subregions to spatial metamemory. Here, we investigated the effects of dorsal versus ventral PFC lesions on two-alternative forced-choice spatial discrimination tasks in male macaque monkeys. Importantly, we were rigorous in approach and applied three independent but complementary indices used to quantify individual animals' metacognitive ability ("Type II sensitivity") by two variants of *meta-d'*/*d'* and phi coefficient (ϕ). Our results were consistent across indices: while neither lesions to superior dorsolateral PFC nor orbitofrontal cortex impaired spatial recognition performance, only monkeys with superior dorsolateral PFC lesions were impaired in meta-accuracy. Together with the observation that the same orbitofrontal cortex lesioned monkeys were impaired in updating rule value in a Wisconsin Card Sorting Test analog, we therefore document a functional double-dissociation between these two PFC regions. Our study presents important causal evidence that other dimensions, namely, domain-specific processing (e.g., spatial vs nonspatial metamemory), also need considerations in understanding the functional specialization in the neural underpinnings of introspection.

Key words: introspection; lesion; macaques; metacognition; prefrontal cortex; recognition memory

Significance Statement

This study demonstrates macaque monkeys' metacognitive capability of introspecting its own memory success is causally dependent on intact superior dorsolateral prefrontal cortices but not the orbitofrontal cortices. Combining neurosurgical techniques on monkeys and state-of-the-art measures of metacognition, we affirm a critical role of the PFC in supporting spatial meta-recognition memory and delineate functional specificity within PFC for distinct elements of metacognition.

Introduction

Metacognition refers to awareness of one's own cognition (e.g., knowledge of one's accuracy, knowing what one knows, or in-

deed knowing when one does not know). Human neuroimaging has generally associated PFC more with meta-judgments based on perceptual as opposed to memory decisions (Fleming et al., 2012; Morales et al., 2018), backed up by structural neuroimaging measures (Fleming et al., 2010), and neuropsychology (Rounis et al., 2010; Shekhar and Rahnev, 2018). Nonetheless, two recent studies showed that neural activations in dorsal PFC and anterior PFC in the macaque brain are associated with metamemory of

Received Feb. 8, 2019; revised May 6, 2019; accepted May 7, 2019.

Author contributions: S.C.K., Y.C., and M.J.B. designed research; S.C.K., Y.C., and M.J.B. performed research; S.C.K. and Y.C. analyzed data; S.C.K. and Y.C. wrote the first draft of the paper; S.C.K., Y.C., and M.J.B. edited the paper; S.C.K. and M.J.B. wrote the paper.

This work was supported by the Ministry of Education of PRC Humanities and Social Sciences Research Grant 16YJC190006, STCSM Natural Science Foundation of Shanghai 16ZR1410200, Fundamental Research Funds for the Central Universities 2018ECNU-HWFW007, NYU Shanghai and NYU-ECNU Institute of Brain and Cognitive Science at NYU Shanghai to S.C.K., and the United Kingdom Medical Research Council Project Grants G0300817 and MR/K005480/1 to M.J.B. We thank Dr. Farshad Mansouri and Dr. Keiji Tanaka for allowing us to analyze the data that they obtained in some animals at the RIKEN Brain Science Institute.

The authors declare no competing financial interests.

Correspondence should be addressed to Sze Chai Kwok at sze-chai.kwok@st-hughs.oxon.org.
<https://doi.org/10.1523/JNEUROSCI.0330-19.2019>

Copyright © 2019 Kwok et al.

This is an open-access article distributed under the terms of the Creative Commons Attribution License Creative Commons Attribution 4.0 International, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

experienced object recognition (Miyamoto et al., 2017) and awareness of ignorance of experience of objects, respectively (Miyamoto et al., 2018), implying a more complex functional architecture. Indeed, pharmacological intervention delineated three subareas supporting metamemory: one for temporally remote items in dorsal area 9 (or 9/46d), a more posterior one for more recent items in area 6, and a third for awareness of ignorance in the most anterior part of PFC, namely, area 10 (frontopolar cortex).

Recent human neuroimaging associates frontopolar cortex with metacognitive control processes, whereas other metacognitive processes underlying decision-making per se are associated with dissociable neural systems more posteriorly within PFC (Qiu et al., 2018). The orbitofrontal cortex (OFC) is one such hub associated with value-based decision-making (Buckley et al., 2009; Noonan et al., 2010; Baltz et al., 2018), inferring the consequences of potential behavior (Schuck et al., 2016), and decision confidence (Kepecs et al., 2008), such that lesions to the OFC might affect decision confidence without affecting first-order task performance (Lak et al., 2014). Since confidence estimation is a fundamental component of decision-making, and the OFC has been implicated in goal-directed decisions that require the evaluation of predicted outcomes (Rudebeck and Murray, 2014), OFC might be causally required to support the computation of some dissociable elements of metacognition (Kepecs et al., 2008; Lak et al., 2014). Moreover, connections differ, for example, posterior parietal areas of the brain involved in egocentric spatial processing have more robust connections with dorsal PFC, whereas temporal lobe areas implicated in object-identity processing have more robust connection with ventral PFC (Yeterian et al., 2012). In light of the above and findings that second-order metacognitive processes could be separated from confidence per se (Dotan et al., 2018), we investigated the causal roles of one dorsal (sdIPFC) and one ventral (OFC) subregion of PFC in spatial recognition memory and hypothesized that the sdIPFC but not OFC will be causally required for accurate spatial memory introspection. Specifically, we contrasted the first-order memory and second-order metamemory performances of macaques with superior dorsolateral PFC (sdIPFC) lesion (i.e., lateral area 9) ($n = 3$), or with OFC-lesioned monkeys ($n = 3$), to unoperated controls ($n = 7$) (Fig. 1) in two delayed-matching-to-position (DMP) spatial recognition tasks (Fig. 2). On the basis of a wide-ranging PFC lesion study literature review in macaque monkeys, we expect that neither lesion would impair first-order spatial recognition per se; for example, OFC lesions do not impair spatial delayed response (Meunier et al., 1997), and unlike lesions in the region of the principal sulcus, lesions more dorsal to area 9 and 8B do not impair spatial delayed response (Goldman et al., 1971; Levy and Goldman-Rakic, 1999).

Cognizant of confidence reporting not being available in our task design, we used response time as a proxy (Dotan et al., 2018) for confidence; but as this is a rough proxy, we accordingly adopted a rigorous analytical approach and applied three quite different, but complementary, indices used to quantify individual animals' metacognitive ability (Type II sensitivity) by two variants of $meta-d'/d'$ and phi coefficient (ϕ). We found that sdIPFC, but not OFC, lesioned monkeys were impaired in meta-accuracy in a high spatial memory demand task variant without showing any impairments in spatial recognition performance itself. We further established that these putative metacognitive deficits were specific to spatial recognition memory rather than to other confounds, such as rule learning, reward evaluation, or general representation of task information using analyses of existing data from the same sdIPFC-lesioned mon-

keys when previously tested on a Wisconsin Card Sorting Test (WCST) analog (Buckley et al., 2009).

Materials and Methods

Animals. Thirteen adult male macaque monkeys provided data to the main experiment (8 *Macaca mulatta*, 3 *Macaca fuscata*, and 2 *Macaca fascicularis*): Three monkeys had orbitofrontal lesion (OFC, consisting of two *M. mulatta* and one *M. fuscata*), another 3 monkeys had superior dorsolateral prefrontal cortex lesion (sdIPFC, consisting of one *M. mulatta* and two *M. fuscata*), and 7 served as unoperated controls (CON, consisting of 5 *M. mulatta* and 2 *M. fascicularis*) for the main tasks. Data from 3 further CON (all *M. fuscata*) were included only for the WCST analog analysis. In a total of 16 animals, 10 macaque monkeys were trained, operated, and tested in Oxford, UK, and 6 in RIKEN Brain Science Institute, Wako, Japan. All animal training, surgery, and experimental procedures were the same in both laboratories. Those conducted in the United Kingdom were licensed in compliance with the UK Animals (Scientific Procedures) Act 1986, and those in Japan were done in accordance with the guidelines of the Japanese Physiological Society and approved by RIKEN's Animal Experiment Committee.

Surgery. The operations were performed in sterile conditions (under gaseous general anesthesia, and with preoperative, perioperative, and postoperative analgesia) with the aid of an operating microscope, and the same surgeon performed all operations in both laboratories. Detailed description of the surgical procedures has been reported previously (Buckley et al., 2009). The intended extent of the sdIPFC lesion was designed to include the cortex on the dorsolateral aspect of the PFC extending up to midline (i.e., lateral area 9 and the dorsal portions of areas 46 and 9/46) but excluding ventrally situated dlPFC cortex; the lesion excluded posteriorly located premotor areas 8A, 8Bd, and 8Bv, nor did it extend anteriorly into area 10. The intended extent of the OFC lesion included, at its lateral extent, the cortex in the medial bank of the lateral orbital sulcus; the lesion included all of the cortex between the medial and lateral orbital sulci, and also extended medially until the lateral bank of the rostral sulcus. The anterior extent of the lesion was an imaginary line drawn between the anterior tips of the lateral and medial orbital sulci, and the posterior extent was an imaginary line drawn just anterior to the posterior tips of these two sulci. The intended lesion therefore included areas 11, 13, and 14 of the orbital surface and did not extend posteriorly into the agranular insula.

Experimental design and statistical analysis. Metamemory quantification: Three different but complimentary indices were used to quantify individual animals' metacognitive ability ("Type II sensitivity"), as defined as the ability to accurately link confidence with performance. Here, we calculated the $meta-d'/d'$, a metric for estimating the metacognitive efficiency (level of metacognition given a particular level of performance or signal processing capacity), which enables a model-based approach to the computation of Type II sensitivity that is independent of response bias and Type I sensitivity (d') on the primary task (Rounis et al., 2010). We also computed metacognitive efficiency using a hierarchical Bayesian estimation method, which can avoid edge-correction confounds and enhance statistical power (Fleming and Daw, 2017). Both $meta-d'$ and d' measures assume that the variance of the internal response takes a Gaussian distribution, and that the distributions associated with the two Type I responses, respectively, are of equal variance. To ensure our results were not due to any idiosyncratic violation of the assumptions of signal detection theory (SDT), we additionally calculated the ϕ , which does not make these parametric assumptions (Fleming and Lau, 2014), and supplemented the statistical analyses using nonparametric tests (Kruskal–Wallis).

