# The geometry of domain-general performance monitoring representations in the human medial frontal cortex

Zhongzheng Fu[1,2], Danielle Beam[1], Jeffrey M. Chung[3], Chrystal M. Reed[3], Adam N. Mamelak[1], Ralph Adolphs[2,4], Ueli Rutishauser[1,3,4,5,*]

[1] Department of Neurosurgery, Cedars-Sinai Medical Center, Los Angeles, CA, USA

[2] Division of Humanities and Social Sciences, California Institute of Technology, Pasadena, CA, USA

[3] Department of Neurology, Cedars-Sinai Medical Center, Los Angeles, CA, USA

[4] Division of Biology and Bioengineering, California Institute of Technology, Pasadena, CA, USA

[5] Center for Neural Science and Medicine, Department of Biomedical Sciences, Cedars-Sinai Medical Center, Los Angeles, CA, USA

*Correspondence: rutishauseru@csmc.edu

36 **Abstract**

37

38 Flexibly adapting behavior to achieve a desired goal depends on the ability to monitor one's own

39 performance. A key open question is how performance monitoring can be both highly flexible to

40 support multiple tasks and specialized to support specific tasks. We characterized performance

41 monitoring representations by recording single neurons in the human medial frontal cortex

42 (MFC). Subjects performed two tasks that involve three types of cognitive conflict. Neural

43 population representations of conflict, error and control demand generalized across tasks and

44 time while at the same time also encoding task specialization. This arose from a combination of

45 single neurons whose responses were task-invariant and non-linearly mixed. Neurons encoding

46 conflict ex-post served to iteratively update internal estimates of control demand as predicted

47 by a Bayesian model. These findings reveal how the MFC representation of evaluative signals are

48 both abstract and specific, suggesting a mechanism for computing and maintaining control

49 demand estimates across trials and tasks.

50

51

## Introduction

Successful goal-directed behavior in uncertain environments depends critically on continual evaluation of one's own performance (Ullsperger, 2017; Ullsperger et al., 2014). We constantly evaluate whether we made an error, experienced conflict, received reward, or responded fast or slow. Information about past and present performance is in turn used by various downstream processes for cognitive control, affective responses and autonomic homeostasis. The resulting behavioral and physiological adaptations encompass task-specific attentional modulation of perception (Danielmeier et al., 2011; Egner and Hirsch, 2005; King et al., 2010; Purcell and Kiani, 2016), estimation of control demand (Darlington et al., 2018; Jiang et al., 2015; Shenhav et al., 2013), global modulation of motor system readiness (Aron et al., 2007; Danielmeier et al., 2011; King et al., 2010; Murphy et al., 2016; Niv et al., 2007; Wessel and Aron, 2017), emotional state (Bach and Dayan, 2017; Eldar et al., 2016; Shackman et al., 2011), or arousal levels (Crone et al., 2004; Ebitz and Platt, 2015). The medial frontal cortex (MFC) computes and represents many aspects of performance monitoring (Bonini et al., 2014; Carter et al., 1998; Ebitz and Platt, 2015; Fu et al., 2019; Heilbronner and Hayden, 2016; Ito et al., 2003; Kerns et al., 2004; Pouget et al., 2011; Sajad et al., 2019; Sarafyazd and Jazayeri, 2019; Shenhav et al., 2013; Sheth et al., 2012; Stuphorn et al., 2000; Tang et al., 2016; Ullsperger et al., 2014; Wang et al., 2018a, 2018b), making it a primary substrate for communicating evaluative signals to downstream processes (Miller and Cohen, 2001; Shenhav et al., 2013; Ullsperger, 2017; Ullsperger et al., 2014).

On the one hand, cognitive control involves modulating specific sensory or motor processes involved in the task performed, thus requiring the availability of task-specific information in performance monitoring signals, a form of "credit assignment" (McDougle et al., 2016; Sarafyazd and Jazayeri, 2019). For example, an error made by dialing the wrong number should be distinguishable from an error made by calling an old friend by the wrong name because they require different correction mechanisms. On the other hand, humans excel at performing novel tasks with little prior training – an aspect of flexible behavior that is difficult to study in

81    nonhuman primates, since they cannot be verbally instructed to execute an untrained novel task.

82    This kind of cognitive flexibility requires domain-general mechanisms that abstract from the

83    sensorimotor details of any particular task (Bernardi et al., 2020; Minxha et al., 2020). In a novel

84    setting, errors and conflicts can have unanticipated causes, and generic control mechanisms such

85    as slowing all movement down and increasing arousal are adaptive because they buy time to

86    recruit more resources for domain-specific adaptations to take effect (Ullsperger, 2017;

87    Ullsperger et al., 2014). Downstream processes implementing generic adaptations in arousal,

88    global motor suppression (Danielmeier et al., 2011; King et al., 2010; Wessel and Aron, 2017),

89    and urgency (Cavanagh et al., 2011; Heitz and Schall, 2012; Murphy et al., 2016; Thura and Cisek,

90    2017) also depend on the availability of domain-general performance monitoring signals to avoid

91    the need to re-learn how to interpret them for every task. Together, these requirements raise

92    the critical question of how performance monitoring signals are represented in the MFC so that

93    they are accessible to inform both domain-specific and domain-general downstream processes.

94    Answering this question requires recording from multiple single neurons in order to characterize

95    the population-level structure of representations in MFC, about which little is known.

96

97        Theoretically, a fundamental trade-off exists between representations that support task

98    specialization and generalization (DiCarlo and Cox, 2007; Fusi et al., 2016). Specialization requires

99    that as many different conditions as possible can be differentiated from each other by a

100   downstream process that has access to a large subset of the neurons in the representation

101   ("dichotomies" in the case of pairwise differentiations as implemented, for instance, by a linear

102   classifier). This requirement can be fulfilled by increasing the dimensionality of the

103   representation: if there are as many dimensions as differentiations, all possible dichotomies can

104   be read out, in principle. By contrast, generalization requires low-dimensional representations

105   that abstract away ("disentangle") details specific to a single task (DiCarlo and Cox, 2007; Higgins

106   et al., 2016, 2018). Theoretical work shows that the geometry of neural representations can be

107   configured to accommodate both of these seemingly conflicting needs (Bernardi et al., 2020).

108   Given a representation formatted in this way, linear decoders (which represent conservatively

109   what downstream neurons could read out assuming a feed-forward architecture), depending on

110   how they are trained, can either differentiate between many different conditions (specialization)

111   or generalize across conditions, tasks, and time (abstraction, "cross-condition generalization").

112   In such a geometry, closely related conditions in each task/context are placed at similar locations

113   on low-dimensional manifolds; these manifolds in turn are approximately parallel to each other

114   to allow generalization. On each manifold, different conditions within a task/context are placed

115   sufficiently apart to allow maximal differentiation (Bernardi et al., 2020). To satisfy these

116   conditions at the population level, the constituent single neurons must be tuned to combinations

117   of several cognitive variables at once ("non-linear mixed selectivity") (Rigotti et al., 2013) . This

118   set of clear theoretical predictions has been supported empirically to some extent by rare

119   empirical studies in monkeys (Bernardi et al., 2020) and humans (Minxha et al., 2020), but has so

120   far not been explored for the important topic of performance monitoring. Here, we examine the

121   hypothesis that neuronal populations in the human MFC represent conflict, error, and estimated

122   control demand in such a format, making them accessible to the many downstream domain-

123   specific and domain-general processes we reviewed above.

124

125        Estimation of the statistical likelihood of environmental events is essential for efficient

126   goal-directed behavior (Behrens et al., 2007; Jiang et al., 2015; Shenhav et al., 2013). A key aspect

127   of this process is estimating the probability of encountering a situation where cognitive control

128   will be needed (Jiang et al., 2015; Shenhav et al., 2013). Human participants engage reactive or

129   proactive control depending on whether conflict is likely or not (Braver, 2012; Carter et al., 2000;

130   Logan and Zbrodoff, 1979; Tzelgov et al., 1992). The former strategy is efficient when conflict is

131   rarely encountered, whereas the latter is necessary when conflict occurs often (Braver, 2012).

132   Neuroimaging studies have shown that the MFC encodes contexts that are implicitly defined by

133   conflict probability (Carter et al., 2000), but it remains unknown how knowledge about such

134   implicit contexts is acquired from the 'trial-by-trial' feedback provided by performance

135   monitoring. Motivated by prior results (Behrens et al., 2007; Darlington et al., 2018; Jiang et al.,

136   2015; Shenhav et al., 2013; Sohn et al., 2019), we here examine the hypothesis that human MFC

137   neurons signal and continuously update the probability of encountering a control-demanding

138   situation. Since the type of control triggered by different kinds of conflicts differs, this requires

139    representations that support both domain-specific as well as domain-general readouts.

140    Estimating the expected frequency of each type of trial requires integrating information over the

141    history of trials, offering the unique opportunity to examine the mechanisms whereby

142    representations are maintained and updated over time. We model the trial-by-trial changes of

143    activity of the neurons encoding estimated conflict probability as a Bayesian updating process, in

144    which estimated priors are updated iteratively every time after an action is completed. We show

145    that a novel type of conflict signal appears only after an action is completed and thereby provides

146    the critical information for updating the conflict prior into a posterior.

147

148         We recorded single neurons in the MFC while human epilepsy subjects perform two

149    different cognitive control tasks in blocks: the Multi-source Interference task ("MSIT") (Bush and

150    Shin, 2006) and the color-word Stroop task ("Stroop")(Stroop, 1935). The causes of conflict and

151    thus errors are different in the two tasks - stimulus-response spatial incompatibility ("Simon"

152    effect) and/or stimulus conflict ("Flanker" effect) in the MSIT, and stimulus-response

153    incompatibility due to reading colored words in the Stroop task. This allowed us to study how

154    different sources of conflict are encoded within a single task and across tasks, and how errors are

155    encoded across tasks. Subjects were instructed verbally and performed both tasks with little prior

156    practice, which is critical to examine the underlying neural mechanisms (which might be different

157    for extensively practiced tasks as is typically done in animals).

158

159    **Results**

160

161         _Task and behavior_

162         We recorded well-isolated single units in the dorsal ACC and pre-SMA, which are two

163    areas within MFC associated with different aspects of performance monitoring. Subjects

164    performed two speeded response tasks that require cognitive control: the multi-source

165    interference task (MSIT) and the color-word Stroop task (Stroop) (**Fig. 1a** and **Methods**; Stroop:

166    593 neurons in dACC and 607 neurons in pre-SMA across 32 participants (10 females); MSIT: 326

167    neurons in dACC and 412 neurons in pre-SMA in 12 participants (6 females); some patients only

168    performed one of the tasks due to time constraints, **Table S1**). In the MSIT task (**Fig. 1a**, left),

169    conflict arises due to incompatibility between target identity and target location (Simon conflict,

170    or "si" trials) and/or distracting number identity (Flanker conflict, or "fl" trials; trials with both

171    are referred to as "sf" trials). In the Stroop task (**Fig. 1b**, right), conflict arises due to

172    incompatibility between ink color and semantic meaning of the displayed words. In both tasks,

173    sequences of stimuli were randomized, with each type of trial occurring with a fixed probability

174    (In Stroop, 33% of trials had conflict; In MSIT, 15%, 15% and 30% of trials had si, fl, and sf type of

175    conflict, respectively). Subjects were encouraged to respond quickly by an adaptive response

176    threshold (see **Methods**), ensuring maximal task engagement. The different stimulus-response

177    mappings lead to different goal-relevant and irrelevant stimulus features and thus to different

178    kinds of cognitive conflict and reasons for committing errors.

179        Reaction times (RT) were significantly prolonged in the presence of conflicts,

180    demonstrating the Simon and Flanker effect in the MSIT task (**Fig. 1b**, left; average RT of 0.76s,

181    0.86s, 0.93s, 1.03s for non-conflict, si, fl, and sf, respectively) and the Stroop effect in the Stroop

182    task (**Fig. 1b**, right; Stroop: 0.76s vs 0.97s for non-conflict and conflict, respectively). We analyzed

183    participants' sequential performance (RT and accuracy) with a Bayesian online learning

184    framework, building on existing models (Behrens et al., 2007; Jiang et al., 2014, 2015). Our

185    models assume that participants iteratively estimate of how likely it is to encounter a certain type

186    of conflict on the next trial. We refer to this variable as the prior for conflict probability (a real

187    number between 0-1 referred to as 'conflict prior'). Since trial sequences were randomized,

188    subjects could not predict with certainty whether the upcoming trial involved conflict or not.

189    However, they could estimate the conflict probability, which is a task parameter set by the

190    experimenter whose value is unknown to the subject a priori. For MSIT, our models estimated

191    two conflict probabilities (one for si, one for fl) at the same time, based on the finding that both

192    conflicts influenced RT (**Fig. 1b**, left). The trial horizon by which past trials ("conflict history")

193    informed the current estimate was dynamically adjusted by a learning rate parameter, which was

194    also estimated online from the data. In order to obtain an individual conflict prior for every

195    subject (even if the trial sequence was identical), we tuned the iterative estimation model by

196    incorporating RT information, using the expectation-maximization procedure described in prior

197    work (Friston, 2002; Jiang et al., 2015). We modelled the RT generation process as a drift-

198    diffusion process (DDM), where the decision variable represents the *difference* in evidence for

199    the target and distractor response; one bound thus represents the correct outcome whereas the

200    other represents the erroneous outcome. This DDM likelihood function for RT is specified with

201    three hyperparameters: decision bound, drift rates, and drift rate bias (Navarro and Fuss, 2009).

202    Conflict prior entered the decision process by biasing the drift rate (Urai et al., 2019) (**Fig. 1c**,

203    right**;** term $v_{bias}$): the bias term was the scaled Stroop prior for the Stroop model and the sum of

204    independently scaled Simon prior and the Flanker prior for the MSIT model. The same scaling

205    parameters were used regardless of whether the trial had conflict or not because, by definition,

206    the effect of conflict prior started before the trial congruency was revealed.

207          We estimated the hyperparameters of this model using an expectation-maximization

208    algorithm (see **Methods** for details). Conflict probability was estimated iteratively by updating

209    the current prior with the observed conflict type on each trial using the Bayes' law; the updated

210    conflict posterior then served as the prior for the next trial. This online nature of the model

211    captured how human subjects learned about the statistics of conflict trials as they were

212    experienced sequentially. In the following analyses, we refer to the means of the prior and

213    posterior distributions as conflict "prior" (before stimulus onset) or "posterior" (after action

214    completion; **Fig. 1d** shows an example MSIT session). We considered two alternative classes of

215    models with additional free parameters: 1) models estimating conflict probability (Stroop, Simon

216    or Flanker) using all data at once instead of trial-by-trial updating; 2) reinforcement learning

217    models that perform trial-by-trial updating using a constant learning rate. All of these alternative

218    models required offline fitting using all data. Our RT-tuned Bayesian learning model performed

219    significantly better than either class of alternative models in terms of explaining RT and the

220    conflict sequence (**Tables S2 and S3** for a summary of model comparisons). Additionally, RT

221    tuning significantly improved the Bayesian model in terms of explaining RT (**Table S2**, compare

222    columns "RT tuned" and "no RT tuned"; MSIT delta BIC = -348.5; Stroop delta BIC = -508) and the

223    trial congruency sequence (**Table S3**, compare columns "RT tuned" and "no RT tuned"; MSIT delta

224    BIC = -157; Stroop delta BIC = -232). We thus used the RT-tuned Bayesian model for all neural

225    analyses.

226       We next examined what aspects of behavior were related to the model-derived

227    regressors (see **Fig. S2** for values of derived hyperparameters). First, in addition to current trial

228    conflict, the estimated conflict prior had a significantly positive (i.e., an increase in RT) main effect

229    on RT in both tasks (**Fig. 1e;** $\chi^2(1) = 6.75$, p = 0.009 for Simon and $\chi^2(1) = 6.79$, p = 0.009 for

230    Flanker in MSIT; $\chi^2(1) = 28.1$, p < 0.001 for Stroop. Likelihood ratio test). The extent to which RT

231    varied with the conflict prior depended on the type of conflict (in the case of MSIT, Simon and

232    Flanker separately), as indicated by a significant negative interaction term ($\chi^2(1) = 12.94$ for

233    Simon and $\chi^2(1) = 14.2$ for Flanker in MSIT; $\chi^2(1) = 33.3$ for Stroop. p < 0.001 for all conflict types.

234    Likelihood ratio test). This relation between conflict prior and RT remained significant when trial

235    ID was added as a nuisance variable (**Fig. S1a**), or when the conflict prior was estimated without

236    RT tuning (**Fig. S1b**). These RT effects were replicated by online participants collected using

237    Amazon mTurk as behavioral controls (**Fig. S1d**). Conflict prior was systematically related to

238    errors: when conflict was likely, subjects were *less* likely to commit an error on this trial,

239    suggesting that more control was engaged (**Fig. S1c**; for MSIT we only considered "sf" trials where

240    most errors occurred; significant main effect $\chi^2(1) = 6.81$, p = 0.009 for MSIT; significant

241    interaction with non-significant main effect $\chi^2(1) = 18.59$, p < 0.001 for Stroop. Likelihood ratio

242    test). Prior work analyzes the influence conflict on the immediately preceding trial has on RT as

243    a signature of cognitive control (Egner and Hirsch, 2005; Kerns et al., 2004). However, the

244    robustness and generality of the conflict adaptation effect varies between studies (Duthoo et al.,

245    2014; Egner, 2007; Schmidt and De Houwer, 2011) and is not the focus of our study.  Rather, we

246    here consider conflict learning effects that occur over the span of many trials and that exists

247    independent of conflict adaptation, as shown in prior work (Jiang et al., 2015). In fact, in our data,

248    conflict on the immediately preceding trial provided a poor estimate of conflict probability;

249    compared to our model, previous trial conflict alone explained significantly less variance in RT

250    (MSIT delta BIC = -365.6; Stroop delta BIC = -298.5; **Table S2**, compare columns "RT tuned" and

251    "Prev conflict"), suggesting that our participants incorporated conflict information from multiple

252    trials back. Collectively, these behavioral data from two tasks demonstrate that our models that

253    estimate conflict probability online explained variance in RT and error likelihood, demonstrating

254    a proactive engagement of control.

255

256         *Neuronal correlates of performance monitoring signals*

257

258         We focused on three types of epochs for analyses (**Fig. 3a**): baseline before stimulus onset,

259 a 500ms epoch centered at the mid-point between 100ms after stimulus onset and button

260 presses ("ex-ante"), and epochs immediately following button presses ("ex-post"). To assess

261 whether signals relevant for performance monitoring are represented in each epoch, we

262 classified neurons by cognitive variables important for performance monitoring. We identified

263 neurons selective for prior mean or prior variance in the baseline period, for conflict in the ex-

264 ante and ex-post period, and for error, surprise, posterior, and posterior variance in the ex-post

265 period (see single-unit examples in **Fig. 2**; schematic of analysis epochs in **Fig. 3a**; and a summary

266 of overall cell counts in **Fig. 3b**). In MSIT, in order to isolate effects related to the Simon conflict,

267 we refer to the union of "si" and "sf" trials as "Simon trials" and the union of "fl" and no-conflict

268 trials as "non-Simon trials". Similarly, to isolate the effect of Flanker conflict we refer to the union

269 of "fl" and "sf" trials as "Flanker trials", and the union of "si" and no-conflict trials as "non-Flanker

270 trials". Except when noted otherwise, we pooled neurons across dACC and pre-SMA because

271 neuronal responses were similar across areas (**Fig. S3a-b**).

272         Single units tracked aspects of performance monitoring in both tasks (single-unit

273 examples in **Fig. 2**; summary in **Fig. 3b**). During the baseline epoch, a significant proportion of

274 neurons encoded the mean or the variance of the prior distribution for conflict probability (**Fig.**

275 **3b**, blue). In the ex-ante epoch, a significant proportion of neurons encoded conflict (15% in MSIT

276 and 12% in Stroop; Fig. 3b, green), consistent with previous reports (Fu et al., 2019; Sheth et al.,

277 2012). In the ex-post epoch (**Fig. 3b**, yellow), neurons encoded conflict (20% in MSIT; 17% in

278 Stroop), conflict surprise (19% in MSIT; 10% in Stroop), occurrence of errors (22% in MSIT; 19%

279 in Stroop), and the mean and variance of posterior distribution of conflict probability (14/26% in

280 MSIT; 20/12% in Stroop). The signal we refer to as conflict surprise is an unsigned conflict

281 prediction error generated by the experienced conflict given the current prior estimate, a critical

282 component in computing the posterior from the prior (see below). The percentage of units

283 selective for a given variable were similar between the two tasks (**Fig. 3b**).

284        Many identified neurons showed selectivity for more than one cognitive variable (**Fig. 3c-**

285    **d**), suggesting a role in bridging different types of information. Approximately 30% of conflict

286    neurons were active *exclusively* in either the ex-ante, early (0-0.5s after button presses) or late

287    (0.5-1.5s after button presses) ex-post epochs (**Fig. 3c**), with some (~12%) active throughout the

288    trial after stimulus onset ("extended"). The distribution of conflict signals across time was

289    strikingly similar between MSIT and Stroop (**Fig. 3c,** compare left and right). We were particularly

290    intrigued by the prominence of neurons signaling conflict ex-post (15-20% of neurons in both

291    tasks; **Fig. 2c** shows an example), which has not been reported before. This conflict signal, whose

292    timing was too late to be useful for within-trial cognitive control, was more prominent compared

293    to the one found in the ex-ante epoch in both tasks (15% vs 20%, $\chi^2(1)$ = 5.08, p = 0.024 for MSIT;

294    12% vs 16%, $\chi^2(1)$ = 9.19, p = 0.0024 for Stroop, chi-squared test). We note that signaling conflict

295    "after the fact" is predicted by our Bayesian conflict learning framework, in which this ex-post

296    conflict serves as the conflict "outcome" signal indicating that the trial was not only *correct* but

297    also with or without conflict, information necessary for computing the conflict posterior from the

298    prior. We found that many conflict neurons also signaled errors, surprise, posterior, or

299    combinations of these variables (for example, signaling conflict, error, and posterior at the same

300    time) (**Fig. 3d**). This multiplexing of signals depended on the timing of conflict signals. The

301    proportion of conflict neurons that also carried information about the posterior (light green bars)

302    increased significantly towards the end of the ex-post epoch, when updating would be most

303    complete and thus the conflict posterior was computed (compare proportion of conflict neurons

304    that multiplexed posterior information in the late ex-post epoch with those that do so in other

305    epochs; $\chi^2(1)$ = 6.14, p = 0.01 for MSIT; $\chi^2(1)$ = 6.22, p = 0.01 for Stroop, chi-squared test).

306    Consistent with this idea, the group of neurons signaling conflict exclusively in the ex-ante epoch

307    ("ex-ante conflict only") showed the least multiplexing, indicating a primary role in monitoring

308    conflict during action production (proportion of "pure" conflict neurons active only during the

309    ex-ante epoch vs. those that are active in other epochs; $\chi^2(1)$ = 5.31, p = 0.02 for MSIT; $\chi^2(1)$ =

310    8.78, p = 0.003 for Stroop, chi-squared test). Additional evidence for a differential role of ex-ante

311    and ex-post conflict signals is provided by comparing the point in time when these signals were

312    first available in each brain area. Here, we extracted for each *conflict* trial the point in time when

313    spike train was first significantly modulated for ex-ante and ex-post conflict neurons (using a

314    Poisson spike train statistics-based approach (Hanes et al., 1995)). By this measure, ex-ante

315    conflict information was first available in dACC, followed by pre-SMA (**Fig. 3e**; median difference

316    = 138ms; $p < 0.001$, Wilcoxon rank sum test). By contrast, ex-post conflict information was

317    available first in pre-SMA, followed by dACC (**Fig. 3e**; median difference = 161ms; $p = 0.002$,

318    Wilcoxon rank sum test). This pattern is consistent with a leading role of pre-SMA in post-action

319    performance monitoring (Fu et al., 2019), and a leading role of dACC in conflict monitoring during

320    action production. Collectively, these ex-post neuronal responses appeared to reflect the process

321    of updating internal estimates of conflict probability based on present trial outcome as signaled

322    by conflict and error neurons. We next tested this hypothesis.

323          Posterior neurons demonstrated the greatest degree of multiplexing (**Fig. 3f**). Only ~18%

324    of posterior neurons signaled posterior exclusively, with the remainder in addition also signaling

325    prior, conflict, surprise, or a mixture of these. This extensive overlap between posterior signals

326    and each of these ex-post constituents might reflect the computation of the conflict posterior,

327    which would involve all these variables. We next tested whether prior neurons (which are

328    selected during the baseline period) changed their spike rates to reflect the updating process in

329    the ex-post epoch (1s after button press). If a neuron correlates with prior on a trial-by-trial basis

330    and the prior is updated into the posterior after each action, the spike rates of this neuron should

331    reflect this updating. As a neural measure of updating, we used the difference of mean-removed

332    firing rates in two epochs: the early ex-post epoch (1s after button presses) and the baseline. As

333    a behavioral measure of updating, we used the numerical difference between posterior and prior

334    means as estimated by the Bayesian models. We then correlated these two trial-by-trial

335    measures for each prior neuron. Across all prior neurons, correlation was significantly positive

336    for all types of conflict priors (**Fig. 3g**, $p < 0.001$, t test against zero. Mean correlations in Simon,

337    Flanker and Stroop are 0.042, 0.032, 0.065, respectively). This result indicates that prior neurons

338    changed their spike rates in the early ex-post epoch, where the conflict outcome was revealed,

339    to reflect the updated posterior. Together, these data demonstrate that a potential role for the

340    ex-post monitoring signals is to update an online estimate of conflict probability.

341

342       *Event-related potentials that reflect activity of prior cells*

343

344       The intracranial EEG data recorded simultaneously with the single units revealed an event

345    related potential following button presses on correct trials (**Fig. 3h**; "CRP", or correct-related

346    potentials). Event-related potentials (ERPs) reflect synchronous postsynaptic potentials of

347    cortical pyramidal neurons within the cortical microcircuitry (Buzsáki et al., 2012; Herrera et al.,

348    2020; Woodman, 2010). Similar to the ex-post neurons we investigated, CRPs on average

349    followed button presses, had larger amplitude on conflict trials compared to non-conflict trials

350    (**Fig. S3c**, $\chi^2(1) = 21.05$, $p < 0.001$, likelihood ratio test) and showed an interaction effect between

351    trial congruency and conflict prior (a measure of conflict surprise; $\chi^2(1) = 8.48$, $p < 0.001$,

352    likelihood ratio test), carrying population-level information important for updating conflict prior.

353    We thus hypothesized that these prominent ERPs might represent inputs for the prior neurons

354    recorded simultaneously. We tested whether variance in the spike counts of prior neurons could

355    be explained by the CRP amplitude for each point in time across the trial (mixed-effect Poisson

356    regression models tested with likelihood ratio test, see methods). We investigated dACC and pre-

357    SMA separately, consistent with our previous work (Fu et al., 2019). We found that the activity

358    of prior neurons in both dACC and pre-SMA around button presses was significantly correlated

359    with the CRP amplitude on a trial-by-trial basis (**Fig. 3i-j**, Poisson mixed-effect regression model,

360    which included RT and prior as nuisance variables. $p < 0.01$ for all time bins marked by black dots

361    on top, likelihood ratio test. Multiple comparisons were corrected for using the false-discovery

362    rate method), but with earlier onset in dACC than in pre-SMA (0s vs 0.325s after button presses).

363    This indicates that the CRP amplitude (which occurs in the ex-post period) predicted the activity

364    of prior neurons around button presses on a trial-by-trial basis, revealing a neuronal correlate for

365    this prominent ERP.

366

367       *Biophysical basis for encoding of prior/posterior*

368

369       Estimating priors/posteriors in our task necessitates the integration and maintenance of

370    information across multiple trials, a non-trivial property of neural circuitry (Wang, 2002). We

371 therefore investigated whether the functional properties of neurons that encoded priors differed

372 from those that did not. The metrics we used was the temporal correlation profile of baseline

373 spike counts across trials and the width of the extracellular waveform. We chose these properties

374 because the autocorrelation of spike counts at rest of PFC neurons predicts a neuron's

375 participation in working memory (Cavanagh et al., 2018) as well as value coding (Cavanagh et al.,

376 2016), and PFC neurons that encode past reward outcomes have narrower waveforms (Kawai et

377 al., 2019). The timescale of autocorrelation we sought to investigate here is over the span of

378 minutes (multiple trials). To this end, we employed Detrended Fluctuation Analysis (DFA) (Peng

379 et al., 1994) to quantify the self-similarity of baseline spike counts for each neuron, treating the

380 trial-by-trial baseline spike counts as time series data. DFA provides a measure of self-similarity

381 closely related to the slope of the power spectrum (and thus the autocorrelation), but without

382 assuming stationarity. A DFA $\alpha$ value greater than 0.5 indicates a positively correlated process,

383 whereas $\alpha$ = 0.5 indicates an uncorrelated process. We found that in both tasks, neurons

384 representing priors had significantly higher DFA $\alpha$ values compared to other categories of

385 neurons (**Fig. 4a-b**, left panels; $p < 0.001$, ANOVA), with DFA $\alpha$ positively correlated with the

386 strength of prior information carried by a particular neuron (**Fig. 4a-b**, right panels; $p < 0.001$, r =

387 0.24 for MSIT, $p < 0.001$, r = 0.21 for Stroop, Spearman's rank correlation. Separate data were

388 used for computing these two metrics, see **Methods** for detail).

