

Research Articles: Behavioral/Cognitive

Partially overlapping neural correlates of metacognitive monitoring and metacognitive control

https://doi.org/10.1523/JNEUROSCI.1326-21.2022

Cite as: J. Neurosci 2022; 10.1523/JNEUROSCI.1326-21.2022

Received: 25 June 2021 Revised: 10 January 2022 Accepted: 12 January 2022

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.jneurosci.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

1	Regular Manuscript
2 3	Partially overlapping neural correlates of metacognitive monitoring and metacognitive
4	control
5	Abbreviated title: Metacognitive monitoring and control
6	
7	Annika Boldt ¹ & Sam J Gilbert ¹
8	¹ Institute of Cognitive Neuroscience, University College London, Alexandra House, 17-19
9	Queen Square, London WC1N 3AZ, UK
10	Correspondence should be addressed to Annika Boldt: a.boldt@ucl.ac.uk
11 12	10 magazi 6 firayyari 5 tahlar
13	49 pages; 6 figures; 5 tables Word counts: Abstract 150; Introduction 1173; Discussion 1652
13	word counts. Abstract 130, Introduction 11/3, Discussion 1032
14 15	Conflict of interests statement: The authors declare that they have no competing financial or non-financial interests.
15	or non-financial interests.
15 16	or non-financial interests. Acknowledgements: This research was funded by the Wellcome Trust, who awarded a Sir
15 16 17 18 19	or non-financial interests. Acknowledgements: This research was funded by the Wellcome Trust, who awarded a Sir Henry Wellcome Postdoctoral Fellowship (206480/Z/17/Z) to AB, and the Economic & Social Research Council (ESRC), who awarded a Research Grant (ES/N018621/1) to SJG. Neither of these funding bodies played a role in the conceptualization, design, data collection
15 16 17 18 19 20	or non-financial interests. Acknowledgements: This research was funded by the Wellcome Trust, who awarded a Sir Henry Wellcome Postdoctoral Fellowship (206480/Z/17/Z) to AB, and the Economic & Social Research Council (ESRC), who awarded a Research Grant (ES/N018621/1) to SJG. Neither of these funding bodies played a role in the conceptualization, design, data collection analysis, decision to publish, or preparation of the manuscript. For the purpose of Open
15 16 17 18 19 20 21	or non-financial interests. Acknowledgements: This research was funded by the Wellcome Trust, who awarded a Sir Henry Wellcome Postdoctoral Fellowship (206480/Z/17/Z) to AB, and the Economic & Social Research Council (ESRC), who awarded a Research Grant (ES/N018621/1) to SJG. Neither of these funding bodies played a role in the conceptualization, design, data collection analysis, decision to publish, or preparation of the manuscript. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted
15 16 17 18 19 20 21 22	or non-financial interests. Acknowledgements: This research was funded by the Wellcome Trust, who awarded a Sir Henry Wellcome Postdoctoral Fellowship (206480/Z/17/Z) to AB, and the Economic & Social Research Council (ESRC), who awarded a Research Grant (ES/N018621/1) to SJG. Neither of these funding bodies played a role in the conceptualization, design, data collection analysis, decision to publish, or preparation of the manuscript. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission. The authors would like to thank the
15 16 17 18 19 20 21 22 23	or non-financial interests. Acknowledgements: This research was funded by the Wellcome Trust, who awarded a Sir Henry Wellcome Postdoctoral Fellowship (206480/Z/17/Z) to AB, and the Economic & Social Research Council (ESRC), who awarded a Research Grant (ES/N018621/1) to SJG. Neither of these funding bodies played a role in the conceptualization, design, data collection analysis, decision to publish, or preparation of the manuscript. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission. The authors would like to thank the MetaOffloading lab for help with the MRI data collection, Pei-Chun Tsai for help with
15 16 17 18 19 20 21 22	or non-financial interests. Acknowledgements: This research was funded by the Wellcome Trust, who awarded a Sir Henry Wellcome Postdoctoral Fellowship (206480/Z/17/Z) to AB, and the Economic & Social Research Council (ESRC), who awarded a Research Grant (ES/N018621/1) to SJG. Neither of these funding bodies played a role in the conceptualization, design, data collection analysis, decision to publish, or preparation of the manuscript. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission. The authors would like to thank the

Abstract

Metacognition describes the process of monitoring one's own mental states, often for the purpose of cognitive control. Previous research has investigated how metacognitive signals are generated (metacognitive monitoring), for example when people (both f/m) judge their confidence in their decisions and memories. Research has also investigated how metacognitive signals are used to influence behavior (metacognitive control), for example setting a reminder (i.e. *cognitive offloading*) for something you are not confident you will remember. However, the mapping between metacognitive monitoring and metacognitive control needs further study on a neural level. We used fMRI to investigate a delayed-intentions task with a reminder element, allowing human participants to use their metacognitive insight to engage metacognitive control. Using multivariate pattern analysis, we found that we could separately decode both monitoring and control, and, to a lesser extent, cross-classify between them. Therefore, brain patterns associated with monitoring and control are partially, but not fully, overlapping.

Significance Statement

Models of metacognition commonly distinguish between monitoring (how metacognition is formed) and control (how metacognition is used for behavioural regulation). Research into these facets of metacognition has often happened in isolation. Here, we provide a study which directly investigates the mapping between metacognitive monitoring and metacognitive control at a neural level. We applied multivariate pattern analysis to fMRI data from a novel task in which participants separately rated their confidence (metacognitive monitoring) and how much they would like to use a reminder (metacognitive control). We find support for the notion that the two aspects of metacognition overlap partially but not

- 50 fully. We argue that future research should focus on how different metacognitive signals are
- 51 selected for control.

Introduction

Our brains possess a remarkable ability to monitor performance and to then use
metacognition to control future behavior. For example, if you have low confidence that you
will remember a delayed intention (metacognitive monitoring; MetaM) like regular
medication intake, you might set a reminder on your phone (metacognitive control; MetaC).
This distinction between monitoring and control is found in the seminal metamemory
framework by Nelson and Narens (1990; see also Flavell, 1976; Kluwe, 1982; Brown, 1987;
Efklides, 2008; Shea et al., 2014; Yeung et al., 2004; Fleming & Daw, 2017; Fletcher &
Carruthers, 2012; Fleur, Bredeweg, & van den Bos, 2021), which proposes that cognition
functions at two distinct levels: the object and the meta level (Figure 1A). Information at the
object level about decisions, memories, attention, action and so forth is re-represented at the
meta level via a process of MetaM. Meanwhile, information at the meta level controls
processing at the object level (MetaC). Shimamura's (2000) dynamic filtering theory extends
the framework by Nelson & Narens (1990), ascribing the role of the object level to posterior
cortical regions and the role of the meta level to prefrontal cortex (PFC). The information
flow between these regions forms the basis of MetaM and MetaC.
We are only slowly beginning to understand the neural mapping between MetaM and
MetaC. This mapping or link describes the relationship that exists between MetaM and
MetaC on a functional level – are these labels describing the identical process or two
different computations with different inputs? This question is important because one rationale
for studing MetaM is that it can provide insight into MetaC (e.g. Boldt & Yeung, 2015;
Wokke et al., 2020; Masset et al., 2020; Gherman & Philiastides, 2018; Miyamoto et al.,
2018; Odegaard et al., 2018; Bang & Fleming, 2018; Ye et al., 2018; Shekhar & Rahnev,
2018). This would be strengthened if the mapping between the two were better understood.
Furthermore dissociations have been found between MetaM and MetaC. For example in

//	some circumstances, young children (Redshaw et al., 2018), OCD patients (Vaghi et al.,
78	2017) older adults (Dunlosky & Connor, 1997), and individuals with Autism Spectrum
79	Conditions (Grainger et al., 2016) have a diminished mapping between MetaM and MetaC,
80	which could lead to suboptimal behavioral regulation. However, the potential neural
81	substrates for this variability are unknown.
82	One of the reasons why the MetaM-MetaC mapping has received little attention is
83	that the two aspects of metacognition are usually studied in isolation (though see Koriat et al.,
84	2006, 2014; Mei et al., 2020; Son & Schwartz, 2009; Schulz, Fleming & Dayan, 2021; Qiu et
85	al, 2018). Studies on MetaM commonly explore the variables that affect how confident
86	people feel and the associated neural correlates. For example, neuroimaging studies have
87	identified a widespread network of involved regions, including the rostrolateral prefrontal
88	cortex (rlPFC; Yokoyama et al., 2010; Fleming, Huijgen, & Dolan, 2012; Allen et al., 2017)
89	and also the precuneus specifically for metamemory studies (e.g. McCurdy et al., 2013;
90	Baird, Smallwood, Gorgolewski, & Margulies, 2013; Ye et al., 2018). Moreover, machine-
91	learning techniques have been used to "decode" brain patterns associated with low versus
92	high confidence, using both fMRI (Hebart et al., 2014; Cortese et al., 2016; Morales, Lau &
93	Fleming, 2018) and EEG (Boldt & Yeung, 2015). Research on MetaC, on the other hand, has
94	focused on situations in which metacognitive experiences are utilized for learning,
95	communication, or speed-accuracy tradeoff, to name a few (e.g. Metcalfe & Finn, 2008;
96	Guggenmos et al., 2016; Lak et al., 2020; Shea et al., 2014; Bahrami et al., 2010; Desender et
97	al., 2019; Frömer, Nassar, Bruckner, Stürmer, Sommer, & Yeung, 2021).
98	Most of what we know about the link between monitoring and control comes from the
99	field of cognitive control and error monitoring. Electrophysiological correlates have been
100	found that signal not only when an error has been committed but are also sensitive to correct-
101	trial performance fluctuations (Allain et al. 2004; Venna Rotvinick & Cohen 2004) Such

monitoring of errors often results in lower response speed immediately after a mistake, a
robust and often-replicated phenomenon termed post-error slowing (Rabbitt, 1966;
Danielmeier & Ullsperger, 2011; Notebaert et al., 2009). In addition to errors, conflict signals
appear to be monitored by the posterior medial frontal cortex (pMFC) including the dorsal
anterior cingulate cortex (dACC). The lateral prefrontal cortex (laPFC) is thought to receive
this input and implement cognitive control (Ridderinkhof, Ullsperger, Crone, Nieuwenhuis,
2004). It should be noted that participants are often not aware of such errors or response
conflicts and that these studies are not directly measuring metacognitive signals.
Nevertheless, evidence from this domain suggests that similar brain regions support
metacognitive monitoring and control. Qiu and colleagues (2018) conducted four elegant
fMRI experiments, using a decision-redecision paradigm: Participants were presented twice
in a row with each stimulus and rated both their response and confidence for each
presentation. They reasoned that participants would engage metacognitive monitoring for
their initial response and use metacognitive control to revise and improve decisions in the
redecision phase. Their analyses revealed an involvement of dACC in the first response and
IFPC in the second. However, because the order of the decision-redecision phases was always
the same, it is impossible to conclude whether the redecision phase really triggered more
MetaC or whether the signal observed in IPFC was instead a 'late' monitoring one. Another
open question is whether MetaM and MetaC rely on similar representations.
In order to address these questions, it is necessary to study both aspects of
metacognition in a single paradigm, which we did using a cognitive offloading task.
Cognitive offloading is the use of physical action to reduce cognitive demand, e.g. setting
external reminders rather than relying on internal memory. Previous research has
demonstrated a MetaM-MetaC link whereby individuals are more likely to set reminders

(MetaC) when they have low confidence in their memory abilities over and above the

influence of their actual memory performance (MetaM; Risko & Gilbert, 2016; Hu et al., 2019; Dunn & Risko, 2016). This finding is a robust pattern that can even be observed when reminder setting is not explicitly instructed (Boldt & Gilbert, 2019) or when confidence was measured in an unrelated perceptual task (Gilbert, 2015). Here, we use a decoding approach to examine this link at a neural level.