Signal detection theoretic and hierarchical Bayesian estimation meta-index ($meta-d'/d'$). Using a Type II SDT toolbox (Maniscalco and Lau, 2012), which has been extensively used for evaluation of metacognitive ability (Baird et al., 2013; Fleming et al., 2014), it is possible to compute a measure of metacognitive accuracy that is unconfounded by Type I performance directly from the empirical Type II receiver operating characteristic curve. The Type II receiver operating characteristic curve reflects the relationship between the accuracy of Type I and the observer's confidence rating. This approach exploits the link between Type I and Type II SDT models to express observed Type II sensitivity at the level of the

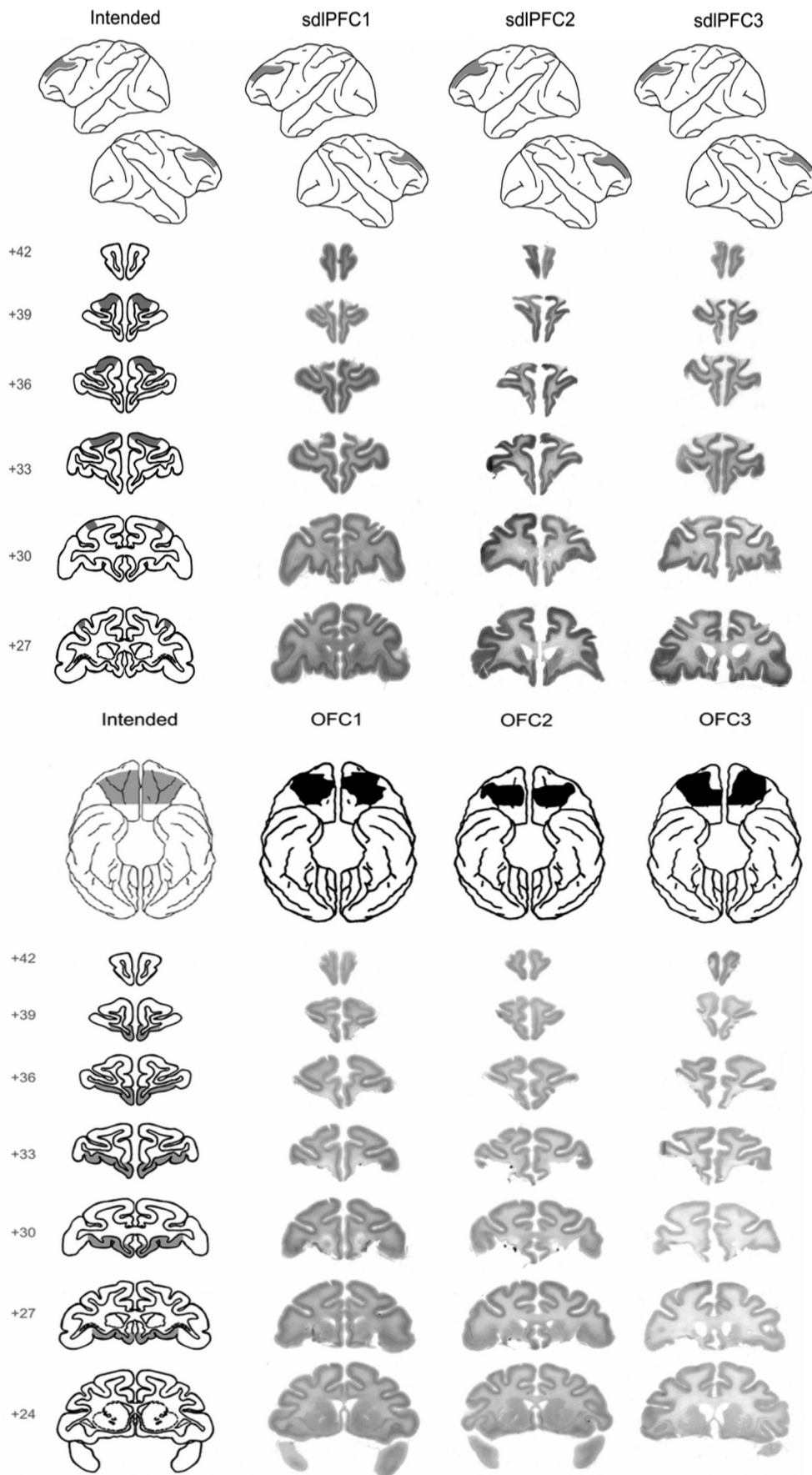


Figure 1. Histology. Top, Photomicrographs of stained coronal sections through the area of the intended lesion in the 3 animals with sdIPFC lesions (sdIPFC1 to sdIPFC3) alongside drawings of the intended extent of the lesions on drawings of representative coronal sections (left). Bottom, Photomicrographs of stained coronal sections through the area of the (Figure legend continues.)

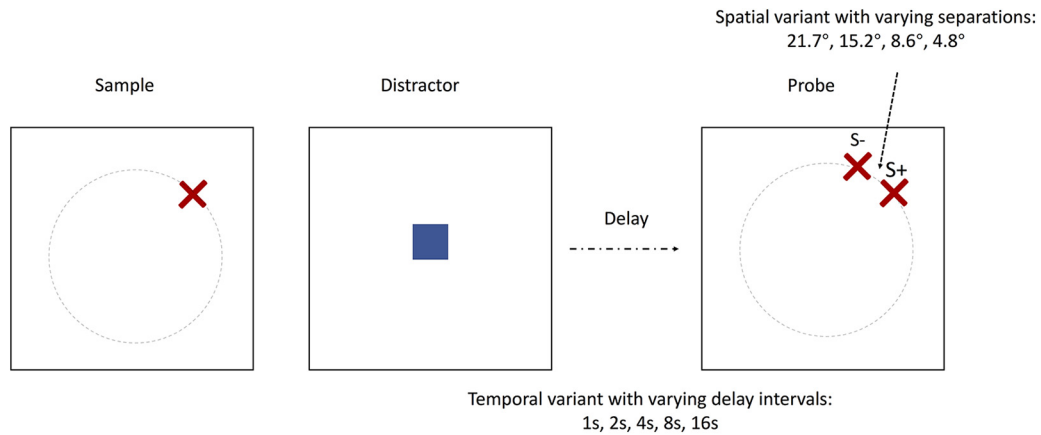


Figure 2. DMP tasks. Each trial consisted of a sample (red cross) stage, a distractor (blue square) stage, a delay, and then a probe/choice (2 red crosses) stage. Temporally taxing DMP: five levels of delay interval between distractor and probes (1, 2, 4, 8, or 16 s); spatially taxing variant DMP: four levels of separation between two red crosses in probe (visual angles of either 4.8°, 8.6°, 15.2°, or 21.7°, which are equivalent to 23, 16, 9, and 5 cm, respectively, on screen; all delay fixed at 1 s). S⁺, Target; S⁻, foil. Gray dotted circle in the figure is invisible to the animal and is just shown to illustrate the two choices are always equidistant to the distractor to obviate proximity bias.

Type I SDT model (termed *meta-d'*). Maximum likelihood estimation is used to determine the parameter values of the Type I SDT model that provide the best fit to the observed Type II data. A measure of metacognitive ability that controls for differences in Type I sensitivity is then calculated by taking the ratio of *meta-d'* and the Type I sensitivity parameter *d'*: meta-efficiency, computed as *meta-d'/d'*. The most straightforward approach to computing meta-efficiency involves an equal variance SDT model in which the variances of internal distributions of evidence for “target” and “foil” in the Type I model are assumed to be equal. We thus quantified metacognitive sensitivity with the SDT-based measure *meta-d'*. Based on Type II signal detection theory, meta-efficiency (in terms of *meta-d'/d'*) reflects how much information, in signal-to-noise units, provides a response-bias free measure of how well confidence ratings track task accuracy. Of note, the standard Type II SDT toolbox is designed for 2AFC tasks, in which S1 and S2 are always constant in the left or right, but our target and foil are randomly presented at the screen, and we only recorded the separation of the two probes and did not track the specific position of the target and the foil. Since the target and the foil were presented randomly on the screen, and the SDT algorithm only requires the distribution of those four kinds of trials, we divided the number of those trials equally to S1 trials and S2 trials in a random manner considering the animals would not have any preference to any given side/location of the screen. In addition, we have also replicated the analyses using a variant of metacognitive efficiency (*H-meta-d'*) with a hierarchical Bayesian estimation method, which can avoid edge-correction confounds and enhance statistical power (Fleming and Daw, 2017).

ϕ . To ensure our results were not due to any idiosyncratic violation of the parametric assumptions of SDT, we additionally calculated the ϕ coefficient index, which does not make the SDT assumptions. The ϕ coefficient is a contingency index of preference for optimal choice (Kornell et al., 2007; Middlebrooks and Sommer, 2011) and was calculated according to the following formula using the number of trials classified in each case [n(case)]:

phi coefficient(ϕ)

$$= \frac{n(\text{Correct High}) \times n(\text{Incorrect Low}) - n(\text{Correct Low}) \times n(\text{Incorrect High})}{\sqrt{n(\text{Correct}) \times n(\text{Incorrect}) \times n(\text{High}) \times n(\text{Low})}}$$

←

(Figure legend continued.) intended lesion in the 3 animals with OFC lesions (OFC1 to OFC3) alongside drawings of the intended extent of the lesions on drawings of representative coronal sections (left). Top row of each panel represents reconstructions of the area of cortex lesioned on drawings of representative lateral and ventral surfaces. Numerals indicate distance in millimeters from the interaural plane. This figure is based on data from, and a supplementary figure in, Buckley et al. (2009).

The ϕ coefficient evaluates how optimally each trial was assigned for high or low confidence based on performance in the preceding cognitive judgment, reflecting the correlation between the two binary variables. Despite differences in their mathematical assumptions, the three metacognitive metrics are highly correlated with each other.

For the computation for SDT *meta-d'/d'*, hierarchical-model *meta-d'/d'*, and ϕ , four types of trials and their distribution are required. The computation performed here are based on the premise that confidence is computed in a retrospective manner (Pleskac and Busemeyer, 2010). Using a summary of the decision process, trials that are responded fast are judged as more of higher certainty (Kiani et al., 2014). Following this logic, we accordingly used trial-specific reaction times (RTs) as a proxy for confidence (Dotan et al., 2018). We collapsed all trials per monkey and classified the trials within-monkeys by the median of all RT crossing with correct/incorrect responses into four kinds: correct/high confidence (fast RT), incorrect/low confidence (slow RT), correct/low confidence (slow RT), incorrect/high confidence (fast RT). For each of the final analyses, each monkey had one single value for the measurement of meta-ability. Admittedly, measures of association, such as ϕ , are prone to metacognitive bias. In contrast, theoretically, measures based on signal detection theory, such as *meta-d'* which was used here, are bias-free. By using standard SDT, Type 2 *d'* is argued to be independent from metacognitive bias (Fleming and Lau, 2014). Therefore, our SDT model *meta-d'* should help alleviate this issue. It should be noted that all the statistical analyses on meta-cognition were performed within-subjects; thus, the numerally (but not statistically significant) faster RT for the sdlPFC monkeys would not affect our main findings based on confidence for trial classification.

Preanalysis. For the formal analyses of the main experiments, trials with RT longer than 20 s and shorter than 100 ms in the memory judgment were discarded (<0.5% of all trials). We also set a stringent selection of good trials on which the monkeys were attentive, crucial for the metacognition analysis. To this end, we set a requirement for the monkeys to touch the distractor before the memory task so as to ensure they were not distracted and/or less willing/ready before initiating the memory judgment. Trials with touch-distractor times longer than 1000 ms (15.0% and 20.2% trials discarded, respectively, for temporal-variant task and spatial-variant task) were not included for analysis. There were no differences in touch-distractor times between the groups in either of the tasks following this trial removal procedure, with a one-way ANOVA on temporal-variant ($F_{(2,10)} = 1.27, p = 0.322$) and on spatial-variant ($F_{(2,9)} = 0.800, p = 0.479$). The results did not differ if we chose to use other touch-distractor times cutoff criteria of 800, 900, 1100, or 1200 ms.