389 We next investigated the relation between a neuron's tendency for long-term

390 maintenance of information (as indicated by $\alpha > 0.5$) and its spike width, a biophysical measure

391 that differs between different types of cells (Bean, 2007; Mosher et al., 2020). The relation

392 between DFA $\alpha$ value, autocorrelation, slope of the power spectrum, and spike width can be seen

393 in the two example neurons shown in **Fig. 4c-f**. The orange neuron, which had a $\alpha$ = 0.91, had a

394 narrower spike waveform, larger autocorrelation and steeper power spectrum slope than the

395 gray neuron with $\alpha$ = 0.54. Across all recorded neurons in both tasks, DFA $\alpha$ values were

396 negatively correlated with spike width (**Fig. 4g-h**; r = -0.19 in MSIT, r = -0.12 in Stroop, $p < 0.001$

397 in both cases, Spearman's rank correlation). Neurons encoding conflict prior/posterior, which

398 requires long-term maintenance, in either task had significantly narrower spike waveforms than

399 all other recorded neurons (**Fig. 4g-h**, right; $p < 0.001$, Wilcoxon's rank sum test). Taken together,

400    these data establish that the long-range temporal correlation of baseline spike counts is an

401    intrinsic firing property of neurons that was predictive of the neuron's spike width as well as the

402    encoding strength of conflict prior/posterior. Neurons that represent conflict priors/posteriors

403    appear to be biophysically distinct and of a different cell type from those that do not code such

404    information, due to their systematically different extracellular waveform and firing properties.

405

406

407    *Temporal progression of performance monitoring signals*

408

409    Given the diversity of firing dynamics seen at the single neuron level (**Figs. 2** and **3**), we

410    next examined the temporal dynamics and robustness of performance monitoring signals at the

411    population level using decoding. We used linear classifiers that had access to *all* recorded

412    neurons, which represents a conservative measure of information available to downstream

413    neurons (Fusi et al., 2016). We first focused on within-time decoding (i.e., training and testing a

414    decoder using data collected in the same epoch). Error, conflict, prior and posterior could be

415    decoded reliably on single trials with high cross-validated accuracy (**Fig. 5a-g** and **Fig. S5a-d**,

416    dotted traces or dotted square shows within-time decoding accuracy). In both tasks, the error

417    signal was decodable with high accuracy throughout the whole ex-post epoch, consistent with

418    our previous report of its role in mediating post-error RT adjustments (Fu et al., 2019) (**Fig. 5a-b**,

419    dotted line). Decoding performance of conflict in both tasks peaked first in the ex-ante epoch

420    (**Fig. 5c-e**, dotted line in the green shading) and then again in the early ex-post epoch (**Fig. 5c-e**,

421    dotted line in the orange shading), before gradually decreasing towards the end of the trial. This

422    time course is consistent with a putative role in estimating conflict probability: conflict is first

423    monitored before committing to a response (ex-ante), followed by a representation of the

424    detected conflict as an outcome signal after button press (ex-post).

425    We next investigated whether the neural code changed over time by using cross-temporal

426    generalization analysis (i.e., training and testing a decoder in different periods of time). We tested

427    the temporal generalization performance of error and conflict decoders trained using data from

428    three defined ROIs: the ex-ante epoch (0.5s; green shading), the early ex-post epoch (0-0.5s after

429    button presses; orange shading) and the late ex-post epoch (0.75-1.25s after button presses; blue

430    shading). In terms of error coding, the early ex-post decoder generalized poorly to later periods

431    (**Fig. 5a-b**; orange line) whereas the late ex-post decoder generalized well across early and late

432    ex-post epochs (**Fig. 5a-b**; blue line). This interesting asymmetry in generalization suggests that

433    there are two groups of error neurons, one signaling errors strongly but transiently, and one

434    signaling error persistently throughout the entire ex-post epochs. In terms of conflict coding,

435    population decoders performed well only within the training epochs but generalized poorly to

436    other epochs for both tasks and all types of conflict (**Fig. 5c-e**), suggesting a dynamic coding

437    patterns for conflict that changed rapidly as the trial unfolds. In particular, the ex-ante decoder

438    did not decode conflict above chance in the ex-post epochs (green traces in **Fig. 5c-e**). This

439    confirms our single neuron findings that the ex-post conflict signals were not simply a

440    continuation of ex-ante conflict signals, but rather signals carried by different groups of neurons

441    at different points in time.

442         For the population coding patterns of prior/posterior, we took a region of interest (ROI)

443    approach given their slow-varying nature, using the baseline and ex-post epochs for prior and

444    posterior, respectively. Since the conflict priors are continuously valued and differed between

445    sessions, we binned trials using quartiles of conflict prior to aggregate data across sessions

446    (labelling trials by four prior levels). We also binned the trials by quartiles of posterior for

447    posterior-related analyses (labelling trials by four posterior levels). We then trained a linear

448    decoder to differentiate between priors/posteriors of two different quantiles. The prior and

449    posterior decoder could differentiate between all pairs of prior/posterior quantiles with high

450    accuracy (**Fig. 5f-g** and **Fig. S5a-d**; within-ROI decoding, upper or lower triangular matrices

451    enclosed by dotted boxes), with accuracy scaling with the distance between pairs of quartiles

452    (i.e., higher accuracy for differentiating $1^{st}$ vs. $4^{th}$ than for $1^{st}$ vs. $2^{nd}$ levels). The prior decoders

453    are able to decode all pairs of posterior levels with high accuracy and vice-versa (**Fig. 5f-g** and **Fig.**

454    **S5a-d**, plots *not* enclosed by dotted boxes), indicating that the representation of prior/posterior

455    is stable across time.

456         Decoding performance for error, conflict and prior/posterior had similar temporal profiles

457    in both dACC and pre-SMA, but with higher decoding accuracy in pre-SMA (**Fig. S4a-e**). Notably,

458    conflict in the immediately preceding trial could be decoded only weakly (lower than 60% in

459    accuracy) in the baseline, as expected, and in the early ex-post epoch for Stroop conflict (**Fig. S5e-**

460    **g**). The weak representation of previous conflict is consistent with our observation that the

461    previous conflict alone was a poor predictor of RT compared to the conflict prior (**Table S2**).

462    Together, these data demonstrate that error, conflict and prior/posterior information can be

463    read out trial-by-trial from the MFC population with high accuracy, with dynamic coding patterns

464    for conflict and error and static coding patterns for prior/posterior.

465

466        *State-space representation of conflict*

467

468        We have shown robust encoding of each of the four conflict types involved in both tasks

469    (sf, fi, sf, Stroop) separately, leaving open the question of how the different encoding schemes

470    are related to each other. Are the different types of conflict encoded along a common 'conflict'

471    axis or are they encoded separately with no generalization between the types of conflict? We

472    tested this question in the MSIT task, which has three types of conflict (Simon, Flanker, both). We

473    took as a putative common conflict coding dimension the line that, in neural state space,

474    connects the neural state during "sf" and "none" trials (both conflict vs. no conflict; **Fig. 6a-b**,

475    dotted lines). Projecting left-out single trials from all four trial types onto this coding dimension

476    allowed differentiation between all pairs of conflict conditions in the ex-ante epoch (**Fig. 6b**, left).

477    and all pairs but one (si vs. sf) in the ex-post epochs (**Fig. 6b**, right). Importantly, this result holds

478    even when RT was equalized across the four conflict conditions (**Fig. S6c**; si, fl, sf, non-conflict.

479    See **Methods** for RT equalization procedure), suggesting that this conflict coding dimension was

480    independent of trial difficulty for which RT is a proxy (Gratton et al., 1992). We next investigated

481    whether Simon and Flanker conflict encoding is related to each other by projecting the activity

482    of single trials onto the coding dimension formed by connecting, in the neural state space, the

483    mean of Simon (si+sf) with the mean of non-Simon (fl+none) separately for each time bin (and

484    vice-versa for Flanker (fl+sf) vs non-Flanker (si+none)). Data for testing were held out (not used

485    for constructing the coding dimensions). Coding dimensions for one type of conflict allowed

486    decoding of the other type of conflict with high accuracy (**Fig. 6c**; black trace, coding dimension

487     of Flanker tested with Simon vs. non-Simon; gray trace, coding dimension of Simon tested with

488     Flanker vs. non-Flanker). Together, these data demonstrate that within a single cognitive task,

489     the MFC population formed a conflict representation that generalized across two types of conflict

490     while at the same time also allowed maximal separation between the different types of conflict,

491     a geometry that supports both abstraction as well as task-specific specialization (Bernardi et al.,

492     2020).

493          When two types of conflict coincide on a trial ("sf" trials in MSIT), is the neural state

494     occupied equal to the sum of the two states occupied by the components ("si" and "fl")? In other

495     words, is the Simon and Flanker representation compositional? Perfect compositionality implies

496     that the vectors for the four trial types ("si", "fl", "none", "sf") are coplanar and form a

497     parallelogram, with the "sf" vector being the diagonal and the opposite sides being parallel to

498     each other (**Fig. 6a-b**). We tested this prediction of parallelism (**Fig. 6a-b**, orange edges and blue

499     edges, respectively) using decoding. If the opposing sides are parallel, a decoder trained to

500     differentiate the two classes connected by one edge should be able to decode two classes

501     connected by the opposite edge (and vice-versa). We found that this was largely the case for both

502     ex-ante and ex-post conflict representation: a decoder trained to differentiate "sf" from "fl" trials,

503     which is simply the axis connecting "sf" and "fl" (orange edge in **Fig. 6a**), was able to differentiate

504     "si" from non-conflict trials projected to this axis above chance, and vice versa (**Fig. 6d**, p < 0.001

505     for both the ex-ante and ex-post data, permutation test). The same was true for the other pair

506     of parallel edges (**Fig. 6d**, testing blue edges in **Fig. 6a**; p < 0.001 for both the ex-ante and ex-post

507     data, permutation test). The parallelism was not perfect because the decoding accuracy, while

508     above chance, was relatively low (< 70%) compared to the performance reached when decoding

509     individual types of conflict (**Fig. 6c**). This structure of the representation was disrupted on error

510     trials, in which generalization performance dropped significantly in the ex-ante (**Fig. S6d**; for both

511     edges, 68% and 58% on correct trials vs. 56% and 47% on error trials) as well as the ex-post epoch

512     (**Fig. S6d**; for both edges, 55% and 66% on correct trials vs. 51% and 59% on error trials) on error

513     trials. Lastly, we examined which neurons contributed to the deviation that keeps the axes from

514     being perfectly parallel and thus perfectly compositional, which was assessed by the mismatch

515     between the actual location of "sf" and the predicted location by vector addition of fl + si.

516    Neurons that encoded Simon and Flanker non-linearly (as measured by the F statistic of the

517    interaction term between Simon and Flanker derived from an ANOVA model) contributed the

518    most to the deviation from linear additivity at the population level (**Fig. 6e**, r = 0.74, p < 0.001,

519    for ex-post data; **Fig. S6e**, r = 0.75, p < 0.001, for ex-ante data; Spearman's rank correlation).

520    Collectively these data suggested that in the MSIT task, neural representations of conflict were

521    structured in a compositional way that separated the four conflict conditions in a parallelogram.

522    This geometry was disrupted on error trials, indicating that this representation was behaviorally

523    relevant.

524

525          *State-space representation of prior/posterior*

526

527          We reasoned that the conflict prior can be viewed as a state (an initial condition) that is

528    present before stimulus and to which the population returns after completing a trial. To test this

529    idea, we again binned trials using quartiles of prior/posterior of each trial (labelling trials with 4

530    prior or posterior levels) and aggregated data across the population. Plotting the neural dynamics

531    in a low-dimensional space spanned by three principal components with largest variance

532    explained (PCA is unsupervised and has no access to the ordinal relation between prior/posterior

533    levels) revealed that the variability across different levels of prior/posterior (~8% of variance) was

534    captured mostly by a single axis (PC3s in **Fig. 6f-h**; green dots mark trial start, red dot trial end),

535    which was orthogonal to most of the time-dependent state changes (captured by PC1s and PC2s,

536    ~68% of variance). During the baseline period (green to cyan dots in **Fig. 6f-h**) neural state

537    changed with low speed (**Fig. 6i** and **Fig. S6f-g**), whereas the speed of changed increased

538    significantly after stimulus onset (**Fig. 6i** and **Fig. S6f-g**, red; p < 0.001, paired t-test), eventually

539    returning to baseline near the starting position (red dots in **Fig. 6f-h**). The distance between the

540    four trajectories was kept approximately constant at all time (**Fig. 6i** and **Fig. S6f-g**), consistent

541    with the levels of prior/posterior being states stably maintained at the individual trial level.

542    Remarkably, the state-space trajectories are not only stable but also preserves the ordinal

543    relation between prior/posterior levels: projection values onto the PC that captured the most

544    variance across prior/posterior levels were arranged in an order consistent with the

545    prior/posterior levels, even though PCA did not have access to such ordinal information (**Fig. 6j-**

546    **l**, see **Legends** for statistics of the multinomial logistic regression). Taken together, the MFC

547    representation of the conflict prior/posterior information is low-dimensional, stable across time

548    and parametric, consistent with the dynamics of line attractors.

549

550    *Domain-general performance monitoring signals at the population*

551

552    What was the relationship between the performance monitoring signals we documented

553    separately in the MSIT and Stroop tasks? Did neurons specialize in encoding a given signal in only

554    a given task or did neurons form a domain-general representation across tasks? If the latter, was

555    this representation abstract in the sense that information about the task identity or task-specific

556    conditions (e.g., the different types of conflicts in MSIT) were no longer available? To answer

557    these questions, we next analyzed a subset of data in which we tracked the same neuron in both

558    Stroop and MSIT (see **Table S1** for a tally of recordings). Note that participants had no knowledge

559    of the second experiment they were going to perform when they performed the first (which was

560    either MIST or Stroop), thereby allowing us to ask how two novel tasks with no prior practice

561    engaged MFC. We used demixed PCA (dPCA) (Kobak et al., 2016) on the neural activity recorded

562    across both tasks to identify coding dimensions for error, conflict and conflict prior/posterior

563    which were stable across time and with task information maximally marginalized out. Namely,

564    the goal is to factorize data into performance monitoring signals, non-specific temporal dynamics

565    as well as signals related to task sets (see **Methods** for details). To match the number of

566    conditions between tasks, we picked non-conflict and "sf" trials in MSIT and non-conflict and

567    Stroop trials in the Stroop task to construct the task-invariant conflict dimension. To assess

568    whether the extracted coding dimensions were meaningful statistically, we used them to decode

569    left-out data that were not used to construct these dimensions. To test generalization across

570    tasks, we first projected both left-out training and testing data onto a dPCA dimension, and then

571    classified the testing data in task using training data from the other task.

572    We found that the dPCA task-invariant coding dimensions identified this way explained

573    between 9-12% of the variance and allowed training of a decoder in one task and testing it in the

574    other with high accuracy (**Fig. 7a-c**; > 80% accuracy for error and conflicts). To quantify the extent

575    of "demixing" of performance monitoring signals from task set information, we computed the

576    angle between the task invariant coding dimensions and their corresponding task dimension. The

577    error dimension supported task-invariant decoding of error throughout the ex-post epoch (**Fig.**

578    **7a**; significant clusters see horizonal bars; p < 0.001, cluster-based permutation tests). The angle

579    between the error coding dimension and the task dimension was 94.47° and did not deviate

580    significantly from orthogonality (p = 0.53, tau = -0.032, Kendall rank correlation). The conflict

581    dimension supported task-invariant decoding in the ex-ante and the early ex-post epoch (**Fig. 7b-**

582    **c**; significant clusters see horizonal bars; p < 0.001, cluster-based permutation tests). The angles

583    between Stroop-Simon conflict coding dimensions and the task dimension were 81.13°, which

584    did not differ significantly from orthogonality (p=0.19, tau = 0.048, Kendall rank correlation). The

585    angles between Stroop-Flanker conflict coding dimensions and the task dimension were 78.6°,

586    which deviated weakly from orthogonality (p=0.02, tau = 0.086, Kendall rank correlation). This

587    task generalizability did not compromise the capacity of this coding dimension to separate

588    different kinds of conflict within MSIT: classifiers could differentiate between 5 out of 6 pairs of

589    conflict conditions with high accuracy (60% - 90%) in both the ex-ante and ex-post periods based

590    on data projected onto the task-invariant dPCA conflict axes (**Fig. 7d**; p values see figure legend,

591    permutation tests). As a control, we repeated all above dPCA decoding analyses using trial

592    conditions equalized by RT (e.g., selecting conflict and non-conflict trials that had similar RTs) and

593    obtained very similar findings (**Fig. S7a-d**). These results suggests that task-invariant

594    representation of error and conflict did not result from the coincidental condition differences in

595    difficulty for which RT is a proxy (Gratton et al., 1992).

596          Similarly, the representation of conflict priors and posteriors also allowed task-invariant

597    decoding of this information while at the same time differentiation between 94% of pairs of

598    prior/posterior levels in both tasks (**Fig. 7e-f**; p values see figure; permutation tests). None of the

599    coding dimensions for conflict prior or posterior were significantly non-orthogonal with the task

600    dimension (angles and Kendall's tau values see **Legend**; p > 0.2 for all, Kendall rank correlation),

601    suggesting that the coding dimensions identified for conflict prior or posterior were significantly

602    "demixed" from the task set dimension, a complete factorization. Consistent with this, the task

603    dimension support decoding of which task a trial was from with very high accuracy (>90% leave-

604    one-out accuracy, p < 0.001 by permutation tests). Together, these data demonstrate that the

605    neural representation of performance monitoring signals in MFC is configured in such a way that

606    it supports generalizability between two different tasks while at the same time also allowing the

607    readout of task-specific information.

608

609            *Domain-general performance monitoring signals at the single-neuron level*

610            What gave rise to the flexible coding scheme that supported both task-invariant and task-

611    specific readouts as revealed above? Was this a population level phenomenon or did individual

612    neurons encode a given variable reliably in both tasks? To answer this question, we quantified

613    cross-task coding stability for each neuron using linear regression (see **Methods**). To do so, we

614    pooled data from both tasks and regressed firing rates against a performance-monitoring

615    variable (error, conflict, prior or posterior), a task indicator, and an interaction term

616    (performance monitoring x task). The statistical significance for each regressors were determined

617    by an F test. We refer to neurons that had a significant main performance-monitoring effect but

618    non-significant interaction as "task invariant", and to neurons that had a significant "performance

619    monitoring x task" interaction as "task dependent" neurons. We selected neurons whose

620    response signaled conflict in ex-ante and ex-post epochs, error in the ex-post epoch, prior in the

621    baseline epoch, and posterior in the ex-post epoch.

622            Out of the selected neurons of each kind (error, conflict, prior), 33-68% were classified as

623    task invariant (**Fig. 7g-j** and **Fig. S7g-k**, red in pie charts). The extent to which a given neuron

624    encoded a performance-monitoring variable by itself (assessed by t statistic for the main effect)

625    and as part of a population (as derived by weight assigned to the neuron by the identified

626    common dPCA coding dimensions) correlated significantly (**Fig. 7g-j**, scatter; see **Legend** for

627    statistics), with the signs of these measures agreeing with each other in most cases (**Fig. 7g-j** and

628    **Fig. S7 g-k**; for cases where signs differed, see **Fig. S7i-j**). Our dPCA analyses marginalized out

629    information about task and time. As a result, neurons selective during either the ex-ante and/or

630    ex-post epochs contributed to the identified common axis and their contribution were thus

631    analyzed separately. On average, neurons identified as "task invariant" or "task dependent" were

632    assigned significantly larger absolute dPCA weights than non-selective neurons ("others" neurons)

633    (**Fig. 7g-j** and **Fig.S7g-k**, dot density plots on the right). While "task-invariant" neurons in many

634    cases had numerically larger dPCA weights than "task-dependent" neurons, on average they did

635    not contribute significantly more to the task-invariant coding dimensions (**Fig. 7g-j** and **Fig.S7 g-**

636    **k**, dot density plots on the right; see **Fig. S7g-k** for the case where task-invariant neurons

637    contributed significantly more than task-dependent neurons). This illustrates how the diversity

638    of encoding schemes at the single neuron level gave rise to a population-level one-dimensional

639    coding dimension that supports a robust domain-general readout of performance monitoring

640    signals with high accuracy in both tasks (**Fig. 7a-f** and **Fig. S7a-f**). A downstream neuron receiving

641    input from MFC could in theory derive domain-general performance monitoring signals simply

642    by taking a linear sum and thresholding, with the weights equal to those assigned by dPCA.

## Discussion

We found that neurons in the human MFC represent estimated control demand (conflict prior) during the baseline, conflict during action production, and conflict outcome and error immediately after an action was performed in both the Stroop and MSIT tasks. A Bayesian conflict learning model that updated conflict probability iteratively after every trial predicted the existence of a novel kind of outcome signal (ex-post conflict signals) and neurons signaling estimated conflict probability, and we identified neurons encoding both of these cognitive variables. Single neurons encoded a diverse array of variables: some neurons encoded conflict, error and conflict prior/posterior in a task-invariant way, some encoded these variables exclusively in one task but not the other, and some multiplexed task information to varying degrees. Neurons with different encoding profiles were randomly and homogenously sampled within MFC with no apparent clustering. Such interdigitated representation patterns would be difficult to detect with fMRI because even one voxel aggregates the activity from hundreds of thousands of neurons. This complexity at the level of single neuron responses precludes a clear interpretation of domain generality or domain specificity and instead requires an analysis of population-level representations in a high-dimensional state-space. The key insight in this study is that the representational geometry of performance monitoring takes advantage of this complex pattern of single neuron responses. Neuronal populations in the human MFC represented performance monitoring signals in a geometry that allows domain-general readout across tasks, while simultaneously also allowing linear decoders to extract task-specific details. Achieving this tradeoff, in turn, requires that the constituent single neurons multiplex information about the different variables required (i.e., show mixed selectivity).

**Domain-general performance monitoring**

The evidence for domain generality of performance monitoring from neuroimaging studies is mixed. Some studies conclude that representations of conflict are domain-general from spatially overlapping BOLD activation maps, but the extent of such overlap could depend on the

672  statistical thresholding used (Fan et al., 2003; Liu et al., 2004). Using multi-voxel pattern analysis

673  (MVPA) that avoids statistical thresholding, one study finds voxel clusters that simultaneously

674  encode different types of conflict in the superior frontal gyrus (Jiang and Egner, 2014). However,

675  the majority of conflict-encoding voxel clusters, notably those found in dACC, are specific to one

676  particular type of conflict; domain-general and domain-specific clusters are distributed in distinct

677  anatomical locations. Direct comparisons between single-unit recordings and BOLD-fMRI have

678  demonstrated that the former can exhibit multivariate representations that cannot be detected

679  with fMRI (Dubois et al., 2015). Additional difficulty for neuroimaging studies is the variation in

680  human cingulate anatomy, which reduces the overall signal resolution when registering to a

681  common template brain (Crosson et al., 1999; Vogt et al., 1995) (a constraint our work does not

682  suffer from since we mapped anatomy in individual brains). In the present study we find that it

683  is the *same* group of neurons within MFC that form a geometry allowing the readout of both

684  domain-general and domain-specific conflict signals. The conflict representation we discovered

685  not only generalizes between two types of conflict involved within a single task (which was found

686  with MVPA-fMRI as well (Jiang and Egner, 2014)), but also between two tasks with completely

687  different stimuli, response requirements, and task rules. A key component of performance

688  monitoring is the ability to detect action errors without relying on external feedback. This type

689  of self-monitoring is a central component of metacognition (Yeung and Summerfield, 2012). In

690  the case of confidence judgments, which are also metacognitive, fMRI results indicate that the

691  same parts of MFC are involved across different cognitive domains (perception or memory,

692  (Morales et al., 2018)), but it remains unknown whether this also holds for error monitoring. We

693  show that a subset of neurons signal errors in a domain general manner across the two tasks and

694  all types of conflict. At the population level, these domain-general error neurons enabled

695  domain-general readouts of self-monitored error across both tasks (**Fig. 7g**). Future work is

696  needed to demonstrate whether this is also the case for metacognitive signals other than errors.

697

698  **Domain-specific performance monitoring**

699

700        We found error neurons independently in both the Stroop task (as previously reported

701        (Fu et al., 2019)) and the MSIT task (a new finding), with a subset of these neurons signaling error

702        only in one or the other task. This would be expected given that what causes an error in the two

703        tasks differs substantially: distraction by the prepotent tendency to read in the case of the Stroop

704        task, and distraction by the tendency to respond to different keys (Simon) or distraction by the

705        flanker number (Flanker) in the case of the MSIT. Importantly, these call for different

706        compensatory mechanisms: for example, Stroop errors can be compensated for by suppressing

707        attention to the word meaning, Simon errors by suppressing attention to the spatial location of

708        the target, and Flanker errors by suppressing attention to the flanking distractors. The same

709        argument applies to the conflict signal; together, these considerations highlight the fact that

710        cognitive control requires information about performance specific to the task performed.

711        Consistent with this requirement, the MFC neural states varied both along a generic conflict

712        dimension that generalized across tasks as well as along all six types of specialized conflict

713        dimensions. By changing connection weights, downstream neuronal processes can flexibly access

714        performance monitoring signals at both extreme as well as intermediate levels of abstraction and

715        drive behavioral adaptations accordingly. Prior work demonstrate the MFC's causal role in credit

716        assignment within a single task by showing that macaque MFC neurons are involved in attributing

717        the cause of an error to either low-level perceptual noise or exogenous changes in the response

718        rule (high-level) (Sarafyazd and Jazayeri, 2019). We found that both at the level of single neuron

719        and population activity, error and conflict signaling offered a robust readout of the task in which

720        these performance disturbances were experienced, even when the subject had no prior task

721        exposure. This task specificity at the population level is supported by the fact that some MFC

722        neurons which did not signal errors or conflict in a previous task started to do so in a novel task

723        with no prior training (**Fig. 7g-I** and **Fig. S7g-h**). These results are broadly consistent with the

724        MFC's role in credit assignment within a task and demonstrate the remarkable flexibility of the

725        MFC performance monitoring circuitry.

726

727        **Compositionality of conflict representation**

728

729     Further insight into how representations can be both general and specific is offered by

730     examining the activity during "sf" conflict trials (both Simon and Flanker conflict) in the MSIT task.

731     The potential compositionality of the representation of two kinds of conflict can be formulated

732     as a generalization problem: if Simon and Flanker conflict are linearly additive, decoders trained

733     to identify the presence of only Simon or Flanker conflict should generalize to the "novel"

734     situation where both types of conflict are present. We found that this was the case, with the

735     neural state approximately equal to the linear vector sum of the two neural states when the two

736     types of conflict are present individually. This suggests that conflict representations are additive

737     to a large extent (with the extent of deviation predicted by the degree of nonlinear mixing

738     present). The (approximate) factorization of conflict representation is important for both

739     domain-specific and domain-general adaptation: when different types of conflict occur

740     simultaneously and the representation can be factorized, downstream processes responsible for

741     resolving each type of conflict can all be initiated. On the other hand, domain-general processes

742     can also read out the representation as a sum and initiate domain-general adaptations.

743

744     **Estimating control demand enabled by ex-post conflict neurons**

745

746     We found that subjects' performance in two conflict tasks was best explained by models

747     that iteratively estimate how likely the next trial is to contain each possible type of conflict. These

748     models integrate information across many trials and outperform the conflict adaptation model

749     whose horizon only includes one trial back in terms of estimating and predicting control demand.