Participants performed a delayed intention task where in separate blocks they engaged in MetaM (how confident am I that I will remember?) or MetaC (how much would I like a reminder?). This allowed us to answer two questions: 1) Do similar brain patterns characterize MetaM and MetaC? If so, 2) Can the neural patterns that characterize specific acts of MetaC be exhaustively characterized in terms of their associated processes of MetaM? We answered these questions by examining cross-classification between MetaM and MetaC: the extent to which a classifier trained on one judgement can decode the other. Insofar as this is possible, this implies a shared neural code for MetaM and MetaC. But if cross-classification is weaker than decoding MetaM and MetaC individually, this implies that their neural bases do not overlap fully.

Materials and Methods

Participants

We trained 29 participants in a behavioral task during a first session. After reviewing their training data, 22 participants returned to the lab for a second MRI session 1 to 21 days later, excluding 7 participants (2 unsuited for MRI due to safety regulations, 2 had extreme staircase values, 3 were unavailable for a second session). Another participant was excluded after scanning due to excessive movement in the scanner. This resulted in a final sample of 21 participants, out of which 15 were female and 6 were male. While we determined our sample size based on practical constraints and on available resources, the final sample size of

N=21 is nevertheless in accordance with previous MRI studies using similar methods (Morales et al., 2018, Qiu et al., 2018; Hebart et al., 2014). Participants were 20.3 years on average (18 – 26 years and paid £36 for their participation in both sessions (about 90 and 150 minutes). All participants were right-handed, had intact color vision, no uncorrected visual impairments and had not been diagnosed with any psychiatric or neurological disorders. All testing was approved by the local ethics committee and participants gave informed consent prior to taking part in the study.

Experimental Design

In order to investigate the extent to which neural patterns associated with MetaM and MetaC are similar or distinct we had to study both aspects of metacognition within a single paradigm. Participants underwent short miniblocks of ongoing shape discrimination trials. For this ongoing task, participants had to quickly and accurately decide whether an array of colored shapes grouped around a fixation dot looked on average more like a circle or a square (De Gardelle & Summerfield, 2011) by pressing one of two buttons. The response categories were equally likely. During some of these miniblocks, participants also had to maintain a delayed intention to press a different button if the stimulus appeared in a target color (Figure 1B). Participants were allowed to use reminders (cognitive offloading) to support their prospective memory in approximately half of the miniblocks, which meant that the central fixation dot of the stimulus took on the target color for the duration of the miniblock. Instead of having to rely on their memory, participants could then simply wait for the color of the shapes to match the color of the fixation dot, making the fulfilment of the delayed intention much easier. There were 12 colors, placed equidistant in RGB space. Within each miniblock, colors were drawn without replacement. There was only one target color per miniblock,

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

presented at the beginning of the miniblock, and its occurrence during the ongoing-task trials always terminated the miniblock.

The task comprised three within-subject experimental conditions (20% Baseline, 40% MetaM, and 40% MetaC; see Table 1) each structured into miniblocks. A miniblock comprised presentation of a target color (except for in the Baseline condition which had no prospective memory element), a single metacognitive rating or cursor placement, followed by 3-7 ongoing-task trials. The number of trials per miniblock was drawn from an exponential distribution with a mean of $\mu = 1.1$, in other words shorter miniblocks were more frequent than longer miniblocks. Each of the eight blocks consisted of 94 shape trials spread unevenly across 40 miniblocks (see Figure 1C). The critical difference between our two key conditions was the metacognitive rating given about the target color before each miniblock. In the MetaC condition, participants reported how much they would want to set a reminder to help them remember this target color. The higher the rating given by the participant, the greater the likelihood of receiving a reminder, which occurred on approximately 50% of miniblocks. More specifically, ratings larger than the moving median of the past 20 MetaC ratings were assigned a reminder, whereas ratings below this cutoff had to be solved using only unaided memory. In the MetaM condition, participants reported their prospective confidence in remembering the target color. However, this had no influence on the likelihood of receiving a reminder, which occurred on a randomly selected 50% of miniblocks. In other words, the two conditions also differed in the relationship between participants' ratings and the provision of reminders. In the MetaM condition, participants' ratings had no influence on whether or not they received a reminder. In the MetaC condition, on the other hand, which miniblock contained a reminder was largely determined by participant's ratings. Therefore, in the MetaM condition participants engaged in metacognitive monitoring but did not exercise metacognitive control. In the MetaC condition they exercised control to make a decision

which is known to be guided by metacognition (Boldt & Gilbert, 2019; Gilbert, 2015). However, they were not explicitly asked to make a direct metacognitive judgment.

In the Baseline condition, there was no target color and thus no prospective-memory component (and no need for a reminder). The rating participants were asked to give was thus an 'empty' one, that is a scale without labels but with a cursor was presented on screen together with two little markers indicating where the cursor should be placed on the scale. Participants then had to move the cursor to the indicated position. In all three conditions, participants were instructed to move the cursor at least once to submit a rating.

Each block was comprised of only two out of the three conditions, the Baseline condition together with either the MetaM or the MetaC condition and alternated between the two. Within each block, conditions were predictable, that is they always followed the order of one Baseline miniblock followed by four other miniblocks. We determined the optimal order of conditions using simulations, allowing us to maximize the efficiency of our design. The main analysis window was the initial 7 seconds of the task (presentation of target color and rating). At the time of these prospective ratings, participants were still unaware whether or not they would receive a reminder, keeping our key contrast free of confounds, which would have been unavoidable had we chosen a retrospective confidence judgement as is more commonly used in the field. To increase the number of instances this analysis window was shown we therefore included partial miniblocks, that is half of the time (20 miniblocks per run), the miniblock ended immediately after the rating without the need to perform any shape classification trials or search for the target.

The study comprised two sessions. The purpose of the first session was assessment of MRI safety, completion of a pre-study questionnaire on how much participants liked the 12 colors used in the task, and training in the behavioral task (presented in MATLAB using Psychtoolbox3; Kleiner, Brainard, Pelli, Ingling, Murray, & Broussard, 2007). Participants

first completed eight practice blocks, each introducing them to a new aspect of the paradigm. They then completed four experimental runs that were identical to the task they would have to complete whilst in the scanner, each lasting ~9 minutes. During the second session, participants first underwent two practice blocks outside of the scanner (each lasting ~5 minutes) to remind them of the task before they completed eight runs in the scanner, with a 6-minute T1 scan between the fourth and fifth run. One participant only completed six blocks due to feeling unwell inside the scanner. Due to the unbalanced design, we decided to exclude this participant from all multivariate analyses.

At the end of the second session, participants were furthermore asked to fill in a post-experiment questionnaire, asking them to rate the liking of all colors again, together with how difficult they found them and several additional questions to determine whether they perceived the MetaM and MetaC conditions as similar, how much control they felt during these conditions, how they used the reminders depending on whether or not they asked for them, and how they approached each rating. The orientation of the rating scales was flipped halfway through the experiment to avoid confounding visuomotor processes with low versus high ratings. The order of scale orientations, response keys for the shape task, and the order of the conditions were counterbalanced across participants.

MRI Data Collection and Preprocessing

We used a 1.5T Siemens Avanto scanner with a 32-channel head coil and MRI-safe button boxes. We acquired both T1-weighted structural images, as well as T2*-weighted echoplanar images (EPI; 64 x 64; 3.2x3.2x3.2 mm voxels) with blood oxygen level-dependent (BOLD) contrast. We used a multiband acquisition sequence with acceleration factor = 3, TE = 54.8 ms, flip angle = 75°, to record 39 interleaved, axial slices (3.2mm thick, oriented approximately to the anterior commissure - posterior commissure plane). This

allowed us to cover most of the brain with an effective repetition time of 1.3s per volume. Encoding phase direction was anterior to posterior. Functional scans were acquired in eight runs, each comprising 410 volumes (~9 min). The first five volumes in each session were discarded to allow for T1 equilibration effects. Between the fourth and fifth functional scans, an approximately 6 min T1-weighted MPRAGE structural scan was collected.

All preprocessing was done using SPM12

(https://www.fil.ion.ucl.ac.uk/spm/software/spm12/). The T1-weighted images were skull stripped and their origin was set to the anterior commissure. We then realigned the EPI volumes and normalized them into 3 mm cubic voxels with fourth-degree B-spline interpolation using normalization parameters derived from segmentation of the co-registered structural scan, then smoothed with an isotropic 8 mm full-width half-maximum Gaussian kernel.

Statistical Analysis

Analyses of behavioral data were conducted using R version 3.6.0 ("Planting of a Tree") with the additional packages plyr, plotrix, Hmisc, R.matlab, viridis, effsize, raincloudplots, ggplot2, grid, gridExtra, and Rmisc. Statistical tests were conducted two-sided if not stated otherwise. For t-tests we reported effect sizes as Cohen's d, and for ANOVAs as partial eta square η^2_p . For the fMRI analyses, the volumes acquired during the eight sessions were treated as separate time series. For each time series, the variance in the BOLD signal was decomposed with a set of regressors in a general linear model. Three regressors were generated to code for the target color presentation and the rating as a 7s boxcar, separately for miniblock and rating conditions (Baseline in MetaM blocks, low MetaM rating, high MetaM rating in MetaM blocks and Baseline in MetaC blocks, low MetaC rating, high MetaC rating in MetaC blocks). Six additional regressors were generated

that represented effects of no interest, specifically, stimulus presentation as a stick function,
separately for targets and non-targets, the ongoing task spanning from the onset of the first to
the last shape stimulus of the miniblock, separately for whether there was a prospective-
memory requirement (Baseline vs. MetaM and MetaC) and the time when the computer
revealed to the participant whether they were allowed to use a reminder as a stick function,
separately for Reminder and Own Memory miniblocks. All regressors were convolved with a
canonical hemodynamic response function. The regressors outlined above, along with six
regressors representing residual movement-related artefacts and the mean over scans
comprised the full model for each session. The data and model were high-pass filtered at a
cutoff of 1/128 Hz. Parameter estimates for each regressor were calculated from the least
mean squares fit of the model to the data. Effects of interest were assessed in a random-effect
analysis by first forming subject-specific contrasts subtracting the Baseline from the other
two conditions. The resulting contrast images were entered into a repeated-measures
ANOVA using nonsphericity correction (Friston, Glaser, Henson, Kiebel, Phillips, &
Ashburner, 2002), representing a condition agnostic selection contrast to identify a network
of regions active in the rating task. Results are reported applying a height threshold of $p <$
0.001 uncorrected in conjunction with an extent threshold determined by SPM12 to achieve p
< 0.05 familywise error correction for multiple comparisons across the whole brain volume.
Region of interest (ROI) analyses were conducted by extracting subject-specific contrast
estimates from the resulting ROIs with the toolbox MarsBaR (Brett, Anton, Valabregue,
Poline, 2002), then entering the resulting data into an ANOVA in R using the same correction
procedure described above.
The logic behind the key analysis of our study was the following: Replicating and
extending previous findings (Hebart et al., 2014; Cortese et al., 2016; Morales et al., 2018) we

first trained separate classifiers to detect A) whether participants were in a high or low

confidence state (MetaM), and B) whether they had high or low desire for a reminder (MetaC). These classifiers could then also be combined in a cross-classification analysis, that is whether a classifier trained on MetaM ratings can also predict MetaC ratings (and vice versa). Insofar as this cross-classification is possible, this suggests shared brain representations for both aspects of metacognition. Going one step further, we then compared within-category classification to cross-classification accuracy to distinguish between two possible patterns of results: If MetaM and MetaC are based on the exact same representational code, there should be no difference in classification accuracy. If, on the other hand, MetaM and MetaC share partially overlapping patterns, we should find significantly higher classification accuracy for within- than across-category classification, but significantly-different-from-zero accuracy for cross-classification.