Behavioral tasks. Spatial recognition tasks (DMP). A temporally demanding DMP task and a spatially demanding DMP task were performed by the monkeys. In both variants, each trial consisted of an encoding

phase in which a spatial position (“sample”) was indicated by a red cross. After the monkey touched the sample, a blue square (“distractor”) appeared in the center of an imaginary (i.e., invisible) circle whose circumference transected the center of the red cross. A touch to the blue square initiated a variable delay interval (i.e., the manipulation of delay for the temporal-variant) and then a choice phase consisting of two identical red crosses in different positions, both located on the circumference of the aforementioned imaginary circle (hence equidistant from the blue square distractor just touched), albeit one positioned in the same (i.e., spatial match) position as the first red cross and the other (i.e., non-match) positioned some angle (with respect to the center of the imaginary circle) away from the spatial match along the invisible circumference. From trial to trial, we could vary the angle of separation along the circumference to allow for easy trials (i.e., large angle and accordingly large spatial separation) and harder trials (i.e., smaller angles and accordingly smaller spatial separation) (compare the manipulation of separation between two probes for the spatial-variant). As mentioned above, one of the crosses appeared in the same position as the sample (target; S^+) and the other one in a different position (foil; S^-). A touch to the S^+ resulted in a delivery of a banana-flavored pellet as reward (see Apparatus), removed the S^- , and the S^+ remained alone for a further 1 s for positive feedback. The screen would then be blanked for an ITI of 6 s before the next trial. A touch to the S^- removed both S^+ and S^- from the screen, and the screen would be blanked for an ITI of 12 s. There was no time constraint imposed on responses made to the choices; therefore, there were no missed trials. No repetition correction routines were implemented following an error response; each trial was new and independent of the outcome of the preceding trial. In terms of sizes of visual stimuli, the sample subtended a visual angle of 9° in task acquisition and the temporal-variant task, or 6.8° in the spatial-variant task; the distractor subtended a visual angle of 4.6° in all tasks.

In the temporal-variant DMP, there were five trial types with differing delay intervals (1, 2, 4, 8, or 16 s) between the distractor and probes. Trials within a session were divided into five trial types with differing intervals of delay between the distractor and probes. The trial-type order was randomized within each successive set of five trials (with one trial of each trial-type per set) so that the delay changed unpredictably from one trial to another. The two probe choices were separated by a visual angle of 21.7° . In the spatial-variant DMP, the separation between two red crosses varied across trials; there were four different trial types with differing spatial separations (visual angles of 4.8° , 8.6° , 15.2° , or 21.7° ; equivalent to 5, 9, 16, or 23 cm, respectively, on screen) between probe choices. Delays were fixed at 1 s for the spatial-variant DMP. In the final testing, each animal accrued 200 rewards to complete the temporal-variant (across two daily sessions) and 150 rewards (one session) to complete the spatial-variant. Since the animals accrued varying numbers of errors to complete the tasks, the mean total numbers of trials were 271.2 and 209.8 trials (averaged across groups), respectively, for the two tasks. The analysis on metamemory was done by collapsing all trials across sessions; thus, the learning trend of these data was not considered. One CON did not complete the spatial-variant so only 12 monkeys were analyzed in the spatial-variant task.

WCST analog. Given that the DMP tasks involve multiple processes, which might confound our main results, we therefore analyzed extant data obtained from a WCST analog, which is a validated rule-guided task taxing multiprocesses, such as perception (involved in matching stimuli), memory and acquisition of abstract rules, and reward value evaluation (Buckley et al., 2009; Mansouri et al., 2015). We accordingly used some WCST data to rule out that the putative meta-deficits observed here were not attributable to these perceptual and reward value evaluation processes. The WCST analog paradigm is summarized as follows: on each trial, a randomly selected sample (a square, a circle, or a cross of different colors) is displayed alone on the center of the touch screen; and when the sample is touched, three additional choice items immediately appear (one matching in color, one matching in shape, and one not matching in either dimension), with their positions randomly chosen. If the animal's choice is correct (i.e., the animal selects the choice item that matches according to the currently reinforced rule, which changes unannounced every time the animal attains 85% in 20 consecutive trials), then a reward

pellet is delivered, and the correct choice remains on the screen for 1 s to provide visual feedback; if the animal makes an incorrect choice, then no reward is given, and the stimuli are removed and replaced by an error signal (white circle), which is presented on the screen for 1 s instead.

We analyzed WCST data from 12 monkey data points (9 CON vs 3 sdlPFC). Six of the 9 CON monkey data points here were from the pre-lesion data of the 6 lesioned monkeys (3 sdlPFC and 3 OFC). We included 3000 trials (acquired from 10 300-trial daily sessions) per monkey data point. We collapsed all trials and classified the trials into four types of trials for the computation for the meta-efficiency and ϕ . Since the Type II SDT toolbox was designed for 2AFC tasks, and the WCST task contained three stimuli, we ran three separate sets of computation, each one discarding only either the bottom, left, or right choice, for each of the three meta-indices. For each monkey, we then computed the mean of these three values as his meta-score to enter into the meta-indices calculation. In this analysis, we did not include the OFC monkeys because the OFC monkeys were severely impaired in the WCST Type I task, thus making any analyses on meta-ability invalid (their chance level implies they did not know how to make correct judgments, violating the prerequisite for the meta-assessment of their judgment). It is possible that the estimated metacognitive indices of each individual animal vary across tasks. However, since the most important contrast was on comparing between groups within the same tasks, we have ensured the numbers of trials included in each experiment to be comparable.

Preliminary training. All monkeys completed preliminary training and task acquisition before performing the two main tasks and WCST described above. We conducted the spatial-variant task immediately after the temporal-variant task without any additional training. The monkeys performed one session per day, 6–7 d per week. For the lesioned animals, the task was administered postoperatively (on average 22 months after lesion). For the two DMP tasks, during task acquisition the monkeys were trained until they reached $\geq 90\%$ performance level within a 100-reward session. All trials in this stage consisted of a short delay interval (1 s) and a wide separation between choice positions (21.7° , or 23 cm) to make the trials “easy” to acquire. Upon reaching criterion, the three groups were not different in the number of errors accrued ($F < 1$) and number of rewards received ($F < 1$), indicating that the groups of lesioned monkeys learned to perform these spatial recognition problems as well as controls.

Apparatus. The tasks were performed in an automated test apparatus. The subject sat, unrestrained, in a wheeled transport cage fixed in position in front of a touch-sensitive screen on which the stimuli could be displayed. The animals could reach out between the horizontal or vertical bars (spaced ~ 45 mm apart) at the front of the transport cage to touch the screen. An automated pellet delivery system delivered banana-flavored pellets (190 mg supplied by Noyes and Neuroscience) into a food well (~ 80 mm in diameter) positioned beneath and to one side of the screen, in response to correct choices made by the subject to the touch screen. Pellet delivery was accompanied by an audible click. A spring-loaded lunchbox (length 200 mm, width 100 mm, height 100 mm) was positioned beneath and to one side of the subject; this opened immediately with a loud crack on completion of the testing session and contained the subject's daily diet of wet monkey chow, primate pellets, nuts, raisins, and a slice of apple, banana, and orange (water was provided in the home cage *ad libitum*). An infrared camera allowed the subject to be observed while it was engaged in the task. The entire apparatus was housed in an experimental cubicle that was dark apart from the background illumination from the touch screen. A computer, with a millisecond accuracy timer-card to record RTs, controlled the experiment and data acquisition. Identical software controlled the tasks in both laboratories to ensure that the tasks were replicated exactly.

Histology and analysis on fiber tract damage. After the conclusion of the experiments, the animals with ablations were sedated, deeply anesthetized, and then perfused through the heart with saline solution (0.9%), which was followed by formol saline solution (10% formalin in 0.9% saline solution). The brains were blocked in the coronal stereotaxic plane posterior to the lunate sulcus, removed from the skull, allowed to sink in sucrose formalin solution (30% sucrose, 10% formalin), and sectioned coronally at $50 \mu\text{m}$ on a freezing 10 microtomes. Every 10th section

through the temporal lobe was stained with cresyl violet and mounted. When referring to cytoarchitecturally defined regions in the lesion description below, we have adopted the nomenclature and conventions of Petrides and Pandya (1999) and have reconstructed lesion extents on standard drawings based upon those provided by the Laboratory of Neuropsychology at National Institute of Mental Health. All three of the sdlPFC lesions were as intended. None of the OFC lesioned animals sustained any bilateral damage outside the area of the intended region; 2 animals sustained extremely slight unilateral damage beyond the intended lateral boundary of the lesion OFC2 and OFC3; in all 3 animals, the lesions did not extend as far medially as intended.

The lesion method in this study was careful aspiration, so it is important to consider, given recent observations of different effects of aspiration versus neurotoxic lesion (Izquierdo et al., 2004; Rudebeck et al., 2013), the possibility that white matter fibers may have been damaged by the sdlPFC lesion (either inadvertent damage to fiber bundles in underlying white matter proximal to the lesion, or damage to local fibers of passage), and we need to consider whether such damage may have contributed to the behavioral deficit seen after sdlPFC lesions. We will first comment briefly on the extent to which different white matter fibers in the frontal lobes may have been compromised by the sdlPFC aspiration lesion, and then follow-up with a summary conclusion:

Cingulum bundle (CB): fibers of CB do not course in the white matter proximal (i.e., just below) the gray matter removed in our aspirative lesion of sdlPFC, except for at the most anterior extent of the lesion (e.g., see Thiebaut de Schotten et al., 2012, their Fig. 1; Schmahmann and Pandya, 2006, their Fig. 10.1). These fibers are unlikely to have been significantly damaged by the sdlPFC lesion as the expert surgical approach we adopt with aspiration lesions, with the help of a high-magnification binocular operating microscope, is to stop immediately upon the underlying white matter starting to become visible as the gray matter is gradually thinned/removed; so while we cannot rule out that some underlying fibers proximal to the sdlPFC gray matter may have been inadvertently damaged, we expect any such damage to be both slight and also asymmetrical across hemispheres (lessening its likely causal impact upon behavior). White matter fibers of the CB also extend into the gray matter per se but only at the most anterior extent of the area of the sdlPFC lesion (see, e.g., Thiebaut de Schotten et al., 2012, their Fig. 1; Schmahmann and Pandya, 2006); these fibers will certainly have been removed with the sdlPFC gray matter aspiration, and so will have primarily affected CB connections to sdlPFC itself (which was lesioned in any case) and any that extend to frontal polar cortex (FPC) anterior to our sdlPFC lesion.

Superior longitudinal fasciculus (SLF): of the three branches of the SLF, the most dorsal branch runs in the white matter proximal (i.e., just below) the gray matter removed in our aspirative lesion of sdlPFC, albeit only at the most posterior extent of this lesion. The same is true for the SLF branch coursing more ventral to that aforementioned branch (e.g., see Thiebaut de Schotten et al., 2012, their Fig. 2; Schmahmann and Pandya, 2006). For the same reasons as above, there is unlikely to be significant or symmetrical damage to SLF. White matter fibers of these two SLF branches also extend into the gray matter per se, albeit only at the most posterior extent of the area of the sdlPFC lesion (see, e.g., Thiebaut de Schotten et al., 2012, their Fig. 2; Schmahmann and Pandya, 2006). These fibers will likely have been removed with the gray matter but primarily affect connections to sdlPFC itself (which was lesioned in any case) as these fibers do not extend more anteriorly into FPC.

Arcuate fasciculus (AF): few fibers of AF run in the white matter proximal (i.e., just below) the gray matter removed in our sdlPFC lesion (see, e.g., Thiebaut de Schotten et al., 2012, their Fig. 3; Schmahmann and Pandya, 2006). As explained above, however, there is unlikely to be significant symmetrical inadvertent damage to AF. A small proportion of AF fibers extend into the gray matter per se within the area of the sdlPFC lesion (see, e.g., Thiebaut de Schotten et al., 2012, their Fig. 3; Schmahmann and Pandya, 2006); these fibers will likely have been removed with the gray matter, and so will have primarily affected connections to sdlPFC itself (which was lesioned in any case) as these fibers do not extend more anteriorly into FPC.

