# The online metacognitive control of decisions

*Juliette Bénon[1], Douglas Lee[2,3], William Hopper[1], Morgan Verdeil[1], Mathias Pessiglione[1], Fabien Vinckier[1], Sebastien Bouret[1], Marion Rouault[1], Raphael Lebouc[1], Giovanni Pezzulo[3], Christiane Schreiweis[1], Eric Burguière[1], Jean Daunizeau[1]*

[1] Paris Brain Institute, Paris, France

[2] Tel-Aviv University, Tel-Aviv, Israel

[3] Institute of Cognitive Sciences and Technologies, Rome, Italy

Address for correspondence:

Jean Daunizeau

Motivation, Brain and Behavior Group

Paris Brain Institute (ICM)

47, bd de l'Hôpital, 75013, Paris, France.

Tel: +33 1 57 27 43 26

Mail: jean.daunizeau@icm-institute.org

**Abstract**

Difficult decisions typically involve mental effort, which scales with the deployment of cognitive (e.g., mnesic, attentional) resources engaged in processing decision-relevant information. But how does the brain regulate mental effort? A possibility is that the brain optimizes a resource allocation problem, whereby the amount of invested resources balances its expected cost (i.e. effort) and benefit. Our working assumption is that subjective decision confidence serves as the benefit term of the resource allocation problem, hence the "metacognitive" nature of decision control. Here, we present a computational model for the *online metacognitive control of decisions* or oMCD. Formally, oMCD is a Markov Decision Process that optimally solves the ensuing resource allocation problem under agnostic assumptions about the inner workings of the underlying decision system. We demonstrate how this makes oMCD a quasi-optimal control policy for a broad class of decision processes, including -but not limited to- *progressive attribute integration*. We disclose oMCD's main properties (in terms of choice, confidence and response time), and show that they reproduce most established empirical results in the field of value-based decision making. Finally, we discuss the possible connections between oMCD and most prominent neurocognitive theories about decision control and mental effort regulation.

41    **Introduction**

42

43    There is no such thing as a free lunch: obtaining reward typically requires investing effort. This

44    holds even for mental tasks, which may involve mental effort for achieving success (in terms

45    of, e.g., mnesic or attentional performance). Nevertheless, we sometimes invest very little

46    mental effort, eventually rushing decisions and falling for all sorts of cognitive biases [1]. So how

47    does the brain regulate mental effort? Recent theoretical neuroscience work proposes to view

48    mental effort regulation as a resource allocation problem: namely, identifying the amount of

49    cognitive resources that optimizes a cost/benefit tradeoff [2–4]. In this context, mental effort

50    signals the subjective cost of investing resources, the aversiveness of which is balanced by

51    the anticipated benefit. In conjunction with simple optimality principles, this idea has proven

52    fruitful for understanding the relationship between mental effort and peoples' performance in

53    various cognitive tasks, in particular those that involve cognitive control [5,6]. Recently, it was

54    adapted to the specific case of value-based decision making, and framed as a self-contained

55    computational model: the Metacognitive Control of Decisions or MCD [7].

56    The working assumption here is that decision confidence serves as the main benefit term of

57    the resource allocation problem [8,9], hence the "metacognitive" nature of decision control. On

58    the one hand, this formalizes the regulating role of confidence in decision making, which has

59    recently been empirically demonstrated in the context of perceptual evidence accumulation

60    [10,11]. On the other hand, this apparently contrasts with standard treatments of value-based

61    decision making, which insists on equating the benefit of value-based decisions with the value

62    of the chosen option [12–14]. This notion is a priori appealing, because the purpose of investing

63    resources   into   decisions   is   reducible   to   approaching   reward   and/or   avoiding

64    losses/punishments. Nevertheless, the benefit of such resource investments may be detached

65    from the subjective evaluation of alternative options [15]. This is partly because the brain attaches

66    subjective value to acquiring information about future rewards. In fact, this holds even when

67    this information cannot be used to influence decision outcomes [16–18]. Recall that, in Marr's

68    sense, any type of decision induces the same computational problem, i.e. the comparison of

69  alternative options. In this view, *evidence-based* and *value-based* decisions simply differ w.r.t.

70  to the underlying comparison criterion: the former relies on truthfulness judgments while the

71  latter involves idiosyncratic preferences [19]. Hence, in both cases, the benefit of allocating

72  resources to decisions is to raise the chance of identifying the best option, i.e. confidence. In

73  other words, if resource allocation aims at comparing alternative options, then decision

74  confidence can be viewed as a probe for goal achievement. This is essentially a simplifying

75  assumption, in the sense that it enables a unique computational architecture to control

76  resource allocations, irrespective of the nature of the underlying decision-relevant

77  computations.

78  In value-based decision making, confidence derives from the discriminability of uncertain value

79  representations, which evolve over decision time as the brain processes more value-relevant

80  information. Low confidence then induces a latent demand for mental effort: the brain refines

81  uncertain value representations by deploying cognitive resources, until they reach an optimal

82  confidence/effort trade-off. Interestingly, this mechanism was shown to explain the -otherwise

83  surprising- phenomenon of choice-induced preference change [7]. More importantly, the MCD

84  model makes quantitative out-of-sample predictions about many features of value-based

85  decisions, including decision time, subjective feeling of effort, choice confidence and changes

86  of mind. These predictions have already been tested -and validated- in a systematic manner,

87  using a dedicated behavioral paradigm (Lee and Daunizeau, 2021). Despite its remarkable

88  prediction accuracy, the original derivation of the model suffers from one main simplifying but

89  limiting approximation: it assumes that MCD operates in a purely *prospective* manner, i.e., the

90  MCD controller commits to a level of mental effort investment identified prior to the decision.

91  In principle, this early commitment would follow from anticipating the prospective benefit (in

92  terms of confidence gain) and cost of effort, given a prior or default representation of option

93  values that would rely on fast/automatic/effortless processes [20]. The issue here, is twofold.

94  First, it cannot explain variations in decision features (e.g., response time, choice confidence,

95  etc.) that occur in the absence of changes in default preferences. Second, it is somehow

96  suboptimal, as it neglects *reactive* processes, which enable the MCD controller to re-evaluate

97   – and improve on- the decision to stop or continue allocating resources, as new information is

98   processed and value representations are updated. The current work addresses these

99   limitations, effectively proposing an "online" variant of MCD which we coin oMCD.