For the multivariate-pattern analyses, we used The Decoding Toolbox (TDT; Hebart, Görgen, & Haynes, 2015), based on the beta images resulting from the previously described general linear models (except that the models were re-fit to unsmoothed, unnormalized data and the MetaM and MetaC boxcar regressors were split into two regressors each using a median split on the respective metacognitive rating). When we ran our four separate decoding analyses, two drew the training and testing data from the same condition (low vs. high ratings for the MetaM and MetaC conditions respectively; defined by block-, condition- and subject-wise median splits), whereas the other two cross-classified (train on low vs. high MetaM ratings and test on high vs. low MetaC ratings and vice versa; note that the rating scale had to be flipped for MetaC as low confidence implies high desire for a reminder). For each of these analyses, a linear support vector machine (SVM) was trained to discriminate between low versus high ratings given the patterns of BOLD activity across voxels. Given the alternating block design and the fact that the orientation of the scale was flipped halfway through the

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

study, we had two low and two high rating images available for each training or testing fold, resulting in a 2-fold procedure (see Figure 1D).

We used a whole-brain searchlight approach (Kriegeskorte, Goebel, & Bandettini, 2006), meaning that for each voxel a separate SVM was built, fitted to the beta values within a sphere with a radius of 3 voxels (9.6 mm). This resulted in three-dimension decoding accuracy maps in native space for each participant and analyses. Decoding accuracy is calculated relative to chance level (subtracted by 50%, so a 5% accuracy corresponds to 55%). These maps were then normalized into Montreal Neurological Institute (MNI) space (using the same normalization parameters as the univariate analyses) and smoothed using a Gaussian kernel (full-width half-maximum, 4mm). Please note that this kernel was half of the one used for the univariate analyses. This was done to avoid excessive smoothing, given that the searchlight analysis already imposes spatial smoothing on the data. The resulting images were entered into a one-sample t-test using SPM12. This allowed assessment of voxels showing consistently higher decoding accuracy in a random-effect analysis. We note that the suitability of second-level t-tests has been challenged for information-like measures such as classification accuracy, where classifier performance can meaningfully be above, but not below, chance levels (Allefeld, Görgen & Haynes, 2016; Hirose, 2020). However, this characteristic does not apply to our two key hypothesis-testing analyses. For the crossclassification between MetaM and MetaC, high MetaM could either predict higher or lower MetaC. For the comparison between within- and cross-classification, accuracy for one classification could be higher or lower than the other. Therefore, in both cases our statistical tests are valid because they are performed on data that could meaningfully take values both above and below zero.

Along with the main MVPA analyses described above, we conducted an additional analysis. Here, we used a similar approach to the univariate ROI analysis described above by

defining a condition agnostic contrast (the mean of all four decoding analyses), extracting ROIs with significantly above-chance decoding accuracy and then entering the resulting classification accuracies into a repeated-measures ANOVA with factors ROI, Training condition (MetaM/MetaC), and Classification type (within-condition/cross-classification).

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

350

351

352

353

Results

Behavioral Results

Our sample included 22 participants, one of which was excluded due to excessive motion in the MRI scanner, see Methods for more details. Participants performed the tasks with a high level of accuracy (mean shape-discrimination accuracy = 93.4%, SEM = 0.84%; non-significant shape bias, t(20) = 1.2, p = 0.25, d = 0.26; mean target-detection rate = 88.1%, SEM = 3.37%; NB chance target-detection accuracy would be 8.3%; false alarm rate = 0.8%). With our design, we decided against using a direct manipulation of difficulty (such as spacing some colors closer to each other in color space) as this would have made it difficult to interpret any effect of confidence due to its inherent confound with a difficulty manipulation. Instead, we relied upon natural fluctuations in confidence, caused for example by individual preferences for colors or fatigue. Figure 2A shows that average, unaided memory performance varied across colors with some colors (e.g. the 4th color, a shade of green) being associated with lower accuracy when participants had to remember this target color unaided by a reminder. Moreover, this figure shows that not all participants had the same inherent color-difficulty profile and that instead some participants perceived particular colors as more difficult than others. Performance in the Baseline condition was high. Here, an indicator of compliance with instructions is participants' placement of the cursor between two thin lines marked on the scale (Figure 2B). Participants reported that those lines were difficult to see in the scanning session. Nevertheless, their cursor locations peaked around the

marked location and landed within the marked positions on approximately half the trials, $\delta = 47.6\%$.

We next established that the reminders aided participants in their fulfilment of the delayed intentions by comparing target-detection error rates for miniblocks in which participants had to use their own memory (fixation dot stayed white) to miniblocks in which they were allowed to use a reminder (fixation dot took on target color), shown in Figure 2C. In both conditions, error rates were reduced when reminders could be used (F(1,20) = 20.5, p < 0.001, $\eta^2_p = 0.51$; ts > 3.4, ps < 0.01, ds > 0.37 when tested separately for the MetaM and MetaC conditions). Error rates did not differ significantly between conditions, F < 1, nor was there an interaction between the two factors, F < 1.

When asked explicitly after the experiment how similar they perceived the two conditions, participants rated the conditions as similar but not identical (M = 0.68 on a scale from 0 = 'totally different' to 1 = 'exactly the same'; min = 0.28; max = 0.98). In fact, we found that participants' perception of the two conditions differed in how much control participants felt they had over the reminders. On a scale ranging from 0 = no control to 1 = full control, participants rated the MetaM condition with a mean of M = 0.32 (min = 0.00; max = 0.88) and the MetaC condition with a mean of M = 0.80 (min = 0.06; max = 0.98). This difference was significant, t(20) = 6.4; p < 0.001, d = 1.94. This shows that participants were able to grasp the key difference that distinguished the two conditions.

We furthermore aimed to rule out that any condition differences found in the pattern classification analyses could be caused by behavioral differences in how the different ratings were approached. Firstly, Figure 3A shows that the average ratings participants gave for each individual color were almost indistinguishable whether they were giving a metacognitive-monitoring or metacognitive-control rating. In fact, if we correlated the average ratings for each color for each individual participant, there was an average relationship of r = 0.76 with

19 out of 21 participants showing a significant, positive relationship between the MetaM and the MetaC rating for different colors. Relatedly, participants' rating and rating RT distributions for the two types of ratings were closely matched (Figures 3B and 3C). It is important to note that participants did not receive any instructions to use these scales in the same way (except for being asked to use the entire range of the scale in both cases).

Futhermore, neither of the metacognitive rating conditions showed a systematic relationship between confidence and accuracy: For retrospective confidence judgements, it is commonly found that these correlate, that is participants express lower confidence on errors than on correct trials (confidence resolution or type-II sensitivity). In the MetaC condition on the other hand, participants' ratings triggered reminders, so we would expect to see the opposite pattern: Trials for which they expressed a high need for a reminder should naturally be the ones on which they were allowed to offload and error rates should therefore be lower. However, we found no significant difference between correct- and error-trial ratings in any of the four conditions (MetaC reminder, t(14) = 0.2, p = 0.88, d = 0.04; MetaC own memory, t(19) = 0.1, p = 0.95, d = 0.02; MetaM reminder, t(17) = 1.0, p = 0.34, d = 0.30; MetaM own memory, t(20) = 1.1, p = 0.28, d = 0.23; participants with missing data excluded from the respective analysis). We furthermore correlated the dichotomous accuracy vector with our continuous confidence measure for all four data cells, separately for each participant. The distributions of these correlations are shown in the right panels of Figure 3D. None were significantly different from zero, ts < 1.0, ps > 0.32. Taken together, both the prospective nature of the ratings in the present task (i.e. participants might have felt they needed to invest more into trials in which they felt less confident or wanted a reminder more) and our unique offloading design could potentially have led to a reduced confidence resolution, but this was the case for both rating conditions.

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

Univariate fMRI Results

We first performed univariate analyses to identify brain regions activated by the
requirement to encode new intentions and make metacognitive judgements about them. We
therefore averaged across the two metacognition conditions (MetaM and MetaC) and
compared them to the Baseline condition, allowing us to find regions of interest (ROIs)
activated by our task. After family-wise error correcting for multiple comparisons, this
contrast revealed seven regions showing increased BOLD signal in the metacognitive
conditions (see Table 2 and Figures 4A and 4B).
Within the seven ROIs, activity was then compared between the metacognition
conditions. More specifically, activity was extracted in two separate contrasts (MetaM >
Baseline and MetaC > Baseline) and then compared. Note that this comparison is orthogonal
to the initial selection contrast and therefore unbiased (Kriegeskorte, Simmons, Bellgowan, &
Baker, 2009). BOLD signal was higher for the MetaC than the MetaM condition in all seven
ROIs (Figure 4C) and this main effect was significant when examined in a ROI (7) x
Condition (2: MetaC/MetaM) repeated-measures ANOVA, $F(1,20) = 8.1$, $p = 0.01$, $\eta_p^2 =$
0.29. There was furthermore a reliable main effect of ROI, $F(7,140) = 7.8$, $p < 0.001$, $\eta_p^2 =$
0.28, as well as a significant interaction of the two factors, $F(7,140) = 3.4$, $p < 0.01$, $\eta_p^2 =$
0.14, reflecting that the absolute signal change and also the difference in signal change was
larger in some ROIs compared to others. Taken together, these results show that regions
which respond to the conditions requiring delayed intentions and metacognitive judgments
showed higher activity when participants rated how much they would like a reminder
(MetaC) compared to how confident they were (MetaM).
We repeated the univariate analyses for deactivations, revealing six "task-negative"
regions showing decreased signal in the conditions requiring delayed intentions and

metacognitive judgments compared with baseline (Figure 5 and Table 3). These regions

included the cingulate and paracingulate cortices, supplementary motor area, supramarginal gyrus, middle and inferior temporal gyri, occipital gyri, and anterior cingulate gyrus. Within these task-negative ROIs, there was more deactivation when participants rated how confident they were (MetaM) compared to how much they would like a reminder (MetaC), however, BOLD signal did not differ significantly between the MetaC and the MetaM condition, F(1,20) = 1.3, p = 0.26, $\eta_p^2 = 0.06$. There was a reliable main effect of ROI, F(5,100) = 18.2, p < 0.001, $\eta_p^2 = 0.48$. The interaction was not significant, F < 1.