Extreme capsule (EC): many fibers from branches of the EC run in the white matter proximal (i.e., just below) the gray matter removed in our

sdlPFC lesion (see, e.g., see Thiebaut de Schotten et al., 2012, their Fig. 4; Schmahmann and Pandya, 2006). Again, there is unlikely to be significant or symmetrical inadvertent damage to EC. Some EC fibers extend into the gray matter per se within the area of the sdlPFC lesion (see, e.g., Thiebaut de Schotten et al., 2012, their Fig. 4; Schmahmann and Pandya, 2006); these fibers will have been removed along with the sdlPFC gray matter; but as the majority of those fibers do not extend more anteriorly into FPC anterior to our lesion, this will have primarily affected connections to sdlPFC itself. Some EC fibers that do connect to more anterior FPC run in a different and more ventrally situated branch of the EC that would not have been transected by our cortical aspirative lesion.

Uncinate fasciculus: this white matter tract connects FPC (anterior to our sdlPFC lesions) to temporal lobe cortical regions; but as uncinate fasciculus fibers do not course through the sdlPFC nor in its underlying region, the uncinate fasciculus could not have been affected by the sdlPFC aspiration lesion.

Frontal aslant tract: the frontal aslant tract connects dorsal to ventral frontal areas; but as the frontal aslant tract courses just posterior to the posterior extent of our sdlPFC lesion, it would not have been affected by our sdlPFC lesion.

Orbito-polar tract: the orbito-polar tract connects OFC to FPC; but as the orbito-polar tract courses below the cortex just anterior to the anterior extent of our sdlPFC lesion, it would not have been affected by our sdlPFC lesion.

In conclusion, while it cannot be ruled out that the aspiration lesions of sdlPFC caused some inadvertent damage to underlying white matter tracts, analyses of the topography of frontal lobe fibers above led us to expect that any such inadvertent damage would be minimal; indeed, photomicrographs of stained coronal sections through the regions of our sdlPFC aspirative lesions indicate little white matter damage (Fig. 1). In addition to considering white matter fiber tracts proximal to the gray matter, one has to consider fibers of passage within the gray matter. Although no major fiber bundles course in the gray matter per se through the region of the intended sdlPFC lesion (Thiebaut de Schotten et al., 2012; Schmahmann and Pandya, 2006), some local fibers exist in the gray matter at certain locations as reviewed above, and these will have been damaged; most of these connect to sdlPFC and only a minority innervate FPC, which is anteriorly adjacent to the sdlPFC and the majority of connections to FPC would remain intact. This in no way minimizes the impact of the present findings as this study now opens up that possibility in the future, and indeed, some researchers have commenced studies of the behavioral effects of FPC lesions that have not existed in the literature until very recently (Boschin et al., 2015; Mansouri et al., 2015, 2017), albeit not yet in the context of metamemory behavior across memory domains, so this would be a key area of future research interest.

Results

Meta-deficits in sdlPFC-lesioned group in spatial recognition

Consistently with the two meta-indices, we revealed a significant main effect of group in the spatial-variant task with a one-way ANOVA on SDT $meta-d'/d'$ ($F_{(2,9)} = 5.464$, $p = 0.028$, $\eta^2 = 0.548$); and given our specific predictions for an impairment in the sdlPFC group, we ran *post hoc* tests for the comparisons (CON vs sdlPFC, Dunnett, $p = 0.034$; CON vs OFC, $p = 0.968$, one-tailed; and on hierarchical-model $meta-d'/d'$: $F_{(2,9)} = 6.524$, $p = 0.018$, $\eta^2 = 0.594$, *post hoc* test: CON vs sdlPFC, Dunnett, $p = 0.020$, CON vs OFC, $p = 0.964$, one-tailed).

Given the relatively small sample size, we additionally ran nonparametric tests and revealed the same pattern for both SDT $meta-d'/d'$, Kruskal–Wallis test across monkeys in three groups ($\chi^2_{(2)} = 5.154$, $p = 0.076$, Dunn's *post hoc* test: CON vs sdlPFC, $p = 0.075$, CON vs OFC, $p = 0.1193$, one-tailed; for hierarchical-model $meta-d'/d'$, Kruskal–Wallis test across monkeys; $\chi^2_{(2)} = 6.436$, $p = 0.040$, Dunn's *post hoc* test: CON vs sdlPFC, $p = 0.036$, CON vs OFC, $p = 0.1478$, one-tailed). The sdlPFC monkeys were impaired in meta-accuracy in spatially demanding recognition, whereas the OFC group did not show any meta-deficit in either of

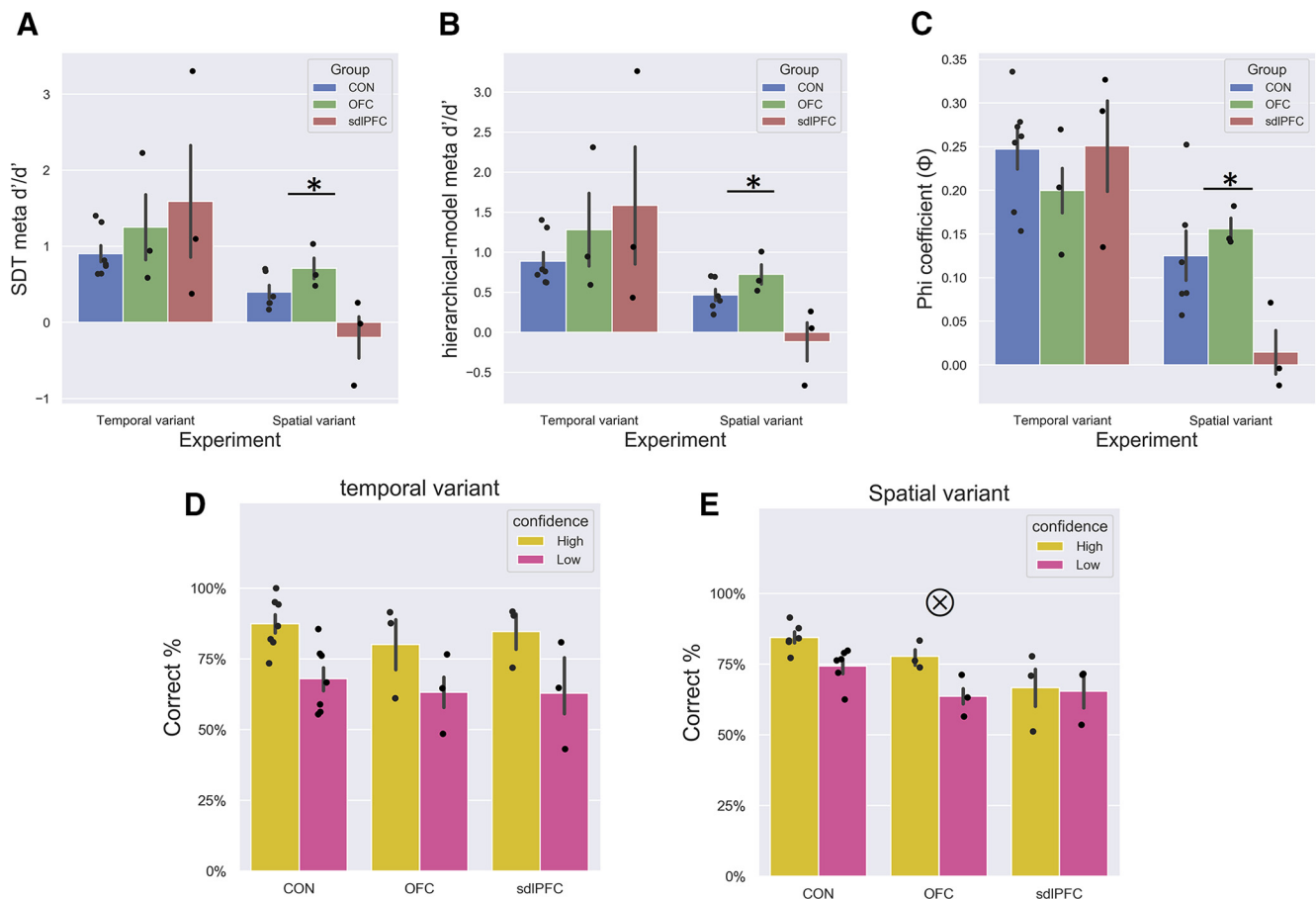


Figure 3. Differential deficits in meta-indices and accuracy confidence interaction in sdIPFC group in spatially demanding recognition. Meta-performance for the three monkey groups (OFC, sdIPFC, CON). Metacognitive accuracy in sdIPFC group was lower than CON group for spatial-variant task, but not for the temporal-variant task: (A) SDT *meta-d'/d'*; (B) hierarchical-model *meta-d'/d'*; and (C). Horizontal axes represent the two spatial recognition tasks (temporal-variant; spatial-variant). Vertical axes represent the three meta-indices. D, E, Accuracy in high confidence trials is usually higher than low confidence trials (for both CON and OFC monkeys in both tasks), but such effects were disrupted in the sdIPFC monkeys, especially in the spatial-variant task. “×” indicates significant group × confidence interaction ($p < 0.05$). * $p < 0.05$. Colored dots represent individuals. Error bars indicate SEM.

the tasks. In contrast, in the temporal-variant task, we found no main effect of group in meta-accuracy, *meta-d'/d'* ($F_{(2,10)} = 0.773$, $p = 0.487$, $\eta^2 = 0.134$, Kruskal–Wallis test across monkeys; $\chi^2_{(2)} = 0.195$, $p = 0.907$; hierarchical-model *meta-d'/d'*: $F_{(2,10)} = 0.500$, $p = 0.621$, $\eta^2 = 0.141$, Kruskal–Wallis test across monkeys; $\chi^2_{(2)} = 1.92$, $p = 0.902$).

To further validate these results, we replicated these findings with the ϕ in both tasks, that is, in the spatial-variant, group effect ($F_{(2,9)} = 4.904$, $p = 0.036$, $\eta^2 = 0.521$, *post hoc* test: CON vs sdIPFC, Dunnett, $p = 0.026$, CON vs OFC, $p = 0.904$, one-tailed) and in the temporal-variant, group effect ($F_{(2,10)} = 0.822$, $p = 0.467$, $\eta^2 = 0.091$). The nonparametric tests confirm the same pattern for the spatial-variant, Kruskal–Wallis test across monkeys in three groups ($\chi^2_{(2)} = 6.064$, $p = 0.048$; Dunn’s *post hoc* test: CON vs sdIPFC, $p = 0.029$, CON vs OFC, $p = 0.198$, one-tailed), but not the temporal-variant, Kruskal–Wallis test across monkeys ($\chi^2_{(2)} = 1.513$, $p = 0.469$). These results using three indices convergently revealed severe impairment in metamemory of recognition in the sdIPFC lesioned monkeys (but not in OFC group), confirming that metacognitive ability was impaired in the spatially demanding spatial recognition task, but not the temporally demanding task (Fig. 3A–C).

These meta-indices in principle refer to how meaningful a subject’s confidence is in distinguishing between correct and incorrect responses. We plotted the distributions of RTs in histograms for the

correct versus incorrect trials separately for individual animals (Fig. 4). By indicating the position of median RT within these histograms, we show that the sdIPFC monkey could not distinguish between correct trials and incorrect trials as well as the unimpaired monkeys in the spatial variant experiment; that is, sdIPFC monkeys have more overlapping area between correct trial RT and incorrect trial RT than CON and OFC monkeys.