100  As we will see, oMCD reduces to identifying the optimal policy for a specific instance of a

101  known class of stochastic control problems: namely, "optimal stopping" [21]. This kind of problem

102  can be solved using Markov Decision Processes or MDPs [22], under assumptions regarding the

103  (stochastic) dynamics of costs and/or benefits. Although less concerned with the notion of

104  mental effort, a similar MDP has already been derived for a specific type of "ideal" value-based

105  decisions [14,23,24]. The underlying assumption here is threefold: (i) the system that computes

106  option values is progressively "denoising" -in a Bayesian manner- its input value signals, (ii),

107  the system that monitors and controls the decision knows how the underlying value

108  computation system works, and (iii) the net benefit of decisions (i.e. the benefit discounted by

109  decision time) is the estimated reward rate. The ensuing MDP is very similar to so-called Drift-

110  Diffusion decision models [25,26], whereby the decision stops whenever the current estimate of

111  option value differences reaches a threshold. Interestingly, the authors show that the

112  assumptions (i), (ii) and (iii) imply that the optimal threshold is a decreasing function of time.

113  This is not innocuous, since this predicts that decision confidence necessarily decreases with

114  decision time, which is not always verified empirically [27]. In retrospect, these assumptions may

115  thus be deemed too restrictive. In this work, we intend to generalize this kind of approaches

116  by relaxing these three assumptions.

117  In particular, we will consider that the decision control system (i.e. the system that decides

118  when to stop deliberating) has only limited information regarding the inner workings of the

119  system that computes option values. We will show how decision confidence can serve both as

120  an efficient titration for the benefit of resource investments and as a shortcut summary statistic

121  for (hidden) value computations. That is, we will show that confidence monitoring is sufficient

122  to operate quasi-optimal decision control for a wide class of value-based decision processes.

123  We demonstrate the generalizability of the ensuing oMCD policy on two distinct decision

124     scenarios. In the above "*Bayesian value denoising*" case, it replicates existing MDPs and

125     extends their repertoire of confidence/RT relationships. We also consider the case of value

126     computation by *progressive attribute integration* [28–33]. As we will see, the latter scenario cannot

127     be reduced to the *Bayesian value denoising* case. This is because the main source of

128     uncertainty in value representations derive (as is the case for, e.g., forward planning) from the

129     arbitrary incompleteness of value computations. We demonstrate that, for both decision

130     scenarios, oMCD's control policy provides a close approximation to the ideal control policy,

131     which requires complete knowledge of the underlying value computations. We also identify

132     testable properties of oMCD control policies under both types of value computations, and show

133     that they are reminiscent of empirical value-based decisions.

134

135

**Methods**

As we will see below, deriving an optimal reactive variant of MCD requires specific mathematical developments, which falls under the frame of Markov decision processes [22]. But before we describe the oMCD model, let us first recall the prospective variant of MCD [7].

Note on ethics (see data re-analysis in the Results section): This work complies with all relevant ethical regulations and received formal approval from the INSERM Ethics Committee (CEEI-IRB00003888, decision no 16–333). All participants gave informed consent.

1.  The prospective MCD model

Note: this section is a summary of the mathematical derivation of the MCD model, which has already been published [7].

Let $z$ be the amount of cognitive (e.g., executive, mnemonic, or attentional) resources that serve to process value-relevant information. Allocating these resources will be associated with both a benefit $B(z)$, and a cost $C(z)$. As we will see, both are increasing functions of $z$: $B(z)$ derives from the refinement of internal representations of subjective values of alternative options or actions that compose the choice set, and $C(z)$ quantifies how aversive engaging cognitive resources is (mental effort). In line with the framework of *expected value of control* [2,4], we assume that the brain chooses to allocate the amount of resources $\hat{z}$ that optimizes the following cost-benefit trade-off:

$$\hat{z} = \arg \max_{z} E\left[ B(z) - C(z) \right] \tag{1}$$

where the expectation accounts for the anticipated impact of allocating resources into decision deliberation (this will be clarified below). Here, the benefit term is simply given by $B(z) = R \times P_c(z)$, where $P_c(z)$ is choice confidence and its weight $R$ quantifies the

161   importance of making a confident decision. As we will see, $P_c(z)$ plays a pivotal role in the

162   model, in that it captures the efficacy of allocating resources for processing value-relevant

163   information. So, how do we define choice confidence?

164   We assume that the subjective evaluation of alternative options in the choice set is uncertain.

165   In other words, the internal representations of values of alternative options are probabilistic.

166   Such a probabilistic representation of value can be understood in terms of, for example, an

167   uncertain prediction regarding the to-be-experienced value of a given option. In what follows,

168   the probabilistic representation of option value $V_i$ takes the form of Gaussian probability

169   density functions $p(V_i) = N(\mu_i, \sigma_i)$, where $\mu_i$ and $\sigma_i$ are the mode and the variance of the

170   probabilistic value representation, respectively (and $i$ indexes alternative options in the choice

171   set). This allows us to define choice confidence $P_c$ as the probability that the (predicted)

172   experienced value of the (to be) chosen item is higher than that of the (to be) unchosen item.

173   When the choice set is composed of two alternatives, $P_c$ is given by:

$$P_c \approx s\left( \frac{\pi |\Delta\mu|}{\sqrt{3(\sigma_1 + \sigma_2)}} \right) \qquad (2)$$

175   where $s(x) = 1/1 + e^{-x}$ is the standard sigmoid mapping, and we assume that the choice

176   follows the sign of the preference $\Delta\mu = \mu_1 - \mu_2$. Equation (2) simply derives from a moment-

177   matching approximation to the Gaussian cumulative density function [34]. Note that Equation 2

178   implicitly assumes that the option with the highest value estimate is chosen. This satisfies the

179   same formal criteria as for choice confidence in the context of evidence-based decisions [35].

180   We assume that the brain valuation system may, in some contexts, automatically generate

181   uncertain estimates of options' value [36,37], before cognitive effort is invested in decision making.

182   In what follows, $\mu_i^0$ and $\sigma_i^0$ are the mode and variance of the ensuing prior value

183   representations. They yield an initial confidence level $P_c^0$. Importantly, this prior or default

184    preference neglects existing value-relevant information that would require cognitive effort to

185    be retrieved and processed [20].

186    Now, how can a decision control system anticipate the benefit of allocating resources to the

187    decision process without knowing the details of the underlying value computations? Recall that

188    the purpose of allocating resources is to process (yet unavailable) value-relevant information.

189    The critical issue is thus to predict how both the uncertainty $\sigma_i$ and the modes $\mu_i$ of value

190    representations will eventually change, before having actually allocated the resources (i.e.,

191    without having processed the information). In brief, allocating resources essentially has two

192    impacts: (i) it decreases the uncertainty $\sigma_i$, and (ii) it perturbs the modes $\mu_i$ in a stochastic

193    manner.