750     We modelled the abstract decision process (which predicts RT and correct/incorrect, but not the

751     actual choices) as a drift diffusion process with different drift rates on a conflict and non-conflict

752     trial. Conflict prior was incorporated by adding a bias to the drift rates. The choice of DDM is

753     informed by prior work demonstrating that sequential effects in perceptual decisions are

754     modelled best with drift rate biases (rather than biases in other factors such as starting point or

755     boundary) (Urai et al., 2019). We chose a Bayesian DDM framework because: 1) it estimates

756     control demand (conflict probability) iteratively as our participants did and is thus neurally

757     feasible; 2) similar models have found success in explaining behavior in cognitive control tasks

758    (Ide et al., 2013; Jiang et al., 2015); 3) it provides trial-by-trial regressors for neural analyses.

759    Notably, this model performed significantly better than one assuming that subjects try to

760    estimate a fixed conflict probability (which is how the task is designed). This shows that subjects

761    were sensitive to the random variability in the trial congruency sequence to adjust their response

762    strategy as reflected in their response times.

763

764        A prediction of our Bayesian learning framework is that during the ex-post period,

765    neurons would signal whether the just completed trial was a conflict or not. This is because this

766    kind of "after the fact" conflict signal is needed to compute the posterior from the prior.

767    Confirming this prediction, we found two kinds of conflict signals: one that occurred after action

768    completion as predicted ("ex-post"), and one that occurred during action production ("ex-ante")

769    as expected (Botvinick et al., 2001; Sheth et al., 2012). Separate groups of neurons gave rise to

770    these two types of conflict signals. To the best of our knowledge, this ex-post coding of conflict

771    is a novel kind of conflict signal not previously documented. We posit that the ex-post conflict

772    signal is an outcome signal (Shenhav et al., 2013) that is used for updating slowly varying

773    representations of estimated conflict probability. Interestingly, there is significant overlap

774    between error neurons and these ex-post conflict neurons. Confirming this, we found that a

775    common coding axis exists that supports decoding of both error and conflict, though the

776    decoding accuracy is significantly lower for conflict than for error (**Fig. S6a-b**). These results

777    suggest the origins of ex-post conflict and error signals may be similar: a putative prediction error

778    computed based on an efference copy (Lo and Wang, 2006). Future work is needed to test this

779    new hypothesis. We also revealed a direct neural correlate of the updating step: firing rate

780    changes of prior neurons during the ex-post period are positively correlated with prior-posterior

781    differences in the model. Altogether, the neuronal responses we found fit remarkably well to the

782    parameters of a Bayesian model that used trial-wise updating to estimate upcoming conflict – a

783    critical ingredient in the control of flexible behavior in changing environments.

784

785    **Prior neurons as a substrate for proactive control**

786

787        The dual mechanisms of control (DMC) framework (Braver, 2012) distinguishes between

788    proactive and reactive control and suggests that proactive control involves sustained anticipatory

789    activity. We posit that the prior-signaling neurons reported in this study can inform proactive

790    control processes about the estimated demand for control in the upcoming trial. Compatible with

791    this view, responses of these prior/posterior neurons in the MFC were sustained and stable

792    across time. At the population level, the prior/posterior representations were parametric with

793    dynamics indicative of a line attractor, from which neural activity departs and then returns to

794    when completing a trial. The distance in neural state space between neural trajectories remained

795    stable across the trial, as expected from a signal that only varies slowly and that reflects a learning

796    process that occurs over multiple trials.

797

798        A likely contribution to the temporal stability of prior coding are the biophysical

799    properties of prior neurons, which differed from other neurons in two ways. First, these neurons

800    exhibited long-range autocorrelations in their baseline spike counts. While spike-count

801    autocorrelations in the range of seconds are known to differ between brain regions (Bernacchia

802    et al., 2011) and between neurons with different tuning (Cavanagh et al., 2016, 2018), here we

803    examined long-range temporal correlation on a substantially longer timescale (minutes) by

804    assessing the self-similarity of trial-by-trial spike counts during the baseline using Detrended

805    Fluctuation Analysis (DFA) (Hardstone et al., 2012). Our finding that DFA values for prior-encoding

806    neurons are high (indicating long-range temporal autocorrelations) suggests that they are ideal

807    substrates for representing a slowly varying internal state. In addition, these neurons tended to

808    have shorter extracellular waveforms. While this relationship is complex (Vigneswaran et al.,

809    2011), neurons with narrower spikes are more likely to be interneurons and the functional role

810    of thin- and broad- spike MFC neurons is often different (Bean, 2007; Sajad et al., 2019). Notably,

811    in macaques, thinner-spike neurons are more likely to have long autocorrelations making them

812    ideal to carry slowly changing information across trials (Kawai et al., 2019). Our findings tie

813    together the biophysical properties of single neurons with their tuning, indicating that conflict

814    prior neurons are ideally suited to carry slowly changing information across trials.
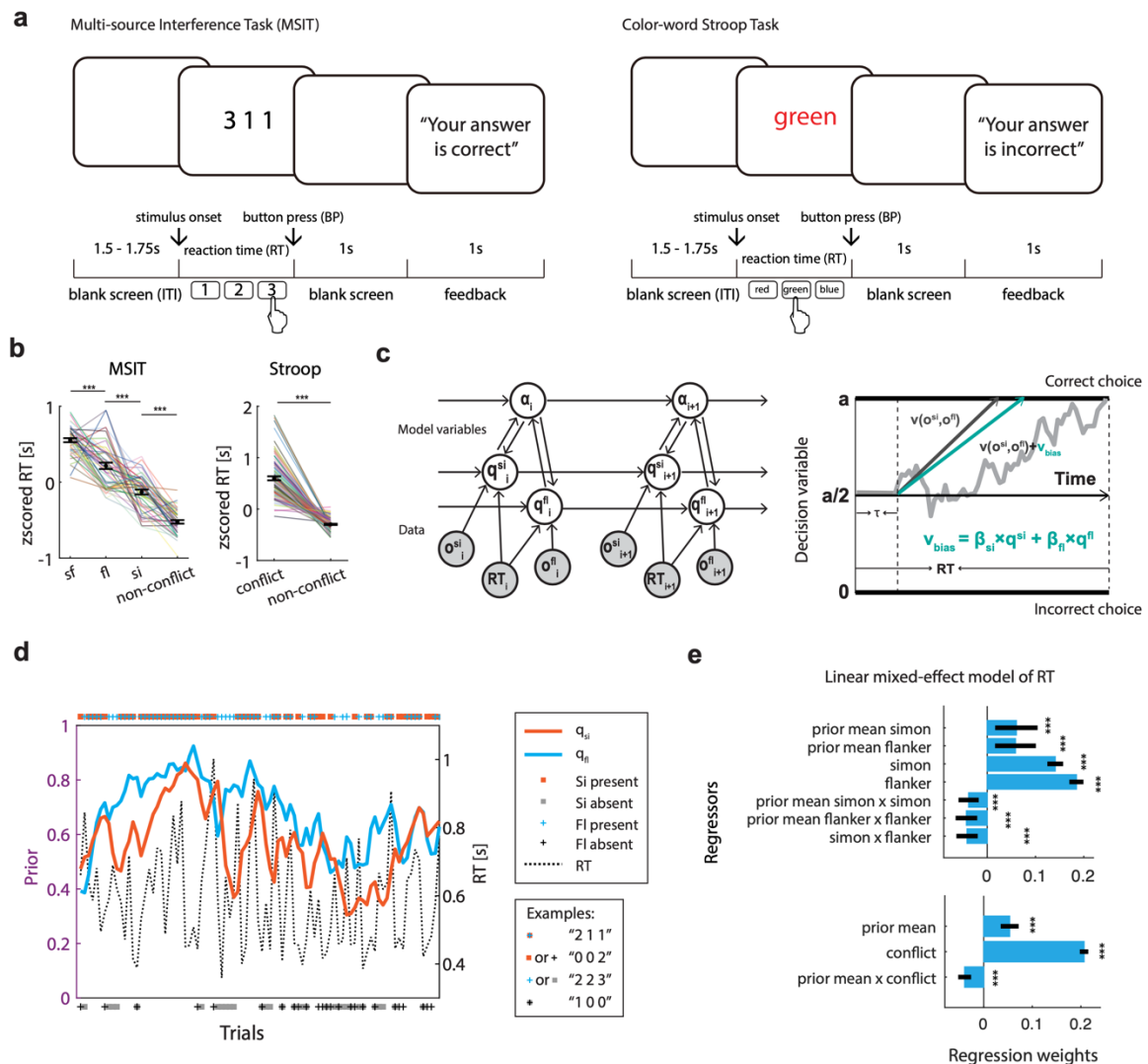
815

816

826 **Figures**
827

828 **Figure 1. Tasks, model, and behavioral results.**
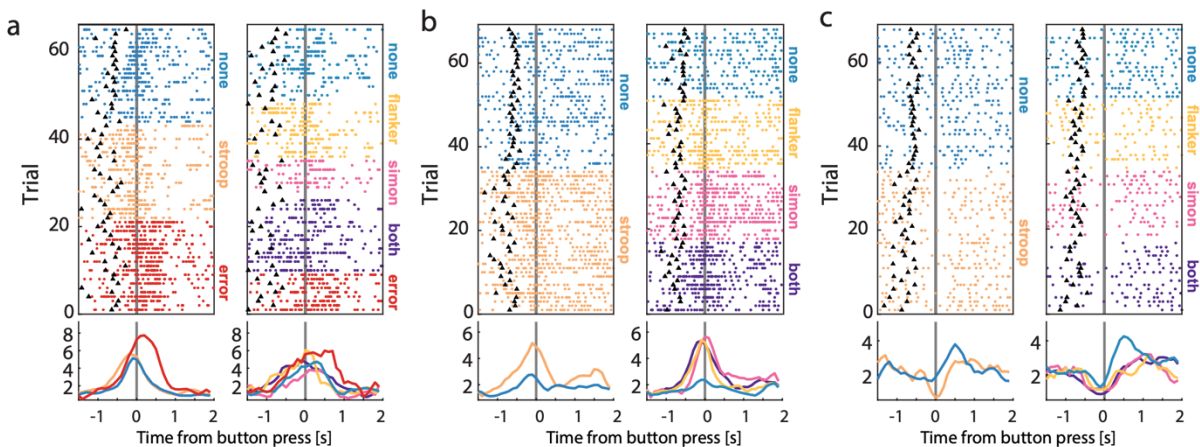829



Figure 1

830

831    (a) Tasks. MSIT (left) and Stroop (right). Participants indicated the identity (1,2, or 3) of the unique
832    number in MSIT, and the ink color of a colored word in Stroop. Feedback followed 1s after
833    responses and indicated trial outcome.
834    (b) RTs were significantly prolonged by conflict in MSIT (left) and Stroop (right), showing the
835    Simon, Flanker and Stroop effects. Each line is a session (N = 41 and 82, respectively).
836    (c) Graphical representation of the updating process (left) and the decision process modelled as
837    a drift diffusion process (right). Incorporating RT likelihood function (DDM) allows the tuning of
838    model estimate for conflict probability. Shown is the MSIT model, which has the five variables
839    volatility ($\alpha$), predicted Simon conflict ($q_{si}$), predicted Flanker conflict ($q_{fl}$), observed Simon
840    conflict ($o_{si}$), observed Flanker conflict ($o_{fl}$), and RT. Observables (conflict type, RT, and outcome)
841    are shown in gray, internal variables in white. Arrows indicate information flow. As the trial
842    started, the volatility variable is updated first, and then both predicted conflicts are updated by
843    the respective observed conflicts (Bernoulli likelihood) and RT (DDM likelihood). A linear
844    combination of the Simon conflict and Flanker conflict priors on each trial are entered as a drift
845    rate bias. The hyperparameters for RT tuning included the two linear coefficients before conflict
846    priors ($\beta_{si}$, $\beta_{fl}$), boundary separation (a) and the four separate drift rates for the four conflict
847    conditions (si only, fl only, si+fl, no-conflict). The conflict priors and posteriors were used as
848    regressors for subsequent behavioral and neural analyses.
849    (d) Estimated mean of the prior for Simon probability (orange) and Flanker probability (blue) from
850    an example MSIT session. Markers placed on the top indicated that either Simon conflict (orange
851    square) or Flanker conflict (blue cross), or both, was *present* on a trial. As is shown here, the
852    priors increase when there is a run of conflict (left part of the graph, both blue and orange traces
853    go up).
854    (e) Regression analyses of RT using linear mixed-effect models. Blue bars show regression
855    coefficients; black bars show confidence interval. All regressors explained significant variance as
856    determined by the likelihood ratio test (see **Methods**). Conflict prior positively predicted RT in
857    MSIT and Stroop.
858    *p < 0.05, ** p < 0.01, *** p < 0.001, n.s., not significant (p > 0.05).
859
860    **Figure 2. Example neurons in Stroop (left) and MSIT (right).**



861

862

863 Shown here are raster plots and peri-stimulus time histograms for three example neurons in both
864 tasks. Left panel shows data from Stroop, right show data from MSIT. Data are aligned to button
865 presses (t=0). The black triangles mark stimulus onset.
866 (a) Neuron signaling action errors.
867 (b) Neuron signaling conflict by firing rate increase around button presses.
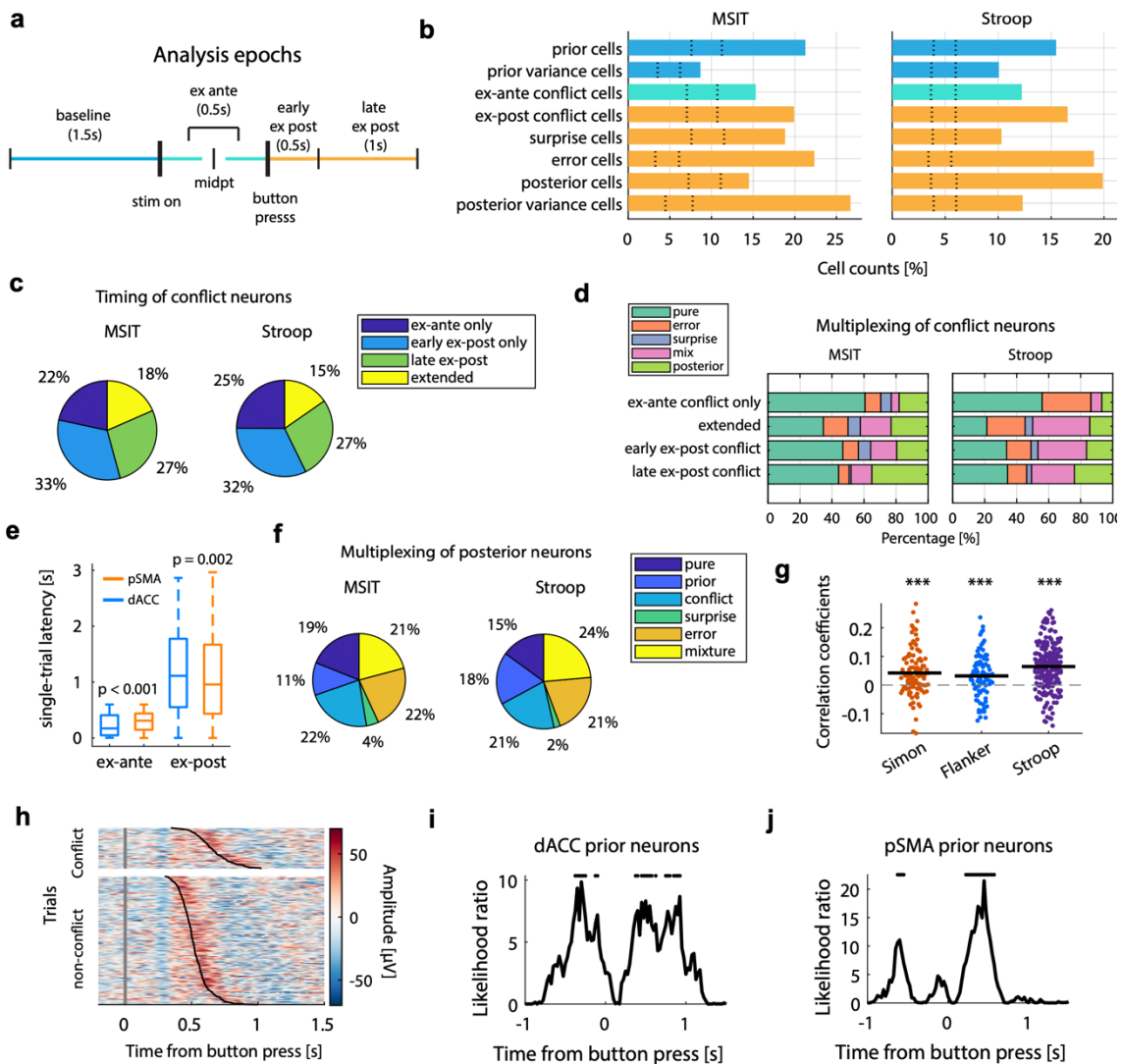868 (c) Neuron signaling conflict by a firing rate decrease around button presses.
869 Trial types are marked by colored words on the right side of the box. These example neurons
870 show similar responses dynamics in both tasks. Trials were re-sorted into groups for display
871 purposes only.
872

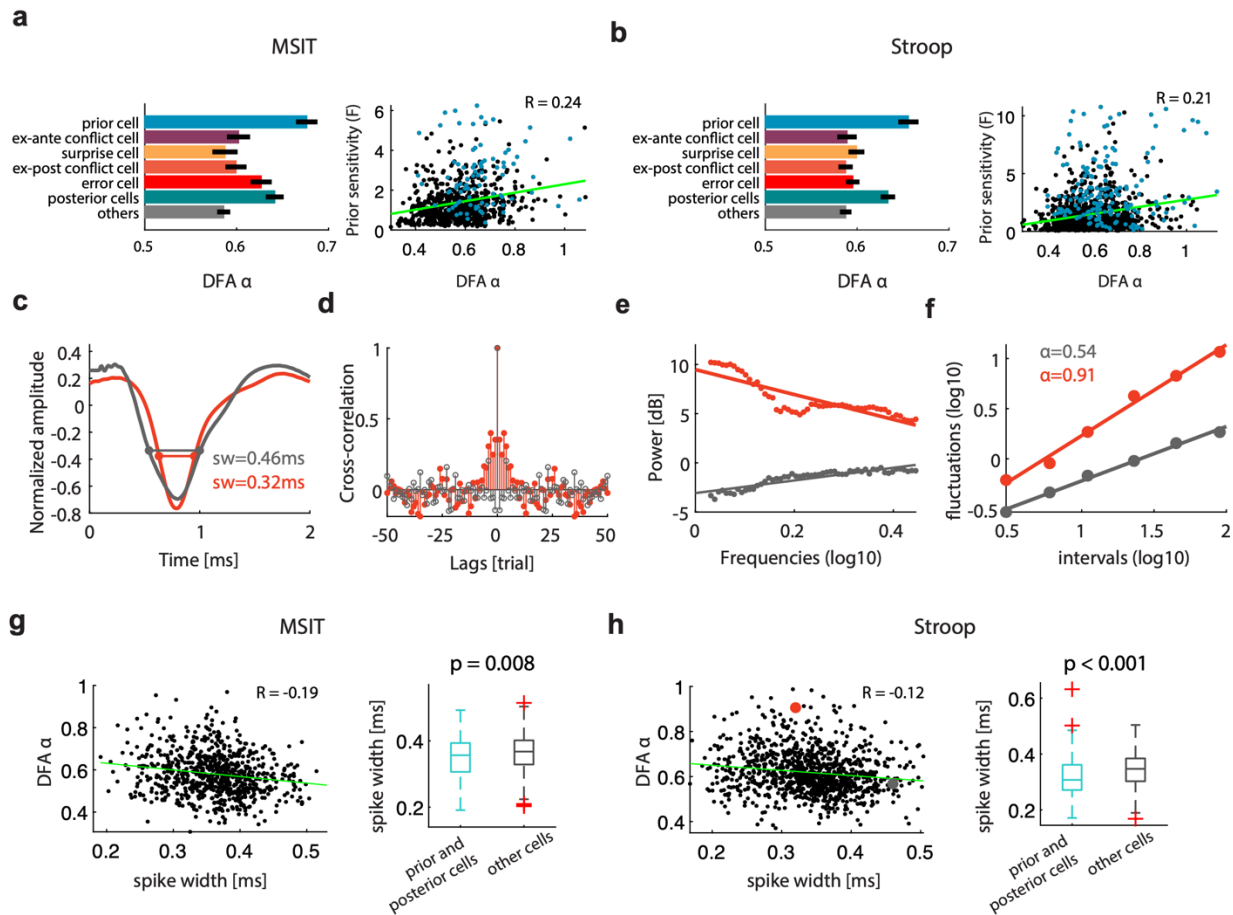873 **Figure 3. Neuronal selection and ERP analysis.**

874



875
876

877    (a) Analyses epochs used in neural analyses. The ex-ante epoch is defined as 0.5s epoch centered
878    at the midpoint between 100ms after stimulus onset and button presses. This was the moment
879    when conflict should achieve its maximum. Ex-post epochs are defined as epochs following
880    button presses. Epochs are colored by blue (baseline), cyan (ex-ante) and orange (ex-post). The
881    thickened vertical bars represent actual physical events, the slim vertical bars demarcate epochs.
882    Note that we also sometimes use 1s after button press as an analysis window.
883    (b) Percentage of neurons that encode task variables in MSIT (left) and Stroop (right). The color
884    code represents the epochs used to select these neurons (see (a)). Rows are arranged from top
885    to bottom in temporal order. Dotted lines represent $2.5^{th}$ and $97.5^{th}$ percentiles of the null
886    distribution obtained from permutation. For all groups shown, $p < 0.001$.
887    (c) Percentage of conflict neurons that are active exclusively in four groups: ex-ante, early ex-
888    post, late ex-post, and throughout ex-ante and ex-post epochs.
889    (d) Percentage of conflict neurons that were also selective for error, surprise, posterior, or any
890    combination of these factors ("mix"). Substantial proportions of conflict neurons multiplex error
891    and posterior information. The intersection between conflict and posterior increases towards the
892    later part of the trial.  Rows are arranged from top to bottom in temporal order.
893    (e) Comparison of single-trial neuronal response latency of conflict neurons in dACC and pre-SMA.
894    Ex-ante conflict neurons become active earlier in dACC than in pre-SMA, whereas ex-post conflict
895    neurons become active earlier in pre-SMA than in dACC. Only correct conflict trials are used in
896    this analysis.
897    (f) Percentage of posterior neurons that intersect with prior, conflict, surprise and error signaling.
898    We hypothesize that the extensive overlap between these groups reflects posterior computation.
899    (g) Neuronal signature of updating conflict prior based on the posterior. Correlation is computed
900    between the difference between prior and posterior (behavioral update) and the difference
901    between demeaned $FR_{ex-post}$ and $FR_{baseline}$ (neural update) for all prior neurons. On average, the
902    correlation is significantly positive, suggesting that the change in firing rates is commensurate
903    with the extent of updating derived from the behavioral model.
904    (h) An example session of intracranial EEG in Stroop, aligned to stimulus onset (grey vertical bars)
905    and sorted by RT (black lines). Color code represents amplitude in micro volt. An event-related
906    potential, named correct-trial potential (CRP), occurs shortly after button presses that is present
907    on both conflict and non-conflict trials.
908    (i) Relation between CRP amplitude and spiking activity of prior neurons. Both data were
909    simultaneously recorded in dACC. Likelihood ratio computed by comparing the full Poisson
910    regression model with CRP as a fixed effect with a reduced model without the CRP term, and is
911    plotted as a function of time. Black dots on top mark significant time bins, corrected for multiple
912    comparisons using the false discovery rate (FDR) method.
913    (j) Same as in (i), but for pre-SMA data.
914
915    **Figure 4. Long-range temporal correlation autocorrelation of spiking of prior neurons.**
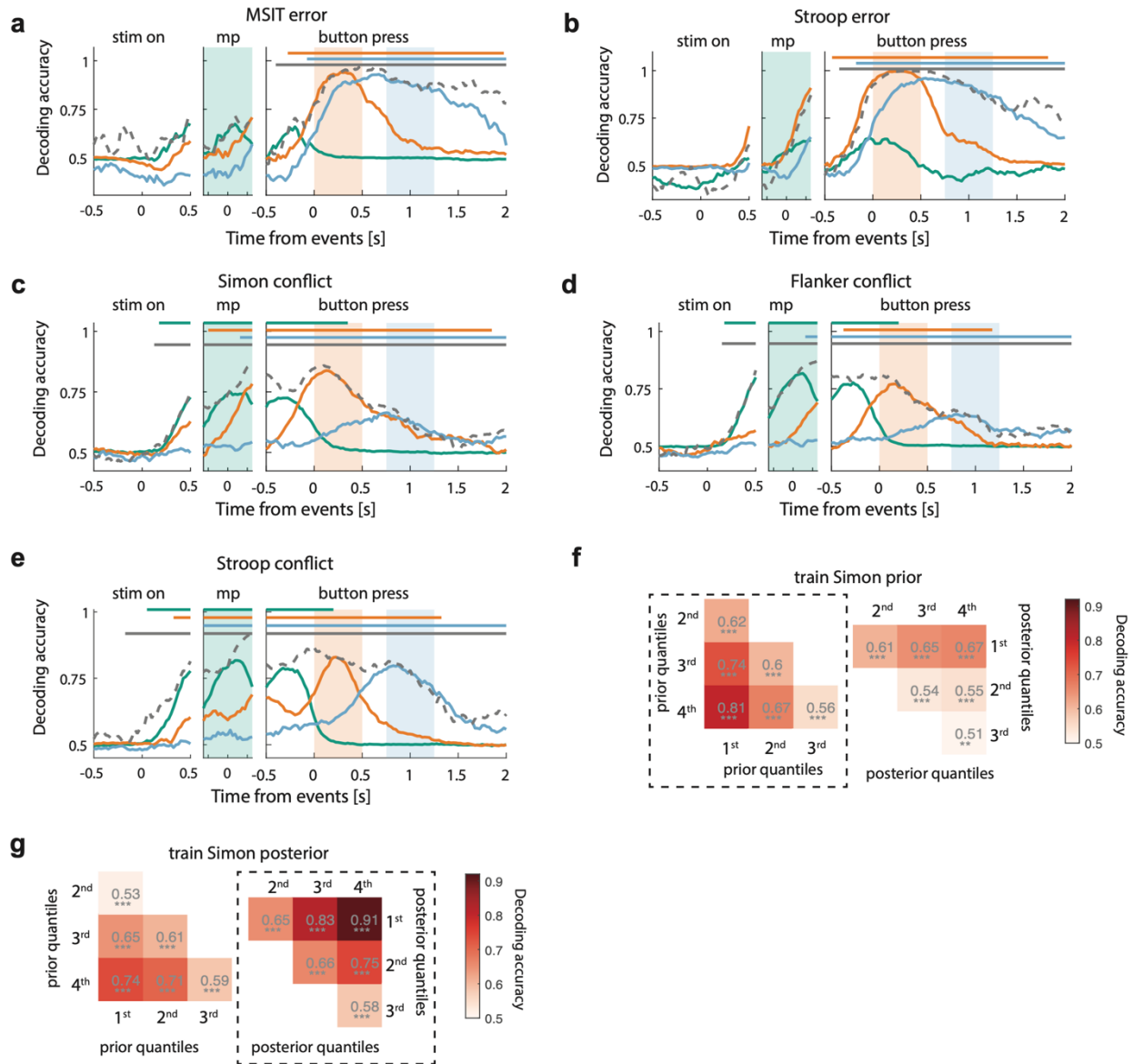916

Figure 4



917
918　(a) Mean baseline DFA $\alpha$ value for different groups of neurons (left), and correlation between
919　baseline DFA $\alpha$ value and the coding strength of prior, which is assessed by the F statistic
920　computed from regressing prior against baseline spike counts (right) for MSIT. Separate data
921　were used to compute $\alpha$ value and prior coding strength to avoid selection bias. Prior neurons
922　have significant larger baseline DFA $\alpha$ value than any other groups (p < 0.001, ANOVA). The
923　coding strength of prior was correlated strongly with the baseline $\alpha$ value (r = 0.24, p < 0.001).
924　(b) same as in (a) but computed for the Stroop task.
925　(c-f) Two example neurons, showing (c) waveforms, (d) autocorrelation, (e), power spectrum,
926　and (f) fluctuations as a function of time intervals used to compute DFA $\alpha$ value (slope). The
927　neuron with narrower spike width has higher DFA $\alpha$ value (r = -0.19, p < 0.001).
928　(g) DFA $\alpha$ value is negatively correlated with spike width for MFC neurons (left) in MSIT. Prior and
929　posterior neurons as a group have significantly narrower spikes than other neurons (right).
930　(h) same as in (g) but for the Stroop task.
931
932　**Figure 5. Temporal dynamics and cross-time generalization.**