We then also ran two mixed-design, repeated-measures ANOVAs on percentage correct with group as a between-subject variable and confidence as a within-subjects variable for the two tasks separately and obtained a significant interaction with the spatial-variant task ($F_{(2,9)} = 5.416$, $p = 0.029$), but not with the temporal-variant task ($F_{(2,10)} = 0.355$, $p = 0.710$). Percentage correct in high confidence trials is usually higher than low confidence trials ($p < 0.01$ for both CON and OFC monkeys in both tasks), but such effects were disrupted in the sdIPFC monkeys in the spatial-variant task ($p = 0.696$), indicating that the sdIPFC monkeys were unable to keep track of the efficacy of confidence during memory judgment (Fig. 3D,E). Correspondingly, one-way ANOVAs having group as a between-subjects variable on *meta-d'* (a sensitivity measure quantifying the ability to discriminate between correct and incorrect judgments) for the two tasks separately also revealed that a significant main effect of group in the spatial-variant task on SDT *meta-d'* ($F_{(2,9)} = 5.701$, $p = 0.025$, $\eta^2 = 0.559$, *post hoc* test: CON vs sdIPFC, Dunnett, $p = 0.015$,

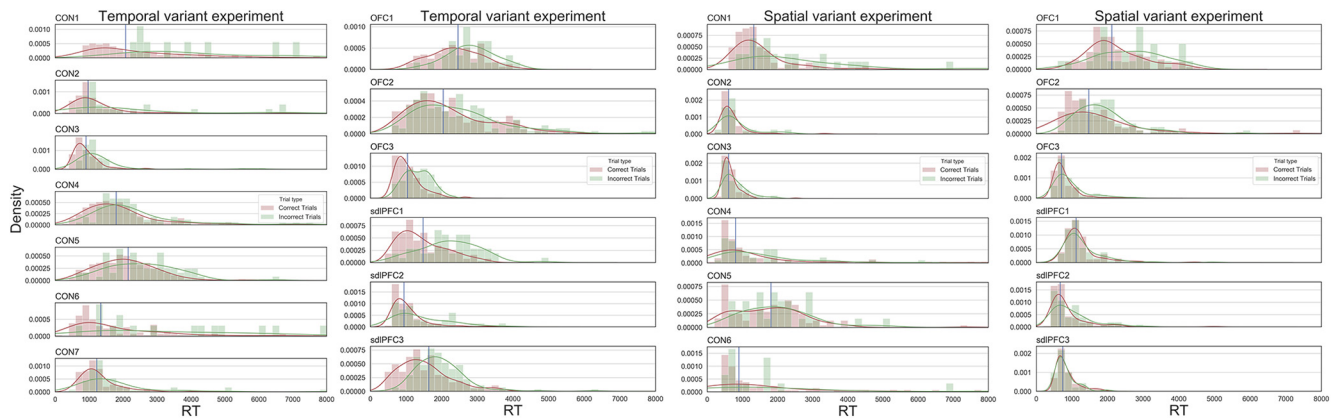


Figure 4. RT distributions for correct and incorrect responses of individual monkeys. x axis, RT; y axis, density. Green and red lines indicate the kernel density estimation. Blue lines indicate the medial RT of each monkey. Bin size for the histograms is set at 200 ms.

CON vs OFC, $p = 0.867$, one-tailed), but not in the temporal-variant ($F_{(2,10)} = 0.558$, $p = 0.589$, $\eta^2 = 0.100$). Given the relatively small sample per group, we additionally ran the statistics using nonparametric tests. These again revealed a significant group effect in the spatial-variant, Kruskal–Wallis test across all monkeys ($\chi^2_{(2)} = 5.154$, $p = 0.025$, Dunn’s *post hoc* test: CON vs sdIPFC, $p = 0.025$, CON vs OFC, $p = 0.257$, one-tailed), but not in the temporal-variant ($\chi^2_{(2)} = 1.438$, $p = 0.487$). To ascertain that these lesion effects were task-specific (task: spatial-variant/temporal-variant), we ran three separate mixed-design, repeated-measures ANOVAs considering only the CON and sdIPFC groups with group as a between-subjects variable and task as a within-subjects variable and confirmed a marginally significant task \times group interaction on SDT *meta-d’/d’* ($F_{(1,7)} = 4.194$, $p = 0.080$) and on hierarchical-model *meta-d’/d’* ($F_{(1,7)} = 4.599$, $p = 0.069$), as well as a slightly weaker effect on ϕ ($F_{(1,7)} = 1.939$, $p = 0.206$).

To address the potential confounds that the memory decision confidence estimated by RTs would be biased by condition difficulty, we calculated the meta-index for each condition for each subject using each animal’s median RT as the cutoff. The results using this within-condition classification replicated our main results. We replicated the significant main effect of group in the spatial-variant task with ϕ ($F_{(2,9)} = 8.847$, $p = 0.008$, $\eta^2 = 0.663$, *post hoc* test: CON vs sdIPFC, Dunnett, $p = 0.011$, CON vs OFC, $p = 0.133$, one-tailed). The results on SDT *meta-d’/d’* showed a trend toward significance ($F_{(2,9)} = 3.161$, $p = 0.091$, $\eta^2 = 0.413$, *post hoc* test: CON vs sdIPFC, Dunnett, $p = 0.152$, CON vs OFC, $p = 0.625$, one-tailed) and on hierarchical-model *meta-d’/d’* ($F_{(2,9)} = 2.477$, $p = 0.139$, $\eta^2 = 0.355$). Additionally, we also ran nonparametric tests and revealed the same pattern for both SDT *meta-d’/d’*, Kruskal–Wallis test across the monkeys ($\chi^2_{(2)} = 6.385$, $p = 0.027$, Dunn’s *post hoc* test: CON vs sdIPFC, $p = 0.021$, CON vs OFC, $p = 0.163$, one-tailed), and ϕ , Kruskal–Wallis test across monkeys ($\chi^2_{(2)} = 7.6154$, $p = 0.022$, Dunn’s *post hoc* test: CON vs sdIPFC, $p = 0.025$, CON vs OFC, $p = 0.120$, one-tailed) and on hierarchical-model *meta-d’/d’* ($\chi^2_{(2)} = 2.477$, $p = 0.093$, Dunn’s *post hoc* test: CON vs sdIPFC, $p = 0.066$, CON vs OFC, $p = 0.163$, one-tailed). Given the relatively low trial numbers per condition, the meta-index shall be more accurate considering all trials within whole sessions. We therefore propose that the analyses based on within-condition RT classification shall be considered as additional tests for confirmatory purpose.

Meta-ability observed in unimpaired animals

It is paramount to establish that the unimpaired monkeys were performing the tasks with certain metacognitive ability (above chance). Considering the relatively small sample sizes (e.g., $n = 3$), Student’s t tests testing the meta-scores against zero might not be the most statistically valid approach. We accordingly performed a series of subject-based distribution simulation, which entail shuffling randomly all the pairings between “responses” (correct/incorrect) and their corresponding “confidence level” (high/low) within each subject. This procedure was then repeated 1000 times, thereby generating 1000 new (random) pairings for each animal. Based on these new pairings, we then computed simulated meta-scores per animal for each of the three metrics (1000 values for each meta-metric) assuming they had not used their meta-ability to perform the tasks. These values are essentially centered on a mean of zero for each virtual animal, indicating their negligible meta-ability. But importantly, we can now test these with the animals’ actual meta-scores using a minimum statistical method (Nichols et al., 2005). The results show that the animals performed significantly above chance in all tasks in which no impairment was found (all p values < 0.005 ; Table 1).

sdIPFC and OFC lesions did not result in recognition impairment

Given that metacognition is quantified by the correspondence between confidence and Type I task performance, it is theoretically important to establish that the task (first-order) performances were matched between the groups to argue for the presence of a true difference in metacognition caused by the sdIPFC lesion. Despite the deficits in metamemory accuracy in the sdIPFC group, importantly, we established that there were not any memory deficits in their Type I performance. In two mixed-design, repeated-measures ANOVAs, we entered the percentage correct or RT with one between-subjects factor group and one within-subjects factor condition and found neither a main effect nor interaction effects with group in the two tasks, temporal-variant: percentage correct: $F_{(2,10)} = 0.284$, $p = 0.759$, RT: $F_{(2,10)} = 0.932$, $p = 0.425$, no group \times condition interaction, all p values > 0.05 ; spatial-variant: percentage correct: $F_{(2,9)} = 3.868$, $p = 0.061$, RT: $F_{(2,9)} = 0.794$, $p = 0.481$, no group \times condition interactions, all p values > 0.05 . Even if we compared only the sdIPFC with the CON group in RTs, the sdIPFC group was not significantly faster in responding (spatial: $F_{(1,7)} = 1.20$, $p = 0.309$;

Table 1. Percentiles of each monkey's actual meta-scores compared with the simulated data for all three tasks in which no impairment was found^a

Group ID	Temporal variant ϕ	Temporal variant meta d'/d'	Temporal variant H-model meta d'/d'	Spatial variant ϕ	Spatial variant meta d'/d'	Spatial variant H-model meta d'/d'	WCST ϕ	WCST meta d'/d'	WCST H-model meta d'/d'
CON 1	99.5	99.5	99.5	99.5	99	99	99.5	99.5	99.5
CON 2	99.5	99.5	99.5	89	89	93	99.5	99.5	99
CON 3	99.5	99.5	99.5	76.5	76.5	83	99.5	91	85
CON 4	99	99	99	96.5	96.5	96	99.5	99.5	99
CON 5	99	99	99	83	83	88	99.5	99.5	99.5
CON 6	99.5	99.5	99.5	75.5	75.5	84	99.5	99.5	97
CON 7	99.5	99.5	99.5						
Statistics for CON	$0.01^7 < 0.001$	$0.01^7 < 0.001$	$0.01^7 < 0.001$	$0.245^6 < 0.001$	$0.245^6 < 0.001$	$0.17^6 < 0.001$	$0.005^6 < 0.001$	$0.09^6 < 0.001$	$0.15^6 < 0.001$
OFC 1	99	93	93	94	94	84	99.5	99.5	99.5
OFC 2	93	99.5	99.5	94	99	99	99.5	99.5	99.5
OFC 3	99.5	99	99	89	98	98	99.5	99.5	99.5
Statistics for OFC	$0.07^3 < 0.001$	$0.07^3 < 0.001$	$0.07^3 < 0.001$	$0.11^3 < 0.001$	$0.06^3 < 0.001$	$0.16^3 = 0.0040$	$0.005^3 < 0.001$	$0.005^3 < 0.001$	$0.005^3 < 0.001$
sdlPFC 1	99.5	99	99	48 ^b	1 ^b	1 ^b	99.5	99.5	99.5
sdlPFC 2	96	81	89	68 ^b	83 ^b	83 ^b	99.5	99.5	99
sdlPFC 3	99.5	99.5	99.5	48 ^b	48 ^b	48 ^b	99.5	99	99
Statistics for sdlPFC	$0.04^3 < 0.001$	$0.19^3 = 0.006$	$0.11^3 = 0.0013$	$0.52^3 = 0.140^b$	$0.99^3 = 0.970^b$	$0.99^3 = 0.970^b$	$0.005^3 < 0.001$	$0.01^3 < 0.001$	$0.01^3 < 0.001$

^aThe inferential statistics are performed using a minimum statistics method (Nichols et al., 2005), showing that the real unimpaired monkeys' meta-scores are all significantly higher than chance level.

^bNot reaching significance.

temporal: $F_{(1,7)} = 2.005$, $p = 0.194$). The ANOVAs also showed that performance decreased with delay, percentage correct: $F_{(4,40)} = 26.964$, $p < 0.001$, RT: $F_{(4,40)} = 9.918$, $p < 0.001$, and with separation, percentage correct: $F_{(3,27)} = 33.960$, $p < 0.001$, RT: $F_{(3,27)} = 0.121$, $p = 0.105$ for all three groups (Fig. 5). These analyses point to the fact that neither sdlPFC nor OFC lesion resulted in any Type 1 recognition memory impairment.