194    The former impact (i) derives from assuming that the amount of information that will be

195    processed increases with the amount of allocated resources. This implies that the precision

196    $1/\sigma_i(z)$ of a given probabilistic value representation necessarily increases with the amount

197    of allocated resources, i.e.:

198    $$1/\sigma_i(z) = 1/\sigma_i^0 + \beta z \tag{3}$$

199    where $1/\sigma_i^0$ is the prior precision of the representation (before any effort has been allocated),

200    and $\beta$ controls the efficacy with which resources increase the precision of the value

201    representation. More precisely, $\beta$ is the precision increase that follows from allocating a

202    unitary amount of resources $z$. In what follows, we will refer to $\beta$ as "*type #1 effort efficacy*".

203    Note that if $\beta = 0$, then mental effort brings no improvement in the precision of value

204    representations.

205    The latter impact (ii) follows from acknowledging the fact that the control system cannot know

206    how processing more value-relevant information will affect its preference before having

207    allocated the corresponding resources. Let $\delta_i$ be the change in the position of the mode of

208      the $i$ th value representation, having allocated an amount $z$ of resources. The direction of the

209      mode's perturbation $\delta_i$ cannot be predicted because it is tied to the information that is yet to

210      be processed. However, a tenable assumption is to consider that the magnitude of the

211      perturbation increases with the amount of information that will be processed. This reduces to

212      stating that the variance of $\delta_i$ increases with $z$, i.e.:

213

$$\mu_i(z) = \mu_i^0 + \delta_i$$
$$\delta_i \sim N(0, \gamma z)$$

(4)

214      where $\mu_i^0$ is the mode of the value representation before any effort has been allocated, and

215      $\gamma$ controls the relationship between the amount of allocated resources and the variance of the

216      perturbation term $\delta$. The higher $\gamma$, the greater the expected perturbation of the mode for a

217      given amount of allocated resources. In what follows, we will refer to $\gamma$ as "*type #2 effort*

218      *efficacy*". Note that Equation 4 treats the impact of future information processing as some form

219      of random perturbation on the mode of the prior value representation. Importantly, Equation 4

220      is not specific to the type of value computations that eventually perturbs the value modes. Our

221      justification for this assumption is twofold: it is simple, and it captures the idea that the MCD

222      controller is agnostic about how the allocated resources will be used by the underlying

223      valuation/decision system. We will see that, in spite of this, the MCD controller can still make

224      quasi-optimal predictions regarding the expected benefit of allocating resources, under very

225      different value computation schemes.

226      Now, predicting the net effect of resource investment onto choice confidence (from Equations

227      (3) and (4)) is not entirely trivial. On the one hand, allocating effort will increase the precision

228      of value representations, which mechanically increases choice confidence, all other things

229      being equal. On the other hand, allocating effort can either increase or decrease the absolute

230      difference $|\Delta\mu(z)|$ between the modes (and hence increase or decrease choice confidence).

231      This depends upon the direction of the perturbation term $\delta$, which is a priori unknown. Having

232    said this, it is possible to derive the *expected* absolute mode difference (as well as its variance)

233    that would follow from allocating an amount $z$ of resources:

234
$$\begin{cases} E\Big[\big|\Delta\mu(z)\big|\Big] = 2\sqrt{\dfrac{\gamma z}{\pi}}\exp\left(-\dfrac{\big|\Delta\mu^0\big|^2}{4\gamma z}\right) + \Delta\mu^0\left(2\times s\left(\dfrac{\pi\,\Delta\mu^0}{\sqrt{6\gamma z}}\right)-1\right) \\[4mm] V\Big[\big|\Delta\mu(z)\big|\Big] = 2\gamma z + \big|\Delta\mu^0\big|^2 - E\Big[\big|\Delta\mu(z)\big|\Big]^2 \end{cases}$$
(5)

235    where we have used the expression for the first-order moment of the so-called "folded normal

236    distribution". Importantly, $E\Big[\big|\Delta\mu(z)\big|\Big]$ is always greater than $\big|\Delta\mu^0\big|$ and increases

237    monotonically with $z$ - as is $V\Big[\big|\Delta\mu(z)\big|\Big]$. In other words, allocating resources is expected

238    to increase the value difference, even though the impact of the perturbation term can go either

239    way.

240    Equation 5 now enables us to derive the expected confidence level $\bar{P}_c(z) \square E[P_c]$ that

241    would result from allocating the amount of resource $z$ :

242
$$\bar{P}_c(z) \approx s\left(\frac{\lambda E\Big[\big|\Delta\mu(z)\big|\Big]}{\sqrt{1+\dfrac{1}{2}\left(\lambda^2 V\Big[\big|\Delta\mu(z)\big|\Big]\right)^{\frac{3}{4}}}}\right)$$
(6)

243    where $\lambda = 1\Big/\sqrt{3\big(\sigma_1(z)+\sigma_2(z)\big)}$. Of course, $\bar{P}_c(0) = P_c^0$, i.e., investing no resources yields no

244    confidence gain. Moreover, the expected choice confidence $\bar{P}_c(z)$ always increase with $z$ ,

245    irrespective of the efficacy parameters, as long as $\beta \neq 0$ or $\gamma \neq 0$. Equation 6 is important,

246    because it quantifies the expected benefit of resource allocation, before having processed the

247    ensuing value-relevant information.

248    To complete the cost-benefit model, we simply assume that the cost of allocating resources to

249    the decision process increases monotonically with the amount of resources, i.e.:

250    $C(z) = \alpha z^{\nu}$            (7)

251    where $\alpha$ determines the effort cost of allocating a unitary amount of resources $z$ (we refer to

252    $\alpha$ as the "unitary effort cost"), and $\nu$ effectively controls the range of resource investments

253    that result in noticeable cost variations (we refer to $\nu$ as the "cost power").

254    Finally, the MCD-optimal resource allocation $\hat{z}$ is identified by replacing Equations (5), (6) and

255    (7) into Equation (1). This can be done before any resource has been invested, hence the

256    *prospective* nature of metacognitive control, here.

257

258

259       2.   Online MCD: optimal control policy

260    We now augment this model, by assuming that the MCD controller re-evaluates the decision

261    to stop or continue allocating resources, as value representations are being updated and online

262    confidence is changing. This makes the ensuing *oMCD* model a *reactive* extension of the

263    above "purely prospective" MCD model, which relieves the system from the constraint of effort

264    investment pre-commitment.

265    Let $t$ be the current time within a decision. For simplicity, we assume that there is a linear

266    relationship between deliberation time and resource investment, i.e.: $z = \kappa t$, where $\kappa$ is the

267    amount of resources that is spent per unit of time. We refer to $\kappa$ as "effort intensity". By

268    convention, the maximal decision time $T$ (the so-called *temporal horizon*) corresponds to the

269    exhaustion of all available resources. This implies that $T = 1/\kappa$ because we consider

270    normalized resources amounts.