Multivariate fMRI Results

The multivariate analyses allowed us to address our two key questions: 1) Do the brain patterns of different metacognitive experiences also distinguish different acts of control? and 2) Can the neural patterns that characterize specific acts of metacognitive control be exhaustively characterized in terms of their associated metacognitive experiences? In a first analysis, we attempted to decode confidence (MetaM). Figures 6A and 6B show the resulting decoding accuracy maps corrected for chance level and multiple comparisons, resulting in nine clusters that contained meaningful information when predicting whether the brain was currently in a low or high confidence state including the anterior cingulate gyrus, parietal occipital sulcus, central sulcus, superior parietal lobule, superior occipital gyrus, cuneus, precuneus, supplementary motor area, occipital fusiform gyrus, calcarine cortex; superior corona radiata, and precentral gyrus (Table 4).

We then repeated the equivalent analysis for the MetaC condition, again successfully decoding whether participants gave a low or high rating (i.e. desire for a reminder) from five clusters including the occipital pole, lateral occipital cortex, superior parietal lobule, superior frontal gyrus (medial segment), middle temporal gyrus (see Table 4). Together these analyses show that the neuroimaging data contains meaningful patterns that distinguish both different

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

metacognitive experiences (low vs. high confidence) and different acts of metacognitive control (low vs. high desire for a reminder).

Having established the existence of meaningful patterns across the brain that distinguish different levels both of MetaM and MetaC, we could then ask whether it was possible to cross-classify the two aspects of metacognition. More specifically, we trained classifiers to distinguish low from high confidence beta images (MetaM) and tested them to predict high versus low MetaC ratings. Note that an inverse relationship is expected between MetaM and MetaC ratings, i.e. low confidence predicts high desire for reminder and vice versa. Therefore one of the scales was inverted in order to perform this analysis. Abovechance classification accuracy can be interpreted as overlapping patterns encoding both MetaM and MetaC. The same analysis was then applied to the opposite direction (train on MetaC, test on MetaM). Importantly, we found overlapping patterns that encode these different types of metacognitive ratings. However, only for the latter analysis direction (train on MetaC, test on MetaM) did we find above-chance classification accuracy after correcting for multiple comparisons. The surviving cluster was located in the left superior and middle frontal gyri. These findings show that brain patterns associated with different metacognitive experiences (low vs. high confidence) also distinguish different acts of metacognitive control (low vs. high desire for a reminder).

To address our second key question, we compared classification accuracy resulting from the two different types of classification analyses described above: within-category (test on MetaM and train on MetaM; test on MetaC and train on MetaC) versus across-category classification (i.e. cross-classification: test on MetaM and train on MetaC; test on MetaC and train on MetaM). We first performed a condition-blind analysis by averaging across all four decoding analyses. This identified ROIs that contain information in one or more of the analyses in an unbiased manner, yielding significant effects in the occipital pole, middle

occipital gyrus, parietal cortex (superior parietal lobule, precuneus), superior frontal gyrus, middle frontal gyrus; precentral gyrus (see Table 5 and Figures 6C and 6D). Within the resulting ROIs, classification accuracies in the four analyses could then be compared (see Figure 6E) to address the question whether decoding accuracy differed significantly between the within-condition classification and the cross-classification analyses. Taking an analogous approach to our univariate analysis, these comparisons were unbiased because they were orthogonal to the analysis used to define the ROIs. We entered the classification accuracies from these regions into a repeated-measures ANOVA with factors ROI, Training condition (MetaM/MetaC) and Classification type (within-condition/between-condition crossclassification). There was a significant main effect of Classification type, F(1,19) = 6.2, p =0.02, $\eta_p^2 = 0.25$, with higher classification accuracy for within-condition classifications than between-condition cross-classifications. This finding can be interpreted as partially overlapping neural representations between MetaM and MetaC as opposed to perfect overlap between the patterns associated with the two aspects of metacognition. Moreover, there was no effect of the conditions on which the classifier was trained or which ROI was analyzed, Fs < 1. We found a significant interaction between ROI and category (within vs. between classification), F(6,114) = 2.4, p = 0.03, $\eta_p^2 = 0.11$, reflecting that the difference between within-condition and across-condition decoding analyses was larger in some ROIs compared to others. No other interactions were significant, Fs < 1. In sum, while our results demonstrate overlapping patterns between metacognitive monitoring and control, they also suggest that patterns of metacognitive control cannot exhaustively be characterized by associated patterns of metacognitive monitoring when participants report their confidence.

522

523

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

Discussion

Metacognitive monitoring is only valuable insofar as it can subsequently influence control. And metacognitive control can only occur if there are metacognitive representations to begin with, which can then be utilized to adjust future behaviour. The two processes must therefore be intimately related, yet the mapping between them requires further study, especially on a neural level. Here we report three main findings: 1) we can separately decode metacognitive monitoring and metacognitive control; 2) brain patterns of different levels of metacognition monitoring (low vs. high confidence) also distinguish different acts of metacognitive control (low vs. high desire for a reminder); and 3) this overlap in patterns while significant is only partial. These findings suggest that patterns of brain activity corresponding to specific acts of metacognitive control are partially, but not fully, characterized by associated acts of metacognitive monitoring.

Our cross classification analysis revealed involvement of the left superior and middle frontal gyri, which form part of the lateral prefrontal cortex (laPFC) in both metacognitive monitoring and control. The role of the laPFC in metacognition has already been highlighted by previous studies, suggesting a role in domain-general metacognition (Morales et al., 2018; see also Vaccaro & Fleming, 2018), in the readout of sensory information as an input for confidence signals (Shekhar & Rahnev, 2018), and more broadly in a mediating role of more rostral parts of laPFC in metacognitive accuracy (Fleming, Weil, Nagy, Dolan, & Rees, 2010; Rounis, Maniscalco, Rothwell, Passingham, & Lau, 2010). Crucially, the laPFC has also been implied in metacognitive control (Qiu et al, 2018; for reviews see Seow, Rouault, Gillan, & Fleming, 2021; Shimamura, 2000; Fleming & Dolan, 2014) matching its more general proposed involvement in cognitive control (MacDonald, Cohen, Stenger, & Carter, 2000; Ridderinkhof et al., 2004). Our study therefore extends this growing body of research that implies an involvement of the lateral prefrontal cortex in metacognition and cognitive control.

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

Seeing as MetaC could not be characterized exhaustively in terms of the MetaM judgments we investigated, this raises the obvious question of which other signals might contribute to MetaC. We consider two main possibilities. The first possibility is that nonmetacognitive signals also play a role in influencing MetaC. A wide variety of signals may be relevant here, such as motivation, the costs and rewards associated with different levels of performance, serial dependencies, fatigue, states of interoceptive and bodily awareness reflecting endogenous signals like arousal (Allen et al., 2016; Hauser et al., 2017; Rouault, McWilliams, Allen, & Fleming, 2018) and so on. This influence of non-metacognitive signals on metacognitive control was already acknowledged in the seminal paper by Nelson and Narens (1990) introducing their metamemory framework. The influence of a wide variety of signals on control is also central to an influential model from the cognitive control literature, the Expected Value of Control model (EVC; Shenhav, Botvinick, & Cohen, 2013). This model emphasizes the flexibility with which different control signals are selected, based on the costs and benefits associated with these signals. The model proposes that the dorsal anterior cingulate cortex integrates both costs and benefits to form the expected value of control. Seeing as MetaC may involve the integration of multiple relevant signals, including the products of MetaM and additional non-metacognitive signals as well, this could potentially explain the greater univariate signal we observed for the MetaC than the MetaM condition. This suggests the incorporation of additional processes into the MetaC judgement beyond those involved in MetaM. We also note that the factor of within-versus crossclassification interacted significantly with region, even though there was no main effect of region. This suggests that the overlap between MetaM and MetaC is greater in some regions than others.

A second possible contribution to the MetaC condition is the integration of additional metacognitive signals, beyond the confidence judgement required by the MetaM condition. In

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

our task, for instance, participants' desire for reminder might have been influenced not only by confidence in their prospective memory but also confidence in their perceptual judgements. Consistent with this, behavioral evidence suggests that confidence judgements are influenced by a variety of domain-general and domain-specific signals (Gilbert, 2015; Kantner, Solinger, Grybinas, Dobbins, 2018; Rouault et al., 2018). Confidence can be regarded as an explicit representation of uncertainty, and uncertainty exists at multiple levels throughout the brain (as noted by the Bayesian brain hypothesis; Knill & Pouget, 2004). Therefore, the metacognitive signals measured in the MetaM condition probably form only a subset of the metacognitive signals which may have contributed to MetaC judgements.

Our paradigm involved measurement of only a single MetaC judgment, which may have been influenced by multiple MetaM signals. In reality, there are multiple types of both MetaM and MetaC. Take for example the situation of a foreign language student studying for a test at her desk during the early evening hours. The student reads a word on a flashcard and we can assume she has access to two relevant metacognitive signals: On the one hand there is the certainty with which the word is perceived in the waning light, the other is the certainty with which the word is recognized from memory. The former confidence should guide her decision whether or not to switch on her desk lamp. The latter confidence should guide her decision whether or not to place the flashcard on the pile marked as 'restudy'. Similarly, the same confidence signal could lead to opposite consequences depending on the situation as shown by Carlebach & Yeung (2021). The authors report that low confidence leads to adviceseeking when the quality of the advice is known and high. However, when the quality of the advice is unknown, people tend to seek advice especially when they have high confidence to test the accuracy of the advisor. How does the brain then 'harvest' these various confidence signals and route them to the appropriate act(s) of metacognitive control? How does it flexibly switch to a different set of signals when required to do so? How are metacognitive

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

signals weighted by past rewards, and how do such weightings shift when our goals change?

Questions such as these could potentially be addressed by adapting the present paradigm to a situation involving two or more forms of metacognitive monitoring and control.

The key finding of our study was the cross-classification between MetaM and MetaC. At a whole-brain corrected threshold, this analysis produced a significant effect in only one direction (train on MetaC and test on MetaM). It is not clear whether this reflects an asymmetry in cross-classification, or simply a thresholding artefact. This could be an interesting question to investigate in future work. Our finding of successful cross classification is in line with the notion that metacognition should be regarded as a cornerstone of cognitive control. Twenty years ago, this point was made prominently by Fernandez-Duque, Baird, & Posner (2000), who drew parallels between metacognitive and executive control functions. Similarly, Yeung & Summerfield (2012, 2014) have suggested that error monitoring, as it is commonly studied in the cognitive-control literature, constitutes an inverse, binary measure of graded confidence. It is therefore not surprising that decision confidence is tracked by a well-established electrophysiological marker of error monitoring, the error positivity (Pe; Boldt & Yeung, 2015). Other empirical examples of links between metacognition and cognitive control are the findings that metacognitive efficiency correlates with cognitive control ability (Drescher, Van Den Bussche, & Desender, 2018) and that confidence modulates the speed accuracy tradeoff on a trial-by-trial basis with participants prioritizing accuracy over response speed after a previous low-confidence decision (Desender et al., 2019). The latter effect is reminiscent of post-error slowing (Rabbitt, 1966; Danielmeier & Ullsperger, 2011; Jentzsch & Dudschig, 2009), one of the most extensively studied effects of the cognitive control literature.