Meta-deficits could not be explained away by speed-and-accuracy trade-off

Here we have used RT as a proxy for memory decision confidence (Dotan et al., 2018). The meta-deficit effects observed here might be confounded by some speed-and-accuracy trade-off strategy adopted by the monkeys toward maximizing the time spent per unit of reward (correct) ratio. Speed-and-accuracy trade-off taps into the monitoring of the current state of mind with regard to the uncertainty properties of the judgment (Bogacz et al., 2010), whereas RT-indexed metacognition, defined as an introspective evaluation process, taps into the higher-level function. We thus analyzed the ratio between percentage correct and RT for each monkey for both tasks. In two mixed-design, repeated-measures ANOVAs, we entered the inverse efficiency (RT/% correct) with one between-subjects factor "group" and one within-subjects factor "condition." We found neither a main effect nor interaction effects with "group" in the two tasks (all p values > 0.05 ; Fig. 6). We conclude that the putative lesion-related meta-deficits were not resultant from any speed-and-accuracy trade-off.

No meta-memory deficit following sdlPFC lesion in short-term abstract rule memory in WCST

While the sdlPFC monkeys were impaired in metamemory for spatial recognition, we have not been able to ascribe the effects specifically to spatial recognition per se. Is this deficit uniquely ascribable to the metamemory in temporospatial recognition, or more generally to the metamemory of learning abstract rules, or other higher cognitive processes? Considering performance supporting WCST demands multiprocesses, such as memory and acquisition of abstract rules, as well as reward value evaluation, we thus analyzed some extant data obtained from WCST to test for metacognitive deficits specifically in the sdlPFC monkeys. In contrast to the spatial recognition task, no meta-deficits were

found with WCST in the sdlPFC group in comparison with the CON group in one-way ANOVAs on SDT meta- d'/d' : $F_{(1,10)} = 0.677$, $p = 0.430$, $\eta^2 = 0.063$, Kruskal–Wallis test across monkeys; $\chi^2_{(2)} = 0.419$, $p = 0.518$; hierarchical-model meta- d'/d' : $F_{(1,10)} = 0.018$, $p = 0.896$, $\eta^2 = 0.002$, Kruskal–Wallis test across monkeys; $\chi^2_{(2)} = 0.077$, $p = 0.782$; and ϕ : $F_{(1,10)} = 1.132$, $p = 0.312$, $\eta^2 = 0.102$, Kruskal–Wallis test across monkeys; $\chi^2_{(2)} = 0.009$, $p = 0.926$ (Fig. 7). These analyses confirm that meta-deficits caused by sdlPFC lesion were highly specific for spatial recognition and ruled out the explanation that such meta-deficits were attributable to processes involved in the maintenance of abstract rules or general representation of knowledge.

Discussion

As expected, neither sdlPFC nor OFC aspirative lesions impaired first-order spatial recognition memory performance; yet importantly, and consistent with our hypotheses, sdlPFC lesions selectively impaired second-order meta-recognition within a recognition memory paradigm taxing recent spatial memory. No such change in metacognitive ability was observed after OFC lesions both affirming a critical functional role of the sdlPFC in supporting spatial metamemory and showing evidence for functional specificity within PFC for elements of metacognition. Our findings are robust because we assessed multiple measures of metamemory (both SDT and hierarchical model meta- d'/d' and ϕ) and found consistent significance across measures. The correlations among the three metacognitive metrics are very high (Fig. 8). Our findings are important because they provide causal evidence, of which there is currently very little, toward refining our understanding of functional specialization within primate PFC underpinning introspection during memory recognition.

Critical role of sdlPFC in meta-recognition memory

While the role of dlPFC in meta-evaluation of visual perception has been relatively well established, by evidence observed in functional activation during postdecision evaluation (Ochsner et al., 2005; Wan et al., 2016), structural and connectivity profiles (Fleming et al., 2010; Fleming and Dolan, 2012), and neuro-modulation studies (Meiron and Lavidor, 2014), the neural basis of meta-recognition memory within this region remains largely unknown. Previously, combined lesions to mid-dlPFC lesion (ar-

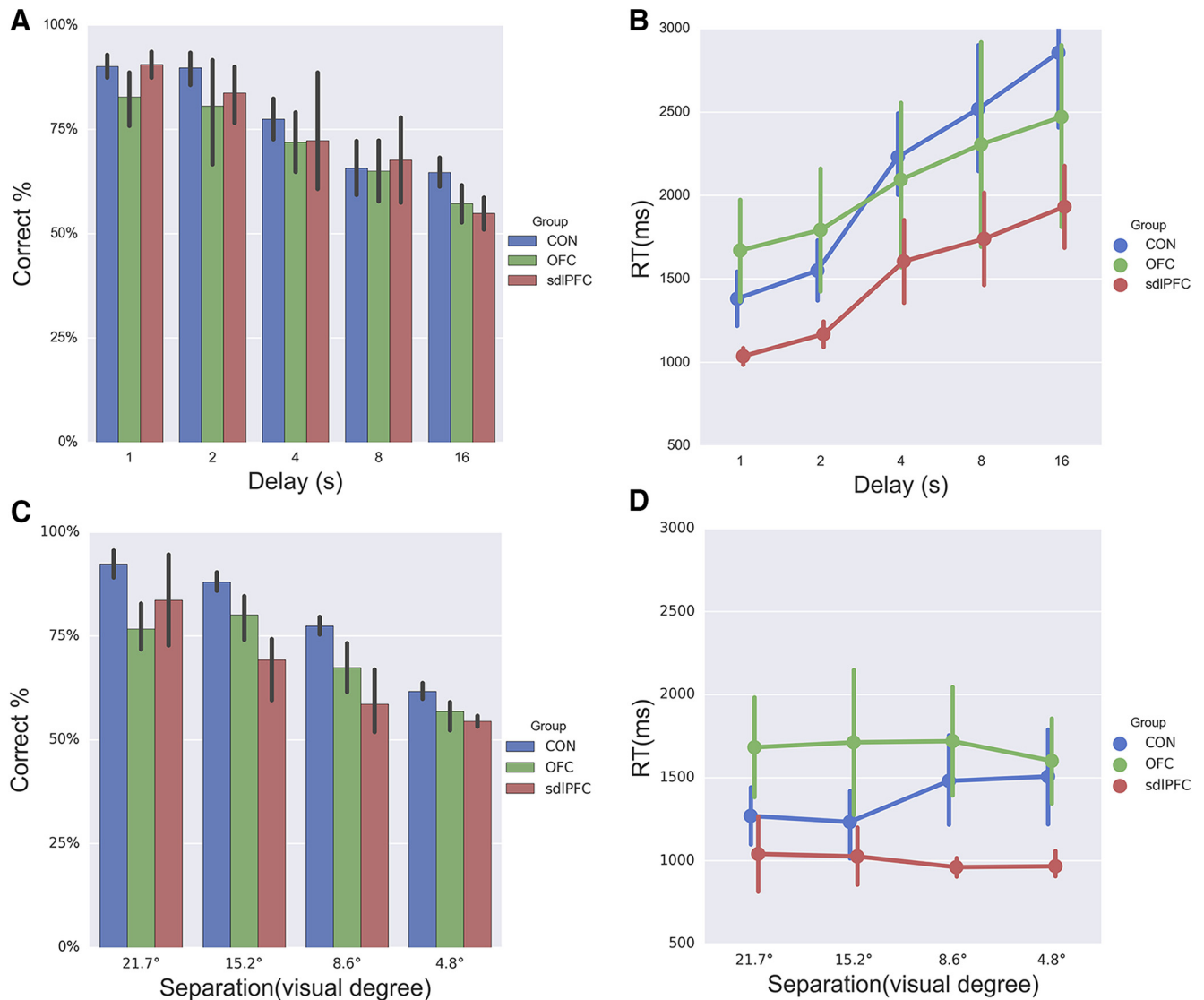


Figure 5. sdIPFC and OFC lesions did not result in recognition impairment. Memory task performance was intact in both tasks: (A) temporal-variant percentage correct; (B) temporal-variant RT; (C) spatial-variant percentage correct; and (D) spatial-variant RT. Error bars indicate SEM.

eas 46 and 9/46) and superior part of the mid-dIPFC (lateral area 9) were found not to affect standard first-order recognition memory maintenance for recently presented objects (in contrast to lesions to ventrolateral PFC and OFC, which do impair first-order object recognition) (Bachevalier and Mishkin, 1986; Kowalska et al., 1991). The mid dIPFC lesions did, however, severely impair executive processes of monitoring visual working memory information in a self-ordered version of the task (Petrides, 2000); moreover, more restricted lesions within lateral area 9 (similar to our sdIPFC lesions) were sufficient to impair the self-ordered task. Together with findings that patients with more diffuse frontal lobe pathology exhibit impaired feeling-of-knowing in the absence of amnesia (Janowsky et al., 1989), our pattern of results corroborate the extant evidence that a dissociation between Type 1 versus Type 2 performance in recognition is associated with dIPFC damage. This dissociation also aligns with the recent distinction stipulating the dIPFC’s putative role in metacognitive control as opposed to decision-making per se (Qiu et al., 2018).

Importantly, and in light of recent evidence from the macaque showing that there exist dissociated networks within PFC underlying metacognition (Miyamoto et al., 2017, 2018), our results

significantly extend the causal evidence for understanding the functional neuroanatomy of metacognition in PFC in several key ways. First, we extend causal evidence for metacognition into the spatial domain. Second, we dissociate the functional neuroanatomy of first-order spatial recognition memory from second-order spatial meta-recognition within PFC. Circumscribed muscimol injections to different localities within mid-dIPFC area impair visually guided saccades in a visuospatial working memory task in a topographical manner (Sawaguchi and Iba, 2001); further evidence that this dIPFC region, but not the more superior dIPFC region, impairs first-order spatial recognition comes from previous surgical lesions in monkeys (Levy and Goldman-Rakic, 1999). Our sdIPFC lesion fails to encroach much on this first-order region and accordingly do not impair first-order spatial recognition; at the same time, our sdIPFC lesions impair meta-recognition in the same task, therefore demonstrating functional dissociation between first-order and second-order spatial recognition within PFC. Together with Miyamoto et al. (2017, 2018), our study confirms that sdIPFC neither contributes exclusively to object nor spatial meta-recognition; rather, it contributes to both. Interestingly, we found dissociation between meta-

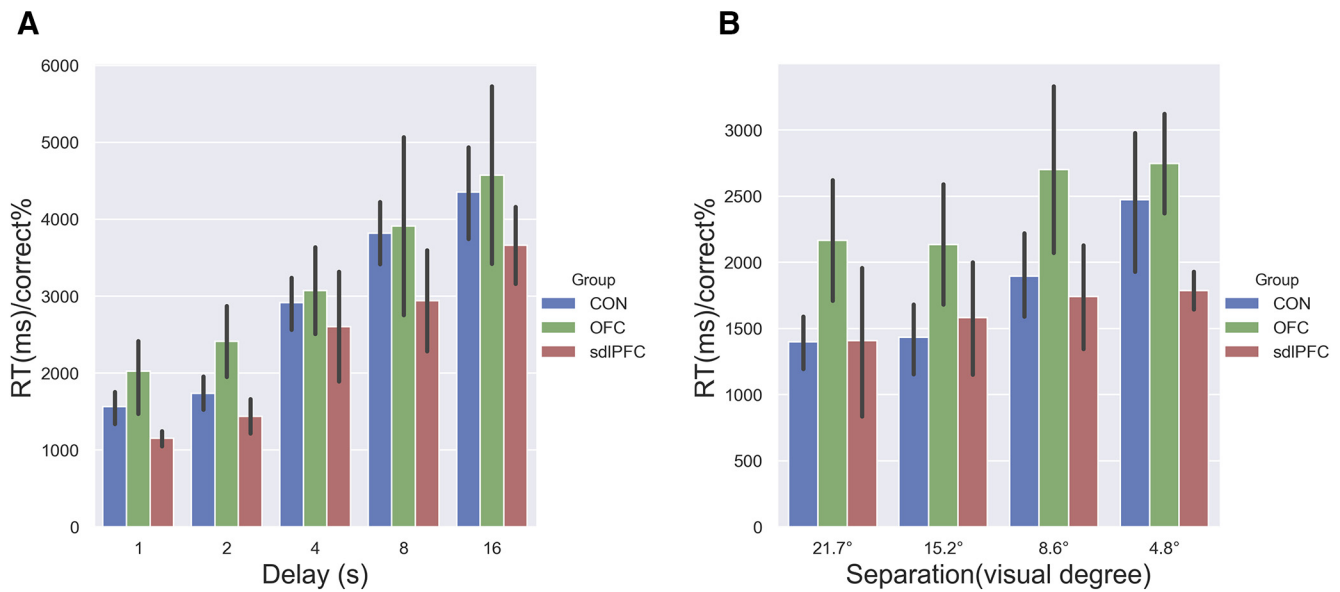


Figure 6. Meta-memory deficits could not be explained away by speed-and-accuracy trade-off. Inverse efficiency (RT/percentage correct) shows no main effect of group in (A) temporal-variant task and (B) spatial-variant task. Error bars indicate SEM.