271    Now, at time $t$, the system holds probabilistic value representations with modes $\mu(t)$ and

272    variance $\sigma(t)$. This yields the confidence level $P_c(\Delta\mu(t))$ given in Equation 2 above, where

273    we have made confidence an explicit function of $\Delta\mu(t)$ for mathematical convenience (see

274    below).

275 This confidence level can be greater or smaller than the initial confidence level $P_c^0$, because

276 new information regarding option values has been assimilated since the start of the

277 deliberation. Of course, the system will anticipate that investing additional resources will

278 increase its confidence (on average). But this may not always overcompensate the cost of

279 spending more resources on the decision. Thus, how should the system determine whether to

280 stop or to continue, in order to maximize the expected cost-benefit tradeoff? It turns out that

281 this problem is one of *optimal stopping*, which is a special case of Markov Decision Processes

282 [22,38]. As we will see, it can be solved recursively (backward in time) using Bellman's optimality

283 principle [39].

284 Let $a(t) \in \{0,1\}$ be the action that is taken at time $t$, where $a(t)=0$ (resp. $a(t)=1$) means

285 that the system stops (resp. continues) deliberating. Let $Q(a(t), \Delta\mu(t))$ be the net benefit

286 that the decision system would obtain at time $t$:

287
$$Q(a(t), \Delta\mu(t)) = \begin{cases} \underbrace{R \times P_c(\Delta\mu(t))}_{B(z)} - \underbrace{\alpha(\kappa t)^\nu}_{C(z)} & \text{if } a(t)=0 \\ 0 & \text{otherwise} \end{cases} \tag{8}$$

288 where both the benefits $B(z)$ and costs $C(z)$ of resource investments have been rewritten

289 in terms of decision time. Without loss of generality, Equation 9 states that the net benefit of

290 resource allocation is only realized when the system decides to stop $(a(t)=0)$. Note that

291 $Q(a(t), \Delta\mu(t))$ is also a function of time (through the precision of value representations and

292 effort cost), but we have ignored this dependency for the sake of notational conciseness.

293 A time $t$, the optimal control policy derives from a comparison between the net benefit of

294 stopping now - i.e., $Q(0, \Delta\mu(t))$ - and some -yet undefined- threshold $\omega(t)$, which may

295 depend upon time. Let $\pi_\omega(t)$ be the control policy (i.e., the temporal sequence of continue/stop

296 decisions) that is induced by the threshold $\omega(t)$:

297  $$\pi_\omega(t) = \begin{cases} 0 \text{ if } Q(0, \Delta\mu(t)) \geq \omega(t) \\ 1 \text{ otherwise} \end{cases} \tag{9}$$

298  Finding the optimal control policy $\pi_\omega^*(t)$ thus reduces to finding the optimal threshold $\omega^*(t)$.

299  By definition, at $t = T$, the system stops deliberating irrespective of its current net benefit

300  $Q(0, \Delta\mu(T))$. By convention, the optimal threshold $\omega^*(T)$ can thus be written as:

301
$$\begin{aligned} \omega^*(T) &= \min_{\Delta\mu(T)} Q(0, \Delta\mu(T)) \\ &= Q(0, 0, T) \\ &= R/2 - \alpha(\kappa T)^\nu \end{aligned} \tag{10}$$

302  Now, at $t = T - 1$, the net benefit $Q(0, \Delta\mu(T-1))$ of stopping now can be compared to the

303  expected net benefit $E\left[Q(0, \Delta\mu(T)) \middle| \Delta\mu(T-1)\right]$ of stopping at time $t = T$, conditional on the

304  current value mode difference $\Delta\mu(T-1)$:

305  $$E\left[Q(0, \Delta\mu(T)) \middle| \Delta\mu(T-1)\right] = R \times E\left[P_c(\Delta\mu(T)) \middle| \Delta\mu(T-1)\right] - \alpha(\kappa T)^\nu \tag{11}$$

306  where the expectation is taken under the transition probability density $p(\Delta\mu(T) | \Delta\mu(T-1))$

307  of the value mode difference for a unitary time increment $(\Delta t = 1 \Leftrightarrow \Delta z = \kappa)$. This density

308  derives from rewriting Equation 4 in terms of the instantaneous change in the moments of the

309  value representations. It is trivial to show that the corresponding first- and second-order

310  moments are $E\left[\mu_i(t) - \mu_i(t-1)\right] = 0$ and $E\left[\left(\mu_i(t) - \mu_i(t-1)\right)^2\right] = \gamma\kappa$, respectively. It

311  follows that the transition probability density of the value mode difference is stationary (i.e. it

312  does not depend upon time) and is given by:

313  $$p(\Delta\mu(t) | \Delta\mu(t-1)) = N(\Delta\mu(t-1), 2\gamma\kappa) \quad \forall t > 1 \tag{12}$$

314  which is of course valid for $t = T$.

315    The optimal policy is to stop if $Q\big(0,\Delta\mu(T-1)\big)\geq E\big[Q\big(0,\Delta\mu(T)\big)\big|\Delta\mu(T-1)\big]$, and to continue

316    otherwise. Note that both $Q\big(0,\Delta\mu(T-1)\big)$ and $E\big[Q\big(0,\Delta\mu(T)\big)\big|\Delta\mu(T-1)\big]$ are deterministic

317    functions of $\Delta\mu(T-1)$. More precisely, they are both monotonically increasing with $\Delta\mu(T-1)$

318    (see Figure 1 below), because current confidence and expected future confidence

319    monotonically increase with $\Delta\mu(T-1)$. Critically, these functions have a different offset, i.e.:

320    $Q(0,0) < E\big[Q\big(0,\Delta\mu(T)\big)\big|\Delta\mu(T-1)=0\big]$ as long as $\gamma > 0$. In addition, they eventually reach

321    a different plateau, i.e.: $\lim\limits_{\Delta\mu(T-1)\to\infty} Q\big(0,\Delta\mu(T-1)\big) > \lim\limits_{\Delta\mu(T-1)\to\infty} E\big[Q\big(0,\Delta\mu(T-1)\big)\big|\Delta\mu(T-1)\big]$ as

322    long as $\alpha > 0$. This is important, because this implies that there exists a critical value mode

323    difference $\Delta\mu^{*}(T-1)$ such that $Q\big(0,\Delta\mu^{*}(T-1)\big) = E\big[Q\big(0,\Delta\mu(T)\big)\big|\Delta\mu^{*}(T-1)\big]$. The net

324    benefit at that critical point is the optimal threshold at $t=T-1$, i.e.:

325    $\omega^{*}(T-1) = Q\big(0,\Delta\mu^{*}(T-1)\big)$. This is exemplified in Figure 1 below.