Our findings bear some interesting parallels to another recent decoding study: Mei and colleagues (2020) reported the results from two behavioral experiments, each focused on

a different type of prospective decision (belief of successfully classifying a visual stimulus vs. deciding whether or not to attend to the stimulus during the upcoming trial). The authors found that it was possible to use the data from one experiment (awareness ratings, confidence ratings and accuracy in previous trials) to predict the prospective decision from the respective other experiment and vice versa. This cross-classification analysis therefore highlights similarities of metacognitive monitoring (in this case: beliefs of successfully classifying the upcoming stimulus) and metacognitive control (in this case: decision to attend), showing that both aspects of metacognition appear in the context of the same behavioral precursors.

Despite the theoretical distinction between two binary facets of metacognition and the two different labels assigned to the conditions, the conceptual distinction between the two is not as straightforward as it may seem. For example, our MetaM condition might still be considered to involve an act of metacognitive control in the sense that participants need to use their metacognitive knowledge to control the act of placing the cursor on the scale to indicate low versus high confidence. We suggest that the key distinction between the conditions is that metacognitive monitoring involves relatively direct read-out of metacognitive (e.g. confidence) signals, whereas metacognitive control involves the use of the signals to inform more complex behaviors rather than report the metacognitive experience itself. However, seeing as metacognitive reports are, at least to some degree, inferential in nature (Koriat, 1993), metacognitive monitoring and control might be seen as extreme points on a continuum rather than dichotomous processes.

In sum, our study delineates the similarities and divisions between neural correlates of metacognitive monitoring and control. Ultimately, understanding the link between monitoring and control could inform interventions such as metacognitive training in conditions including brain injury (Fleming et al., 2017), schizophrenia (Moritz & Woodward, 2007) and OCD (Fisher & Wells, 2008). We propose that a full understanding of the

relationship between monitoring and control will require a focus on the ways in which
distinct metacognitive signals are integrated and selectively routed to appropriate acts of
metacognitive control.

653	References
654	Allain, S., Carbonnell, L., Falkenstein, M., Burle, B., & Vidal, F. (2004). The modulation of
655	the Ne-like wave on correct responses foreshadows errors. Neuroscience Letters,
656	372(1-2), 161-166. https://doi.org/10.1016/j.neulet.2004.09.036
657	Allefeld, C., Görgen, K., & Haynes, J. D. (2016). Valid population inference for information
658	based imaging: From the second-level t-test to prevalence inference. NeuroImage, 141,
659	378–392. https://doi.org/10.1016/j.neuroimage.2016.07.040
660	Allen, M., Frank, D., Schwarzkopf, D. S., Fardo, F., Winston, J. S., Hauser, T. U., & Rees, G.
661	(2016). Unexpected arousal modulates the influence of sensory noise on confidence.
662	ELife, 5, 1–17. https://doi.org/10.7554/eLife.18103
663	Allen, M., Glen, J. C., Müllensiefen, D., Schwarzkopf, D. S., Fardo, F., Frank, D., Rees,
664	G. (2017). Metacognitive ability correlates with hippocampal and prefrontal
665	microstructure. NeuroImage, 149(February), 415-423.
666	https://doi.org/10.1016/j.neuroimage.2017.02.008
667	Allen, M., Poggiali, D., Whitaker, K., Marshall, T. R., & Kievit, R. A. (2019). Raincloud
668	plots: a multi-platform tool for robust data visualization. Wellcome Open Research, 4,
669	63. https://doi.org/10.12688/wellcomeopenres.15191.1
670	Bahrami, B., Olsen, K., Latham, P. E., Roepstorff, A., Rees, G., & Frith, C. D. (2010).
671	Optimally interacting minds. Science, 329(5995), 1081–1085.
672	https://doi.org/10.1126/science.1185718
673	Baird, B., Smallwood, J., Gorgolewski, K. J., & Margulies, D. S. (2013). Medial and lateral
674	networks in anterior prefrontal cortex support metacognitive ability for memory and

675	perception. The Journal of Neuroscience: The Official Journal of the Society for
676	Neuroscience, 33(42), 16657–16665. https://doi.org/10.1523/JNEUROSCI.0786-
677	<u>13.2013</u>
678	Bang, D., & Fleming, S. M. (2018). Distinct encoding of decision confidence in human
679	medial prefrontal cortex. Proceedings of the National Academy of Sciences, 115(23),
680	6082–6087. https://doi.org/10.1073/pnas.1800795115
681	Boldt, A., & Gilbert, S. J. (2019). Confidence guides spontaneous cognitive offloading.
682	Cognitive Research: Principles and Implications, 4(1), 45.
683	https://doi.org/10.1186/s41235-019-0195-y
684	Boldt, A., & Yeung, N. (2015). Shared Neural Markers of Decision Confidence and Error
685	Detection. Journal of Neuroscience, 35(8), 3478–3484.
686	https://doi.org/10.1523/JNEUROSCI.0797-14.2015
687	Brett, M., Anton, JL., Valabregue, R., & Poline, JB. (2002). Region of interest analysis
688	using an SPM toolbox. In 8th International Conference on Functional Mapping of the
689	Human Brain. Sendai, Japan.
690	Brown, A. L. (1987). Metacognition, executive control, self-regulation, and other more
691	mysterious mechanisms. In F. E. Weinert & R. H. Kluwe (Eds.), Metacognition,
692	motivation, and understanding (pp. 95–116). Hillsdale, NJ: Erlbaum.
693	Carlebach, N., & Yeung, N. (2020, October 25). Flexible use of confidence to guide advice
694	requests. PsyArXiv. https://doi.org/10.31234/osf.io/ctyqp
695	Cortese, A., Amano, K., Koizumi, A., Kawato, M., & Lau, H. (2016). Multivoxel
696	neurofeedback selectively modulates confidence without changing percentual

697	performance. <i>Nature Communications</i> , /(1), 13669.
698	https://doi.org/10.1038/ncomms13669
699	Danielmeier, C., & Ullsperger, M. (2011). Post-error adjustments. Frontiers in Psychology,
700	2(September), 233. https://doi.org/10.3389/fpsyg.2011.00233
701	De Gardelle, V., & Summerfield, C. (2011). Robust averaging during perceptual judgment.
702	Proceedings of the National Academy of Sciences of the United States of America,
703	108(32), 13341–13346. https://doi.org/10.1073/pnas.1104517108
704	Desender, K., Boldt, A., Verguts, T., & Donner, T. H. (2019). Confidence predicts speed-
705	accuracy tradeoff for subsequent decisions. ELife, 8. https://doi.org/10.7554/eLife.43499
706	Desender, K., Van Opstal, F., & Van Den Bussche, E. (2017). Subjective experience of
707	difficulty depends on multiple cues. Scientific Reports, 7(March), 1–14.
708	https://doi.org/10.1038/srep44222
709	Drescher, L. H., Van den Bussche, E., & Desender, K. (2018). Absence without leave or
710	leave without absence: Examining the interrelations among mind wandering,
711	metacognition and cognitive control. PLoS ONE, 13(2), 1-18.
712	https://doi.org/10.1371/journal.pone.0191639
713	Dunlosky, J., & Connor, L. T. (1997). Age differences in the allocation of study time account
714	for age differences in memory performance. Memory and Cognition, 25(5), 691-700.
715	https://doi.org/10.3758/BF03211311
716	Dunn, T. L., & Risko, E. F. (2016). Toward a Metacognitive Account of Cognitive
717	Offloading. Cognitive Science, 40(5), 1080–1127. https://doi.org/10.1111/cogs.12273

718	Efklides, A. (2008). Metacognition. European Psychologist, 13(4), 277–287.
719	https://doi.org/10.1027/1016-9040.13.4.277
720	Fernandez-Duque, D., Baird, J. A., & Posner, M. I. (2000). Executive attention and
721	metacognitive regulation. Consciousness and Cognition, 9, 288-307.
722	https://doi.org/10.1006/ccog.2000.0447
723	Fisher, P. L., & Wells, A. (2008). Metacognitive therapy for obsessive-compulsive disorder:
724	A case series. Journal of Behavior Therapy and Experimental Psychiatry, 39(2), 117-
725	132. https://doi.org/10.1016/j.jbtep.2006.12.001
726	Flavell, J. H. (1976). Metacognitive aspects of problem solving. In L. B. Resnick (Ed.), <i>The</i>
727	nature of intelligence (pp. 231–236). Hillsdale, NJ: Erlbaum.
728	Fleming, S. M., & Daw, N. D. (2017). Self-evaluation of decision-making: A general
729	Bayesian framework for metacognitive computation. Psychological Review, 124(1), 91-
730	114. https://doi.org/10.1037/rev0000045
731	Fleming, S. M., & Dolan, R. J. (2014). The neural basis of metacognitive ability. In <i>The</i>
732	Cognitive Neuroscience of Metacognition (pp. 245–265). Berlin, Heidelberg: Springer.
733	https://doi.org/10.1007/978-3-642-45190-4_11
734	Fleming, S. M., Huijgen, J., & Dolan, R. J. (2012). Prefrontal contributions to metacognition
735	in perceptual decision making. The Journal of Neuroscience, 32(18), 6117-6125.
736	https://doi.org/10.1523/JNEUROSCI.6489-11.2012
737	Fleming, J., Ownsworth, T., Doig, E., Hutton, L., Griffin, J., Kendall, M., & Shum, D. H. K.
738	(2017). The efficacy of prospective memory rehabilitation plus metacognitive skills

739	training for adults with traumatic brain injury: Study protocol for a randomized
740	controlled trial. Trials, 18(1), 1-11. doi:10.1186/s13063-016-1758-6
741	Fleming, S. M., Weil, R. S., Nagy, Z., Dolan, R. J., & Rees, G. (2010). Relating introspective
742	accuracy to individual differences in brain structure. Science, 329(5998), 1541–1543.
743	https://doi.org/10.1126/science.1191883
744	Fletcher, L., & Carruthers, P. (2012). Metacognition and reasoning. <i>Philosophical</i>
745	Transactions of the Royal Society of London. Series B, Biological Sciences, 367(1594),
746	1366–1378. https://doi.org/10.1098/rstb.2011.0413
747	Fleur, D. S., Bredeweg, B., & van den Bos, W. (2021). Metacognition: ideas and insights
748	from neuro- and educational sciences. Npj Science of Learning, 6(1), 13.
749	https://doi.org/10.1038/s41539-021-00089-5
750	Friston, K. J., Glaser, D. E., Henson, R. N. A., Kiebel, S., Phillips, C., & Ashburner, J.
751	(2002). Classical and Bayesian inference in neuroimaging: Applications. NeuroImage,
752	16(2), 484–512. https://doi.org/10.1006/nimg.2002.1091
753	Frömer, R., Nassar, M. R., Bruckner, R., Stürmer, B., Sommer, W., & Yeung, N. (2021).
754	Response-based outcome predictions and confidence regulate feedback processing and
755	learning. ELife, 10, 1–29. https://doi.org/10.7554/ELIFE.62825
756	Gherman, S., & Philiastides, M. (2018). Human VMPFC encodes early signatures of
757	confidence in perceptual decisions. <i>ELife</i> , 7, 1–28. <u>https://doi.org/10.7554/eLife.38293</u>
758	Gilbert, S. J. (2015). Strategic use of reminders: Influence of both domain-general and task-
759	specific metacognitive confidence, independent of objective memory ability.