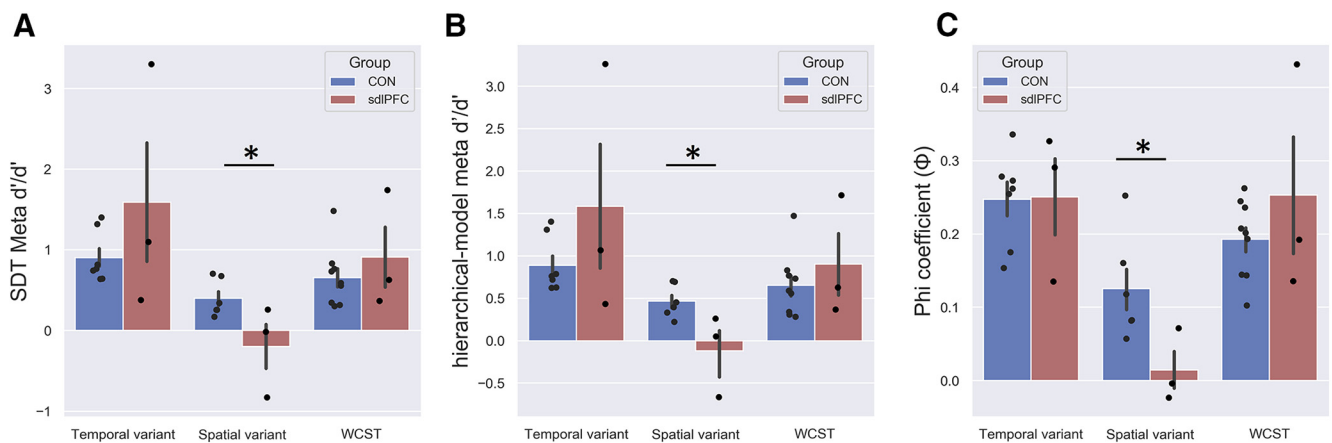


Figure 7. No meta-memory deficit following dlPFC lesion in WCST analog. Meta-performance for the two monkey groups (sdIPFC, CON) in the two spatial recognition tasks (temporal-variant; spatial-variant) and WCST analog. Vertical axes represent the three meta-indices: (A) SDT $meta-d'/d'$; (B) hierarchical-model $meta-d'/d'$; and (C) Metacognitive accuracy in the sdIPFC group was lower than CON group for spatial-variant task (see also Fig. 4A–C), but not for WCST analog or temporal-variant task. Error bars indicate SEM. * $p < 0.05$. Colored dots represent individual monkeys.

recognition deficits after sdIPFC lesions in a spatially demanding but not temporally demanding version of our task. The main difference between these two versions is that in the former the spatial difficulty is intentionally varied between trials; given it is a spatial recognition task, we postulate efficacious metacognitive monitoring of performance will therefore be in flux and continually challenged necessitating a dynamic signal from sdIPFC input into the wider neural system in support of metacognitive computation. By contrast, in the temporal-variant, the spatial difficulty is constant across trials so any metacognitive evaluation signal from sdIPFC would likely be a less important parameter to the system and either absent or less liable to be disrupted by the sdIPFC lesion accordingly, in accordance with our behavioral observations.

OFC represents value/confidence but does not support self-introspection

At first glance, OFC is an interpreter of specific values especially in terms of sorting and representation of inferred information (Wallis, 2007; Jones et al., 2012). Could it be the cornerstone of

meta-appraisal toward one’s own memory performance? Despite proposals that the OFC is a key part of continuous decision-making under uncertainty (Kennerley et al., 2011), related to the explicit manifestation of decision confidence (Kepecs et al., 2008; Lak et al., 2014) and various aspects of decision-making (Izquierdo, 2017), its role in metacognition in the present experimental context appears to be none as evident in the total absence of meta-impairment in the OFC group. One likely possibility is that value assignment (Padoa-Schioppa and Assad, 2006), valuation of inferred information (Jones et al., 2012), and decision confidence (Lak et al., 2014) per se are fundamentally distinct/dissociable from metacognitive introspection. Indeed, meta-decision processes, as measured by the SDT and hierarchical-model $meta-d'$ here, are in principle “bias-free” and are immune to any bias due to “confidence” (Fleming and Lau, 2014), suggesting that the computation performed by the OFC underlying confidence signals (Kepecs et al., 2008; Lak et al., 2014) does not necessarily equate to the same neurobiological prerequisite for metacognitive computation. Relatedly, in our tasks, there was no

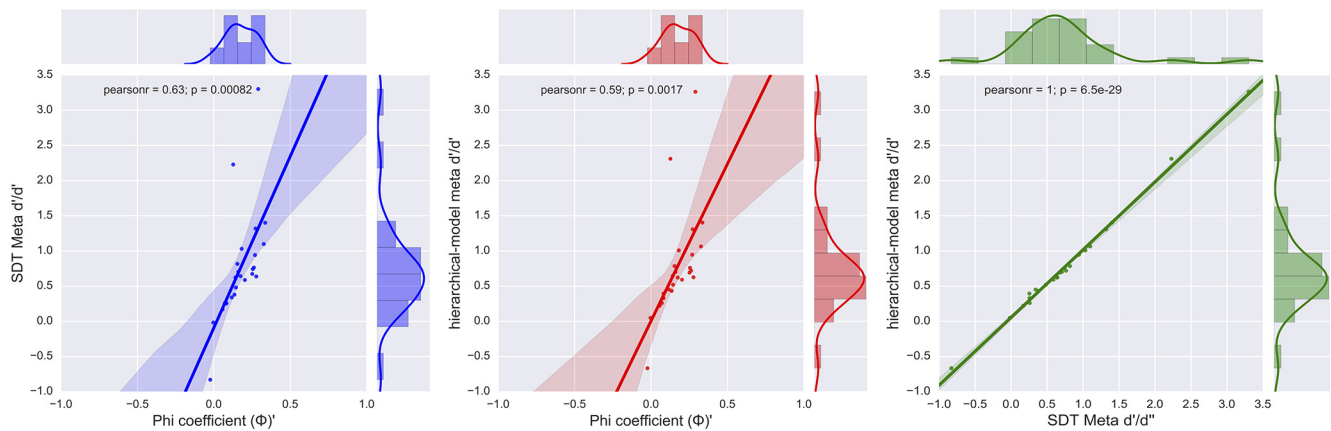


Figure 8. Strong correlations among the three metacognitive metrics. Pearson correlations computed among the three indices were all statistically significant (all p values 0.001). Colored dots represent individual data points collapsed across monkey groups, with each monkey shown twice (temporal and spatial variants).

explicit requirement for reporting confidence, in which case the memory response need not be bound with any explicit valuation processes (Rolls and Grabenhorst, 2008) or reward-based updating (Buckley et al., 2009). The introspection following memory decision was thus based entirely on some self-generated space, without any feedback or input exerted externally on their decision confidence or monitoring of degree of uncertainty (Kennerly et al., 2011). This task feature discrepancy might explain the lack of meta-awareness deficits even when the OFC was absent.

An alternative but not mutually exclusive explanation, given that frontal versus parietal cortex neural correlates of metacognition are somewhat domain-specific (McCurdy et al., 2013), is that OFC's contribution to metacognition might analogously be domain-specific. Indeed, previous studies tapping into the OFC role in meta-related processes were all on perceptual decisions, such as odor discrimination judgment (Kepecs et al., 2008; Lak et al., 2014), whereas at present the tasks in question concern mnemonic decisions. Some recent work on humans has similarly evinced such specificity for perceptual versus memory related metacognition (Fleming et al., 2014; Ye et al., 2018, 2019).

Applications and limitations of RT approach

We categorized trials into fast response trials (high confidence) and slow response trials (low confidence) according to the median of RT (Dotan et al., 2018). A similar method was used in humans, suggesting that fast response and slow response are supported by distinct mechanisms (Novikov et al., 2017) and processes (van de Vijver et al., 2011; Coomans et al., 2016). While recognizing that other possibilities exist (e.g., fast trials and slow trials might also relate to different levels of attention in addition to uncertainty) (Novikov et al., 2017), certainty is nonetheless informed by RT (Kiani et al., 2014) and all considered process are presumably related to self-monitoring. In our approach, we focus on four kinds of trials: correct/high confidence (fast RT), incorrect/low confidence (slow RT), correct/low confidence (slow RT), incorrect/high confidence (fast RT); accordingly, our analyses may be considered to be particularly relevant to implicit metacognitive accuracy of the decision-making system.

In the present experiments, we used RTs as a proxy for memory decision confidence. This might not capture as much meta-component as other paradigms, which explicitly require the animals to provide accurate metacognitive judgments to obtain rewards (Miyamoto et al., 2017, 2018). Moreover, as opposed to some other paradigms (Lak et al., 2014), which can produce a

continuum of confidence ratings, using a sharp, half-half median cutoff might also obscure the true distribution of high versus low confidence pattern. Using this rather implicit measurement for confidence, the RT in Type I tasks may possibly be entangled with other nonmetacognitive components, leaving us with a mixed read-out of the directional certainty (Moreira et al., 2018). However, the advantage of RT-indexed metacognition should not be overlooked too, namely, that it does not suffer from any training-induced associations, such as environmental cue associations, behavioral cue associations, or response competition, which could contaminate true introspection (Hampton, 2009).

A wider theoretical implication afforded by the current study is that this same group of sdIPFC monkeys, despite their impaired memory self-appraisal in the DMP task, were completely intact in all aspects of a rule-guided memory WCST analog (Buckley et al., 2009). This constitutes a stark contrast to the OFC-lesioned monkeys, who were impaired in updating rule value representation in the WCST analog (Buckley et al., 2009), but not in their introspection in the present DMP tasks. These results together constitute a double dissociation within PFC between dorsolateral and ventromedial PFC regions in differentially supporting two related, yet perhaps distinct, higher-order processes, providing compelling evidence suggestive of functional specialization of dual supervisory, self-monitoring abilities between dorsolateral versus ventral parts of the PFC.