326

327    Now, let us move one step backward in time, at $t=T-2$. Here again, the optimal policy is to

328    stop if the current net benefit $Q\big(0,\Delta\mu(T-2)\big)$ is higher than the expected future net benefit

329    $E\big[Q\big(a(T-1),\Delta\mu(T-1)\big)\big|\Delta\mu(T-2)\big]$, conditional on $\Delta\mu(T-2)$. However, the latter now

330    depends upon $a(T-1)$, i.e., whether the system will later decide to stop or to continue:

331    $E\big[Q\big(a(T-1),\Delta\mu(T-1)\big)\big|\Delta\mu(T-2)\big] = \begin{cases} E\big[Q\big(0,\Delta\mu(T-1)\big)\big|\Delta\mu(T-2)\big] & \text{if } a(T-1)=0 \\ E\big[E\big[Q\big(0,\Delta\mu(T)\big)\big|\Delta\mu(T-1)\big]\big|\Delta\mu(T-2)\big] & \text{otherwise} \end{cases}$

332                                                                                                (13)

333    The optimal control policy cannot be directly identified from Equation 13. This is where we

334    resort to Bellman's optimality principle: namely, whatever the current state and action are, the

335    remaining actions of an optimal policy must also constitute an optimal policy with regard to the

336    state resulting from the current action [39]. Practically speaking, the derivation of the optimal

337    threshold at $t = T - 2$ is done under the constraint that oMCD's next action follows the optimal

338    policy, i.e., $a(T-1) = \pi_\omega^*(T-1)$.

339    Let $Q^*(\Delta\mu(t)) \equiv Q(\pi_\omega^*(t), \Delta\mu(t))$ be the net benefit evaluated under the optimal policy at

340    time $t$, which we refer to as the "optimal net benefit". Under Bellman's optimality principle, the

341    optimal policy at $t = T - 2$ is to stop if the current net benefit $Q(0, \Delta\mu(T-2))$ is higher than

342    the expected optimal net benefit $E\left[Q^*(\Delta\mu(T-1)) \big| \Delta\mu(T-2)\right]$, where the expectation is

343    again taken under the transition probability density in Equation 12.

344    Now, at time $t = T - 1$, the optimal net benefit is given by:

345    $$Q^*(\Delta\mu(T-1)) \Box \max\left\{Q(0, \Delta\mu(T-1)), E\left[Q(0, \Delta\mu(T)) \big| \Delta\mu(T-1)\right]\right\} \qquad (14)$$

346    Note that $Q^*(\Delta\mu(T-1))$ is just another function of $\Delta\mu(T-1)$ (cf. dotted green curve in Figure

347    1). This means that the only source of stochasticity in $Q^*(\Delta\mu(T-1))$ comes from $\Delta\mu(T-1)$,

348    which can nonetheless be predicted (with some uncertainty), given the current value mode

349    difference $\Delta\mu(T-2)$. In turn, this makes the expected optimal net benefit

350    $E\left[Q^*(\Delta\mu(T-1)) \big| \Delta\mu(T-2)\right]$ a deterministic function of $\Delta\mu(T-2)$. Again, as long as $\gamma > 0$

351    and $\alpha > 0$, there exists a critical value mode difference $\Delta\mu^*(T-2)$ such that

352    $Q(0, \Delta\mu^*(T-2)) = E\left[Q^*(\Delta\mu(T-1)) \big| \Delta\mu^*(T-2)\right]$. The net benefit at that critical point is the

353    optimal threshold $\omega^*(T-2)$ at $t = T - 2$.

354    In fact, the reasoning is the same for all times $t < T - 1$:

355    First, the expected optimal net benefit obeys the following backward recurrence relationship

356    (Bellman equation for all $t < T - 1$):

357 $$E\left[Q^*\left(\Delta\mu(t)\right)\middle|\Delta\mu(t-1)\right] = E\left[\max\left\{Q\left(0,\Delta\mu(t)\right), E\left[Q^*\left(\Delta\mu(t+1)\right)\middle|\Delta\mu(t)\right]\right\}\middle|\Delta\mu(t-1)\right]$$

358 $$(15)$$

359 This equation is solved recursively backward in time, starting at the expected net benefit at

360 $t = T-1$, as given in Equation 11. Both expectations in Equation 15 are taken under the

361 transition probability density $p\left(\Delta\mu(t)\middle|\Delta\mu(t-1)\right)$ of the value mode difference under a unitary

362 resource investment (cf. Equation 12).

363 Second, the optimal threshold at time $t$ is given by:

364 $$\omega^*(t) = Q\left(0,\Delta\mu^*(t)\right) \tag{16}$$

365 where $\Delta\mu^*(t)$ is the critical value mode difference, i.e., $\Delta\mu^*(t)$ is such that:

366 $$Q\left(0,\Delta\mu^*(t)\right) = E\left[Q^*\left(\Delta\mu(t+1)\right)\middle|\Delta\mu(t) = \Delta\mu^*(t)\right] \tag{17}$$

367 Since the net benefit is a deterministic function of decision confidence, the oMCD-optimal

368 threshold $\omega^*(t)$ for net benefits can be transformed into an oMCD-optimal confidence

369 threshold $\omega_P^*(t)$. Replacing the net benefit with the optimal threshold $\omega^*(t)$ and confidence

370 with $\omega_P^*(t)$ in Equation 9 yields:

371 $$\omega_P^*(t) = \frac{\omega^*(t) + \alpha(\kappa t)^\nu}{R} \tag{18}$$

372 At any point in time, comparing the net benefit $Q\left(0,\Delta\mu(t)\right)$ of resource allocation to $\omega^*(t)$ is

373 exactly equivalent to comparing the current confidence level $P_c(t)$ to $\omega_P^*(t)$. In other terms,

374 the optimal control policy (cf. Equation 10) can be rewritten as:

375 $$\pi_\omega^*(t) = \begin{cases} 0 \text{ if } P_c(t) \geq \omega_P^*(t) \\ 1 \text{ otherwise} \end{cases} \tag{19}$$

376 This highlights the central role of confidence, whose monitoring (during deliberation) is a

377 sufficient condition for operating optimal decision control. In turn, this greatly simplifies the

378     decision control architecture because knowledge about the underlying decision-relevant

379     computations is not required. As we will see later, oMCD is flexible (i.e. it encompasses many

380     kinds of decision processes) and robust to deviations from its working assumptions (i.e. it

381     provides a tight approximation to optimal control under alternative settings of the resource

382     allocation problem).