Figure 5



933
934    (a-e) Decoding accuracy as a function of time for Stroop error, MSIT error, Simon conflict, Flanker
935    conflict and Stroop conflict. The three panels show data aligned to stimulus onset, midpoint
936    between 100ms after stimulus onset and button presses, button press onset. Dotted gray trace
937    represents within-time decoding accuracy, i.e., the data from the same epoch were used to train
938    and test a decoder. Green, blue and orange traces represent decoding accuracy of decoders
939    trained with data from the epochs demarcated by shading with the same colors, which is a test
940    of the temporal generalization of these decoders. Horizontal bars demarcate the extent of
941    significant clusters in time as determined by the cluster-based permutation test (p < 0.5).
942    (f-g) Decoding accuracy for classifying between pairs of Simon prior or posterior levels (binned
943    by quartiles). Color bars show decoding accuracy. Dotted frames mark the within-time decoding
944    results. Decoders that classify prior quantiles are trained using baseline spike counts (1.5s before
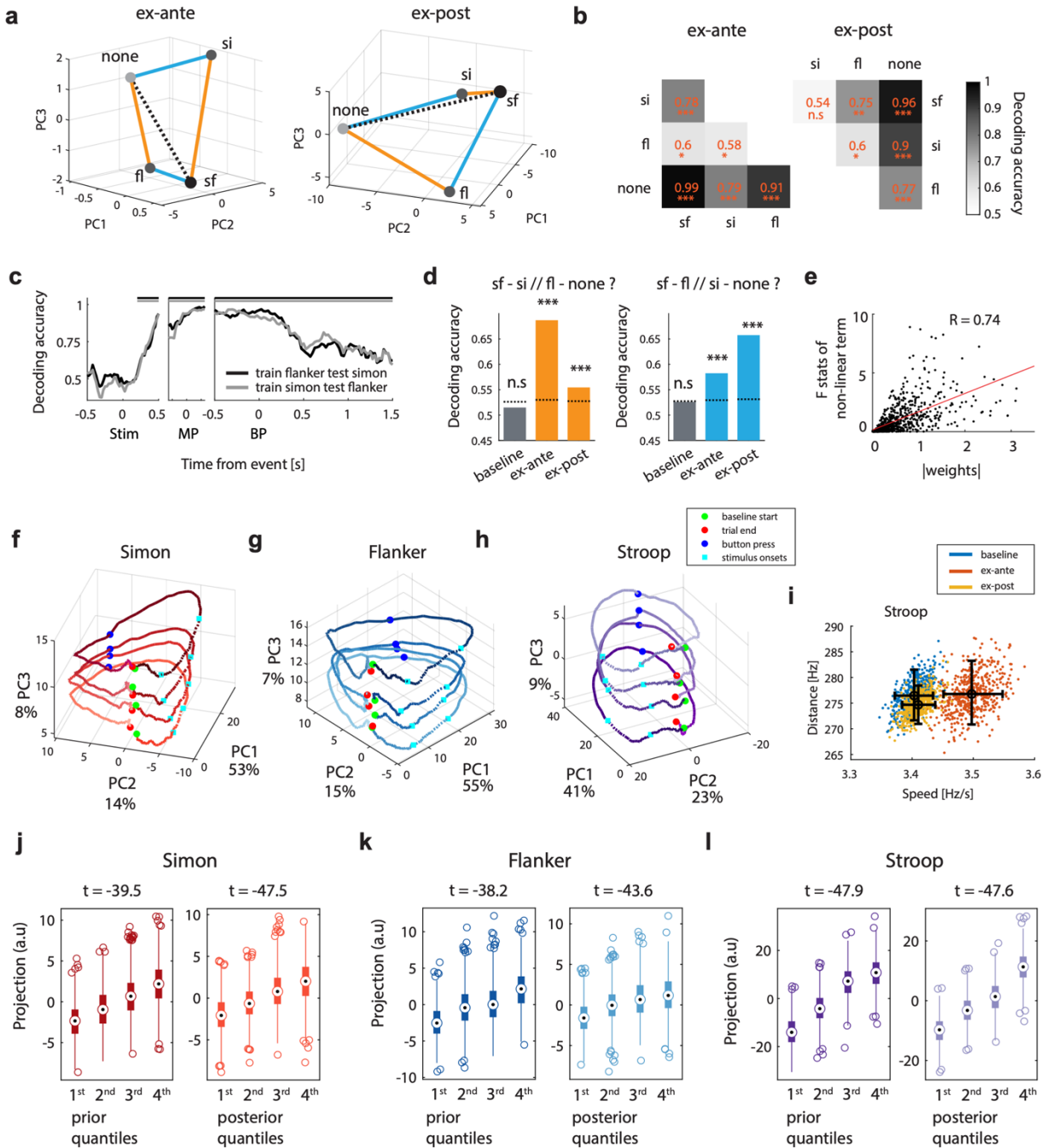
945 stimulus onset), whereas decoders that classify posterior quartiles are trained using ex-post spike
946 counts (2s after button press). To test temporal generalization of these decoders, prior-trained
947 decoders are tested with posterior data and labels, and vice versa. Dashed boxes represent
948 within-time decoding.
949
950 **Figure 6. State-space representation of conflicts, prior, and posterior.**
951



Figure 6

952

953 (a) Visualization of the conflict population representation in MSIT. Trial mean of four MSIT
954 conflict conditions, Simon only ("si"), Flanker only ("fl"), Simon and Flanker both present ("sf"),
955 and non-conflict ("none"), plotted in space spanned by three principal components. Left panel
956 uses ex-ante data. Right panel uses ex-post data. The extent of compositionality of conflict
957 representation is tested by condition generalization of decoding in (d). Dotted line is the vector
958 used to classify pairs of conflict conditions in (b).
959 (b) Decoding accuracy of pairwise classification of conflict conditions. Training data and left-out
960 testing data from the four conflict conditions are projected to the population vector flanked by
961 averages of non-conflict trials and sf trials, shown as the dotted line in (a). Color code represents
962 decoding accuracy. This coding dimension separates the four conflict conditions well.
963 (c) Abstract conflict signal. The three panels show decoding accuracy using data aligned to
964 stimulus onset, midpoint between 100ms after stimulus onset and button press onset, and
965 button press onset. For each time point, a decoder is trained with Simon (union of si and sf) vs.
966 non-Simon (union of fl and none) trials or Flanker (union of fl and sf) vs non-Flanker (union of si
967 and none) trials and used to classify left-out data of Flanker vs non-Flanker trials (grey) or Simon
968 vs. non-Simon trials (black).
969 (d) Testing compositionality of conflict representation with condition generalization of decoding.
970 If compositional, the sum of the representations of fl and si (vectors none -> fl, none -> si) should
971 be equal to the representation of sf (none -> sf), and the four condition means should form a
972 parallelogram. We tested the condition generalization using raw spike count data from the ex-
973 ante and ex-post (1s after button presses) epochs and from the baseline as a control without
974 dimensionality reduction. Data from the means connected by one of the blue edges in (a) were
975 used to train a decoder, which was then tested with left-out data from the means connected by
976 the opposite blue edge, and vice versa. Blue bars show decoding accuracy (baseline 0.51, ex-ante
977 0.69, ex-post 0.55). Same was also tested with data connected by the orange edges. Orange bars
978 show decoding accuracy (baseline 0.53, ex-ante 0.58, ex-post 0.66). Both the blue and orange
979 pairs of opposing edges supported such generalization simultaneously as indicated by the above-
980 chance decoding accuracy ($p < 0.001$), demonstrating parallelism and thus the compositionality
981 of conflict representation. Dotted lines show 97.5th percentile of the null distribution from
982 permutation.
983 (e) Single neuron with nonlinear coding of Simon and Flanker conflict contribute to deviation of
984 conflict representation from perfect linearity. Data used here are from the ex-post epoch.
985 Nonlinear coding of conflict by a single neuron is measured by the F statistic of the interaction
986 term between Simon and Flanker conflict in an ANOVA model with spike counts as the dependent
987 variable. Each neuron's contribution to the deviation from linear additivity in the high
988 dimensional neural space is quantified by the weight of the difference vector between "sf" and
989 "si + fl". Scatter plot shows the relation between these two measures. Red line shows the linear
990 fit.
991 (f-h) Visualization of prior/posterior population representation in MSIT and Stroop. Green dots
992 mark the onset of trial baseline, cyan squares mark the range of possible stimulus onsets, blue
993 dots mark button press and red dots mark end of trial. The range of stimulus onsets (a range
994 because trials are aligned to button press onsets) is shown as broken lines for each prior level.
995 For the portion before button presses (blue dots), the four trajectories correspond to the mean
996 of trials grouped by quartiles of prior. For the portion after button presses, the four trajectories

997 correspond to the mean of trials grouped by quartiles of posterior for after button presses. Trials
998 are aligned to button press onset. The color of trajectories fades as the trial progresses towards
999 the end (dark and light colors correspond to start and end of the trial). Most of the variance
1000 related to prior/posterior is captured by PC3, which is orthogonal to most of the time-dependent
1001 dynamics.
1002 (i) Distance between trajectories and average speed computed from trials grouped by quartiles
1003 of Simon conflict prior in the baseline (blue) and the ex-ante (orange) epoch, and trials grouped
1004 by Simon conflict posterior in the ex-post epoch (yellow). Trajectories are visualized in Figure 6f.
1005 The state space speed stays low during baseline, increases significantly during the ex-ante epoch
1006 and decreases back to a value similar to that during the baseline. Distance between trajectories
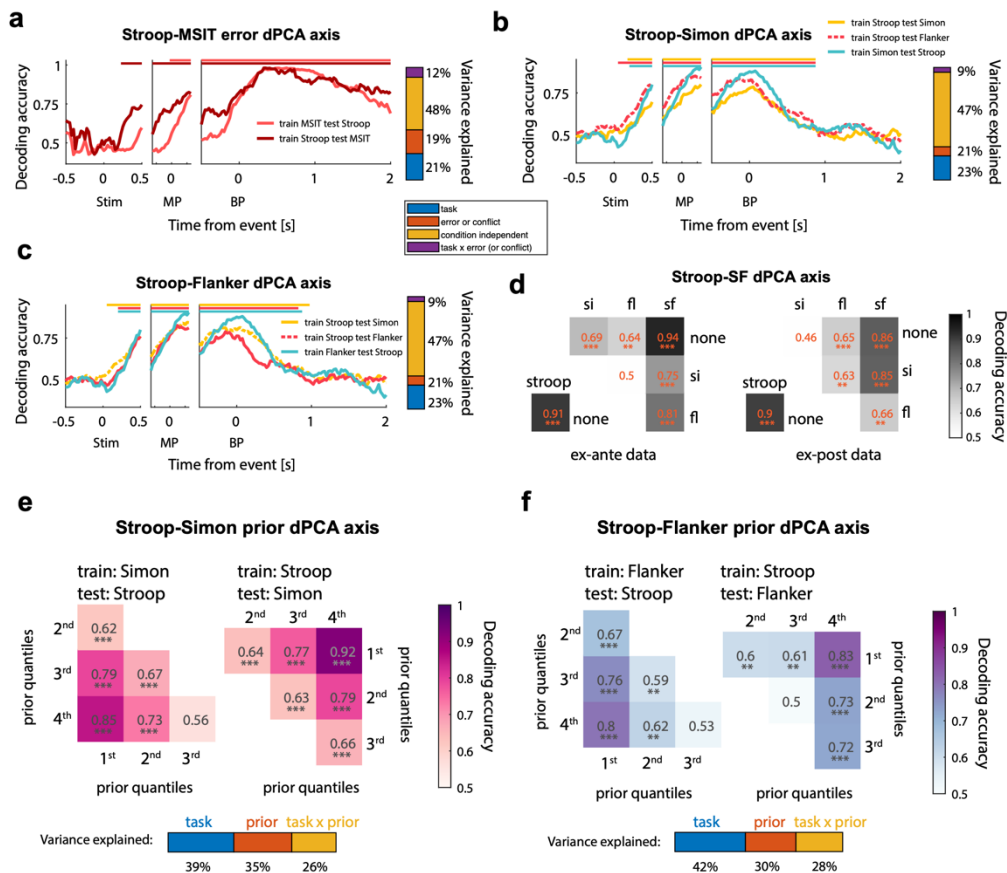1007 is stable across time.
1008 (j-l) Projection values on the coding dimension of Simon (J, p < 0.001 for both prior/posterior,
1009 prior t(11996) = -39.5, posterior t(11996)=-47.5), Flanker (K, p < 0.001 for both prior/posterior,
1010 prior t(11996) = -38.2, posterior t(11996)=-43.6) or Stroop (L, p < 0.001 for both prior/posterior,
1011 prior t(11996) = -47.9, posterior t(11996)=-47.6) prior or posterior (PC3 in f-h). The order of
1012 projection values is on average consistent with the order of the prior/posterior quartiles, even
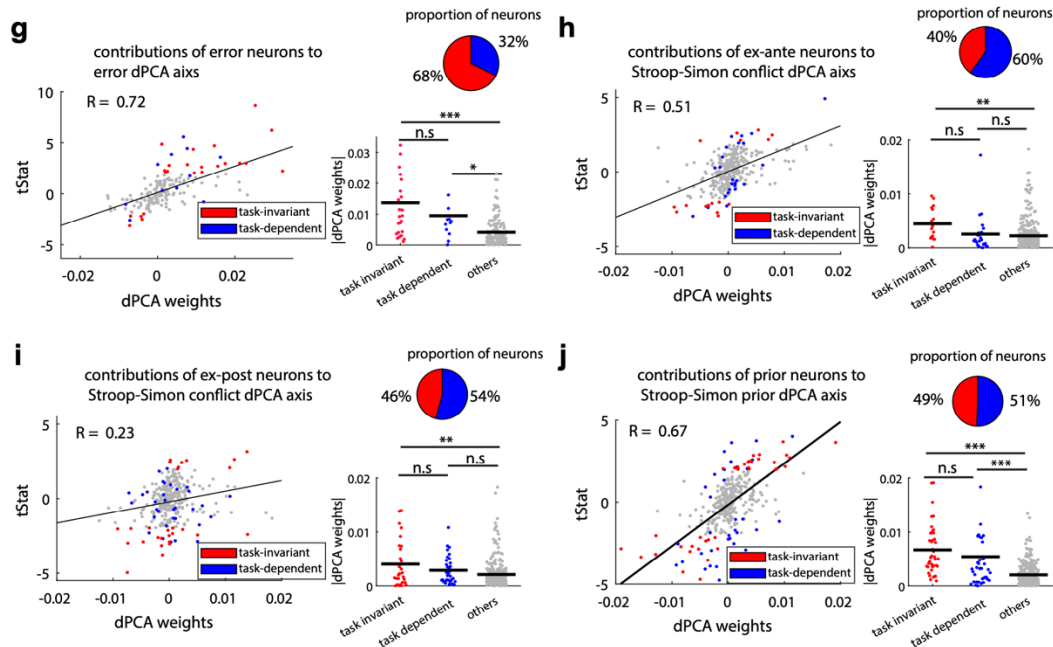1013 though PCA does not have access to the order information.
1014
1015 **Figure 7. Domain-general representation of performance monitoring signals.**



1016

1017
1018     (a) Task-invariant decoding of errors in both MSIT (salmon) and Stroop (crimson). The task-
1019     invariant coding dimension is extracted using dPCA that marginalizes out task information and
1020     time. This dPCA coding dimension is extracted from the error contrast in Stroop (error conflict vs.
1021     correct conflict trials) and the error contrast in MSIT (error "sf" trials and correct "sf" trials). This
1022     controls for trial conflict and isolates effects related only to error. Left, Accuracy for decoding
1023     errors as a function of time. Bar on the right shows the variance explained by the different dPCA
1024     components (color code see figure legend). The angle between the error dPC and the task dPC
1025     drived from dPCA is 94.47° and did not differ significantly from orthogonality (p = 0.53, tau = -
1026     0.032, Kendall rank correlation).
1027     (b-c) Task-invariant decoding of conflict in both MSIT (Simon, yellow; Flanker red) and Stroop
1028     (green). Because MSIT has two conflict conditions, Simon and Flanker, task invariance was
1029     investigated between Stroop/Simon and between Stroop/Flanker conflicts separately. This dPCA
1030     coding dimension is extracted from conflict and non-conflict trials in Stroop and either from
1031     Simon and non-Simon trials (b) or from Flanker and non-Flanker trials (c), by marginalizing out
1032     task and time. The angle between the Stroop-Simon conflict dPC and the task dPC is 81.13° and
1033     did not differ significantly from orthogonality (p = 0.19, tau = 0.048, Kendall rank correlation).
1034     The angle between the Stroop-Flanker conflict dPC and the task dPC is 81.13° and is significantly
1035     but weakly non-orthogonal (p = 0.018, tau = 0.087, Kendall rank correlation). Left-out conflict
1036     trials and non-conflict trials in Stroop, and left-out Simon, non-Simon, Flanker, non-Flanker trials
1037     in MSIT are projected and classified by this coding dimension. Left, decoding accuracy of conflict
1038     as a function of time. The bar on the right represents variance explained by the different dPCA
1039     components (color code see figure legend).
1040     (d) Testing separability of conflict conditions in Stroop and MSIT using data from the ex-ante (left)
1041     and ex-post epochs (right). The dPCA coding dimension used in this analysis is extracted by using
1042     conflict and non-conflict trials in Stroop and sf and non-conflict trials in MSIT by marginalizing out
1043     task information. Because data from ROIs are used, the time dimension is already marginalized

1044     out before entering dPCA algorithm. This coding dimension support classification of 83% of pairs
1045     of MSIT conflict conditions (upper triangle matrices) as well as Stroop conflict (left corner). Color
1046     coding represents decoding accuracy, orange numbers indicate the numerical values of decoding
1047     accuracy of that pair of conflict conditions (e.g., the accuracy is 0.69 for decoding si vs. none).
1048     Conflict monitoring is thus task-invariant but still preserves maximal separability of task-specific
1049     conflict conditions (MSIT).
1050     (e) Task-invariant decoding of all pairs of conflict prior levels in Stroop (lower triangle matrix) and
1051     Simon (upper triangle matrix). The dPCA coding dimension here is extracted by using the Stroop
1052     conflict prior contrast (the $1^{st}$ vs. $4^{th}$ quartiles of Stroop conflict prior) and the Simon conflict prior
1053     contrast (the $1^{st}$ vs. $4^{th}$ quartiles of Simon conflict prior), marginalizing out task information. Color
1054     code represents decoding accuracy. Bar at the bottom shows variance explained of dPCA
1055     components (for decoding, the component labelled as "prior" is used). The angle between
1056     Stroop-Simon conflict prior dPC and the task dPC is 106.42° and did not differ significantly from
1057     orthogonality (p = 0.78, tau = -0.01, Kendall rank correlation). The angle between Stroop-Flanker
1058     conflict prior dPC and the task dPC is 74.42° and did not differ significantly from orthogonality (p
1059     = 0.30, tau = -0.038, Kendall rank correlation).
1060     (f) Same as in (e) but for Flanker prior.
1061     (g-j) Contribution of single neuron coding of error (g), Stroop-Simon ex-ante conflict (h), Stroop-
1062     Simon ex-post conflict (i) and Stroop-Simon prior (j) to the task-invariant population coding of
1063     these variables. Because MSIT has two conflict conditions, Simon and Flanker, task invariance
1064     was tested between Stroop and Simon or between Stroop and Flanker separately. We modelled
1065     each neuron's baseline (j), ex-ante (h) or ex-post (g, i) response using ANOVA. The main effects
1066     are a dummy variable indicating error (g) or Stroop-Simon conflict (h, i) or Stroop-Simon prior (j)
1067     and task ID (Stroop or MSIT), and the interaction term between these two. A significant
1068     interaction suggests that the coding is more prominent in one task than the other. Task-invariant
1069     neurons is defined as having a significant main effect of the variable of interest but an
1070     insignificant interaction with the task ID. Task-dependent neurons is defined as having a
1071     significant interaction term.
1072     (g) Contribution of single neuron ex-post coding of error to task-invariant population coding of
1073     error. Of the 37 of neurons that were selected as signaling error in the ex-post epoch in either
1074     task, 68% did so in a task-invariant way (red) and 32% in a task-dependent (blue) way (pie chart).
1075     There is a strong correlation between the error t-statistic and the dPCA weight of a particular
1076     neuron (scatter plot on the left; R = 0.72). Comparing the mean absolute value of dPCA weights
1077     between task-invariant, task-dependent and uncategorized neurons, both the task-invariant (p <
1078     0.01) and task-dependent (p < 0.05) neurons had significantly larger absolute weights than the
1079     uncategorized neurons.
1080     (h) Contribution of single neuron ex-ante coding of Stroop-Simon conflict to task-invariant
1081     population coding of Stroop-Simon conflict. Of the 40 of neurons that were selected as signaling
1082     conflict in the ex-ante epoch in either task, 40% did so in a task-invariant way (red) and 60% in a
1083     task-dependent (blue) way (pie chart). There is a strong correlation between the error t-statistic
1084     and the dPCA weight of a particular neuron (scatter plot on the left; R = 0.51). Comparing the
1085     mean absolute value of dPCA weights between task-invariant, task-dependent and uncategorized
1086     neurons, only the task-invariant neurons had significantly larger absolute weights than the
1087     uncategorized neurons (p < 0.01).

1088 (i) Contribution of single neuron ex-post coding of Stroop-Simon conflict to task-invariant
1089 population coding of Stroop-Simon conflict. Of the 46 of neurons that were selected as signaling
1090 conflict in the ex-post epoch, 46% did so in a task-invariant way (red) and 54% in a task-
1091 dependent (blue) way (pie chart). There is a strong correlation between the error t-statistic and
1092 the dPCA weight of a particular neuron (scatter plot on the left; R = 0.51). Comparing the mean
1093 absolute value of dPCA weights between task-invariant, task-dependent and uncategorized
1094 neurons, only the task-invariant neurons had significantly larger absolute weights than the
1095 uncategorized neurons (p < 0.01).
1096 (j) Contribution of single neuron baseline coding of Stroop-Simon conflict prior to task-invariant
1097 population coding of Stroop-Simon conflict prior. Of the 75 of neurons that were selected as
1098 signaling prior in the baseline in either task, 49% did so in a task-invariant way (red) and 51% in
1099 a task-dependent (blue) way (pie chart). There is a strong correlation between the error t-statistic
1100 and the dPCA weight of a particular neuron (scatter plot on the left; R = 0.51). Comparing the
1101 mean absolute value of dPCA weights between task-invariant, task-dependent and uncategorized
1102 neurons, both the task-invariant (p < 0.001) and task-dependent (p < 0.001) neurons had
1103 significantly larger absolute weights than the uncategorized neurons.
1104
1105 *p < 0.05, ** p < 0.01, *** p <= 0.001, n.s., not significant (p > 0.05).
1106
1107 **Table S1. Number of sessions and neurons recorded**
1108 Summary of number of neurons recorded in each subject. For some subjects, both the Stroop
1109 task and MSIT were performed.
1110

| Patients ID | Sex | Age | Stroop | | | MSIT | | |
|---|---|---|---|---|---|---|---|---|
| | | | # sessions | dACC | pSMA | # sessions | dACC | pSMA |
| P9HMH | M | 55 | 1 | 9 | 0 | NA | NA | NA |
| P11HMH | M | 16 | 2 | 26 | 0 | NA | NA | NA |
| P14HMH | M | 31 | 2 | 5 | 0 | NA | NA | NA |
| P16HMH | F | 34 | 2 | 22 | 0 | NA | NA | NA |
| P19HMH | M | 34 | 1 | 13 | 0 | NA | NA | NA |
| P21HMH | M | 20 | 2 | 8 | 0 | NA | NA | NA |
| 'P31HMH' | M | 30 | 1 | 3 | 0 | NA | NA | NA |
| 'P41HMH' | M | 19 | 1 | 2 | 0 | NA | NA | NA |
| 'P42HMH' | M | 29 | 1 | 2 | 0 | NA | NA | NA |
| 'P24CS' | F | 47 | 2 | 17 | 46 | NA | NA | NA |
| 'P25CS' | F | 36 | 2 | 32 | 0 | NA | NA | NA |
| 'P26CS' | F | 36 | 1 | 16 | 0 | NA | NA | NA |
| 'P27CS' | M | 45 | 1 | 2 | 2 | NA | NA | NA |
| 'P29CS' | M | 19 | 1 | 9 | 9 | NA | NA | NA |
| 'P31CS' | M | 31 | 2 | 38 | 9 | NA | NA | NA |
| 'P32CS' | M | 19 | 1 | 0 | 5 | NA | NA | NA |
| 'P33CS' | F | 44 | 4 | 66 | 25 | NA | NA | NA |
| 'P34CS' | M | 70 | 5 | 26 | 0 | 1 | 4 | 0 |
| 'P35CS' | M | 63 | 6 | 1 | 47 | 1 | 0 | 9 |
| 'P36CS' | M | 45 | 6 | 8 | 64 | NA | NA | NA |
| 'P37CS' | F | 33 | 11 | 107 | 40 | 5 | 64 | 17 |
| 'P39CS' | M | 26 | 6 | 21 | 96 | NA | NA | NA |
| 'P40CS' | M | 25 | 3 | 7 | 25 | 3 | 20 | 34 |
| 'P42CS' | F | 25 | 5 | 83 | 60 | 8 | 45 | 97 |
| 'P47CS' | M | 33 | 2 | 0 | 18 | NA | NA | NA |
| 'P48CS' | F | 32 | 1 | 20 | 21 | NA | NA | NA |
| 'P44CS' | F | 53 | 2 | 20 | 42 | 2 | 17 | 36 |
| 'P49CS' | F | 24 | NA | NA | NA | 1 | 0 | 4 |
| 'P51CS' | M | 17 | NA | NA | NA | 9 | 107 | 17 |
| 'P55CS' | F | 43 | NA | NA | NA | 4 | 14 | 67 |
| 'P56CS' | M | 48 | 3 | 4 | 15 | NA | NA | NA |
| 'P60CS' | M | 67 | NA | NA | NA | 2 | 29 | 49 |
| 'P61CS' | F | 52 | 4 | 7 | 79 | 4 | 7 | 77 |
| 'P71CS' | M | 40 | 1 | 19 | 4 | 1 | 19 | 5 |

1111

1112
1113 **Table S2. Model comparisons for RT**
1114

Table S2

### BIC of RT

| | RT Tuned | no RT tuning | RL | | | RT Tuned | Prev conflict |
|---|---|---|---|---|---|---|---|
| MSIT | -686.3 | -337.8 | -450.9 | | MSIT | -805.7 | -440.1 |
| Stroop | -1412.1 | -904.1 | -972.5 | | Stroop | -1733.6 | -1435.1 |

1115
1116 Model comparison for RT using BIC. To test whether our RT-tuned Bayesian model explains
1117 variance in RT better than other models, we used linear mixed-effect models that takes into
1118 account subject variability (details of the model see **Methods**) and computed BIC for these
1119 models. The conflict prior is entered as a main fixed-effect and also as a by-session random effect.
1120 Here, conflict priors generated by four models are considered: "RT tuned", Bayesian conflict
1121 learning model with DDM hyperparameters and thus the conflict prior is tuned by RT. "No RT
1122 tuning", Bayesian conflict learning model without incorporating DDM likelihood for RT. "RL", a
1123 reinforcement learning model where the conflict probability is modelled as a "value" function
1124 and updated trial-by-trial by a simple update rule. "Prev conflict", a dummy variable indicating
1125 previous trial conflict. These linear mixed-effect models all have the same number of free
1126 parameters. A separate comparison was done between the RT tuned model with the model that
1127 uses the previous conflict (sub-table on the right) because the number of trials must be kept the
1128 same for the comparison and the "prev conflict" model did not consider the first trial for each
1129 session (there was no "prev conflict" in that case).
1130
1131 **Table S3 Model comparison for trial congruency**

Table S3

### BIC of conflict

| | RT Tuned | no RT tuning | RL | constant prior |
|---|---|---|---|---|
| MSIT | 26517 | 26674 | 26958 | 28271 |
| Stroop | 23807 | 24039 | 24779 | 26498 |