760	Consciousness and Cognition, 33, 245–260.
761	https://doi.org/10.1016/j.concog.2015.01.006
762	Grainger, C., Williams, D. M., & Lind, S. E. (2016). Metacognitive monitoring and control
763	processes in children with autism spectrum disorder: Diminished judgement of
764	confidence accuracy. Consciousness and Cognition, 42, 65-74.
765	https://doi.org/10.1016/j.concog.2016.03.003
766	Guggenmos, M., Wilbertz, G., Hebart, M. N., & Sterzer, P. (2016). Mesolimbic confidence
767	signals guide perceptual learning in the absence of external feedback. <i>ELife</i> , 5, 1–19.
768	https://doi.org/10.7554/eLife.13388
769	Hauser, T. U., Allen, M., Purg, N., Moutoussis, M., Rees, G., & Dolan, R. J. (2017).
770	Noradrenaline blockade specifically enhances metacognitive performance, 1–13.
771	https://doi.org/10.7554/eLife.24901
772	Hebart, M. N., Görgen, K., & Haynes, JD. (2015). The Decoding Toolbox (TDT): a
773	versatile software package for multivariate analyses of functional imaging data.
774	Frontiers in Neuroinformatics, 8(January), 1–18.
775	https://doi.org/10.3389/fninf.2014.00088
776	Hebart, M. N., Schriever, Y., Donner, T. H., & Haynes, J. D. (2014). The Relationship
777	between Perceptual Decision Variables and Confidence in the Human Brain. Cerebral
778	Cortex, 26(1), 118–130. https://doi.org/10.1093/cercor/bhu181
779	Hirose, S. (2020). Valid and powerful group statistics for decoding accuracy: Information
780	Prevalence Inference using the i-th order statistic <i>i</i> -test). <i>BioRxiv</i> , 578930.
781	https://doi.org/10.1101/578930

782	Hu, X., Luo, L., & Fleming, S. M. (2019). A role for metamemory in cognitive offloading.
783	Cognition, 193(June). https://doi.org/10.1016/j.cognition.2019.104012
784	Jentzsch, I., & Dudschig, C. (2009). Why do we slow down after an error? Mechanisms
785	underlying the effects of posterror slowing. Quarterly Journal of Experimental
786	Psychology, 62(2), 209–218. https://doi.org/10.1080/17470210802240655
787	Kantner, J., Solinger, L. A., Grybinas, D., & Dobbins, I. G. (2018). Confidence carryover
788	during interleaved memory and perception judgments. Memory and Cognition.
789	https://doi.org/10.3758/s13421-018-0859-8
790	Kleiner, M., Brainard, D., Pelli, D., Ingling, A., Murray, R., & Broussard, C. (2007). What's
791	new in psychtoolbox-3. <i>Perception</i> , 36(14), 1–16.
792	Kluwe, R. H. (1982). Cognitive knowledge and executive control. In D. Griffin (Ed.), Human
793	mind – animal mind (pp. 201–224). New York: Springer.
794	Knill, D. C., & Pouget, A. (2004). The Bayesian brain: The role of uncertainty in neural
795	coding and computation. Trends in Neurosciences, 27(12), 712-719.
796	https://doi.org/10.1016/j.tins.2004.10.007
797	Koriat, A. (1993). How do we know that we know? The accessibility model of the feeling of
798	knowing. Psychological Review, 100(4), 609-639. https://doi.org/10.1037/0033-
799	<u>295X.100.4.609</u>
800	Koriat, A., Ackerman, R., Adiv, S., Lockl, K., & Schneider, W. (2014). The effects of goal-
801	driven and data-driven regulation on metacognitive monitoring during learning: a
802	developmental perspective. Journal of Experimental Psychology: General, 143(1), 386-
803	403. https://doi.org/10.1037/a0031768

804	Koriat, A., Ma'ayan, H., & Nussinson, R. (2006). The intricate relationships between
805	monitoring and control in metacognition: Lessons for the cause-and-effect relation
806	between subjective experience and behavior. Journal of Experimental Psychology:
807	General, 135(1), 36–69. https://doi.org/10.1037/0096-3445.135.1.36
808	Kriegeskorte, N., Goebel, R., & Bandettini, P. (2006). Information-based functional brain
809	mapping. Proceedings of the National Academy of Sciences of the United States of
810	America, 103(10), 3863–3868. https://doi.org/10.1073/pnas.0600244103
811	Kriegeskorte, N., Simmons, W. K., Bellgowan, P. S., & Baker, C. I. (2009). Circular analysis
812	in systems neuroscience: The dangers of double dipping. Nature Neuroscience, 12(5),
813	535–540. https://doi.org/10.1038/nn.2303
814	Lak, A., Okun, M., Moss, M. M., Kepecs, A., Harris, K. D., Carandini, M., Carandini, M.
815	(2020). Dopaminergic and Prefrontal Basis of Learning from Sensory Confidence and
816	Reward Value Article Dopaminergic and Prefrontal Basis of Learning from Sensory
817	Confidence and Reward Value. Neuron, 1–12.
818	https://doi.org/10.1016/j.neuron.2019.11.018
819	MacDonald, A. W., Cohen, J. D., Stenger, V. A., & Carter, C. S. (2000). Dissociating the
820	Role of the Dorsolateral Prefrontal and Anterior Cingulate Cortex in Cognitive Control.
821	Science, 288(5472), 1835–1838. https://doi.org/10.1126/science.288.5472.1835
822	Masset, P., Ott, T., Lak, A., Hirokawa, J., & Kepecs, A. (2020). Behavior- and Modality-
823	General Representation of Confidence in Orbitofrontal Cortex. Cell, 1–15.
824	https://doi.org/10.1016/j.cell.2020.05.022

825	McCurdy, L. Y., Maniscalco, B., Metcalfe, J., Liu, K. Y., de Lange, F. P., & Lau, H. (2013).
826	Anatomical Coupling between Distinct Metacognitive Systems for Memory and Visual
827	Perception. The Journal of Neuroscience: The Official Journal of the Society for
828	Neuroscience, 33(5), 1897–1906. https://doi.org/10.1523/JNEUROSCI.1890-12.2013
829	Mei, N., Rankine, S., Olafsson, E., & Soto, D. (2020). Similar history biases for distinct
830	prospective decisions of self-performance. Scientific Reports, 10(1), 1-13.
831	https://doi.org/10.1038/s41598-020-62719-z
832	Metcalfe, J., & Finn, B. (2008). Evidence that judgments of learning are causally related to
833	study choice. Psychonomic Bulletin & Review, 15(1), 174-179.
834	https://doi.org/10.3758/PBR.15.1.174
835	Miyamoto, K., Setsuie, R., Osada, T., & Miyashita, Y. (2018). Reversible Silencing of the
836	Frontopolar Cortex Selectively Impairs Metacognitive Judgment on Non-experience in
837	Primates. Neuron, 97(4), 980-989.e6. https://doi.org/10.1016/j.neuron.2017.12.040
838	Morales, J., Lau, H., & Fleming, S. M. (2018). Domain-General and Domain-Specific
839	Patterns of Activity Supporting Metacognition in Human Prefrontal Cortex. The Journal
840	of Neuroscience, 38(14), 2360–17. https://doi.org/10.1523/JNEUROSCI.2360-17.2018
841	Morey, R. D. (2008). Confidence Intervals from Normalized Data: A correction to Cousineau
842	(2005). Tutorials in Quantitative Methods for Psychology, 4(2), 61–64.
843	https://doi.org/10.20982/tqmp.04.2.p061
844	Moritz, S., & Woodward, T. S. (2007). Metacognitive training in schizophrenia: From basic
845	research to knowledge translation and intervention. Current Opinion in Psychiatry,
846	20(6), 619–625. https://doi.org/10.1097/YCO.0b013e3282f0b8ed

847	Nelson, 1. O., & Narens, L. (1990). Metamemory: A theoretical framework and new
848	findings. In G. H. Bower (Ed.), The psychology of learning and motivation: Advances in
849	research and theory (pp. 125-173). San Diego, CA: Academic Press.
850	Notebaert, W., Houtman, F., Opstal, F. Van, Gevers, W., Fias, W., & Verguts, T. (2009).
851	Post-error slowing: An orienting account. Cognition, 111(2), 275–279.
852	https://doi.org/10.1016/j.cognition.2009.02.002
853	Odegaard, B., Grimaldi, P., Cho, S. H., Peters, M. A. K., Lau, H., & Basso, M. A. (2018).
854	Superior colliculus neuronal ensemble activity signals optimal rather than subjective
855	confidence. Proceedings of the National Academy of Sciences, 201711628.
856	https://doi.org/10.1073/pnas.1711628115
857	Qiu, L., Su, J., Ni, Y., Bai, Y., Zhang, X., Li, X., & Wan, X. (2018). The neural system of
858	metacognition accompanying decision-making in the prefrontal cortex. PLoS Biology,
859	16(4), e2004037. https://doi.org/10.1371/journal.pbio.2004037
860	Rabbitt, P. M. (1966). Errors and error correction in choice-response tasks. <i>Journal of</i>
861	Experimental Psychology, 71(2), 264–272. https://doi.org/10.1037/h0022853
862	Redshaw, J., Vandersee, J., Bulley, A., & Gilbert, S. J. (2018). Development of children's use
863	of external reminders for hard-to-remember intentions. Child Development, 89(6),
864	2099–2108. https://doi.org/10.1111/cdev.13040
865	Ridderinkhof, K. R., Ullsperger, M., Crone, E. A., & Nieuwenhuis, S. (2004). The role of the
866	medial frontal cortex in cognitive control. Science, 306(5695), 443-447.
867	https://doi.org/10.1126/science.1100301

868	Risko, E. F., & Gilbert, S. J. (2016). Cognitive Offloading. Trends in Cognitive Sciences,
869	20(9), 676–688. https://doi.org/10.1016/j.tics.2016.07.002
870	Rouault, M., McWilliams, A., Allen, M. G., & Fleming, S. M. (2018). Human Metacognition
871	Across Domains: Insights from Individual Differences and Neuroimaging. Personality
872	Neuroscience, 1, 1–13. https://doi.org/10.1017/pen.2018.16
873	Rounis, E., Maniscalco, B., Rothwell, J., Passingham, R., & Lau, H. (2010). Theta-burst
874	transcranial magnetic stimulation to the prefrontal cortex impairs metacognitive visual
875	awareness. Cognitive Neuroscience, 1(3), 165–175.
876	https://doi.org/10.1080/17588921003632529
877	Schulz, L., Fleming, S. M., & Dayan, P. (2021). Metacognitive computations for information
878	search: Confidence in control. BioRxiv. https://doi.org/10.1101/2021.03.01.433342
879	Seow, T. X. F., Rouault, M., Gillan, C. M., & Fleming, S. M. (2021). How local and global
880	metacognition shape mental health. Biological Psychiatry.
881	https://doi.org/10.1016/j.biopsych.2021.05.013
882	Shea, N., Boldt, A., Bang, D., Yeung, N., Heyes, C., & Frith, C. D. (2014). Supra-personal
883	cognitive control and metacognition. Trends in Cognitive Sciences, 18(4), 186-193.
884	https://doi.org/10.1016/j.tics.2014.01.006
885	Shekhar, M., & Rahnev, D. (2018). Distinguishing the Roles of Dorsolateral and Anterior
886	PFC in Visual Metacognition. The Journal of Neuroscience, 38(22), 5078–5087.
887	https://doi.org/10.1523/JNEUROSCI.3484-17.2018