383     This closes the derivation of oMCD's optimal control policy.

384

385     Although the derivation of oMCD's optimal control policy is agnostic w.r.t. the underlying value

386     computations, it still requires some prior information regarding the upcoming information

387     processing: namely, prior moments of value representations, type #1 and #2 effort efficacies,

388     decision importance, unitary effort cost and cost power. This means that oMCD implicitly

389     includes a *prospective* component, which is used to decide how to optimally *react* to a

390     particular (stochastic) internal state of confidence. In other terms, one can think of oMCD as a

391     mixed prospective/reactive policy, whose prospective component is the shape of the

392     confidence threshold temporal dynamics.

393     Figure 2 below shows a representative instance of oMCD's optimal control policy, from 1000

394     Monte-Carlo simulations (using decision parameters $R=1$, $\alpha=0.2$, $\beta=1$, $\gamma=4$, $\kappa=1/100$, $\upsilon=0.5$, $\sigma_0=1$).

395

396     First, one can see that oMCD's optimal confidence threshold $\omega_P^*(t)$ lies above the average

397     confidence level $\bar{P}_c(t)$ of its prospective variant (cf. Equation 6, whose Monte-Carlo estimate

398     is depicted by the blue line in panel B). This means that oMCD's control policy would, in most

399     cases, demand higher confidence than prospective MCD. Importantly however, oMCD's policy

400     is sensitive to unpredictable fluctuations in the trajectory of value modes, which will induce

401     variations in resource investments (or, equivalently, response times). This enables oMCD to

402     exploit favorable variations in confidence if they eventually reach the threshold sooner than

403     expected.

404    Note that the confidence threshold $\omega_P^*(t)$ is, by construction, the confidence level that the

405    system achieves when committing to its decision. This means that, under oMCD's policy, the

406    relationship between reported confidence levels and response times is entirely determined by

407    the shape of the optimal threshold dynamics. In this example, this relationship will be mostly

408    negative, i.e. reported confidence levels tend to decrease when response times increase. This

409    is despite the fact that average confidence $\bar{P}_c(t)$ always increases as decision time unfolds,

410    as long as effort efficacy parameters are nonzero. In other words, the overt relationship

411    between response times and reported confidence levels (across trials) may be qualitatively

412    different from the covert temporal dynamics of confidence during decision deliberation.

413    So what is the impact of decision parameter on oMCD's confidence threshold dynamics? This

414    is summarized in Figure 3 below, where we systematically vary each parameter in turn (when

415    setting all the others to unity).

416

417    The net effect of increasing effort efficacy (either type #1 or type #2) is to increase the absolute

418    confidence threshold. In other terms, the demand for confidence increases with effort efficacy.

419    In contrast, the demand for confidence decreases with unitary effort cost. Note that the effect

420    of increasing decision importance (not shown) is exactly the same as that of decreasing unitary

421    effort cost. Importantly, the shape of the confidence threshold dynamics is approximately

422    invariant to changes in effort efficacy or unitary effort cost.

423    The only parameter that eventually changes the qualitative dynamics of oMCD's optimal

424    confidence threshold is the effort cost power (panel D). In brief, increasing the cost power

425    tends to decrease the initial slope of oMCD's confidence threshold dynamics. Here, the latter

426    eventually falls below zero (i.e., the confidence threshold decreases with decision time) when

427    the effort cost becomes superlinear ($v>1$). This is because, in this case, late resource

428    investments are comparatively more costly than early ones.

429    Note that, in contrast to effort efficacies, effort cost parameters can be altered without changing

430    the dynamics of expected confidence. In other terms, the shape of the relationship between

431 decision time and confidence is, for the most part, independent from the inner workings of the

432 underlying decision system.

433 Let us now relate the MCD framework to standard decision processes, which differ in terms of

434 their respective value computations.

435

436     3. How does MCD relate to standard decision processes?

437 By itself, the MCD framework does not commit to any specific assumption regarding how value-

438 relevant information is processed. Nevertheless, the properties of decisions that are controlled

439 through MCD actually depend upon how probabilistic value representations change over time.

440 In what follows, we focus on two specific scenarios of value computations, and disclose their

441 connection with MCD.

442

443     • Bayesian value denoising.

444 Let us first consider the *Bayesian value denoising* case, in which value representations are

445 updated Bayesian beliefs on a hidden value signal. Note that, in this case, the optimal control

446 rule - for maximizing expected reward rate - reduces to a specific instance of so-called *drift-*

447 *diffusion decision* models with decaying bounds on the estimated value difference [14,24].

448 Assume that, at each time point, the decision system receives an unreliable copy $y(t)$ of the

449 (hidden) value $V$ of each alternative option. More precisely, $y(t)$ is a noisy input signal that

450 is centered on $V$, i.e.: $y(t) = V + \varepsilon(t)$, where the random noise term $\varepsilon(t)$ is i.i.d. Gaussian

451 with zero mean and variance $\Sigma$ (and we have dropped the option indexing for notational

452 simplicity). One may think of $\Sigma$ as measuring the (lack of) reliability of the input value signal.

453 This induces the following likelihood function for the hidden value: $p(y(t)|V) = N(V, \Sigma)$.

454 Finally, assume that the decision system holds a Gaussian prior belief about the hidden

455    options' value, i.e.: $p(V) = N(\mu_0, \sigma_0)$, where $\mu_0$ and $\sigma_0$ are the corresponding prior mean

456    and variance. At time $t$, a Bayesian observer would assimilate the series of noisy signals to

457    derive a probabilistic (posterior) representation $p(V|y(1),...,y(t)) = N(\mu(t), \sigma(t))$ of

458    hidden options' values with the following mean and variance [40]:

459   
$$\begin{cases} \mu(t) = \mu_0 + \tilde{\delta}(t) \\ \sigma(t) = \dfrac{1}{\dfrac{1}{\sigma_0} + t \times \dfrac{1}{\Sigma}} \end{cases} \qquad (20)$$

460    where the perturbation $\tilde{\delta}$ of the value mode is given by:

461   
$$\tilde{\delta}(t) = \frac{1}{\dfrac{\Sigma}{\sigma_0} + t} \sum_{t'=1}^{t} (y(t) - \mu_0) \qquad (21)$$

462    Equation 21 specifies what the perturbation to the value mode would be, if the underlying value

463    computation was a process of *Bayesian value denoising*, whose outcome is the posterior

464    estimate $\mu(t) = E[V|y(1),...,y(t)]$ of value. In brief, Equation 21 states that the value mode

465    changes in proportion to prediction errors (i.e., $y(t) - \mu_0$), which the Bayesian observer

466    accumulates while sampling more input value signals. The stochasticity of the value mode's

467    perturbation $\tilde{\delta}$ is driven by the random noise term $\varepsilon$ in the incoming noisy value signal.