1132
1133 Model comparison for trial congruency using BIC. We used Bernoulli likelihood when computing
1134 BIC for the conflict sequence. Note that the number of fitted parameters for Bayesian models is
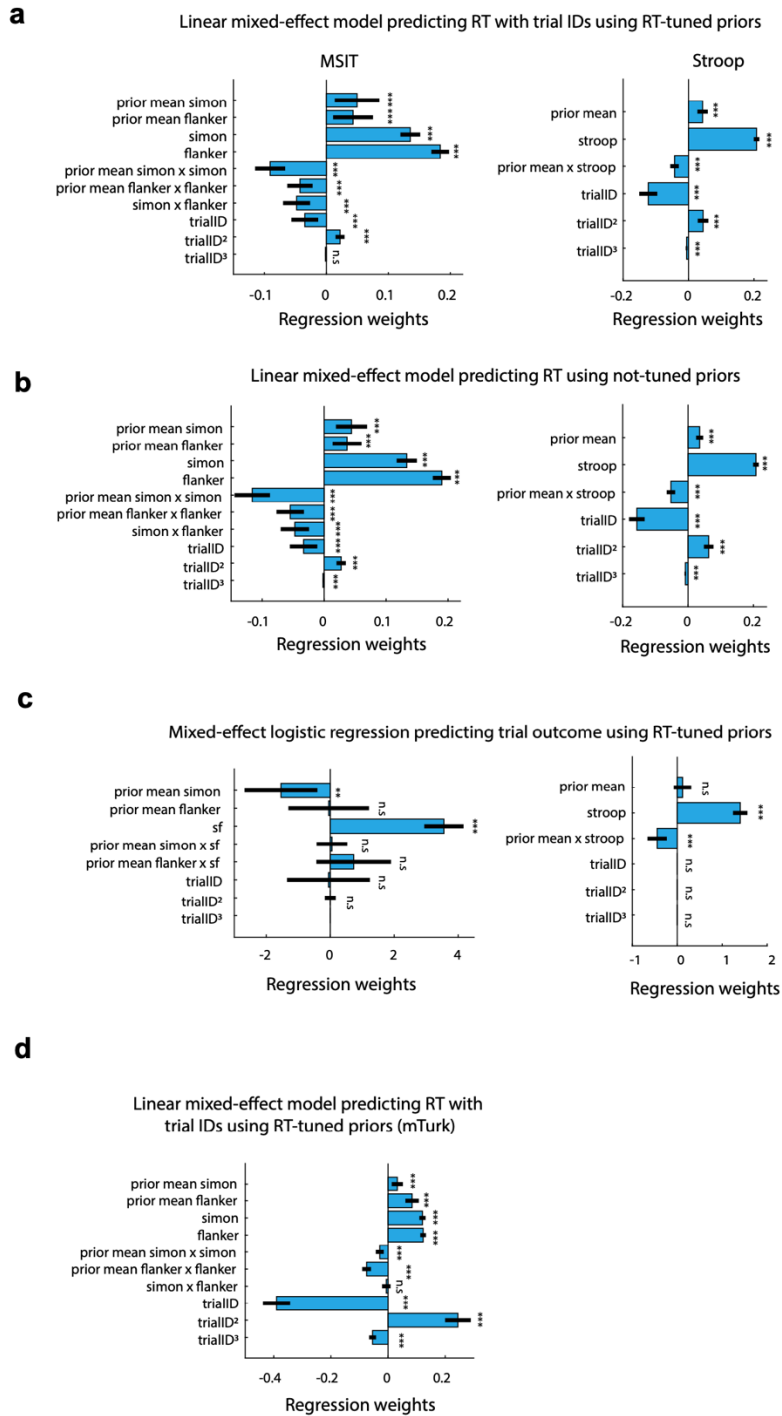
1135   zero, for the "RL" model is one (learning rate), and for "constant prior" model is one (the constant
1136   prior). BIC penalizes free parameters.
1137
1138
1139   **Figure S1 Behavioral models. Related to Figure 1.**

Figure S1



1140

1141 Statistical significance of regressors is determined by comparing the full model and a reduced
1142 model with a particular regressor removed, using a likelihood ratio test.
1143 (a) Linear mixed-effect model for RT that incorporates trial ID regressors for MSIT (left) and
1144 Stroop (right). Conflict priors used are from the Bayesian online learning models with RT tuning.
1145 We added the first, second and third -order trial ID regressors to model putative practice effects.
1146 The main effects of conflict priors, conflict, and their interaction are all significant even in the
1147 presence of trial ID regressors, suggesting these regressors capture behavioral effect that do not
1148 depend on trial ID.
1149 (b) Same as (a), but for the Bayesian online learning models without RT tuning. Thus, in this
1150 instance, conflict prior is estimated based on conflict sequence alone. The main effects of conflict
1151 priors (not tuned by RT), conflict, and their interaction are all significant in the presence of trial
1152 ID regressors. Therefore, RT tuning improves conflict prior (see (a) and Table S2, S3), but this is
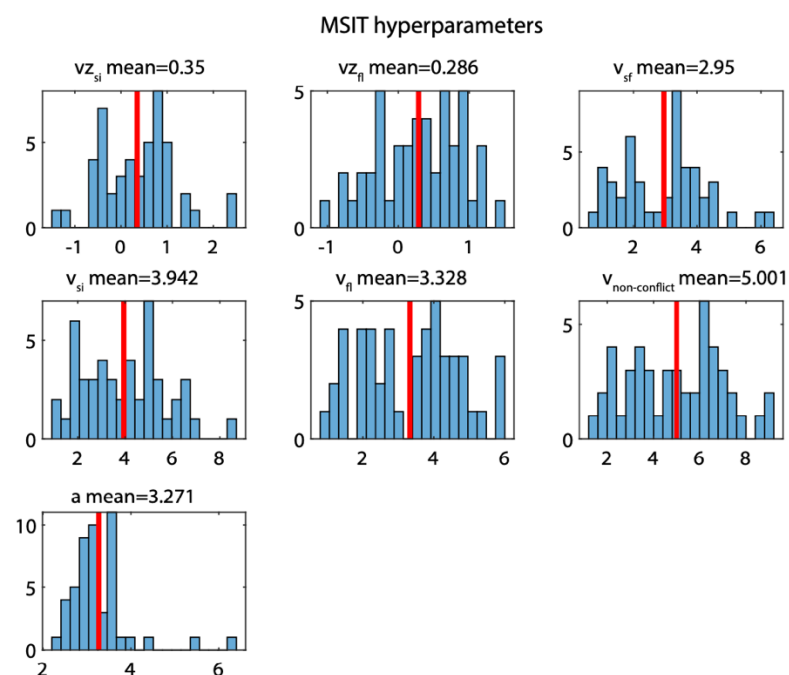1153 not required.
1154 (c) Mixed-effect logistic regression for predicting trial outcome (error or correct) for MSIT (left)
1155 and Stroop (right). Conflict priors used are from the Bayesian online learning models with RT
1156 tuning. For MSIT, we consider only "sf" trials for conflict trials, on which most of errors occur, and
1157 non-conflict trials. Conflict prior reduces error likelihood in both MSIT (significant main effect, p
1158 = 0.009) and Stroop (significant interaction term), p < 0.001).
1159 (d) Linear mixed-effect model for RT that incorporates trial ID regressors for MSIT. Data were
1160 collected from online participants using Amazon Mechanical Turk.
1161
1162 **Figure S2. DDM hyperparameters used in Bayesian conflict learning models. Related to Figure**
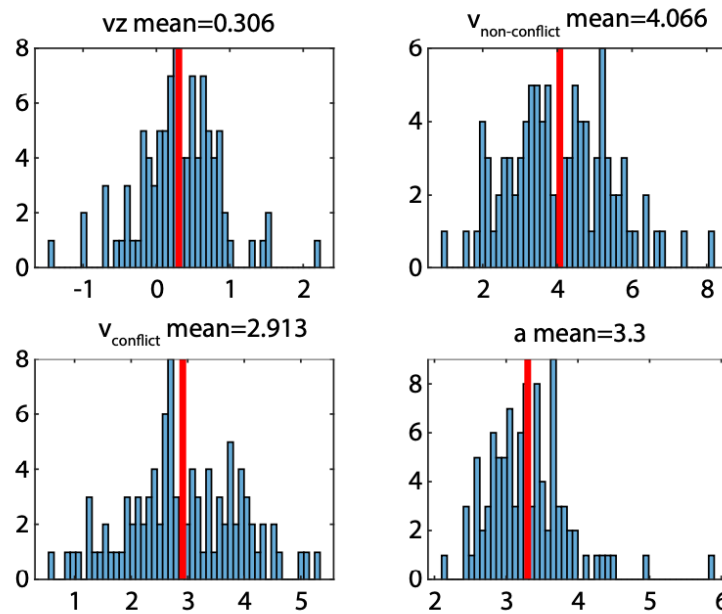1163 **1.**



Figure S2

1164

**b**

### Stroop hyperparameters



1165

1166   (a) Hyperparameters used in the Bayesian conflict learning model for MSIT. $vz_{si}$ and $vz_{fl}$ are
1167   coefficients scaling Simon and Flanker prior. $v_{sf}$, $v_{si}$, $v_{fl}$, $v_{non-conflict}$ are base drift rates in
1168   both Simon and Flanker present ("sf"), Simon-only ("si"), Flanker-only ("fl"), non-conflict trials. $a$
1169   is the boundary separation. The effective drift rate was the sum of the base drift rate and the
1170   scaled conflict prior. The base drift rates were significantly different from each other (p < 0.001,
1171   ANOVA). Post-hoc pairwise testing from a multiple comparison test determined that $v_{si}$ did not
1172   differ significantly from $v_{fl}$; $v_{sf}$ were significantly larger than either $v_{si}$ or $v_{fl}$; both $v_{si}$ and $v_{fl}$
1173   were significantly larger than $v_{non-conflict}$. These
1174   (b) Hyperparameters used in the Bayesian conflict learning model for Stroop. $vz$ is the coefficient
1175   scaling Stroop prior. $v_{conflict}$, $v_{non-conflict}$ are base drift rates in conflict and non-conflict trials.
1176   $a$ is the boundary separation. $v_{non-conflict}$ are significantly larger than $v_{conflict}$ across sessions
1177   (p < 0.001, t test).
1178   Hyperparameters are used in the DDM likelihood function for tuning the prior estimation process
1179   using an expectation-maximization algorithm.
1180
1181   *p < 0.05, ** p < 0.01, *** p < 0.001, n.s., not significant (p > 0.05 or not significant determined
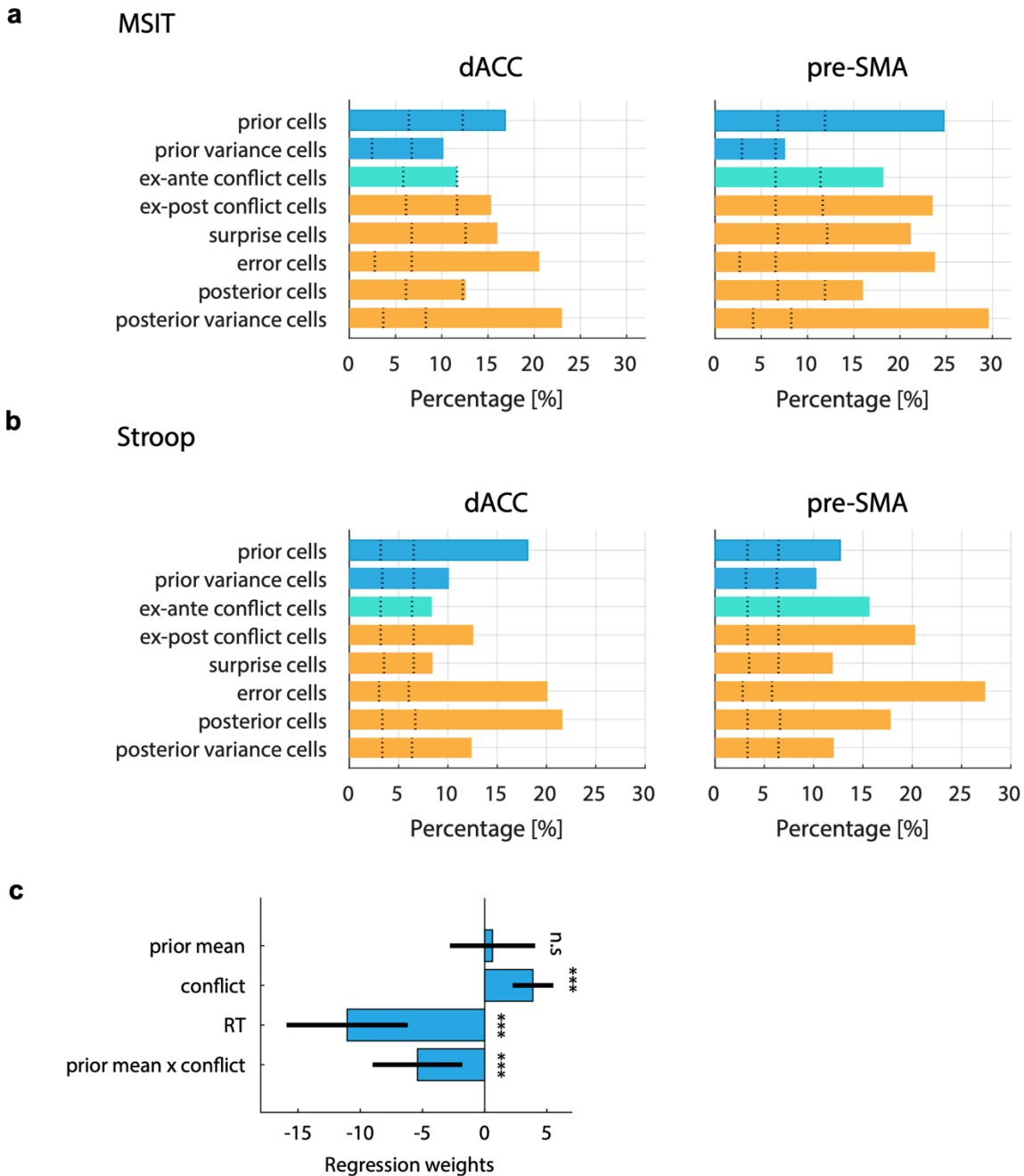1182   using FDR).
1183
1184
1185
1186   **Figure S3. Neuronal selection by areas and ERP analysis. Related to Figure 3.**

Figure S3



1187
1188    (a) Percentages of significant neurons in both dACC (left) and pre-SMA (right) in MSIT.
1189    (b) Percentages of significant neurons in both dACC (left) and pre-SMA (right) in Stroop.
1190    Dotted lines represent 2.5th and 97.5th percentiles of the null distribution obtained from
1191    permutation. For all groups shown, p < 0.001. Patterns of neuronal selection are similar between
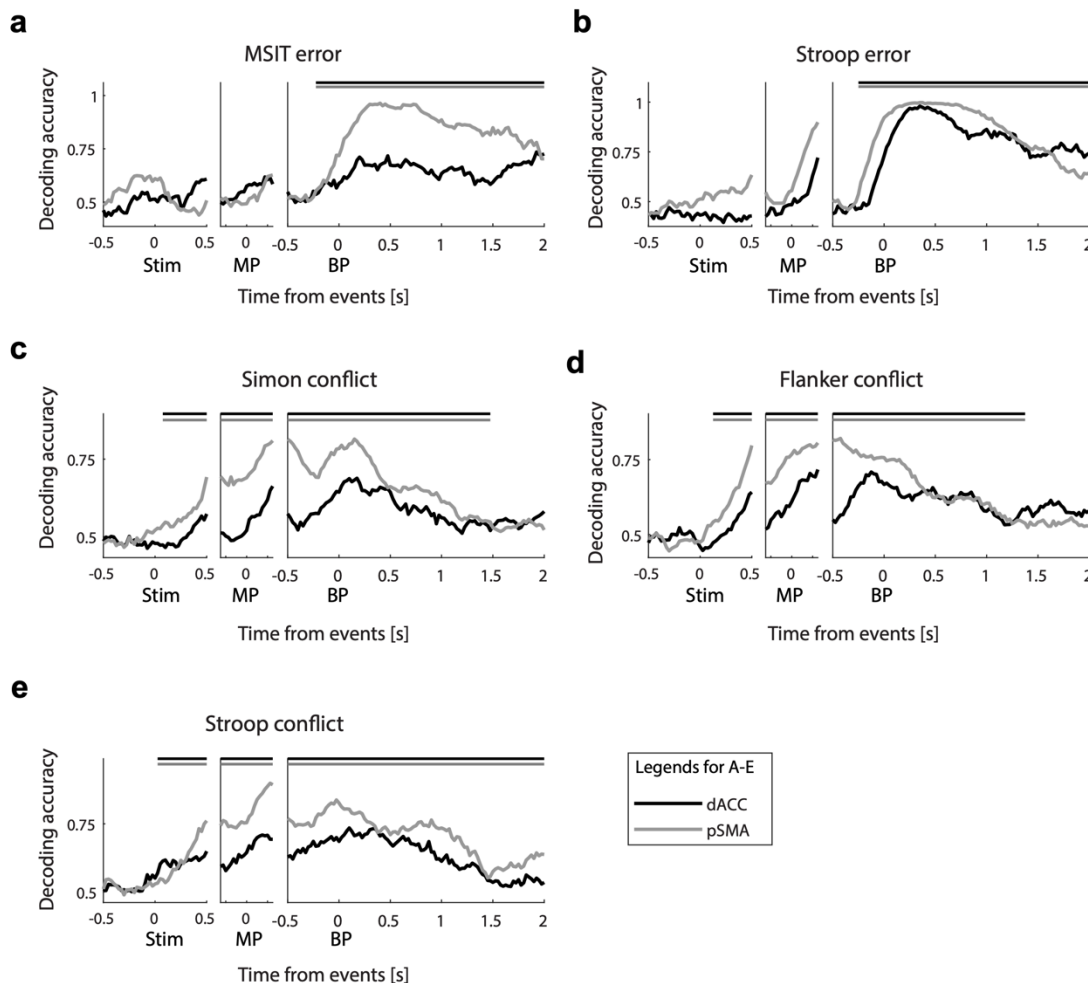1192    dACC and pre-SMA.
1193    (c) Linear mixed-effect model for CRP amplitude in the Stroop task. Conflict priors used are from
1194    the Bayesian online learning models with RT tuning. The main effects of conflict, RT and the

1195    interaction between prior and conflict were all significant. The main effect of conflict prior was
1196    not significant. Statistical significance was determined by a likelihood ratio test (comparing
1197    between the full model and the reduced models with regressors of interest removed).
1198
1199    **Figure S4. Population decoding of error and conflict by areas. Related to Figure 5.**
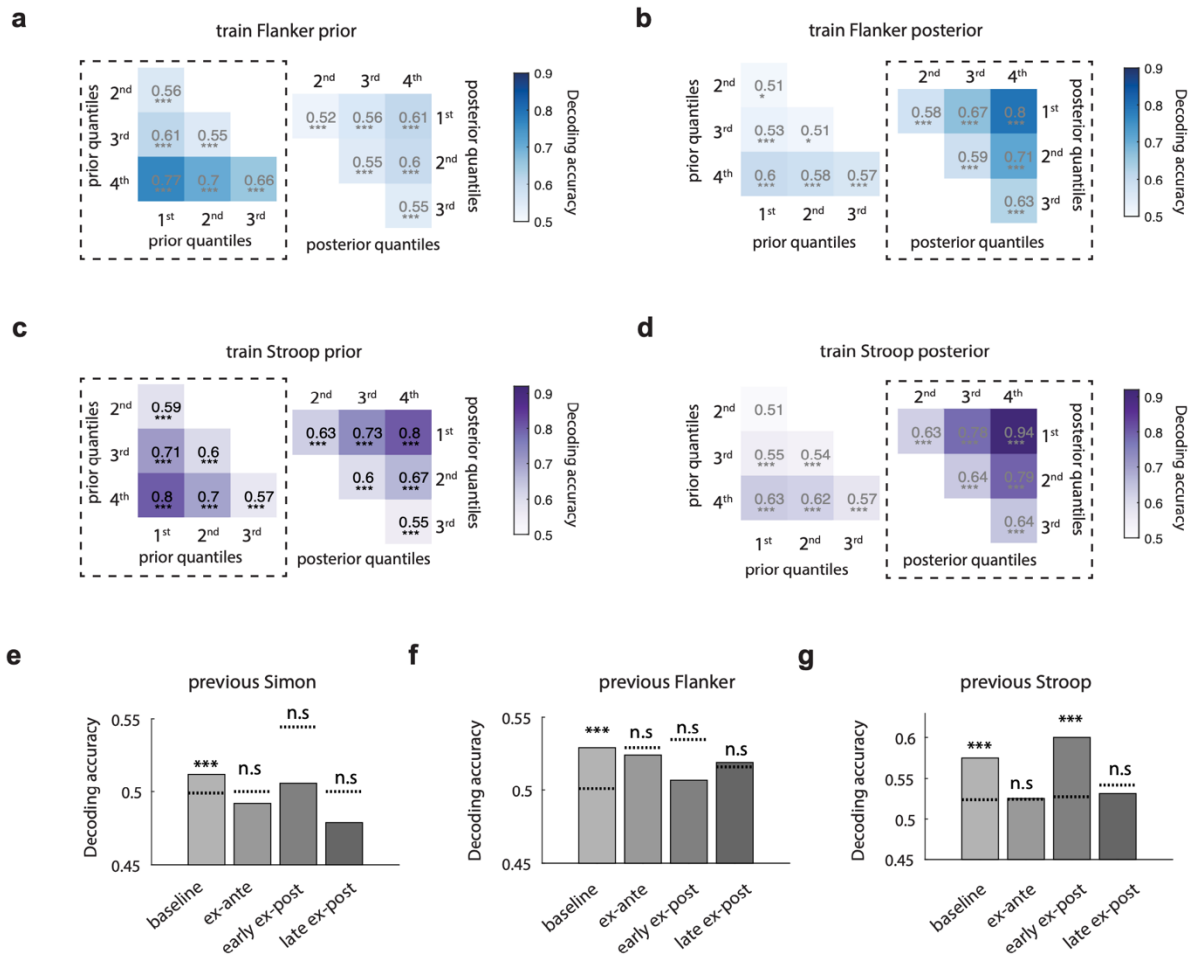1200

Figure S4



1201
1202    (a-e) Population decoding accuracy for MSIT error (a), Stroop error (b), Simon conflict (c), Flanker
1203    conflict (d), Stroop conflict (e).  For (a-e), black traces are from dACC data and grey traces are
1204    from pre-SMA data. Horizontal bars at the top demarcate significant cluster as determined by
1205    the cluster-based permutation test ($p < 0.05$). Overall dynamics are similar between dACC and
1206    pre-SMA, though the decoding accuracy on average is lower in the former.
1207
1208
1209    **Figure S5. Population decoding of prior/posterior (Flanker and Stroop) and past-trial conflict.**
1210    **Related to Figure 5.**

Figure S5



1211
1212 (a-d) Decoding accuracy for classifying between pairs of Flanker (a-b) and Stroop (c-d) prior or
1213 posterior quartiles. Color bars show decoding accuracy. Dotted frames mark the within-time
1214 decoding results. Decoders that classify prior quartiles are trained using baseline spike counts
1215 (1.5s before stimulus onset), whereas decoders that classify posterior quartiles are trained using
1216 ex-post spike counts (2s after button press). To test temporal generalization of these decoders,
1217 prior-trained decoders are tested with posterior data and labels, and vice versa. Dashed boxes
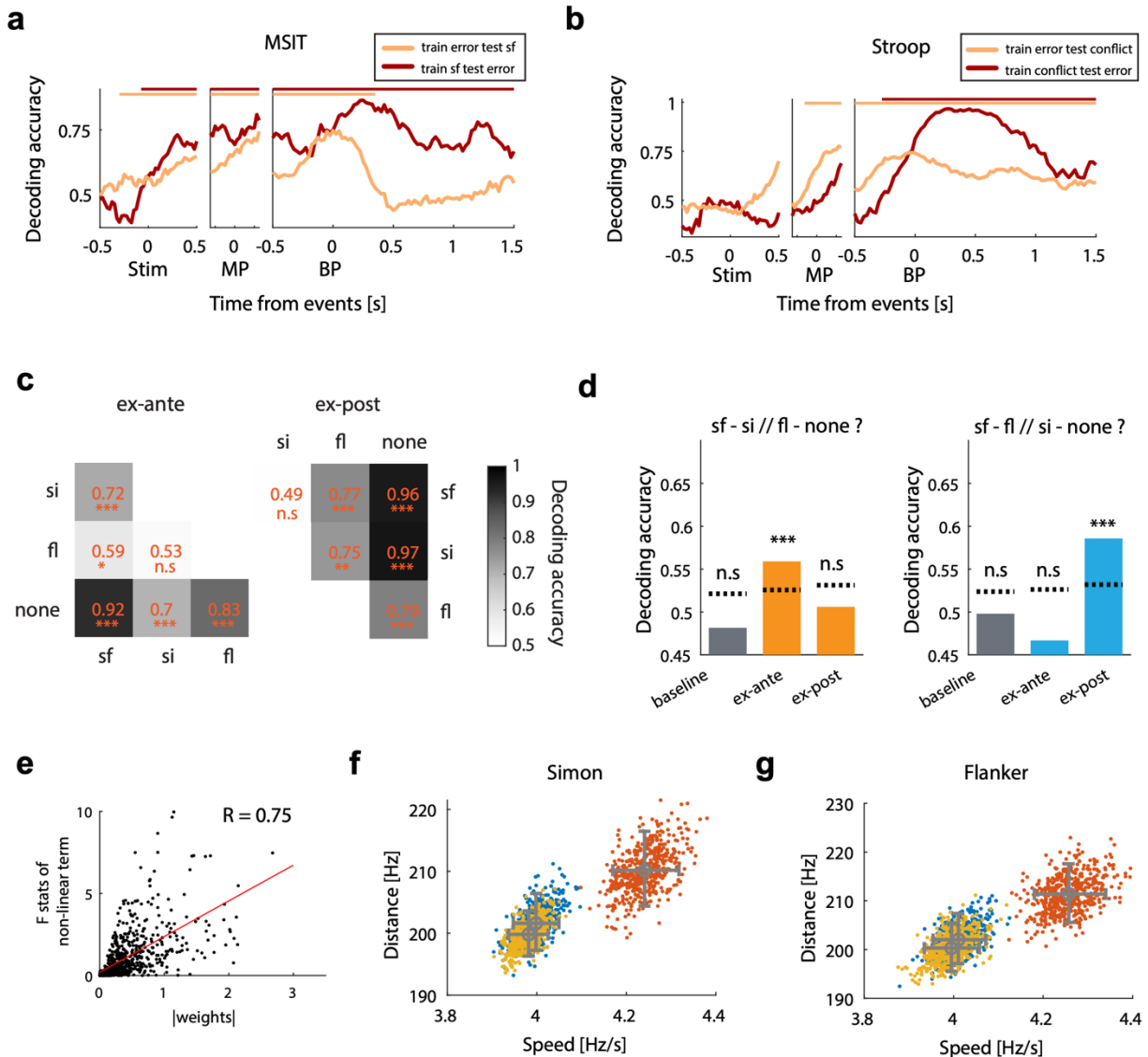1218 represent within-time decoding.
1219 (e-g) Population decoding of Simon (e), Flanker (f), Stroop (g) on the immediately preceding trial
1220 in different epochs. Dotted lines show 97.5% percentile from the null distribution (permutation).
1221 During baseline, there is significant coding of past trial conflict as expected from the persistence
1222 of ex-post conflict signals. Coding of the past trial conflict was non-significant during all other
1223 epochs except for past trial Stroop conflict in the early ex-post epochs, suggesting that this
1224 information in our experimental setup was likely not reliable for cognitive control.
1225
1226 *p < 0.05, ** p < 0.01, *** p < 0.001, n.s., not significant (p > 0.05 or not significant determined
1227 using FDR).
1228

1229

1230 **Figure S6. Within-task state space analyses. Related to Figure 6.**

1231

## Figure S6



1232
1233 (a) A common population coding dimension for error and conflict in MSIT. This coding dimension
1234 is extracted using dPCA, using an error contrast (error "sf" trials vs. "sf" trials) and a conflict
1235 contrast ("sf" trials vs. no conflict). "sf" trials are split into two non-overlapping groups for this.
1236 Plot show the decoding accuracy of both sf (apricot) vs no-conflict trials, and error "sf" vs. "sf"
1237 trials (out-of-sample). Horizontal bars at the top demarcate significant clusters, as determined by
1238 the cluster-based permutation test (p < 0.05).
1239 (b) A common population coding dimension for both error and conflict in Stroop. This coding
1240 dimension is extracted using dPCA, using an error contrast (error conflict trials vs. conflict trials)
1241 and a conflict contrast (conflict trials vs. no conflict) and by marginalizing out the contrast
1242 dimension. Conflict trials are split into two non-overlapping groups for this. Plot show the

1243    decoding accuracy of both sf (apricot) vs no-conflict trials, and error "sf" vs. "sf" trials (out-of-

1244    sample). Horizontal bars at the top demarcate significant clusters, as determined by the cluster-

1245    based permutation test (p < 0.05).