888	Shenhav, A., Botvinick, M. M., & Cohen, J. D. (2013). The expected value of control: an
889	integrative theory of anterior cingulate cortex function. Neuron, 79(2), 217-240.
890	https://doi.org/10.1016/j.neuron.2013.07.007
891	Shimamura, A. P. (2000). The role of the prefrontal cortex in dynamic filtering.
892	Psychobiology, 28(2), 207–218. https://doi.org/10.3758/BF03331979
893	Son, L. K., & Schwartz, B. L. (2009). The relation between metacognitive monitoring and
894	control. Applied Metacognition, 15–38. https://doi.org/10.1017/cbo9780511489976.003
895	Vaccaro, A. G., & Fleming, S. M. (2018). Thinking about thinking: A coordinate-based meta-
896	analysis of neuroimaging studies of metacognitive judgements. Brain and Neuroscience
897	Advances, 2, 239821281881059. https://doi.org/10.1177/2398212818810591
898	Vaghi, M. M., Luyckx, F., Sule, A., Fineberg, N. A., Robbins, T. W., & De Martino, B.
899	(2017). Compulsivity Reveals a Novel Dissociation between Action and Confidence.
900	Neuron, 1–7. https://doi.org/10.1016/j.neuron.2017.09.006
901	Wokke, M. E., Achoui, D., & Cleeremans, A. (2020). Action information contributes to
902	metacognitive decision-making. Scientific Reports, 10(1), 1–15.
903	https://doi.org/10.1038/s41598-020-60382-y
904	Ye, Q., Zou, F., Lau, H., Hu, Y., & Kwok, S. C. (2018). Causal Evidence for Mnemonic
905	Metacognition in Human Precuneus. The Journal of Neuroscience, 38(28), 6379-
906	6387. https://doi.org/10.1523/JNEUROSCI.0660-18.2018
907	Yeung, N., Botvinick, M. M., & Cohen, J. D. (2004). The neural basis of error detection:
908	Conflict monitoring and the error-related negativity. Psychological Review, 111(4),
909	931–959. https://doi.org/10.1037/0033-295X.111.4.939

910	Yeung, N., & Summerfield, C. (2012). Metacognition in human decision-making:
911	Confidence and error monitoring. <i>Philosophical Transactions of the Royal Society B:</i>
912	Biological Sciences, 367(1594), 1310–1321. https://doi.org/10.1098/rstb.2011.0416
913	Yeung, N., & Summerfield, C. (2014). Shared mechanisms for confidence judgements and
914	error detection in human decision making. In S. M. Fleming & C. D. Frith (Eds.), The
915	cognitive neuroscience of metacognition (pp. 147–167). Berlin, Heidelberg: Springer.
916	https://doi.org/10.1007/978-3-642-45190-4_7
917	Yokoyama, O., Miura, N., Watanabe, J., Takemoto, A., Uchida, S., Sugiura, M.,
918	Nakamura, K. (2010). Right frontopolar cortex activity correlates with reliability of
919	retrospective rating of confidence in short-term recognition memory performance.
920	Neuroscience Research, 68(3), 199–206.
921	https://doi.org/10.1016/j.neures.2010.07.2041
922	
923	
924	
925	
926	
927	
928	

	930	Author contributions
	931	AB and SJG designed the research, collected and analyzed the data, and wrote the
Ļ	932	manuscript. AB prepared the figures. Both authors read and approved the final manuscript.
lanuscrip	933	
SC		
<u></u>		
ਲ □		
5		

Figure Legends

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

Figure 1: A) Model of metamemory proposed by Nelson & Narens (1990). The arrows indicate the flow of information. B) Example of a typical MetaM miniblock: Participants were presented with a target color and had to rate how confident they were that they would remember this color later on. It was then revealed whether or not they were allowed to use a reminder for the current miniblock (in this case, the fixation dot took on the target color for the duration of the ongoing task). The ongoing task was a shape discrimination task where participants had to judge whether an array of colored shapes was on average a circle or square. The miniblock ended unpredictably with the target color, which participants were instructed to respond to using a different key. The analysis window for the fMRI analyses is highlighted in yellow. C) Two example blocks showing how participants were alternatingly presented with one Baseline miniblock and then four miniblocks of the current metacognitive rating condition (MetaM or MetaC), shown in different colours. The height of the rectangles reflects the approximate length of the miniblocks, their shading and angle the offloading condition. D) Design matrices for the within-category classifications (first two panels from the right) and cross classifications (second two panels from the right). Lighter colors denote beta images modelling the lower half of ratings in the block in question, whereas darker colors denote higher ratings. Please note that in this example the participant began the experiment with the MetaM condition, but that approximately half of our sample started in the opposite order for balancing reasons. Note also that an inverse relationship between MetaC and MetaM is expected in the cross-classification analysis, hence the ordering of high vs. low MetaC (dark vs. light red) has been flipped in the cross classifications.

956

957

958

Figure 2: A) Target detection accuracy across the range of twelve target colors, placed equidistant in RGB space. Only trials without a reminder were included and the data were

averaged across conditions (metacognitive monitoring and metacognitive control). The thick
black line indicates the sample average, whereas thinner lines represent individual
participants. B) Placement of the cursor in the Baseline condition relative to the marked
position on the scale (shown in red). Each black line indicates the data from a single
participant. C) Target detection accuracy as a function of miniblock condition (Metacognitive
Monitoring and Metacognitive Control) and offloading condition (Own Memory and
Reminder) as a raincloud plot (Allen et al., 2019). The distributions depict the densities of the
participant-wise accuracy averages for each condition. The boxplots show the median and
interquartile range of the data and the lines represent the individual observations.
Figure 3: A) Participants' ratings of Metacognitive Monitoring (low to high confidence) and
Metacognitive Control (high to low need for a reminder; to match orientation of the
confidence scale) showed a similar pattern across the twelve different colors placed
equidistant in RGB space. B and C) Histogram of the B) ratings and C) rating RTs pooled
across all participants as a function of rating type. D) Metacognitive ratings as a function of
objective accuracy and condition shown as participant-wise averages (left panels) and
correlations (right panels). MetaM = metacognitive monitoring; MetaC = metacognitive
control.
Figure 4: A) and B) show the regions of significant signal change in the contrast of (MetaM +
MetaC) > Baseline. A) Significant results were plotted on sagittal $(x = -3)$, coronal $(y = 27)$
and axial ($z = 23$) views of the skull-stripped, mean, normalized structural image. B) Three-
dimensional renderings of results on right hemisphere, left hemisphere, and superior views.
C) Percent signal change for both metacognitive rating conditions in comparison with the

Baseline condition, in regions of interest (ROIs) defined by the contrasts shown in Table 2.

984	MetaM = metacognitive monitoring; MetaC = metacognitive control; r = right; l = left; MFG
985	= middle frontal gyrus; IFG = inferior frontal gyrus; SFG = superior frontal gyrus; SMA =
986	supplementary motor area. Error bars indicate +/- within-subject confidence intervals (95%)
987	according to Morey (2008).
988	
989	Figure 5: A) and B) show the regions of significant signal change in the contrast of Baseline
990	> (MetaM + MetaC). A) Significant results were plotted on sagittal (x = 0), coronal (y = -20)
991	and axial $(z = 0)$ views of the skull-stripped, mean, normalized structural image. B) Three-
992	dimensional renderings of results on right hemisphere, left hemisphere, and superior views.
993	C) Percent signal change for both metacognitive rating conditions in comparison with the
994	Baseline condition, in regions of interest (ROIs) defined by the task-negative contrasts shown
995	in Table 3. MetaM = metacognitive monitoring; MetaC = metacognitive control r = right; 1 =
996	left; SMA = supplementary motor area; MTG = middle temporal gyrus; ITG = inferior
997	temporal gyrus; OcG = occipital gyri. Error bars indicate +/- within-subject confidence
998	intervals (95%) according to Morey (2008).
999	
1000	Figure 6: A) and B) show the above-chance decoding accuracy maps for the condition-
1001	specific classification analyses (blue: train on MetaM, test on MetaM; red: train on MetaC,
1002	test on MetaC; yellow: train on MetaC, test on MetaM). A) Significant results were plotted
1003	on sagittal ($x = -5$), coronal ($y = 7$) and axial ($z = 43$) views of the skull-stripped, mean,
1004	normalized structural image. B) Three-dimensional renderings of results on right hemisphere,
1005	left hemisphere, and superior views. C) and D) show the above-chance decoding accuracy
1006	maps when all four classification analyses were averaged (condition-blind decoding). C)
1007	Significant results were plotted on sagittal ($x = -18$), coronal ($y = 38$) and axial ($z = 3$) views
1002	of the skull-stripped mean normalized structural image. D) Three-dimensional renderings of

resul	ts on right hemisphere, left hemisphere, and superior views. E) Above-chance
class	ification accuracy for all four classification analyses (trained and/or tested on MetaM
and N	MetaC, respectively) in regions of interest (ROIs) defined by a condition-blind selection
contr	east that averaged across all four analyses, listed in Table 5. MetaM = metacognitive
moni	toring; MetaC = metacognitive control, r = right; l = left; SFG = superior frontal gyrus;
MFG	6 (medial frontal gyrus). Error bars indicate +/- within-subject confidence intervals
(95%	according to Morey (2008).

1016 Tables

1017 Table 1: List of experimental conditions.

	Baseline	MetaM	MetaC	
Proportion	20%	40%	40%	
	(32 partial + 32 full	(64 partial + 64 full	(64 partial + 64 full	
	= 64 miniblocks)	olocks) = 128 miniblocks) = 128 mi		
Delayed intention	/	Target color T		
Rating	Cursor placement Very		Sure reminder to	
		very confident		
Reminders	/	50% (random)	based on moving	
			median rating cut-	
			off	

1018

1019 Table 2: Regions of increased signal in the MetaM and MetaC conditions, relative to the

Baseline condition. MetaM = metacognitive monitoring; MetaC = metacognitive control; l =

1021 left; r = right.

Contrast	Label	Laterality	Peak	$\mathbf{k}_{\mathbf{E}}$	p _{FWE}	Z _{max} at
			voxel		cluster-	peak level
			MNI co-		corrected	
			ordinates			
(MetaM +	Occipital and	right and left	30, -55, 5	1338	< 0.001	5.14
MetaC) >	parietal cortex					
Baseline	(calcarine					
	cortex; cuneus;					
	precuneus;					

lateral					
ventricles; al	1				
regions both	1				
and r)					
Middle front	al right	42, 32, 44	345	< 0.001	4.94
gyrus					
Inferior and	left	-42, 20, 26	802	< 0.001	4.68
middle fronta	al				
gyri					
Superior and	right	27, 62, 5	152	0.002	4.49
middle fronts	al				
gyri					
Supplementa	ry left	-6, 23, 44	117	0.009	4.27
motor area					
Angular gyru	ıs left	-57, -55, 44	150	0.003	4.16
Pre- and	right	18, -28, 65	75	0.046	4.08
postcentral g	yri				
Angular gyru	is right	57, -58, 44	87	0.028	4.05

1023 Table 3: Regions of decreased signal in the MetaM and MetaC conditions, relative to the

Baseline condition. MetaM = metacognitive monitoring; MetaC = metacognitive control; l =

1025 left; r = right.