468    Conditioned on the hidden value $V$, it is easy to show that $E[\tilde{\delta}|V] \propto V - \mu_0$. This implies

469    that the random walk in Equation 21 actually has a nonzero drift that is proportional to the

470    hidden value. Importantly however, the Bayesian observer does not know what the hidden

471    value $V$ is. Prior to observing noisy value signals, its expectation is simply that

472    $E[y] = E[V] = \mu_0$ and therefore $E[\tilde{\delta}] = 0$. In fact, this holds true at any time $t$: the Bayesian

473    observer's expectation about the future change in its value belief mode, i.e.

474    $E\left[\mu(t+1)-\mu(t)\middle|y(1),...,y(t)\right]$, is always zero, because its expectation about the next

475    value signal reduces to her current value mode, i.e. $E\left[y(t+1)\middle|y(1),...,y(t)\right]=\mu(t)$. In other

476    words, although the modes' perturbation $\tilde{\delta}$ actually have a nonzero mean (as long as $V$

477    deviates from the mode of the observer's belief), the Bayesian observer's expectation about

478    its future realizations is always zero.

479    Nevertheless, the Bayesian observer can accurately predict how the precision of its belief will

480    change with time. Comparing Equations 3 and 20 suggests that, under the *Bayesian value*

481    *denoising* scenario, type #1 effort efficacy reduces to: $\beta=1/\kappa\Sigma$. This means that type #1 effort

482    efficacy simply increases with the reliability of the input value signal.

483    In addition, although the Bayesian observer cannot anticipate in what direction the to-be-

484    sampled signal $y(t)$ will modify the mode of its posterior belief, it can derive a prediction over

485    the magnitude of the perturbation:

486   
$$E\left[\tilde{\delta}(t)^2\right]=t\times\frac{\Sigma+t\sigma_0}{\left(\dfrac{\Sigma}{\sigma_0}+t\right)^2} \tag{22}$$

487    where the expectation is derived under the agent's prior belief about the hidden value. Now,

488    Equation 4 defines type #2 effort efficacy in terms of the ratio $E\left[\tilde{\delta}(t)^2\right]/\kappa t$ of expected change

489    magnitude over effort investment (where $z=\kappa t$). Note that, under Equation 22, this quantity

490    varies as a function of decision time. Thus, under the *Bayesian value denoising* scenario, type

491    #2 effort efficacy can be approximated as its sample average over all admissible decision

492    times, i.e.: $\gamma\approx 1/T\sum_{t=1}^{T}(\Sigma+t\sigma_0)/(\Sigma/\sigma_0+t)^2\kappa$. This is only an approximation of course, since

493    $E\left[\tilde{\delta}(t)^2\right]$ eventually tails off as time increases, because noisy value signals that are sampled

494    later in time have a smaller effect on the posterior mode. In other words, were the MCD

495     controller to know about the inner computations of the underlying value updating system, it

496     would rely on Equation 22 rather than on Equation 4. The ensuing ideal control policy is

497     summarized in the Supplementary Methods 1 in the Supplementary Information.

498

499     •    The progressive attribute integration case.

500     Second, let us consider another type of value computation, which essentially proceeds from

501     progressively integrating the value-relevant attributes of choice options. This typically happens

502     when choice options can be decomposed into multiple dimensions that may conflict with each

503     other (cf., e.g., tastiness versus healthiness for food items).

504     Let $x_1,...,x_k$ be the set of $k$ such value-relevant attributes, the combination of which is specific

505     to each option. Assume that the decision system constructs the value of alternative options

506     according to a weighted sum of attributes, i.e.: $V = \sum_k w_k \times x_k$ , where the attribute weights

507     $w_k$ are the same for all options. Assume that each attribute is sampled from a Gaussian

508     distribution with mean $\eta_k$ and variance $\varsigma_k$, i.e. $p(x_k) = N(\eta_k, \varsigma_k)$. Finally, assume that

509     attributes are available to the decision system one at a time, i.e. decision time steps co-occur

510     with attribute-disclosing events. For the sake of simplicity, we set the decision's temporal

511     horizon to $T = k$ , i.e. we focus on the decision to stop (potentially prematurely) the integration

512     of all available value-relevant attributes. In what follows, we refer to this scenario as the

513     *progressive attribute integration* model.

514     In the absence of default preferences, the system holds a prior representation about the

515     options' value that is maximally uninformative. This is because, prior to any value computation,

516     any combination of value-relevant attributes is admissible, and the system did not disclose the

517     options' attributes yet. The first two moments of the system's prior value representation

518     $p(V) = N(\mu_0, \sigma_0)$ are thus given by:

519
$$\begin{cases} \mu_0 = \sum_{k'=1}^{k} w_{k'} \times \eta_k \\ \sigma_0 = \sum_{k'=1}^{k} w_{k'}^2 \times \varsigma_k \end{cases}$$
(23)

520    where of $k$ is the number of value-relevant attributes.

521    Now, as time unfolds and the decision system discloses the value-relevant attributes, it

522    progressively removes sources of uncertainty about the value of alternative options. In

523    principle, if the system reaches the temporal horizon, then it knows all the attributes and can

524    evaluate the alternative options with infinite precision. However, as long as some attributes are

525    missing, value representations remain uncertain. Let $K(t)$ be the set of attribute indices that

526    have been available to the decision system up until time $t$. At time $t$, the decision system thus

527    holds an updated probabilistic representation of value $p\left(V \middle| x_{K(t)}\right) = N\left(\mu(t), \sigma(t)\right)$ with the

528    following mean and variance:

529
$$\begin{cases} \mu(t) = \mu_0 + \tilde{\delta}(t) \\ \sigma(t) = \sigma_0 - \sum_{k' \in K(t)} w_{k'}^2 \times \varsigma_{k'} \end{cases}$$
(24)

530    where the change in the value mode is simply given by:

531    $$\tilde{\delta}(t) = \sum_{k' \in K(t)} w_{k'} \times \left(x_{k'} - \eta_{k'}\right)$$
(25)

532    As before, Equation 25 specifies what the perturbation to the value mode would be, if the

533    underlying value computation was a process of progressive *attribution integration*, whose

534    outcome is the value estimate $\mu(t)$. Note that here, variability in mode perturbations does not

535    arise from some form of stochasticity or unreliability of input signals, as is the case for the

536    *Bayesian value denoising* scenario above. Rather, it derives from the arbitrariness of the

537    permutation order with which attributes become available for options' evaluation. However,

538    should the full set of attributes eventually be disclosed, the estimated value would be

539    $\mu(T) = \sum_{k'}^{k} w_{k'} \times x_{k'}$ , with full certainty ($\sigma(T) = 0$).