1246    (c) Decoding accuracy of pairwise classification of conflict conditions after RT was equalized

1247    across conditions. Trials were selected such that RTs on si, fl, sf and non-conflict trials were

1248    equalized (p > 0.1, t test). Training data and left-out testing data from the four conflict conditions

1249    are projected to the population vector flanked by averages of non-conflict trials and sf trials,

1250    shown as the dotted line in Figure 6a. Color code represents decoding accuracy. This coding

1251    dimension separates the four conflict conditions well.

1252    (d) Testing compositionality of conflict representation with condition generalization of decoding

1253    on error trials. We tested the condition generalization using raw spike count data from the ex-

1254    ante and ex-post (1s after button presses) epochs and from the baseline as a control without

1255    dimensionality reduction. We only used data on error trials for this analysis. Data from the means

1256    connected by one of the blue edges in Figure 6a were used to train a decoder, which was then

1257    tested with left-out data from the means connected by the opposite blue edge, and vice versa.

1258    Blue bars show decoding accuracy. Same was also tested with data connected by the orange

1259    edges as shown in Figure 6a. Orange bars show decoding accuracy. Decoding accuracy were

1260    reduced on error trials compared to on correct trials (compare with Figure 6d). Dotted lines show

1261    97.5[th] percentile of the null distribution from permutation.

1262    (e) Single neuron with nonlinear coding of Simon and Flanker conflict contribute to deviation of

1263    conflict representation from perfect linearity. Data used here are from the ex-ante epoch.

1264    Nonlinear coding of conflict by a single neuron is measured by the F statistic of the interaction

1265    term between Simon and Flanker conflict in an ANOVA model with spike counts as the dependent

1266    variable. Each neuron's contribution to the deviation from linear additivity in the high

1267    dimensional neural space is quantified by the weight of the difference vector between "sf" and

1268    "si + fl". Scatter plot shows the relation between these two measures. Red line shows the linear
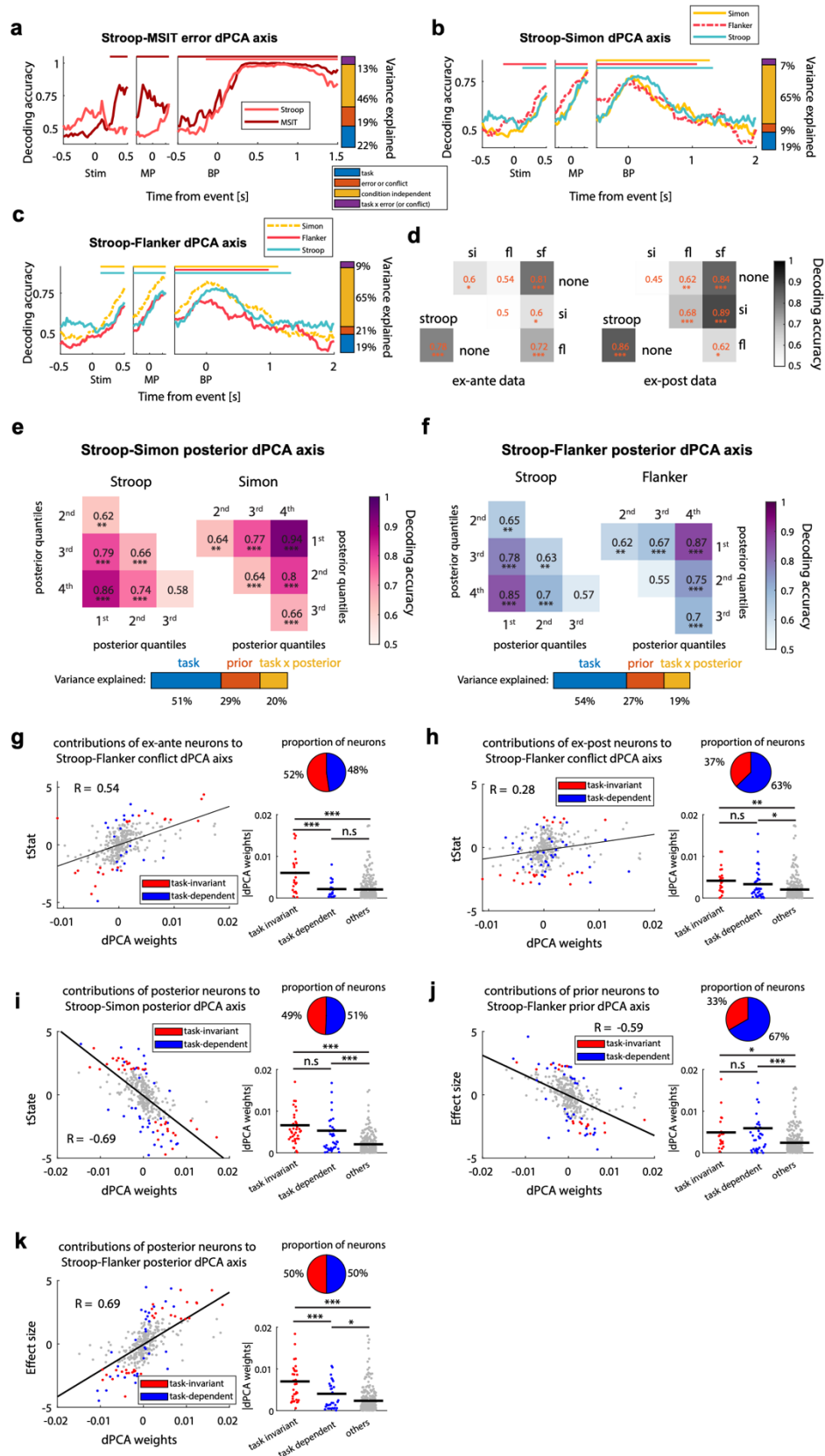
1269    fit.

1270    (f) Distance between trajectories and average speed computed from trials grouped by quartiles

1271    of Simon conflict prior in the baseline (blue) and the ex-ante (orange) epoch, and trials grouped

1272    by Simon conflict posterior in the ex-post epoch (yellow). Trajectories are visualized in Figure 6f.

1273    The state space speed stays low during baseline, increases significantly during the ex-ante epoch

1274    and decreases back to a value similar to that during the baseline. Distance between trajectories

1275    is stable across time.

1276    (g) Same as in (f) but for Flanker conflict prior.

1277

1278    **Figure S7. Domain-general representation of performance monitoring signals. Related to Figure**

1279    **7.**

Figure S7

1281     (a) Task-invariant decoding of errors in both MSIT (salmon) and Stroop (crimson) after RTs were
1282     equalized across conditions. Specifically, trials were selected for Stroop such that the RTs on error
1283     conflict trials did not differ significantly those from correct conflict trials (p > 0.1, t-test). For MSIT,
1284     trials were selected such that RTs on error "sf" trials and correct "sf" trials did not differ
1285     significantly (p > 0.1, t test). The task-invariant coding dimension is extracted using dPCA that
1286     marginalizes out task information and time. This dPCA coding dimension is extracted from the
1287     error contrast in Stroop (error conflict vs. correct conflict trials) and the error contrast in MSIT
1288     (error "sf" trials and correct "sf" trials). This controls for trial conflict and isolates effects related
1289     only to error. Left, Accuracy for decoding errors as a function of time. Bar on the right shows the
1290     variance explained by the different dPCA components (color code see figure legend).
1291     (b-c) Task-invariant decoding of conflict in both MSIT (Simon, yellow; Flanker red) and Stroop
1292     (green) after RTs were equalized across conditions. Because MSIT has two conflict conditions,
1293     Simon and Flanker, task invariance was investigated using Stroop/Simon and Stroop/Flanker
1294     conflicts separately. Specifically, trials were selected for Stroop such that RTs on conflict and non-
1295     conflict trials did not differ significantly (p > 0.1, t test). For MSIT, trials were selected such that
1296     RTs on Simon and non-Simon trials did not differ significantly (p > 0.1) and those on Flanker and
1297     non-Flanker trials did no differ significantly (p > 0.1). This dPCA coding dimension is extracted
1298     from conflict and non-conflict trials in Stroop and either from Simon and non-Simon trials (b) or
1299     from Flanker and non-Flanker trials (c), by marginalizing out task information and time. Left-out
1300     conflict trials and non-conflict trials in Stroop, and left-out Simon, non-Simon, Flanker, non-
1301     Flanker trials in MSIT are projected and classified by this coding dimension. Left, decoding
1302     accuracy of conflict as a function of time. The bar on the right shows variance explained by the
1303     different dPCA components (color code see figure legend).
1304     (d) Testing separability of conflict conditions in Stroop and MSIT using data from the ex-ante (left)
1305     and ex-post epochs (right) after RTs were equalize across conditions. Here, trials were selected
1306     for Stroop such that RTs on conflict and non-conflict trials did not differ significantly (p > 0.1, t
1307     test). For MSIT, trials were selected such that RTs on si, fl,sf, and non-conflict trials did not differ
1308     with each other significantly (p > 0.1, t test). The dPCA coding dimension used in this analysis is
1309     extracted by using conflict and non-conflict trials in Stroop and sf and non-conflict trials in MSIT
1310     by marginalizing out task information. Because data from ROIs are used, temporal information is
1311     already marginalized out before being used by the dPCA algorithm. This coding dimension
1312     support classification of 75% of pairs of MSIT conflict conditions (upper triangle matrices) as well
1313     as Stroop conflict (left corner). Color coding represents decoding accuracy, orange numbers
1314     indicate the numerical values of decoding accuracy of that pair of conflict conditions.
1315     (e) Task-invariant decoding of all pairs of conflict posterior quartiles in Stroop (lower triangle
1316     matrix) and Simon (upper triangle matrix). The dPCA coding dimension here is extracted using
1317     the Stroop conflict posterior contrast (the $1^{st}$ vs. $4^{th}$ quartiles of Stroop conflict posterior) and the
1318     Simon conflict posterior contrast (the $1^{st}$ vs. $4^{th}$ quartiles of Simon conflict posterior),
1319     marginalizing out task information. Color code represents decoding accuracy. Bar at the bottom
1320     shows variance explained of dPCA components (for decoding, the component labelled as
1321     "posterior" is used).
1322     (f) Same as in (e) but for Flanker posterior.
1323     (g-k) Contribution of single neuron coding of Stroop-Flanker ex-ante conflict (g), Stroop-Flanker
1324     ex-post conflict (h), Stroop-Simon posterior (i), Stroop-Flanker prior (j) and Stroop-Flanker

1325     posterior (k) to the task-invariant population coding of these variables. Because MSIT has two
1326     conflict conditions, Simon and Flanker, task invariance was tested between Stroop and Simon or
1327     between Stroop and Flanker separately. We modelled each neuron's baseline (j), ex-ante (g) or
1328     ex-post (h,i,k) response using linear regression. The main effects are a dummy variable indicating
1329     Stroop-Flanker conflict (g,h) or Stroop-Simon posterior (i) or Stroop-Flanker prior (j) or Simon-
1330     Flanker posterior (k) and task ID (Stroop or MSIT), and the interaction term between these two.
1331     A significant interaction suggests that the coding is more prominent in one task than the other.
1332     Task-invariant neurons is defined as having a significant main effect of the variable of interest but
1333     an insignificant interaction with the task ID. Task-dependent neurons is defined as having a
1334     significant interaction term.
1335     (g) Contribution of single neuron ex-ante coding of Stroop-Flanker conflict to task-invariant
1336     population coding of Stroop-Flanker conflict. Of the 42 of neurons that were selected as signaling
1337     conflict in the ex-ante epoch in either task, 52% did so in a task-invariant way (red) and 48% in a
1338     task-dependent (blue) way (pie chart). There is a strong correlation between the error t-statistic
1339     and the dPCA weight of a particular neuron (scatter plot on the left; r = 0.54). Comparing the
1340     mean absolute value of dPCA weights between task-invariant, task-dependent and uncategorized
1341     neurons, the task-invariant neurons had significantly larger absolute weights than the task-
1342     dependent neurons (p < 0.001) and uncategorized neurons (p < 0.001).
1343     (h) Contribution of single neuron ex-post coding of Stroop-Flanker conflict to task-invariant
1344     population coding of Stroop-Flanker conflict. Of the 59 of neurons that were selected as signaling
1345     conflict in the ex-post epoch in either task, 37% did so in a task-invariant way (red) and 67% in a
1346     task-dependent (blue) way (pie chart). There is a strong correlation between the error t-statistic
1347     and the dPCA weight of a particular neuron (scatter plot on the left; r = 0.28). Comparing the
1348     mean absolute value of dPCA weights between task-invariant, task-dependent and uncategorized
1349     neurons, both the task-invariant (p < 0.01) and task-dependent (p < 0.05) neurons had
1350     significantly larger absolute weights than the uncategorized neurons.
1351     (i) Contribution of single neuron ex-post coding of Stroop-Simon conflict posterior to task-
1352     invariant population coding of Stroop-Simon conflict prior. Of the 79 of neurons that were
1353     selected as signaling posterior in the ex-post epoch in either task, 49% did so in a task-invariant
1354     way (red) and 51% in a task-dependent (blue) way (pie chart). There is a strong correlation
1355     between the error t-statistic and the dPCA weight of a particular neuron (scatter plot on the left;
1356     R = 0.51). Comparing the mean absolute value of dPCA weights between task-invariant, task-
1357     dependent and uncategorized neurons, both the task-invariant (p < 0.001) and task-dependent
1358     (p < 0.001) neurons had significantly larger absolute weights than the uncategorized neurons.
1359     (j) Contribution of single neuron baseline coding of Stroop-Flanker conflict prior to task-invariant
1360     population coding of Stroop-Flanker conflict prior. Of the 58 of neurons that were selected as
1361     signaling prior in the baseline in either task, 33% did so in a task-invariant way (red) and 67% in
1362     a task-dependent (blue) way (pie chart). There is a strong correlation between the error t-statistic
1363     and the dPCA weight of a particular neuron (scatter plot on the left; r = -0.59). Comparing the
1364     mean absolute value of dPCA weights between task-invariant, task-dependent and uncategorized
1365     neurons, both the task-invariant (p < 0.05) and task-dependent (p < 0.001) neurons had
1366     significantly larger absolute weights than the uncategorized neurons.
1367     (k) Contribution of single neuron ex-post coding of Stroop-Flanker conflict posterior to task-
1368     invariant population coding of Stroop-Flanker conflict posterior. Of the 66 of neurons that were

1369 selected as signaling posterior in the ex-post epoch in either task, 50% did so in a task-invariant
1370 way (red) and 50% in a task-dependent (blue) way (pie chart). There is a strong correlation
1371 between the error t-statistic and the dPCA weight of a particular neuron (scatter plot on the left;
1372 r = 0.69). Comparing the mean absolute value of dPCA weights between task-invariant, task-
1373 dependent and uncategorized neurons, both the task-invariant (p < 0.001) and task-dependent
1374 (p < 0.05) neurons had significantly larger absolute weights than the uncategorized neurons, and
1375 the task-invariant had significantly larger values than the task-dependent neurons (p < 0.001).
1376
1377
1378
1379
1380 **Methods**
1381
1382 **Tasks**
1383
1384 Subjects performed a speeded version of the Stroop and Multi-Source Interference (MSIT) tasks.
1385 For the Stroop task, subjects were shown a series of randomly intermixed color words ("red",
1386 "green", "blue") printed in either red, green, or blue color (see Figure 1a). Subjects were
1387 instructed to name the color the word stimulus was printed in while ignoring its meaning and to
1388 do so as quickly and accurately as possible. For the MSIT task, subjects were shown an array of
1389 three numbers (0,1,2,3), out of which two were the same and the third of which was different
1390 (target). Subjects were instructed to press the button identical to the target number (which was
1391 unique) regardless of the position at which it was shown. For both tasks, all responses were
1392 recorded as button presses using an external response box (RB-740, Cedrus Corp., San Pedro, CA).
1393 For both tasks, the stimulus disappeared immediately when a button was pressed and was
1394 followed by a blank screen for 1s, followed by a feedback screen, which was shown for 1s. Subject
1395 were given three types of feedback: correct, incorrect or "too slow". 10-15% of trials were rated
1396 as "too slow" based on an adaptive response threshold (see (Fu et al., 2019) for details), which
1397 we used to emphasize the need to respond quickly and thereby resulting in a sufficiently large
1398 error rate (~10% of trials). The inter-trial interval was sampled randomly from a uniform
1399 distribution between 1.5s to 2s. Trial sequences were pseudo-randomized and designed to avoid
1400 back-to-back repetitions of the same stimulus. The proportion of conflict trials in the Stroop task
1401 was 30-40%; For MSIT, the proportions of Simon only ("si"), Flanker only ("fl"), Simon and Flanker
1402 coincident ("sf") trials are 15%, 15%, and 30%, respectively (the remaining 40% of trials have no
1403 conflict). The tasks were implemented in MATLAB (The Mathworks, Inc., Natick, MA) using
1404 Psychtoolbox-3 (Brainard, 1997). The two tasks were performed in sequence, i.e., subjects
1405 finished blocks of one task first and then moved on to blocks of the other task. The order of task
1406 performed was randomized across experimental sessions.
1407
1408 **Behavioral controls**
1409
1410 As a control, we additionally collected behavioral data from N = 51 normal control subjects (24
1411 females; age mean±sd: 44±11) using the Amazon mTurk platform. We implemented the MSIT
1412 task as described above using the jsPsych toolbox (de Leeuw, 2015). These behavioral data were

1413    analyzed using the same methods as documented below. These control subjects exhibited a
1414    robust conflict prior effect like the patients (see Figure S1d).

1415

1416    **Subjects**

1417

1418    34 patients (see Table S1 for age and gender) who were evaluated for possible surgical treatment
1419    of their focal epilepsy using implantation of depth electrodes volunteered for the study and gave
1420    written informed consent. We only included patients with well-isolated single- neuron activity on
1421    at least one electrode in the areas of interest. All research protocols were approved by the
1422    institutional review boards of Cedars-Sinai Medical Center, Huntington Memorial Hospital, and
1423    California Institute of Technology.

1424

1425

1426    **Electrophysiological recordings**

1427

1428    We analyzed data from up to 4 electrodes in each subject (bilateral dACC and pre-SMA) in this
1429    paper. For each depth electrode, there are eight microwires with high impedance microwires at
1430    the tip, and eight macro contacts with low impedance along the shaft (AdTech Medical Inc.). Data
1431    from all microwires and the most medial macro contact (which is placed within dACC or pre-SMA)
1432    are analyzed in this paper. For recordings from microwires, the sampling rate was 32-40kHz and
1433    the raw signal was acquired broadband (0.01Hz-9kHz). One microwire on each depth electrode
1434    was designated as a local reference wire. For intracranial EEG recordings done with macro
1435    contacts, the sampling rate was 2kHz (ATLAS, Neuralynx, Inc., Bozman, MT).

1436

1437    *Electrode localization*

1438

1439    Electrodes were localized using a combination of a pre-operative MRI and a postoperative
1440    MRI/CT using standard procedures we described elsewhere (Fu et al., 2019; Minxha et al., 2017).
1441    Only electrodes that could be clearly localized to the dACC (cingulate gyrus or cingulate sulcus;
1442    for patients with a paracingulate sulcus, electrodes were assigned to the dACC if they were within
1443    the paracingulate sulcus or superior cingulate gyrus) or the pre-SMA (superior frontal gyrus) were
1444    included.

1445

1446    *Spike detection and sorting*

1447

1448    We filtered the raw broadband signal with a zero-phase lag filter in the 300-3000Hz band. Spikes
1449    were detected and sorted using a template-matching algorithm (Rutishauser et al., 2006). Sorting
1450    quality is evaluated using the same metrics reported in (Fu et al., 2019) and only well-isolated
1451    single units are included in this paper. Channels with interictal epileptic events were excluded.

1452

1453    **Quantification and statistical analyses**

1454

1455    *Behavioral modelling and analyses*

1456

1457     We developed a series of Bayesian conflict learning models to infer subjects' internal estimate of
1458     conflict probability (details see below). For this analysis, we concatenated all blocks of an
1459     experiments done in a single session. For trials with unusually long RTs (> 3 sd from the mean of
1460     the whole experiment), we replaced the outlier's RT with the average RT computed from the
1461     neighboring 6 trials. We estimate the Bayesian model parameters using all trials but excluded
1462     error trials (after fitting) for analyses that focuses on conflict and conflict prior. We then analyzed
1463     whether the variance of RT was related to the estimated parameters using linear mixed-effect
1464     models (Aarts et al., 2014). For MSIT, the linear mixed-effect model is specified as follows (all
1465     models are represented in Wilkinson's notation):

1467 $$\log(RT) \sim Simon\ prior * Simon\ conflict + Flanker\ prior * Flanker\ conflict + (1$$
1468 $$+ Simon\ prior + Flanker\ prior | sessionID : subjectID)$$

1470     For Stroop, this is specified as:

1472 $$\log(RT) \sim Stroop\ prior * Stroop\ conflict + (1 + Stroop\ prior | sessionID : subjectID)$$

1474     Here, the fixed effects of Simon, Flanker and Stroop conflicts are dummy variables indicating
1475     whether a particular trial involves conflict (value = 1) or not (value = 0). The fixed effects for priors
1476     are obtained from our Bayesian conflict learning models as detailed below. To test if RT was
1477     affected by conflict on the immediately preceding trial, represented by $Simon\ prevConflict$,
1478     $Flanker\ prevConflict$ and $Stroop\ prevConflict$, we again constructed a linear mixed-effect
1479     models for both MSIT and Stroop. For MSIT, the model is specified as follows:

1481 $$\log(RT) \sim Simon\ prevConflict * Simon\ conflict + Flanker\ prevConflict$$
1482 $$* Flanker\ conflict + (1 | sessionID : subjectID)$$

1484     For Stroop, this is specified as:

1486 $$\log(RT) \sim Stroop\ prevConflict * Stroop\ conflict + (1 | sessionID : subjectID)$$

1488     We investigated the effect of conflict prior on the likelihood of making an error using generalized
1489     linear mixed-effect models. For MSIT, this model is given as

1491 $$Outcome \sim Simon\ prior * SF + Flanker\ prior * SF + (1 + Simon\ prior$$
1492 $$+ Flanker\ prior | sessionID : subjectID)$$

1494     where $SF$ is a dummy variable indicating whether the trial has both Simon and Flanker conflict
1495     (value = 1) or non-conflict (value = 0). We restricted this analysis to sf trials because most errors
1496     occurred on these trials. For Stroop, this model is given as

1498 $$Outcome \sim Stroop\ prior * Stroop\ conflict + (1 + Stroop\ prior | sessionID : subjectID)$$
1499     The response variable $Outcome$ is a categorical variable indicating whether the trial ended in a
1500     correct (0) or incorrect (1) response.

1501    To determine the statistical significance of each fixed effect, we compared the full model with a
1502    reduced model where the fixed effect in question was removed using the likelihood ratio test. To
1503    determine whether RT tuning of the model (see below) is necessary for conflict prior to explain
1504    RT variance, we switched out the conflict prior with the one estimated without RT tuning and
1505    kept all other terms the same. Statistical significance determined this way was indicated by stars
1506    or "n.s" (non-significant) in Figures 1 and S1. To determine whether the conflict prior explains
1507    variance in RT and error likelihood independent of practice, which is assumed to vary with the
1508    trial number, we augmented the aforementioned mixed-effect models by including three
1509    additional trial-ID terms: $trialID, trialID^2, trialID^3$ to capture variance related to  practice
1510    effects.

1511
1512
1513    _Bayesian conflict learning models_
1514
1515    Our models are structurally similar to those used in several previous studies (Behrens et al., 2007;
1516    Jiang et al., 2015). Here, we briefly highlight the modifications we made to extend these existing
1517    models to model behavior in both the Stroop and MSIT tasks, the latter of which has two types
1518    of conflicts that are monitored at the same time. Our models have the following parameters (see
1519    Fig. 1c for a schematic of the model structure): 1) a flexible learning rate $\alpha$, which captures the
1520    subject's belief in the rate of change in control demand in the environment (i.e., a change in
1521    conflict probability), and 2) conflict probability ($q_s$ for Stroop conflict in the model for Stroop, $q_{si}$
1522    for Simon and $q_{fl}$  for Flanker conflicts in the model for MSIT). The models utilize two types of
1523    data (both of which are only available after a trial's response has been made): 1) trial congruency
1524    $o$ (value of 1 indicates an incongruent trial; $o_s$ for Stroop congruency, $o_{si}$ for Simon congruency
1525    and $o_{fl}$ for Flanker congruency; 2) reaction time $RT$, assigning Bernoulli likelihood function to the
1526    former and the drift-diffusion model likelihood function to the latter (details see below). Prior
1527    work (Jiang et al., 2015) uses a Gaussian likelihood function to describe RT generation in a
1528    Bayesian learning framework similar to ours, but we argue that the use of DDM has several
1529    advantages over the Gaussian approach: 1) fewer parameters are used in the DDM, making it
1530    computationally possible to model two types of MSIT conflict at the same time; 2) the DDM
1531    parameters have physiological meaning and thus also provide a clear physiological reasoning for
1532    conflict prior to affect specific components of the decision process; 3) DDM has been widely used
1533    and validated as the generative framework to model RT during decision making (Pedersen et al.,
1534    2017; Wiecki et al., 2013).

1535
1536    Estimating control demand is operationalized as estimating the probability that a certain conflict
1537    (Stroop, Simon or Flanker) would occur in the block. One advantage of our models is that they
1538    estimate both the conflict probability and the rate of change in conflict probability in an online
1539    manner, i.e., the models iteratively update their current estimates after every trial with new
1540    incoming data (of that trial) using Bayes' law. This is consistent with the way human subjects
1541    perform conflict tasks while estimating the associated control demand: they perform and
1542    estimate trial-by-trial. Note that in this study the conflict probability was constant throughout
1543    the experiment (but this was not known to subjects). Nevertheless, we allow the model to infer
1544    the learning rate $\alpha$ online because humans demonstrate inherent bias in believing that

1545 environmental statistics are not stable (Yu and Cohen, 2008). There is therefore no fitting
1546 involved for estimating $\alpha$ and thus the models are not penalized for including this parameter in
1547 model comparisons. Note that we use a single $\alpha$ for both types of conflicts in MSIT. To simplify
1548 model estimation, we made the Markovian assumption that the current estimate of conflict
1549 probability depends only on the current trial congruency and RT, and the estimated conflict
1550 probability on the last trial, but not on the full history of past trial conflict probability (Behrens et
1551 al., 2007). The iterative estimation of conflict probability then involves combining the estimated
1552 conflict probability from the previous trial (prior), transition functions capturing how the current
1553 estimate will change from the previous one (the probability of current estimate conditional on
1554 previous estimate) and the likelihood function.