References

- Bachevalier J, Mishkin M (1986) Visual recognition impairment follows ventromedial but not dorsolateral prefrontal lesions in monkeys. *Behav Brain Res* 20:249–261.
- Baird B, Smallwood J, Gorgolewski KJ, Margulies DS (2013) Medial and lateral networks in anterior prefrontal cortex support metacognitive ability for memory and perception. *J Neurosci* 33:16657–16665.
- Baltz ET, Yalcinbas EA, Renteria R, Gremel CM (2018) Orbital frontal cortex updates state-induced value change for decision-making. *eLife* 7:e35988.
- Bogacz R, Wagenmakers EJ, Forstmann BU, Nieuwenhuis S (2010) The neural basis of the speed–accuracy tradeoff. *Trends Neurosci* 33:10–16.
- Boschin EA, Piekema C, Buckley MJ (2015) Essential functions of primate frontopolar cortex in cognition. *Proc Natl Acad Sci* 112: E1020–E1027.
- Buckley MJ, Mansouri FA, Hoda H, Mahboubi M, Browning PG, Kwok SC, Phillips A, Tanaka K (2009) Dissociable components of rule-guided behavior depend on distinct medial and prefrontal regions. *Science* 325:52–58.
- Coomans F, Hofman A, Brinkhuis M, van der Maas HL, Maris G (2016) Distinguishing fast and slow processes in accuracy response time data. *PLoS One* 11:e0155149.
- Dotan D, Meyniel F, Dehaene S (2018) On-line confidence monitoring during decision making. *Cognition* 171:112–121.

- Fleming SM, Daw ND (2017) Self-evaluation of decision-making: a general Bayesian framework for metacognitive computation. *Psychol Rev* 124:91–114.
- Fleming SM, Dolan RJ (2012) The neural basis of metacognitive ability. *Philos Trans R Soc Lond B Biol Sci* 367:1338–1349.
- Fleming SM, Lau HC (2014) How to measure metacognition. *Front Hum Neurosci* 8:443.
- Fleming SM, Weil RS, Nagy Z, Dolan RJ, Rees G (2010) Relating introspective accuracy to individual differences in brain structure. *Science* 329:1541–1543.
- Fleming SM, Huijgen J, Dolan RJ (2012) Prefrontal contributions to metacognition in perceptual decision making. *J Neurosci* 32:6117–6125.
- Fleming SM, Ryu J, Golfinos JG, Blackmon KE (2014) Domain-specific impairment in metacognitive accuracy following anterior prefrontal lesions. *Brain* 137:2811–2822.
- Goldman PS, Rosvold HE, Vest B, Galkin TW (1971) Analysis of the delayed-alternation deficit produced by dorsolateral prefrontal lesions in the rhesus monkey. *J Comp Physiol Psychol* 77:212–220.
- Hampton RR (2009) Multiple demonstrations of metacognition in nonhumans: converging evidence or multiple mechanisms? *Comp Cogn Behav Rev* 4:17–28.
- Izquierdo A (2017) Functional heterogeneity within rat orbitofrontal cortex in reward learning and decision making. *J Neurosci* 37:10529–10540.
- Izquierdo A, Suda RK, Murray EA (2004) Bilateral orbital prefrontal cortex lesions in rhesus monkeys disrupt choices guided by both reward value and reward contingency. *J Neurosci* 24:7540–7548.
- Janowsky JS, Shimamura AP, Squire LR (1989) Memory and metamemory: comparisons between patients with frontal lobe lesions and amnesic patients. *Psychobiology* 17:3–11.
- Jones JL, Esber GR, McDannald MA, Gruber AJ, Hernandez A, Mirenzi A, Schoenbaum G (2012) Orbitofrontal cortex supports behavior and learning using inferred but not cached values. *Science* 338:953–956.
- Kennerley SW, Behrens TE, Wallis JD (2011) Double dissociation of value computations in orbitofrontal and anterior cingulate neurons. *Nat Neurosci* 14:1581–1589.
- Kepecs A, Uchida N, Zariwala HA, Mainen ZF (2008) Neural correlates, computation and behavioural impact of decision confidence. *Nature* 455:227–231.
- Kiani R, Corthell L, Shadlen MN (2014) Choice certainty is informed by both evidence and decision time. *Neuron* 84:1329–1342.
- Kornell N, Son LK, Terrace HS (2007) Transfer of metacognitive skills and hint seeking in monkeys. *Psychol Sci* 18:64–71.
- Kowalska DM, Bachevalier J, Mishkin M (1991) The role of the inferior prefrontal convexity in performance of delayed nonmatching-to-sample. *Neuropsychologia* 29:583–600.
- Lak A, Costa GM, Romberg E, Koulakov AA, Mainen ZF, Kepecs A (2014) Orbitofrontal cortex is required for optimal waiting based on decision confidence. *Neuron* 84:190–201.
- Levy R, Goldman-Rakic PS (1999) Association of storage and processing functions in the dorsolateral prefrontal cortex of the nonhuman primate. *J Neurosci* 19:5149–5158.
- Maniscalco B, Lau H (2012) A signal detection theoretic approach for estimating metacognitive sensitivity from confidence ratings. *Conscious Cogn* 21:422–430.
- Mansouri FA, Buckley MJ, Mahboubi M, Tanaka K (2015) Behavioral consequences of selective damage to frontal pole and posterior cingulate cortices. *Proc Natl Acad Sci U S A* 112:E3940–E3949.
- Mansouri FA, Egner T, Buckley MJ (2017) Monitoring Demands for Executive Control: Shared Functions between Human and Nonhuman Primates. *Trends in Neurosciences* 40(1):15–27.
- McCurdy LY, Maniscalco B, Metcalfe J, Liu KY, de Lange FP, Lau H (2013) Anatomical coupling between distinct metacognitive systems for memory and visual perception. *J Neurosci* 33:1897–1906.
- Meiron O, Lavidor M (2014) Prefrontal oscillatory stimulation modulates access to cognitive control references in retrospective metacognitive commentary. *Clin Neurophysiol* 125:77–82.
- Meunier M, Bachevalier J, Mishkin M (1997) Effects of orbital frontal and anterior cingulate lesions on object and spatial memory in rhesus monkeys. *Neuropsychologia* 35:999–1015.
- Middlebrooks PG, Sommer MA (2011) Metacognition in monkeys during an oculomotor task. *J Exp Psychol Learn Mem Cogn* 37:325–337.
- Miyamoto K, Osada T, Setsuie R, Takeda M, Tamura K, Adachi Y, Miyashita Y (2017) Causal neural network of metamemory for retrospection in primates. *Science* 355:188–193.
- Miyamoto K, Setsuie R, Osada T, Miyashita Y (2018) Reversible silencing of the frontopolar cortex selectively impairs metacognitive judgment on non-experience in primates. *Neuron* 97:980–989.e6.
- Morales J, Lau H, Fleming SM (2018) Domain-general and domain-specific patterns of activity supporting metacognition in human prefrontal cortex. *J Neurosci* 38:3534–3546.
- Moreira CM, Rollwage M, Kaduk K, Wilke M, Kagan I (2018) Post-decision wagering after perceptual judgments reveals bi-directional certainty readouts. *Cognition* 176:40–52.
- Nichols T, Brett M, Andersson J, Wager T, Poline JB (2005) Valid conjunction inference with the minimum statistic. *Neuroimage* 25:653–660.
- Noonan MP, Walton ME, Behrens TE, Sallet J, Buckley MJ, Rushworth MF (2010) Separate value comparison and learning mechanisms in macaque medial and lateral orbitofrontal cortex. *Proc Natl Acad Sci U S A* 107:20547–20552.
- Novikov NA, Nurislamova YM, Zhzhikashvili NA, Kalenkovich EE, Lapina AA, Chernyshev BV (2017) Slow and fast responses: two mechanisms of trial outcome processing revealed by EEG oscillations. *Front Hum Neurosci* 11:218.
- Ochsner KN, Beer JS, Robertson ER, Cooper JC, Gabrieli JD, Kihlstrom JF, D'Esposito M (2005) The neural correlates of direct and reflected self-knowledge. *Neuroimage* 28:797–814.
- Padoa-Schioppa C, Assad JA (2006) Neurons in the orbitofrontal cortex encode economic value. *Nature* 441:223–226.
- Petrides M (2000) Dissociable roles of mid-dorsolateral prefrontal and anterior inferotemporal cortex in visual working memory. *J Neurosci* 20:7496–7503.
- Petrides M, Pandya DN (1999) Dorsolateral prefrontal cortex: comparative cytoarchitectonic analysis in the human and the macaque brain and corticocortical connection patterns. *Eur J Neurosci* 11:1011–1036.
- Pleskac TJ, Busemeyer JR (2010) Two-stage dynamic signal detection: a theory of choice, decision time, and confidence. *Psychol Rev* 117:864–901.
- Qiu L, Su J, Ni Y, Bai Y, Zhang X, Li X, Wan X (2018) The neural system of metacognition accompanying decision-making in the prefrontal cortex. *PLoS Biol* 16:e2004037.
- Rolls ET, Grabenhorst F (2008) The orbitofrontal cortex and beyond: from affect to decision-making. *Prog Neurobiol* 86:216–244.
- Rounis E, Maniscalco B, Rothwell JC, Passingham RE, Lau H (2010) Theta-burst transcranial magnetic stimulation to the prefrontal cortex impairs metacognitive visual awareness. *Cogn Neurosci* 1:165–175.
- Rudebeck PH, Murray EA (2014) The orbitofrontal oracle: cortical mechanisms for the prediction and evaluation of specific behavioral outcomes. *Neuron* 84:1143–1156.
- Rudebeck PH, Mitz AR, Chacko RV, Murray EA (2013) Effects of amygdala lesions on reward value coding in orbital and medial prefrontal cortex. *Neuron* 80:1519–1531.
- Sawaguchi T, Iba M (2001) Prefrontal cortical representation of visuospatial working memory in monkeys examined by local inactivation with muscimol. *J Neurophysiol* 86:2041–2053.
- Schmahmann JD, Pandya DN (2006) *Fiber pathways of the brain*. Oxford: Oxford University Press.
- Schuck NW, Cai MB, Wilson RC, Niv Y (2016) Human orbitofrontal cortex represents a cognitive map of state space. *Neuron* 91:1402–1412.
- Shekhar M, Rahnev D (2018) Distinguishing the roles of dorsolateral and anterior PFC in visual metacognition. *J Neurosci* 38:5078–5087.
- Thiebaut de Schotten M, Dell'Acqua F, Valabregue R, Catani M (2012) Monkey to human comparative anatomy of the frontal lobe association tracts. *Cortex* 48:82–96.
- van de Vijver I, Ridderinkhof KR, Cohen MX (2011) Frontal oscillatory dynamics predict feedback learning and action adjustment. *J Cogn Neurosci* 23:4106–4121.
- Wallis JD (2007) Orbitofrontal cortex and its contribution to decision-making. *Annu Rev Neurosci* 30:31–56.
- Wan X, Cheng K, Tanaka K (2016) The neural system of postdecision evaluation in rostral frontal cortex during problem-solving tasks. *eNeuro* 3:1–24.
- Ye Q, Zou F, Lau H, Hu Y, Kwok SC (2018) Causal evidence for mnemonic metacognition in human precuneus. *J Neurosci* 38:6379–6387.
- Ye Q, Zou F, Dayan M, Lau H, Hu Y, Kwok SC (2019) Individual susceptibility to TMS affirms the precuneal role in metamemory upon recollection. *Brain Struct Funct*: 1–13.
- Yeterian EH, Pandya DN, Tomaiuolo F, Petrides M (2012) The cortical connectivity of the prefrontal cortex in the monkey brain. *Cortex* 48:58–81.