Contrast	Label	Laterality	Peak	\mathbf{k}_{E}	p _{FWE}	Z _{max} at
			voxel		cluster-	peak level
			MNI co-		corrected	

			ordinates			
Baseline >	Cingulate and	right and	3, 2, 35	311	< 0.001	5.99
(MetaM +	paracingulate	left				
MetaC)	cortices; SMA					
	(supplementary					
	motor area; all					
	regions both r					
	and 1)					
	Supramarginal	right	60, -19, 35	809	< 0.001	5.96
	gyrus					
	Supramarginal	left	-66, -28,	1249	< 0.001	5.86
	gyrus		35			
	MTG and ITG	left	-45, -61, 8	534	< 0.001	4.93
	(middle and					
	inferior					
	temporal gyri);					
	OcG (occipital					
	gyri)					
	MTG and ITG	right	57, -55, -4	508	< 0.001	4.92
	(middle and					
	inferior					
	temporal gyri);					
	OcG (occipital					
	gyri)					
	Anterior	right	3, 32, -4	560	< 0.001	4.61

cingulate gyrus			

Table 4: Clusters of above-chance classification accuracy in the four classification analyses.

1028 MetaM = metacognitive monitoring; MetaC = metacognitive control; l = left; r = right.

MVPA	Label	Laterality	Peak	k _E	PFWE	Z _{max} at
			voxel		cluster-	peak level
			MNI co-		corrected	
			ordinates			
MetaM (low	Anterior	left	-3, 17, 26	58	< 0.001	4.57
vs. high	cingulate gyrus					
confidence)						
	Parietal	left	-18, -85, 41	310	< 0.001	4.51
	occipital sulcus					
	Central sulcus	right	21, -28, 53	90	< 0.001	4.51
	Superior	right	24, -70, 50	404	< 0.001	4.43
	parietal lobule;					
	superior					
	occipital gyrus;					
	cuneus;					
	precuneus					
	Supplementary	right	15, 14, 44	149	< 0.001	4.10
	motor area					
	(both 1 and r)					
	Occipital	left	-33, -67, - 19	80	< 0.001	3.94
	fusiform gyrus					

	Calcarine	left	-6, -67, 14	43	0.003	3.89
	cortex; cuneus					
	(all regions					
	both l and r)					
	Superior	left	-24, -13, 32	38	0.006	3.75
	corona radiata		32			
	Precentral	left	-45, 2, 35	26	0.046	3.69
	gyrus					
MetaC (low	Occipital pole	left	-9, -100, 14	93	< 0.001	4.09
vs. high			14			
need for a						
reminder)						
	Lateral	right	36, -73, 5	28	0.029	4.08
	occipital cortex					
	Superior	left	-12, -67, 53	50	0.001	3.79
	parietal lobule		33			
	Superior		0, 29, 50	31	0.017	3.63
	frontal gyrus					
	(medial					
	segment)					
	Middle	right	57, -52, -4	32	0.014	3.58
	temporal gyrus					
MetaC →	Superior and	left	-21, 11, 62	52	< 0.001	4.14
MetaM	middle frontal					
	gyri					

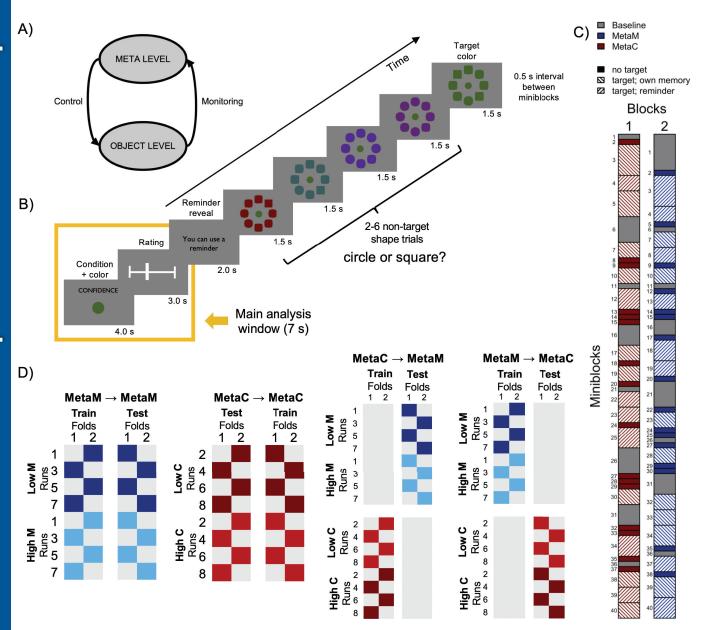
MetaM →	/	/	/	/	/	/
MetaC						

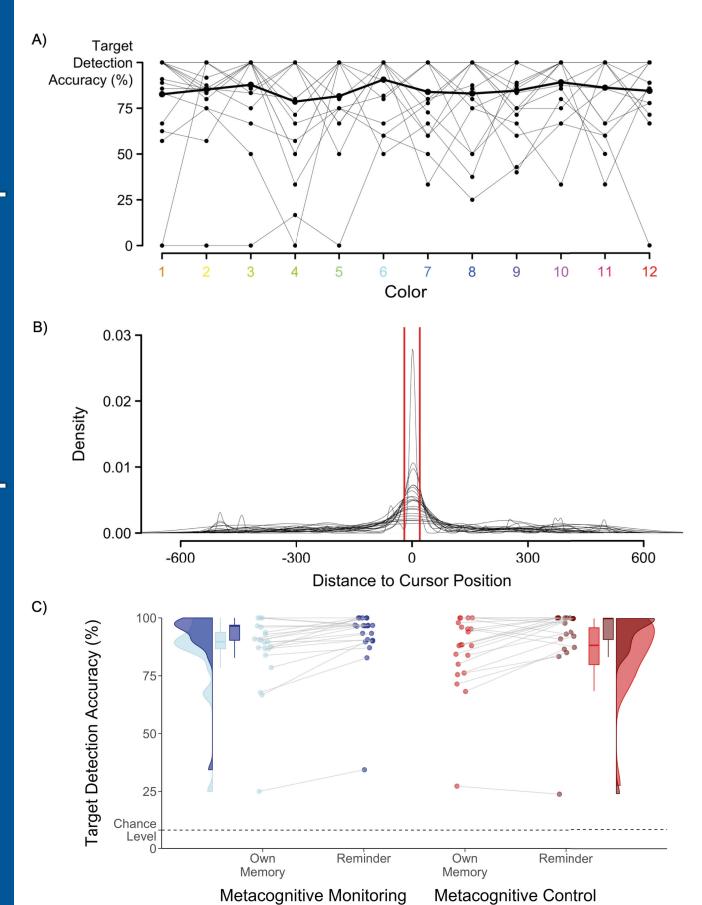
1030 Table 5: Clusters of above-chance classification accuracy in the condition-blind classification

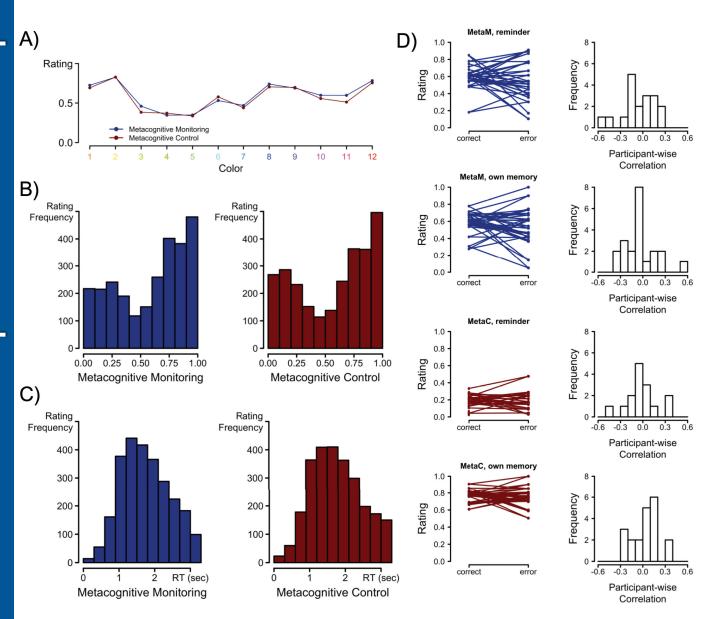
1031 analyses.

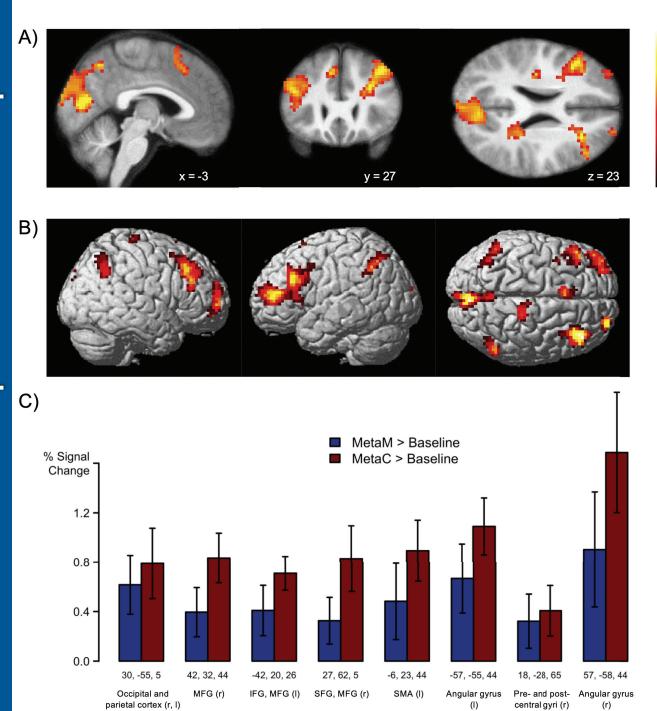
Label	Laterality	Peak voxel	$\mathbf{k}_{\mathbf{E}}$	p _{FWE} cluster-	Z _{max} at peak
		MNI co-		corrected	level
		ordinates			
Occipital pole	right	15, -94, 11	90	0.002	4.54
Occipital pole	left	-24, -91, -1	321	< 0.001	4.51
Middle occipital	left	-30, -70, 26	89	0.002	4.47
gyrus					
Parietal cortex	right	18, -55, 59	115	< 0.001	4.22
(superior parietal					
lobule;					
precuneus)					
Superior frontal	left	-12, 17, 44	114	< 0.001	4.12
gyrus					
Superior and	left	-27, 8, 62	121	< 0.001	3.96
middle frontal					
gyri; precentral					
gyrus					
Parietal cortex	left	-15, -70, 44	91	0.002	3.65
(superior parietal					
lobule;					

precuneus)			









cortex, SMA (I,r)

gyrus (r)

gyrus (I)

OcG (I)

OcG (r)

gyrus (r)