1555

1556 The model starts with a transition function for the learning rate $\alpha$:

1557

1558 $$p(\alpha_{i+1}|\alpha_i) = k\delta(\alpha_{i+1} - \alpha_i) + 1 - k$$

1559

1560 This formulation assumes that the learning rate has a probability k of having the same distribution
1561 as that of the preceding trial but with a probability 1-k of switching to a uniform distribution (over
1562 all possible $\alpha$), because the learning rate is largely stable across time. The transition function for
1563 conflict probability concerning the transition from the current estimate to a future estimate is
1564 computed in two steps. Here, we refer to the current-trial estimate of conflict probability for
1565 Stroop, Simon or Flanker generically as $q_i$, to which we assigned a uniform prior. The transition
1566 function is thus denoted as $p(q_{i+1}|q_i, \alpha_i)$. First, an auxiliary variable $q_{i+0.5}$ is constructed, which
1567 is a beta-distributed random variable with its mode being $q_i$ and the sum of two parameters
1568 being $\frac{1}{\alpha_{i+1}}$:

1569

1570 $$v_{i+1} = \frac{1}{\alpha_{i+1}} - 2$$

1571
1572 $$q_{i+0.5} \sim Beta(q_i v_{i+1} + 1, v_{i+1} - q_i v_{i+1} + 1)$$

1573

1574 The conflict probability transition function is then constructed as

1575

1576 $$q_{i+1} \sim q_{i+0.5} + \alpha_{i+1}(o_i - q_{i+0.5})$$

1577

1578 The transition function adopts a classical update rule used in reinforcement learning models, and
1579 the learning rate controls the balance between past $(q_{i+0.5})$ and current information $(o_i)$. For
1580 the MSIT model, we take the product of the transition functions computed separately for Simon
1581 and Flanker predicted conflict. Finally, we consider the likelihood function. Since the trial
1582 sequences were designed and re-used for different subjects, the estimated conflict probability
1583 would be the same across subjects for the same sequence, but this is inconsistent with the fact
1584 that such individual estimates should be subjective and different between participants. We thus
1585 incorporated RTs from each subject, which are assumed to be generated through a drift-diffusion
1586 process, to estimate a *subjective* conflict probability based on the assumption that a subjects' RT

1587    reflects the extent to which they engaged control. RT is assumed to be generated by a diffusion
1588    process. We used an abstract version of the drift-diffusion model where the two bounds
1589    represent the correct and wrong choice (and not the actual choices). The diffusion process
1590    accumulated the *difference* in the evidence between the target and distractor response (Fig. 1c,
1591    right), which is smaller for conflict trials and thus leads to longer RTs. We refer to the drift-
1592    diffusion likelihood function for RT as

1593    $p_{DDM}(v_{si}, v_{fl}, v_{sf}, v_{non-conflict}, vz_{si}, vz_{fl}, a, q_{si}, q_{fl})$     for     the     MSIT     model,     and
1594    $p_{DDM}(v_{conflict}, v_{stroop\ non-conflict}, vz, a, q_s)$ for the Stroop model (Navarro and Fuss, 2009). The
1595    hyperparameters specifying the DDM are boundary separation ($a$), base drift rates for Simon-
1596    only, Flanker-only, both Simon and Flanker present, and non-conflict trials in MSIT
1597    ($v_{si}, v_{fl}, v_{sf}, v_{non-conflict}$), base drift rates for conflict and non-conflict trials in Stroop
1598    ($v_{conflict}, v_{stroop\ non-conflict}$), and drift rate bias coefficients that scales the conflict probability
1599    of Simon, Flanker and Stroop ($vz_{si}, vz_{fl}, vz$). The effective drift rate is then the sum between the
1600    base drift rate and the drift rate bias (see Fig. 1c). Here we made the assumption that conflict
1601    prior affects RT by biasing *drift rates* based on a previous work investigating the effect of choice
1602    history on RT (Urai et al., 2019). We also assumed that the drift rate diffusion started at the half
1603    point of the boundary separation (i.e., $z = 0.5$). With the Markovian assumption, the updating
1604    process for the MSIT model is thus given by

1605

1606    $p(k, \alpha_{i+1}, q_{si,i+1}, q_{fl,i+1} | o_{si,\leq i+1}, o_{fl,\leq i+1}, RT_{\leq i+1}) \propto p(o_{si,i+1}, o_{fl,i+1}, RT_{i+1} | q_{si,i+1}, q_{fl,i+1})$
1607    $\iint \left[ \int p(k, \alpha_i, q_{si,i}, q_{fl,i} | o_{si,\leq i}, o_{fl,\leq i}, RT_{\leq i}) p(\alpha_{i+1} | \alpha_i, k) d\alpha_i \right] p(q_{si,i+1} | q_{si,i}, \alpha_i) p(q_{fl,i+1} | q_{fl,i}, \alpha_i) dq_{si,i} dq_{fl,i}$

1608

1609    The updating process for the Stroop model is given by

1610

1611    $p(k, \alpha_{i+1}, q_{s,i+1} | o_{s,\leq i+1}, RT_{\leq i+1})$
1612    $\propto p(o_{s,i+1}, RT_{i+1} | q_{s,i+1}) \int \left[ \int p(k, \alpha_i, q_i | o_{s,\leq i}, RT_{\leq i}) p(\alpha_{i+1} | \alpha_i, k) d\alpha_i \right] p(q_{s,i+1} | q_{s,i}, \alpha_i) dq_{s,i}$

1613

1614    The likelihood function is the product of Bernoulli likelihood for trial congruency and DDM
1615    likelihood for RT. For MSIT, the likelihood function is given as:

1616

1617    $p(o_{si,i+1}, o_{fl,i+1}, RT_{i+1} | q_{si,i+1}, q_{fl,i+1}) = (1 - |o_{si,i+1} - q_{si,i+1}|)(1 - |o_{fl,i+1} - q_{fl,i+1}|)p_{DDM}$

1618

1619    For Stroop, the likelihood function is given as:

1620

1621    $p(o_{s,i+1}, RT_{i+1} | q_{s,i+1}) = (1 - |o_{s,i+1} - q_{s,i+1}|)p_{DDM}$

1622

1623    These hyperparameters are estimated using an expectation-maximization (EM) algorithm as
1624    shown in earlier work (Jiang et al., 2015). Briefly, the model parameters were at first estimated
1625    *without* incorporating the DDM likelihood for RT ("E" step). Hyperparameters were then fit by
1626    maximizing the DDM likelihood for the observed RT using the conflict prior (s) obtained ("M"
1627    step). The DDM likelihood function with the fitted hyperparameters were then incorporated into

1628    the Bayesian updating process ("E" step) to generate a new set of conflict prior (s), which were
1629    then used to maximize the DDM likelihood over observed RT again. These steps were repeated
1630    until the convergence of both model parameters and hyperparameters (Euclidean distance
1631    between parameter vectors from successive iterations $< 10^{-5}$).

1632

1633    We considered three alternative classes of models: 1) reinforcement learning (RL) models; 2)
1634    constant model; 3) Bayesian learning model without RT tuning. For the RL model, we constructed
1635    a value function corresponding to the estimated conflict probability, and this estimate is also
1636    updated trial-by-trial using a Rescorla-Wagner rule. For MSIT, the update rule is:

1637

1638
$$q_{si,i+1} = q_{si,i} + \alpha\big(o_{si,i+1} - q_{si,i}\big)$$

1639

1640
$$q_{fl,i+1} = q_{fl,i} + \alpha\big(o_{fl,i+1} - q_{fl,i}\big)$$

1641

1642    For Stroop, the update rule is:

1643

1644
$$q_{s,i+1} = q_{s,i} + \alpha\big(o_{s,i+1} - q_{s,i}\big)$$

1645

1646    The free parameter $\alpha$ in the RL models was fit by maximizing the data likelihood (Bernoulli) for
1647    trial congruency. For the constant model, $q_s$, $q_{si}$, $q_{fl}$ were fit directly by maximizing the data
1648    likelihood for trial congruency. For Bayesian conflict learning models without RT tuning, $q_s$, $q_{si}$,
1649    $q_{fl}$ were estimated online but the likelihood function for RT was not incorporated in the process.

1650

1651    We used the Bayesian Information Criterion (BIC) to compare the RT-tuned Bayesian conflict
1652    learning models with the RL models, constant models and the non-RT tuned Bayesian conflict
1653    learning models. We compared these models separately for their ability to explain trial
1654    congruency and RT. For this analysis, we pooled all data from all sessions and computed the BIC
1655    for each model and for each data type (RT or trial congruency), consistent with a previous study
1656    (Behrens et al., 2007). Results of model comparisons can be found in Tables S2 and S3.

1657

1658

1659    *Selection of neurons*

1660

1661    We defined the epochs of interests according to events in the tasks (see Figure 3a for an
1662    illustration). The baseline epoch starts at 1.5s before stimulus onset and ends at stimulus onset.
1663    This epoch is used to analyze encoding of conflict prior. The ex-ante epoch is anchored to the
1664    midpoint of a period of time starting at 100ms after stimulus onset (to account for the minimal
1665    delay needed for visual information to reach the MFC) and ending at the time of button presses.
1666    We then defined the ex-ante epoch as a 500ms window centered on the midpoint of this period.
1667    The rationale for analyzing conflict signals in this epoch is as follows: at the early stage of stimulus
1668    processing, information about the different response options is not yet fully processed and hence
1669    minimal conflict; the conflict signal should reach its maximum when the different stimulus

1670   dimensions that drive competing responses are fully available; and finally, it should subside after
1671   a response is selected.
1672   We counted the number of spikes in these epochs and regressed the spike counts against the
1673   different regressors (error, conflict, conflict surprise, conflict prior and conflict posterior) using
1674   linear regression. For each regressor, we extracted a p value computed from the F test. A neuron
1675   was deemed selective for this regressor when $p < 0.05$. For MSIT, since there were both Simon
1676   and Flanker conflicts, neurons were selected when regressors related to either Simon or Flanker
1677   conflict were significant (e.g., ex-ante conflict cells in MSIT were the union of neurons selective
1678   for Simon conflict and Flanker conflict during the ex-ante epoch). To assess whether a neuronal
1679   class is significantly present in the population, we derive a null distribution by permuting the
1680   relation between spike counts and the regressor of interest for 1000 times. A p value is computed
1681   by comparing the true proportion of selected neurons against this null distribution. The 95[th]
1682   interval of the null distribution for each neuronal class is plotted as dotted lines in Figure 3b.
1683   To statistically compare the extent of multiplexing between two groups of cells active in different
1684   epochs (Figure 3d), we used the chi-squared test and reported the p-value and effect size of the
1685   test.
1686
1687   *Single-trial spike train latency*
1688
1689   We estimated the onset latency in individual trials using Poisson spike-train analysis (Figure 3e)
1690   (Hanes et al., 1995). This method detects the moments when the observed inter-spike intervals
1691   (ISI) deviate significantly from that assumed by a constant-rate baseline Poisson process. We
1692   used the spike rate averaged across the whole block of experiment as a baseline spike rates for
1693   each neuron. This baseline rate was then used to compute a Poisson surprise metric across the
1694   spike train. We started our detection algorithm from the onset of stimulus for each trial. For the
1695   ex-ante conflict neurons (two columns on the left), we restricted the range in which the detection
1696   algorithm looks for bursts to after stimulus onset and before button presses. This is because, by
1697   their definition, ex-ante conflict neurons carried a conflict signal before action. For the ex-post
1698   conflict neurons (two columns on the right), we restricted the range to 200ms before button
1699   presses and before end of trial. We then extract the latency of the first significant burst. The
1700   statistical threshold for detecting an onset was $p < 0.01$. Repeating the same procedure with a
1701   threshold of $p < 0.001$ did not affect our conclusions. For these analyses, we only used the conflict
1702   trials as we focused on the single-trial conflict response of selected conflict-encoding neurons.
1703
1704   *Correct-related potential (CRP) analyses*
1705
1706   We simultaneously also recorded the intracranial electroencephalography (iEEG) while we
1707   recorded single unit activity. iEEG data were acquired from low-impedance macro contacts
1708   closest to the microwires. We focused on the contacts that were directly placed in dACC and pre-
1709   SMA (as confirmed by post-operative imaging, see (Fu et al., 2019)). To extract the CRP, we
1710   downsampled the iEEG data to a sampling frequency of 100Hz (using MATLAB "resample") and
1711   then bandpass filtered (0.1Hz-10Hz) the data with a finite impulse response filter (MATLAB
1712   function "fir1"). Filtered data were then shifted in time to account for average filter delay

1713   (computed using MATLAB function "grpdelay"). We then computed the CRP amplitude for each
1714   trial by averaging the filtered iEEG data within [0,250ms] after button presses.
1715
1716   To analyze whether the CRP amplitude was related to conflict and/or RT, we used a linear mixed-
1717   effect model, pooling experimental sessions and electrodes. The model in Wilkinson's notation
1718   is given by
1719
1720   $$\text{CRP} \sim Stroop\ prior * Stroop\ conflict + RT$$
1721   $$+ (1 + Stroop\ prior + RT | sessionID:subjectID)$$
1722
1723   The fixed effect of $Stroop\ conflict$ is a dummy variable indicating whether a trial is a conflict
1724   trial (value = 1) or not (value = 0). We analyzed the relation between spike counts of prior neurons
1725   in dACC and pre-SMA and the simultaneously recorded CRP using a Poisson mixed-effect model.
1726   Spike counts were gathered using a 500ms bin swept across the trial in steps of 25ms. The model
1727   was computed for spike counts in each bin. The full model is given by:
1728
1729   $$\text{Spike counts}(t) \sim CRP + Stroop\ prior + RT + (1 | cellID)$$
1730
1731   To determine the statistical significance of each fixed effect, we compare the full model with a
1732   reduced model with the fixed effect of interest removed and tested the likelihood ratio between
1733   the full and reduced model using a likelihood ratio test. For the CRP-spike count relation model,
1734   the reduced model is the same as the full model except that the fixed effect of $CRP$ is removed.
1735   We plotted the likelihood ratio from this model comparison as a function of time in Figure 3I. The
1736   p-values were obtained from the likelihood ratio tests and corrected using false discovery rate
1737   method.
1738
1739   *Detrended fluctuation analysis*
1740
1741   Detrended fluctuation analysis (DFA) was first developed by Peng and colleagues (Hardstone et
1742   al., 2012; Peng et al., 1994) to quantify long-range temporal correlations (LRTC). We use DFA to
1743   quantify the extent of LRTC in baseline spike counts on the scale of *trials.* First, the cumulative
1744   sum of the spike counts during baseline was computed. To be consistent with prior literature we
1745   refer to this cumulative sum as the *signal profile*. A set of trial window sizes were defined
1746   between the lower bound of 4 trials and the upper bound of the block length. For each window
1747   size, we then partitioned the signal profile into a series of data snippets. Partitioning was done
1748   such that two adjacent snippets had an overlap of half the window size. We then removed the
1749   linear trend from each data snipped (using least square regression) and computed the standard
1750   deviation across time. The mean of the standard deviations across all snippets of identical
1751   window size was then computed (y axis of Fig. 4f). Finally, the mean standard deviations were
1752   regressed linearly against the logarithmically scaled time windows and the slope was extracted
1753   as the DFA $\alpha$ value (Fig. 4f shows the fluctuations as a function of logarithmically transformed
1754   trial window sizes for two example neurons).
1755   For Fig 4a-b, we tested the relation between a neuron's baseline DFA $\alpha$ value and its tendency to
1756   encode conflict prior. To avoid selection bias, we split trials into two sets of equal size, with one

1757 half consisting of a consecutive run of trials. This is because DFA is used for time series data and
1758 thus required the data be consecutive and temporally ordered. For randomization, we first
1759 randomly sampled one trial from the first half of the block. Then a consecutive run of trials
1760 starting with this randomly picked trial as the starting point were extracted. The consecutive set
1761 was used to compute DFA $\alpha$ value while the rest of the trials were used to correlate with conflict
1762 prior (Simon, Flanker or Stroop) using Spearman rank correlation.

1763

1764 *Decoding analysis (Support-vector machine)*

1765

1766 Data were aggregated from different experimental sessions to create pseudo-population data
1767 matrices. We constructed for each trial a peri-stimulus time histogram (PSTH) using 500ms bins
1768 in steps of 25ms. For all conflict or conflict prior -related decoding, we used correct trials only.
1769 Since different behavioral sessions had different number of trials (some subjects participated in
1770 less sessions than the others), we subsampled the same number of trials from each condition for
1771 each neuron and repeat this process 50 times. For error decoding, we subsampled 10 error trials
1772 and 10 correct conflict trials for Stroop and 10 correct sf trials for MSIT. Since most errors
1773 occurred on high conflict trials, these contrasts isolate the effect of error while controlling for the
1774 effect of conflict. For conflict decoding, we subsampled 30 trials from each conflict conditions:
1775 conflict and non-conflict trials for Stroop; Simon and non-Simon trials, Flanker and non-Flanker
1776 trials for MSIT. For each time bin, we performed 5-fold cross validation using LIBSVM (Chang and
1777 Lin, 2011). We used the linear kernel and set the $c$ parameter to 1 for all analyses. In brief, trials
1778 were first randomly split into 5 equal parts; each part was used in turn as the testing data while
1779 the rest of the four parts were used as training data. Decoding accuracy was the proportion of
1780 correct classifications among the 250 samples (50 resamples x 5 folds). Note that the resampling
1781 was done once to generate testing and training sets for the whole time series and used for both
1782 within-time and across-time decoding. For within-time decoding (All plots with dotted lines in Fig.
1783 5), the SVM classifier trained using the training data from each time bin was tested using the
1784 testing data from the same time bin. For cross-temporal decoding (temporal generalization; all
1785 plots with solid lines in Fig. 5), the SVM classifier trained using the training data from each time
1786 bin were tested across the trial using testing data gathered from other time bins.

1787

1788

1789 *Reaction times equalization*

1790

1791 For analyses shown in Figures S6c and S7a-d, we selected a subset of trials from each condition
1792 so that the RTs did not differ significantly across conditions (e.g., equalizing RTs between conflict
1793 and non-conflict trials in the Stroop data). Here we detail the RT equalization procedure we used
1794 to create "RT equalized sets". We first selected a condition as the "anchor" condition. We sorted
1795 the RTs of the anchor condition, and for each RT we searched in the target (to-be-equalized)
1796 condition(s) for a trial whose RT did not differ from the anchor RT by more than 0.1s. If all RTs in
1797 the target condition differed from the anchor RT by > 0.1s, the anchor RT was not included in the
1798 RT equalized set. Once selected, the anchor RT and the target RT were both removed from the
1799 original set to ensure that no trials were included twice in any RT equalized trial sets.  This
1800 procedure was repeated until one of the conditions considered were emptied. We confirmed

1801 post hoc that RT equalization was successful by testing whether RTs were not significantly
1802 different using ANOVA (p > 0.5 for all the RT equalized sets).
1803
1804
1805 _Decoding analysis (population activity vectors and demixed PCA)_
1806
1807 Data were aggregated from different experimental sessions to create a pseudo-population. We
1808 randomly selected one trial for each neuron from one condition and concatenate the data from
1809 each neuron to form a single-trial testing data matrix. The rest of the trials were averaged for
1810 each condition and concatenated to form a training data matrix. Coding dimensions were defined
1811 based on the condition-averaged training data. To define the coding dimensions used to decode
1812 conflict conditions within MSIT, we used the population activity vectors (a high dimensional
1813 vector in the raw firing rate space) defined by the difference between the two condition means.
1814 To define coding dimensions for the cross-task decoding problems we used dPCA to extract
1815 demixed principal components (dPC). Details of which trials were used to define the coding
1816 dimensions used to generate Figures 6,7, S6 and S7 are given in the sections to follow. Both
1817 testing and training data were projected onto the identified coding dimensions. The labels for
1818 testing data were assigned according to the label of the nearest neighbor of the training data. To
1819 test condition generalization, we projected the testing data from one pair of conditions to a
1820 coding dimension defined by another pair of conditions (e.g., a Simon trial and non-Simon trial
1821 projected to the population vector flanked by Flanker and non-Flanker trial averages) and
1822 classified using the labels of the nearest projected training data. This decoding procedure was
1823 repeated 1000 times (resulting in 1000 single-trial testing data matrices and the corresponding
1824 training data matrices), and the decoding accuracy was defined by the proportion of correct
1825 classifications among these 1000 repetitions. To determine statistical significance, we permuted
1826 the trial labels for 500 times and for each permutation, we repeated all above steps to generate
1827 a null distribution. A p-value was computed from comparing the true decoding accuracy with the
1828 null distribution.
1829
1830 _Pseudo-population matrices for MSIT analyses_
1831
1832 For Figure 6a, we formed the pseudo-population data matrix by taking the average of spike
1833 counts within the ex-ante or ex-post (1s after button presses) epoch across all Simon-only,
1834 Flanker-only, Simon+ Flanker, and non-conflict trials, respectively. We then used PCA on this
1835 condition-averaged data matrix to extract the three principal components (PC) that explained
1836 most variance to visualize the geometric arrangement of the four conflict types. For Figure 6f-h,
1837 trials were binned by quartiles of prior and posterior into four bins separately. However, because
1838 conflict prior was updated into conflict posterior after each button press, binning priors does not
1839 guarantee that the posteriors would fall into the same bins. This is because updating is specific
1840 to each behavioral session and thus differs between neurons. Averaging trials using only bins
1841 formed by prior quartiles would thus mix trials with different levels of posterior for each neuron.
1842 To avoid this problem, we thus formed the data matrix (which now includes the time dimension
1843 rather than a single ROI; spikes were counted in 500ms bins swept across the whole trial in steps
1844 of 25ms; spike trains were aligned to button presses) by concatenating two submatrices: one that

1845     was constructed by averaging trials within bins defined by prior quartiles using data *before* button
1846     presses, and one that was constructed from averaging trials within bins defined by posterior
1847     quartiles for neural data *after* the button press. We then used PCA to find the three PCs that
1848     explained the most variance for this matrix. The concatenated data matrix was then projected
1849     onto these PCs to generate the visualization of trajectory corresponding to prior/posterior levels.
1850

1851     <u>Vectors in the state space to quantify population geometry within MSIT</u>
1852

1853     We next describe how coding dimensions were defined in each case using population activity
1854     vectors in the raw firing rate space. For Figure 6b, the coding dimension was the population
1855     vector flanked by the trial averages of sf and non-conflict trials (Fig. 6a, dashed lines).
1856     Classifications were carried out between pairs of conflict conditions (e.g., between si and fl trials)
1857     as detailed above. For Figure 6c, we took a bin-wise approach to investigate whether Simon
1858     conflict representation generalize to Flanker representation, and vice versa. For this, we split
1859     trials into four non-overlapping groups: Simon, Flanker, non-Simon, non-Flanker trial sets. We
1860     split sf and non-conflict trials randomly in half. One half of sf trials were pooled with si trials to
1861     form the Simon trial set, and one half of non-conflict trials were pooled with fl trials to form the
1862     non-Simon trial set. The other half of sf trials were then pooled with fl trials to form the Flanker
1863     trial set, and the other half of non-conflict trials were pooled with si trials to form non-Flanker
1864     trial set. Using these trial sets, for each time bin we extracted two coding dimensions from the
1865     training data: one population vector flanked by trial averages of Simon and non-Simon trials
1866     (Simon coding dimension), and one population vector flanked by the trial averages of Flanker and
1867     non-Flanker trials (Flanker coding dimension). We then projected the testing data from
1868     Simon/non-Simon trials onto the Flanker coding dimension and classified the testing data using
1869     the closest projected training data, and vice versa. For details of this classification procedure see
1870     above paragraph. This assesses the extent to which coding of Simon and Flanker conflict is
1871     abstract.
1872

1873     <u>Compositionality of conflict representation</u>
1874

1875     For Figure 6d, the coding dimension were taken to be the blue and orange edges as shown in
1876     Figure 6a. The purpose of this analysis is to assess to what extent the representation of conflict
1877     is compositional (within a task). We assumed that in the neuronal firing rate space, the
1878     representation of Simon/Flanker conflict is a vector pointing from non-conflict trial averages to
1879     the si/fl trial averages. Compositionality of such conflict representation would imply that the sf
1880     representation (vector pointing from non-conflict trial average to the sf trial average) is equal to
1881     the sum of the Simon and Flanker representations. According to the parallelogram law of vector
1882     addition, this then corresponds to the blue and orange edges in Figure 6a forming a parallelogram.
1883     We tested the extent of parallelism in the data using decoding. The coding dimensions here were
1884     defined by the following population vectors using training data: one flanked by non-conflict and
1885     si trial averages (Fig. 6a, blue), one flanked by fl and sf trial averages (blue), one flanked by fl and
1886     non-conflict trial averages (orange) and one flanked by si and sf trial averages (orange). Left-out
1887     testing data from conditions flanking one of the blue or orange pair of edges were then projected
1888     to the other edge in the pair and classified by the training data defining this edge. For example,

1889 single-trial testing data of non-conflict and si trials were projected to the coding dimension
1890 flanked by fl and sf trial averages and were classified by fl or sf trial averages.
1891
1892 <u>Relationship of single neuron tuning with parallelism in geometry</u>
1893
1894 For Figure 6e and Figure S6e, the goal is to investigate the relation between the nonlinearity in
1895 single neuron conflict coding and the deviation from perfect compositionality in state space
1896 representation of conflict. We denote the state-space representation of Simon and Flanker -only
1897 conflict as the population vectors flanked by the trial averages of si and non-conflict and by the
1898 trial averages of fl and non-conflict. We refer to the state space location occupied by the linear
1899 sum of Simon and Flanker representation defined above as "s+f". The deviation from perfect
1900 compositionality is then given by the population vector flanked by "sf" and "s+f". The loading of
1901 "sf" to the "s+f" vector reflects the single neuron contribution to the deviation at the population
1902 level. To quantify nonlinearity of conflict coding for each neuron, we first regressed the spike
1903 counts in the ex-ante or ex-post epoch (1s) against three fixed effects: a Simon effect (dummy
1904 variable indicating the presence or absence of Simon conflict on a trial), a Flanker effect (dummy
1905 variable indicating the presence or absence of Flanker conflict on a trial) and the interaction term
1906 between these two. We extracted the F statistic related to the interaction term, which captures
1907 the effect of nonlinear mixing of Simon and Flanker conflict. We then extracted a population
1908 vector flanked by the sf trial average and and "s+f", the sum of two population vectors one
1909 flanked by trial averages of si trials and non-conflict trials, and one flanked by trial averages of fl
1910 trials and non-conflict trials. We then correlated the loading of "sf" - "s+f" vector and the F
1911 statistics from a particular neuron.
1912
1913 <u>Quantification of state space dynamics</u>
1914
1915 For Figure 6i, we binned spike counts using 250ms bins swept across the trial in steps of 10ms.
1916 The state-space speed was defined to be the Euclidean distance between population vectors of
1917 adjacent time bins divided by the step size. We averaged the state-space speed across time within
1918 an epoch. We also computed the Euclidean distance between pairs of trajectories (1st and 2nd,
1919 2nd and 3rd,3rd and 4th) and averaged this across trajectories and across time bins within an epoch.
1920 State-space speed and the averaged distance between trajectories were plotted against each
1921 other in Figure 6I. Our method for extracting speed in state-space follows prior work (Stokes et
1922 al., 2013).
1923
1924 <u>Testing ordinal relationship of prior/posterior projections</u>
1925
1926 We analyzed the ordinal relation between neural projections of prior/posterior as shown in
1927 Figure 6j-l. PCA axes encoding prior/posterior variance were extracted from spike count data
1928 collected in ROIs (baseline for prior and the ex-post epoch (0-1s after button presses) for
1929 posterior). Since prior/posterior is continuously valued, we created four trial conditions by
1930 binning the trials using quartiles of prior/posterior. For each type of prior or posterior (Simon,
1931 Flanker and Stroop), we projected the left-out trial (not used for computing the PCA axis) onto
1932 the PCA axis for each trial condition and this procedure was repeated 1000 times, yielding 1000

1933    projected values for each trial condition. We then regressed the projected values (concatenated
1934    into a vector) against their trial condition labels (1st,2nd,3rd,4th quartile bins) using a multinomial
1935    logistic regression with the assumption of ordinal relation between trial groups. Essentially, we
1936    were testing whether the out-of-sample project values can reliably predict the trial condition
1937    they belong to assuming that the conditions were ordinal. We reported the p-value and t-statistic
1938    of the effect of projected values.

1939

1940

1941    *Demixed Principal Component Analyses (dPCA)*

1942

1943    We used dPCA to extract task-invariant representation of performance monitoring signals. For
1944    Figures 7 and S7, we investigated task-invariant coding of error, conflict and conflict prior
1945    separately, resulting in three separate optimization problems. For Figure S6, we investigated the
1946    invariance coding of error and conflict. Analyses on conflict and conflict prior used only correct
1947    trials. We used dPCA as described previously (Kobak et al., 2016), with the following adaptions
1948    made for our purposes. The dPCA algorithm first decomposes population neural activity into
1949    marginalized data matrices with respect to the variables of interest. For analyses in Figure 7a-c,
1950    we constructed the marginalized population activity (referred generically as $\overline{X_\phi}$) with respect to
1951    error ($\overline{X_{error}}$ in Fig. 7a) or conflict ($\overline{X_{conflict}}$, Fig. 7b-c) by marginalizing out time and task
1952    dimensions (denoted by "$\langle \cdot \rangle_{task,t}$").
1953    For interpretability, we investigated whether the neural representation is abstract across tasks
1954    separately between Stroop and Simon conflict ("$s \;\&\; si\; conflict$" is the task-invariant dimension
1955    indicating presence of absence of conflict for both tasks) and between Stroop and Flanker
1956    ("$s \;\&\; fl\; conflict$" is the task-invariant dimension indicating presence of absence of conflict for
1957    both tasks). Set up this way, the "task" dimension captures variance related to task set
1958    differences (Stroop vs. MSIT).   To compute marginalized averages, we use N-dimensional
1959    population activity

$$\overline{X_{Stroop\; error\; \&\; MSIT\; error}} = \langle r(error, task, t) - \bar{r}(t) \rangle_{task,t}$$

$$\overline{X_{s\; \&\; si\; conflict}} = \langle r(s \;\&\; si\; conflict, task, t) - \bar{r}(t) \rangle_{task,t}$$

$$\overline{X_{s\; \&\; fl\; conflict}} = \langle r(s \;\&\; fl\; conflict, task, t) - \bar{r}(t) \rangle_{task,t}$$

1967    , where $\bar{r}(t)$ is the firing rate averaged across trials and time bins. For Figure S7a-c, these
1968    definitions are the same except that RT equalized trial sets were used.
1969    For Figure S6a-b, we sought a common coding dimension between error and conflict separately
1970    for MSIT and Stroop, by marginalizing out the information about time and which pair of
1971    conditions were contrasted ("contrast" indicator, for MSIT it indicates whether the contrast
1972    considered is sf vs. non-conflict or error sf vs. correct sf; for Stroop it indicates whether the
1973    contrast considered is correct conflict vs. correct non-conflict or error vs. correct conflict).
1974    we constructed the marginalized population activity with respect to error vs conflict in both
1975    MSIT ($\overline{X_{MSIT\; error/sf}}$ in Fig. S6a, $MSIT\; error/sf$ is the contrast-invariant dimension indicating

1976 presence or absence of errors and presence or absence of sf conflict) and Stroop

1977 ($\overline{X_{Stroop\ error/conflict}}$ in Fig. S6b, $Stroop\ error/conflict$ is the contrast-invariant dimension

1978 indicating presence or absence of errors and presence or absence of sf conflict) as follows:

1979

1980 $$\overline{X_{MSIT\ error/sf}} = \langle r(MSIT\ error/sf, contrast, t) - \bar{r}(t) \rangle_{contrast,t}$$

1981

1982 $$\overline{X_{Stroop\ error/conflict}} = \langle r(Stroop\ error/conflict, contrast, t) - \bar{r}(t) \rangle_{contrast,t}$$

1983

1984 For analyses in Figure 7d-f and Figure S7e-f investigating task-invariant coding of conflict prior,

1985 we used data from a single ROI (ex-ante or ex-epoch) and hence only the task but not time

1986 dimension was marginalized out. Here again for interpretability, we investigated cross-task

1987 representation between Stroop prior and Simon prior and between Stroop prior and Flanker

1988 prior separately, ensuring that the task dimension captures task set difference

1989

1990 $$\overline{X_{s\ \&\ si\ conflict}} = \langle r(s\ \&\ si\ conflict, task) - \bar{r}(t) \rangle_{task}$$

1991

1992 $$\overline{X_{s\ \&\ fl\ conflict}} = \langle r(s\ \&\ fl\ conflict, task) - \bar{r}(t) \rangle_{task}$$

1993

1994 $$\overline{X_{s\ \&\ si\ prior}} = \langle r(s\ \&\ si\ conflict, task) - \bar{r}(t) \rangle_{task}$$

1995

1996 $$\overline{X_{s\ \&\ fl\ prior}} = \langle r(s\ \&\ fl\ conflict, task) - \bar{r}(t) \rangle_{task}$$

1997

1998 For analyses in Figure S7d, these definitions are the same except that RT equalized trial sets were

1999 used.

2000 The algorithm then finds encoding ($F_\phi$) and decoding ($D_\phi$) matrices separately for each

2001 marginalized averages using the regularized reduced-rank regression:

2002

2003 $$L_\phi = \left\| \overline{X_\phi} - F_\phi D_\phi \bar{X} \right\|^2 + \mu \left\| F_\phi D_\phi \right\|^2$$

2004

2005 We assigned a regularization coefficient $\mu$ to avoid overfitting ($\mu = 6e^{-6}$ determined from

2006 results reported in (Kobak et al., 2016)). We used the columns of $D_\phi$ as the demixed principal

2007 components (dPC) and projected N-dimensional data (single-trial data for testing and trial-

2008 averaged data for training) to these dPCs. The numerical values of $D_\phi$ reflects the contribution

2009 for each neuron to task-invariant representation.

2010

2011 To test the statistical significance of coding dimensions, we randomly chose one trial for each

2012 trial type (e.g., one error trial and one correct trial) and constructed a single-trial activity matrix

2013 $X_{test}$. We then used the remaining trials to form the trial-averaged training data $\overline{X_{train}}$, which is

2014 used to find the dPCA coding dimensions. The left-out single-trial data $X_{test}$ is then projected

2015 onto the first coding dimension that captures the most variance computed from $X_{train}$, and

2016 classified according to the closest class mean. We repeated this procedure 1000 times and

2017 determined the decoding accuracy as the proportion of correct classification among the 1000

2018 test trials. We then generated the null distribution by shuffling the trial labels and then repeated

2019 the decoding procedure 500 times. For Figure 7d-f, statistical significance is determined by
2020 comparing the true decoding accuracy with this null distribution. For Figure 7a-c, statistical
2021 significance is determined by the cluster-based permutation test using this null distribution
2022 (Maris and Oostenveld, 2007). The fraction of explained variance (Bars in Figures 7a-c,e-f and
2023 S7a-c,e-f) for each marginalization is given by:

2024
$$R_\phi{}^2 = \frac{\left\|\overline{X_\phi}\right\|^2 - \left\|\overline{X_\phi} - F_\phi D_\phi \bar{X}\right\|^2}{\left\|\overline{X_\phi}\right\|^2}$$

2025
2026 For analyses in Figure 7g-j and Figure S7 g-k, we first quantify for each single neuron its task-
2027 invariant coding strength of error, conflict or conflict prior within a certain ROI. Spike counts
2028 within the baseline, ex-ante or ex-post epochs from MSIT and Stroop were concatenated and
2029 were regressed against three fixed effects: a cognitive effect (trial outcome, trial congruency or
2030 conflict prior), a task effect (a dummy variable with value 1 for MSIT and value 0 for Stroop) and
2031 an interaction between these two. The signed effect size is taken to be the t-statistic computed
2032 from this linear regression. The t-statistic related to the cognitive effect characterizes the
2033 strength of task-invariant coding of cognitive variables (error, conflict and conflict prior). The task
2034 dependency of such coding is captured by the t-statistic of the interaction term (which indicates
2035 that a neuron exhibits non-linear mixing).
2036 Neurons with a significant cognitive effect but a non-significant task effect and a non-significant
2037 interaction term were classified as "task-invariant" neurons. Neurons with a significant
2038 interaction term are classified as "task-dependent" neurons.
2039
2040

2041 **References**

2042

2043 Aarts, E., Verhage, M., Veenvliet, J.V., Dolan, C.V., and van der Sluis, S. (2014). A solution to
2044 dependency: using multilevel analysis to accommodate nested data. Nat. Neurosci. *17*, 491–
2045 496.

2046 Aron, A.R., Behrens, T.E., Smith, S., Frank, M.J., and Poldrack, R.A. (2007). Triangulating a
2047 cognitive control network using diffusion-weighted magnetic resonance imaging (MRI) and
2048 functional MRI. J. Neurosci. Off. J. Soc. Neurosci. *27*, 3743–3752.

2049 Bach, D.R., and Dayan, P. (2017). Algorithms for survival: a comparative perspective on
2050 emotions. Nat. Rev. Neurosci. *18*, 311–319.

2051 Bean, B.P. (2007). The action potential in mammalian central neurons. Nat. Rev. Neurosci. *8*,
2052 451–465.

2053 Behrens, T.E.J., Woolrich, M.W., Walton, M.E., and Rushworth, M.F.S. (2007). Learning the
2054 value of information in an uncertain world. Nat. Neurosci. *10*, 1214–1221.

2055 Bernacchia, A., Seo, H., Lee, D., and Wang, X.-J. (2011). A reservoir of time constants for
2056 memory traces in cortical neurons. Nat. Neurosci. *14*, 366–372.

2057    Bernardi, S., Benna, M.K., Rigotti, M., Munuera, J., Fusi, S., and Salzman, C.D. (2020). The
2058    Geometry of Abstraction in the Hippocampus and Prefrontal Cortex. Cell *183*, 954-967.e21.

2059    Bonini, F., Burle, B., Liégeois-Chauvel, C., Régis, J., Chauvel, P., and Vidal, F. (2014). Action
2060    Monitoring and Medial Frontal Cortex: Leading Role of Supplementary Motor Area. Science
2061    *343*, 888–891.

2062    Botvinick, M.M., Braver, T.S., Barch, D.M., Carter, C.S., and Cohen, J.D. (2001). Conflict
2063    monitoring and cognitive control. Psychol. Rev. *108*, 624–652.

2064    Brainard, D.H. (1997). The Psychophysics Toolbox. Spat. Vis. *10*, 433–436.

2065    Braver, T.S. (2012). The variable nature of cognitive control: a dual mechanisms framework.
2066    Trends Cogn. Sci. *16*, 106–113.

2067    Bush, G., and Shin, L.M. (2006). The Multi-Source Interference Task: an fMRI task that reliably
2068    activates the cingulo-frontal-parietal cognitive/attention network. Nat. Protoc. *1*, 308–313.

2069    Buzsáki, G., Anastassiou, C.A., and Koch, C. (2012). The origin of extracellular fields and currents
2070    — EEG, ECoG, LFP and spikes. Nat. Rev. Neurosci. *13*, 407–420.

2071    Carter, C.S., Braver, T.S., Barch, D.M., Botvinick, M.M., Noll, D., and Cohen, J.D. (1998). Anterior
2072    Cingulate Cortex, Error Detection, and the Online Monitoring of Performance. Science *280*,
2073    747–749.

2074    Carter, C.S., Macdonald, A.M., Botvinick, M., Ross, L.L., Stenger, V.A., Noll, D., and Cohen, J.D.
2075    (2000). Parsing executive processes: Strategic vs. evaluative functions of the anterior cingulate
2076    cortex. Proc. Natl. Acad. Sci. *97*, 1944–1948.

2077    Cavanagh, J.F., Wiecki, T.V., Cohen, M.X., Figueroa, C.M., Samanta, J., Sherman, S.J., and Frank,
2078    M.J. (2011). Subthalamic nucleus stimulation reverses mediofrontal influence over decision
2079    threshold. Nat. Neurosci. *14*, 1462–1467.

2080    Cavanagh, S.E., Wallis, J.D., Kennerley, S.W., and Hunt, L.T. (2016). Autocorrelation structure at
2081    rest predicts value correlates of single neurons during reward-guided choice. ELife *5*, e18937.

2082    Cavanagh, S.E., Towers, J.P., Wallis, J.D., Hunt, L.T., and Kennerley, S.W. (2018). Reconciling
2083    persistent and dynamic hypotheses of working memory coding in prefrontal cortex. Nat.
2084    Commun. *9*, 3498.

2085    Chang, C.-C., and Lin, C.-J. (2011). LIBSVM: A library for support vector machines. ACM Trans.
2086    Intell. Syst. Technol. *2*, 27:1-27:27.

2087    Crone, E.A., Somsen, R.J.M., Beek, B.V., and Molen, M.W.V.D. (2004). Heart rate and skin
2088    conductance analysis of antecedents and consequences of decision making. Psychophysiology
2089    *41*, 531–540.

2090  Crosson, B., Sadek, J.R., Bobholz, J.A., Gökçay, D., Mohr, C.M., Leonard, C.M., Maron, L.,
2091  Auerbach, E.J., Browd, S.R., Freeman, A.J., et al. (1999). Activity in the Paracingulate and
2092  Cingulate Sulci during Word Generation: An fMRI Study of Functional Anatomy. Cereb. Cortex *9*,
2093  307–316.

2094  Danielmeier, C., Eichele, T., Forstmann, B.U., Tittgemeyer, M., and Ullsperger, M. (2011).
2095  Posterior Medial Frontal Cortex Activity Predicts Post-Error Adaptations in Task-Related Visual
2096  and Motor Areas. J. Neurosci. *31*, 1780–1789.

2097  Darlington, T.R., Beck, J.M., and Lisberger, S.G. (2018). Neural implementation of Bayesian
2098  inference in a sensorimotor behavior. Nat. Neurosci. *21*, 1442–1451.

2099  DiCarlo, J.J., and Cox, D.D. (2007). Untangling invariant object recognition. Trends Cogn. Sci. *11*,
2100  333–341.

2101  Dubois, J., de Berker, A.O., and Tsao, D.Y. (2015). Single-unit recordings in the macaque face
2102  patch system reveal limitations of fMRI MVPA. J. Neurosci. Off. J. Soc. Neurosci. *35*, 2791–2802.

2103  Duthoo, W., Abrahamse, E.L., Braem, S., Boehler, C.N., and Notebaert, W. (2014). The
2104  heterogeneous world of congruency sequence effects: an update. Front. Psychol. *5*.

2105  Ebitz, R.B., and Platt, M.L. (2015). Neuronal Activity in Primate Dorsal Anterior Cingulate Cortex
2106  Signals Task Conflict and Predicts Adjustments in Pupil-Linked Arousal. Neuron *85*, 628–640.

2107  Egner, T. (2007). Congruency sequence effects and cognitive control. Cogn. Affect. Behav.
2108  Neurosci. *7*, 380–390.

2109  Egner, T., and Hirsch, J. (2005). Cognitive control mechanisms resolve conflict through cortical
2110  amplification of task-relevant information. Nat. Neurosci. *8*, 1784–1790.

2111  Eldar, E., Rutledge, R.B., Dolan, R.J., and Niv, Y. (2016). Mood as Representation of Momentum.
2112  Trends Cogn. Sci. *20*, 15–24.

2113  Fan, J., Flombaum, J.I., McCandliss, B.D., Thomas, K.M., and Posner, M.I. (2003). Cognitive and
2114  Brain Consequences of Conflict. NeuroImage *18*, 42–57.

2115  Friston, K.J. (2002). Bayesian Estimation of Dynamical Systems: An Application to fMRI.
2116  NeuroImage *16*, 513–530.

2117  Fu, Z., Wu, D.-A.J., Ross, I., Chung, J.M., Mamelak, A.N., Adolphs, R., and Rutishauser, U. (2019).
2118  Single-Neuron Correlates of Error Monitoring and Post-Error Adjustments in Human Medial
2119  Frontal Cortex. Neuron *101*, 165-177.e5.

2120  Fusi, S., Miller, E.K., and Rigotti, M. (2016). Why neurons mix: high dimensionality for higher
2121  cognition. Curr. Opin. Neurobiol. *37*, 66–74.

2122   Gratton, G., Coles, M.G., and Donchin, E. (1992). Optimizing the use of information: strategic
2123   control of activation of responses. J. Exp. Psychol. Gen. *121*, 480–506.

2124   Hanes, D.P., Thompson, K.G., and Schall, J.D. (1995). Relationship of presaccadic activity in
2125   frontal eye field and supplementary eye field to saccade initiation in macaque: Poisson spike
2126   train analysis. Exp. Brain Res. *103*, 85–96.

2127   Hardstone, R., Poil, S.-S., Schiavone, G., Jansen, R., Nikulin, V.V., Mansvelder, H.D., and
2128   Linkenkaer-Hansen, K. (2012). Detrended Fluctuation Analysis: A Scale-Free View on Neuronal
2129   Oscillations. Front. Physiol. *3*.

2130   Heilbronner, S.R., and Hayden, B.Y. (2016). Dorsal Anterior Cingulate Cortex: A Bottom-Up
2131   View. Annu. Rev. Neurosci. *39*, 149–170.

2132   Heitz, R.P., and Schall, J.D. (2012). Neural Mechanisms of Speed-Accuracy Tradeoff. Neuron *76*,
2133   616–628.

2134   Herrera, B., Sajad, A., Woodman, G.F., Schall, J.D., and Riera, J.J. (2020). A Minimal Biophysical
2135   Model of Neocortical Pyramidal Cells: Implications for Frontal Cortex Microcircuitry and Field
2136   Potential Generation. J. Neurosci. *40*, 8513–8529.

2137   Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner,
2138   A. (2016). beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework.

2139   Higgins, I., Amos, D., Pfau, D., Racaniere, S., Matthey, L., Rezende, D., and Lerchner, A. (2018).
2140   Towards a Definition of Disentangled Representations. ArXiv181202230 Cs Stat.

2141   Ide, J.S., Shenoy, P., Yu, A.J., and Li, C.R. (2013). Bayesian Prediction and Evaluation in the
2142   Anterior Cingulate Cortex. J. Neurosci. *33*, 2039–2047.

2143   Ito, S., Stuphorn, V., Brown, J.W., and Schall, J.D. (2003). Performance Monitoring by the
2144   Anterior Cingulate Cortex During Saccade Countermanding. Science *302*, 120–122.

2145   Jiang, J., and Egner, T. (2014). Using Neural Pattern Classifiers to Quantify the Modularity of
2146   Conflict–Control Mechanisms in the Human Brain. Cereb. Cortex *24*, 1793–1805.

2147   Jiang, J., Heller, K., and Egner, T. (2014). Bayesian modeling of flexible cognitive control.
2148   Neurosci. Biobehav. Rev. *46*, 30–43.

2149   Jiang, J., Beck, J., Heller, K., and Egner, T. (2015). An insula-frontostriatal network mediates
2150   flexible cognitive control by adaptively predicting changing control demands. Nat. Commun. *6*,
2151   8165.

2152   Kawai, T., Yamada, H., Sato, N., Takada, M., and Matsumoto, M. (2019). Preferential
2153   Representation of Past Outcome Information and Future Choice Behavior by Putative Inhibitory

2154    Interneurons Rather Than Putative Pyramidal Neurons in the Primate Dorsal Anterior Cingulate
2155    Cortex. Cereb. Cortex *29*, 2339–2352.

2156    Kerns, J.G., Cohen, J.D., MacDonald, A.W., Cho, R.Y., Stenger, V.A., and Carter, C.S. (2004).
2157    Anterior Cingulate Conflict Monitoring and Adjustments in Control. Science *303*, 1023–1026.

2158    King, J.A., Korb, F.M., von Cramon, D.Y., and Ullsperger, M. (2010). Post-Error Behavioral
2159    Adjustments Are Facilitated by Activation and Suppression of Task-Relevant and Task-Irrelevant
2160    Information Processing. J. Neurosci. *30*, 12759–12769.

2161    Kobak, D., Brendel, W., Constantinidis, C., Feierstein, C.E., Kepecs, A., Mainen, Z.F., Qi, X.-L.,
2162    Romo, R., Uchida, N., and Machens, C.K. (2016). Demixed principal component analysis of
2163    neural population data. ELife *5*, e10989.

2164    de Leeuw, J.R. (2015). jsPsych: A JavaScript library for creating behavioral experiments in a Web
2165    browser. Behav. Res. Methods *47*, 1–12.

2166    Liu, X., Banich, M.T., Jacobson, B.L., and Tanabe, J.L. (2004). Common and distinct neural
2167    substrates of attentional control in an integrated Simon and spatial Stroop task as assessed by
2168    event-related fMRI. NeuroImage *22*, 1097–1106.

2169    Lo, C.-C., and Wang, X.-J. (2006). Cortico–basal ganglia circuit mechanism for a decision
2170    threshold in reaction time tasks. Nat. Neurosci. *9*, 956–963.

2171    Logan, G.D., and Zbrodoff, N.J. (1979). When it helps to be misled: Facilitative effects of
2172    increasing the frequency of conflicting stimuli in a Stroop-like task. Mem. Cognit. *7*, 166–174.

2173    Maris, E., and Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data. J.
2174    Neurosci. Methods *164*, 177–190.

2175    McDougle, S.D., Boggess, M.J., Crossley, M.J., Parvin, D., Ivry, R.B., and Taylor, J.A. (2016). Credit
2176    assignment in movement-dependent reinforcement learning. Proc. Natl. Acad. Sci. *113*, 6797–
2177    6802.

2178    Miller, E.K., and Cohen, J.D. (2001). An Integrative Theory of Prefrontal Cortex Function. Annu.
2179    Rev. Neurosci. *24*, 167–202.

2180    Minxha, J., Mosher, C., Morrow, J.K., Mamelak, A.N., Adolphs, R., Gothard, K.M., and
2181    Rutishauser, U. (2017). Fixations Gate Species-Specific Responses to Free Viewing of Faces in
2182    the Human and Macaque Amygdala. Cell Rep. *18*, 878–891.

2183    Minxha, J., Adolphs, R., Fusi, S., Mamelak, A.N., and Rutishauser, U. (2020). Flexible recruitment
2184    of memory-based choice representations by the human medial frontal cortex. Science *368*.

2185    Morales, J., Lau, H., and Fleming, S.M. (2018). Domain-General and Domain-Specific Patterns of
2186    Activity Supporting Metacognition in Human Prefrontal Cortex. J. Neurosci. *38*, 3534–3546.

2187   Mosher, C.P., Wei, Y., Kamiński, J., Nandi, A., Mamelak, A.N., Anastassiou, C.A., and Rutishauser,
2188   U. (2020). Cellular Classes in the Human Brain Revealed In Vivo by Heartbeat-Related
2189   Modulation of the Extracellular Action Potential Waveform. Cell Rep. *30*, 3536-3551.e6.

2190   Murphy, P.R., Boonstra, E., and Nieuwenhuis, S. (2016). Global gain modulation generates time-
2191   dependent urgency during perceptual choice in humans. Nat. Commun. *7*, 13526.

2192   Navarro, D.J., and Fuss, I.G. (2009). Fast and accurate calculations for first-passage times in
2193   Wiener diffusion models. J. Math. Psychol. *53*, 222–230.

2194   Niv, Y., Daw, N.D., Joel, D., and Dayan, P. (2007). Tonic dopamine: opportunity costs and the
2195   control of response vigor. Psychopharmacology (Berl.) *191*, 507–520.

2196   Pedersen, M.L., Frank, M.J., and Biele, G. (2017). The drift diffusion model as the choice rule in
2197   reinforcement learning. Psychon. Bull. Rev. *24*, 1234–1251.

2198   Peng, C.K., Buldyrev, S.V., Havlin, S., Simons, M., Stanley, H.E., and Goldberger, A.L. (1994).
2199   Mosaic organization of DNA nucleotides. Phys. Rev. E Stat. Phys. Plasmas Fluids Relat.
2200   Interdiscip. Top. *49*, 1685–1689.

2201   Pouget, P., Logan, G.D., Palmeri, T.J., Boucher, L., Paré, M., and Schall, J.D. (2011). Neural basis
2202   of adaptive response time adjustment during saccade countermanding. J. Neurosci. Off. J. Soc.
2203   Neurosci. *31*, 12604–12612.

2204   Purcell, B.A., and Kiani, R. (2016). Neural Mechanisms of Post-error Adjustments of Decision
2205   Policy in Parietal Cortex. Neuron *89*, 658–671.

2206   Rigotti, M., Barak, O., Warden, M.R., Wang, X.-J., Daw, N.D., Miller, E.K., and Fusi, S. (2013). The
2207   importance of mixed selectivity in complex cognitive tasks. Nature *497*, 585–590.

2208   Rutishauser, U., Schuman, E.M., and Mamelak, A.N. (2006). Online detection and sorting of
2209   extracellularly recorded action potentials in human medial temporal lobe recordings, in vivo. J.
2210   Neurosci. Methods *154*, 204–224.

2211   Sajad, A., Godlove, D.C., and Schall, J.D. (2019). Cortical microcircuitry of performance
2212   monitoring. Nat. Neurosci. *22*, 265–274.

2213   Sarafyazd, M., and Jazayeri, M. (2019). Hierarchical reasoning by neural circuits in the frontal
2214   cortex. Science *364*.

2215   Schmidt, J.R., and De Houwer, J. (2011). Now you see it, now you don't: Controlling for
2216   contingencies and stimulus repetitions eliminates the Gratton effect. Acta Psychol. (Amst.) *138*,
2217   176–186.

2218   Shackman, A.J., Salomons, T.V., Slagter, H.A., Fox, A.S., Winter, J.J., and Davidson, R.J. (2011).
2219   The integration of negative affect, pain and cognitive control in the cingulate cortex. Nat. Rev.
2220   Neurosci. *12*, 154–167.

2221   Shenhav, A., Botvinick, M.M., and Cohen, J.D. (2013). The Expected Value of Control: An
2222   Integrative Theory of Anterior Cingulate Cortex Function. Neuron *79*, 217–240.

2223   Sheth, S.A., Mian, M.K., Patel, S.R., Asaad, W.F., Williams, Z.M., Dougherty, D.D., Bush, G., and
2224   Eskandar, E.N. (2012). Human dorsal anterior cingulate cortex neurons mediate ongoing
2225   behavioural adaptation. Nature *488*, 218–221.

2226   Sohn, H., Narain, D., Meirhaeghe, N., and Jazayeri, M. (2019). Bayesian Computation through
2227   Cortical Latent Dynamics. Neuron *103*, 934-947.e5.

2228   Stokes, M.G., Kusunoki, M., Sigala, N., Nili, H., Gaffan, D., and Duncan, J. (2013). Dynamic
2229   Coding for Cognitive Control in Prefrontal Cortex. Neuron *78*, 364–375.

2230   Stroop, J.R. (1935). Studies of interference in serial verbal reactions. J. Exp. Psychol. *18*, 643–
2231   662.

2232   Stuphorn, V., Taylor, T.L., and Schall, J.D. (2000). Performance monitoring by the supplementary
2233   eye field. Nature *408*, 857–860.

2234   Tang, H., Yu, H.-Y., Chou, C.-C., Crone, N.E., Madsen, J.R., Anderson, W.S., and Kreiman, G.
2235   (2016). Cascade of neural processing orchestrates cognitive control in human frontal cortex.
2236   ELife *5*, e12352.

2237   Thura, D., and Cisek, P. (2017). The Basal Ganglia Do Not Select Reach Targets but Control the
2238   Urgency of Commitment. Neuron *95*, 1160-1170.e5.

2239   Tzelgov, J., Henik, A., and Berger, J. (1992). Controlling Stroop effects by manipulating
2240   expectations for color words. Mem. Cognit. *20*, 727–735.

2241   Ullsperger, M. (2017). Neural bases of performance monitoring. In The Wiley Handbook of
2242   Cognitive Control, (Wiley Blackwell), pp. 292–313.

2243   Ullsperger, M., Danielmeier, C., and Jocham, G. (2014). Neurophysiology of Performance
2244   Monitoring and Adaptive Behavior. Physiol. Rev. *94*, 35–79.

2245   Urai, A.E., de Gee, J.W., Tsetsos, K., and Donner, T.H. (2019). Choice history biases subsequent
2246   evidence accumulation. ELife *8*, e46331.

2247   Vigneswaran, G., Kraskov, A., and Lemon, R.N. (2011). Large Identified Pyramidal Cells in
2248   Macaque Motor and Premotor Cortex Exhibit "Thin Spikes": Implications for Cell Type
2249   Classification. J. Neurosci. *31*, 14235–14242.

2250   Vogt, B.A., Nimchinsky, E.A., Vogt, L.J., and Hof, P.R. (1995). Human cingulate cortex: Surface
2251   features, flat maps, and cytoarchitecture. J. Comp. Neurol. *359*, 490–506.

2252   Wang, X.-J. (2002). Probabilistic Decision Making by Slow Reverberation in Cortical Circuits.
2253   Neuron *36*, 955–968.

2254   Wang, J., Narain, D., Hosseini, E.A., and Jazayeri, M. (2018a). Flexible timing by temporal scaling
2255   of cortical responses. Nat. Neurosci. *21*, 102–110.

2256   Wang, S., Mamelak, A.N., Adolphs, R., and Rutishauser, U. (2018b). Encoding of Target
2257   Detection during Visual Search by Single Neurons in the Human Brain. Curr. Biol. *28*, 2058-
2258   2069.e4.

2259   Wessel, J.R., and Aron, A.R. (2017). On the Globality of Motor Suppression: Unexpected Events
2260   and Their Influence on Behavior and Cognition. Neuron *93*, 259–280.

2261   Wiecki, T.V., Sofer, I., and Frank, M.J. (2013). HDDM: Hierarchical Bayesian estimation of the
2262   Drift-Diffusion Model in Python. Front. Neuroinformatics *7*.

2263   Woodman, G.F. (2010). A brief introduction to the use of event-related potentials in studies of
2264   perception and attention. Atten. Percept. Psychophys. *72*, 2031–2046.

2265   Yeung, N., and Summerfield, C. (2012). Metacognition in human decision-making: confidence
2266   and error monitoring. Philos. Trans. R. Soc. B Biol. Sci. *367*, 1310–1321.

2267   Yu, A.J., and Cohen, J.D. (2008). Sequential effects: Superstition or rational behavior? Adv.
2268   Neural Inf. Process. Syst. *21*, 1873–1880.

2269