540    Here again, the decision system cannot anticipate in which direction the future value mode will

541    change, i.e. its expectation over future mode changes always is $E\left[\tilde{\delta}(t)\right] = 0$ at any point in

542    time (because $E\left[x_k\right] = \eta_k$). Nevertheless, it can derive a prediction over the magnitude of

543    the perturbation, by averaging over all possible permutation orders:

$$
\begin{aligned}
E\left[\tilde{\delta}(t)^2\right] &= \frac{t}{k}\sum_{k'=1}^{k} w_{k'}^2 \times \varsigma_{k'} \\
&= t\sigma_0
\end{aligned}
$$

544                                                                                  (26)

545    Comparing Equations 4 and 26 suggests that, under the *progressive attribute integration*

546    scenario, type #2 effort efficacy simplifies to: $\gamma = \sigma_0$ . This means that type #2 effort efficacy

547    simply scales with the expected range of attributes' variation. This also implies that, in contrast

548    to the above *value denoising* case, the transition probability density of value modes under the

549    *progressive attribute* integration scenario is stationary and complies with oMCD's assumption

550    (cf. Equation 12).

551    What about type #1 effort efficacy? Note that one cannot directly compare Equation 24 to

552    Equation 4, because of the arbitrariness of the order of attribute-disclosing events. In fact, this

553    arbitrariness implies that the dynamics of value variances is decreasing with time but

554    stochastic. Although oMCD is neglecting this stochasticity, type #1 efficacy can be derived

555    from the first-order moment of value variance dynamics. Accordingly, averaging over all

556    possible    permutations    yields    the    following    expected    change    in    precision:

557    $E\left[1/\sigma(t) - 1/\sigma_0\right] \approx t \times 1/\sigma_0 (k-t)$. Using the same logic as above, this suggests that type

558    #1 effort efficacy can now be approximated as: $\beta \approx 1/(k-1)\sum_{t=1}^{k-1} 1/\kappa\sigma_0(k-t)$. Note that

559    we have removed the time horizon from averaging over admissible decision times, since it

560    induces a singularity (infinite precision).

561    Importantly, the *progressive attribute integration* scenario implies that both first- and second-

562    order moments of value representations follow stochastic dynamics. This means that the ideal

563    control policy does not reduce to a single threshold (on either net benefits or confidence), but

564    rather unfolds onto the bidimensional space spanned by both moments of value

565    representations. This makes the *progressive attribute* integration scenario qualitatively

566    different from the *Bayesian value denoising* case. We refer the interested reader to the

567    Supplementary Methods 2 in the Supplementary Information for details regarding the

568    mathematical derivation of the ideal control policy under *progressive attribute integration*.

569

570    One can see that the definition of type #1 and type #2 effort efficacies depends upon the way

571    in which the decision process perturbs the value representations (the above scenarios are just

572    two examples out of many possible forms of value computations). In principle, optimal control

573    would thus require variants of MCD controllers that are tailored to the underlying decision

574    system. For the sake of completeness, the derivation of such ideal control policies are

575    summarized in Appendices 1 and 2. In this context, the MCD architecture that we propose

576    provides an efficient alternative, which generalizes across decision processes and still

577    operates quasi-optimal decision control (see below). The only requirement here, is to calibrate

578    the MCD controller over a few decision trials to learn effort efficacy parameters. Note that such

579    calibration is expected to be very quick (at the limit: only one decision trial), because effort

580    efficacies can be learned on within-trial dynamics (of value representations). This is effectively

581    what we have done here, in an analytical manner, when deriving approximations for the effort

582    efficacy parameters under distinct decision scenarios.

583

584 **Results**

585

586 In the previous section of this manuscript, we derived the online, dual prospective/reactive

587 variant of MCD (and disclosed its connection with two exemplar decision systems). We now

588 wish to illustrate its properties.

589

590     1.  How do prospective MCD and oMCD differ?

591 Formally speaking, online/reactive and prospective MCD policies are solving the same

592 resource allocation problem, i.e. they both aim at stopping resource investment when its net

593 benefits are maximal. At this point, one may thus ask whether oMCD produces better decisions

594 than prospective MCD, which operates by committing to a predefined resource investment.

595 More precisely, under prospective MCD, the decision stops when the expected net benefit is

596 maximal, which is evaluated at the onset of the decision (this corresponds to the red vertical

597 line in Figure 2). But does oMCD yield higher net benefits than prospective MCD (on average)?

598 To answer this question, we resort to Monte-Carlo simulations. In brief, we simulate a particular

599 decision trial in terms of the stochastic dynamics of value representations, according to

600 Equations (3) and (4), using the same decision parameters as for Figure 2. At each time step,

601 oMCD's policy proceeds by comparing the ensuing confidence level to the optimal confidence

602 threshold. When the confidence threshold is reached, we store the resource investment, as

603 well as the ensuing confidence level and net benefit. We proceed similarly for prospective

604 MCD, except that resource investment is defined according to Equation (1). We then repeat

605 the procedure to evaluate the average confidence levels, amounts of invested resources, and

606 net benefits induced by both MCD variants. These are summarized in Figure 4 below, where

607 the averages are taken over 500 sample path trajectories of value modes. Note: as a reference,

608 we also compare MCD control policies to a so-called "oracle" dummy policy, which

609 retrospectively identifies the net benefit apex, i.e. the time at which the stochastic trajectory of

610    net benefits is maximal. This provides an upper (though unachievable) bound to the expected

611    net benefit of any online control policy.

613    One can see that oMCD tends to invest fewer resources and yet achieves higher confidence

614    than prospective MCD (on average). In turn, the ensuing average net benefit is lower for

615    prospective MCD than for oMCD (which is closer to the oracle). Unsurprisingly, under oMCD,

616    the statistical relationship between resource investments and reported confidence levels

617    unfolds along the dynamics of the optimal confidence threshold. In this setting, decisions that

618    take longer eventually yield lower confidence (although this actually depends upon decision

619    parameters, see Figure 3). For prospective MCD, there is no such relationship because

620    resource investment is fixed once decision parameters are set.

622    So do these observations generalize over decision parameter settings? To answer this

623    question, we repeat the same analysis as above, under 200 random settings of all decision

624    parameters. Figure 5 below summarizes the results of this Monte-Carlo simulations series.

626    One can see that the impact of decision parameters on resource investment and confidence

627    is very similar under both MCD variants. This is important, because this means that the known

628    properties of prospective MCD [7] generalize to oMCD. In addition, oMCD's optimal control policy

629    tends to yield lower resource investments and higher confidence levels than prospective MCD.

630    Both effects almost compensate each other, but oMCD tends to provide a small but systematic

631    improvement on the ensuing net benefit, which typically increases with type #2 effort efficacy

632    ($\gamma$). This is because increasing $\gamma$ increases the stochasticity of value mode dynamics, which

633    provides oMCD with more opportunities to exploit favorable variations in confidence (cf. panel

634    B).

636    Now, when compared to prospective MCD, oMCD possesses a unique feature: the potentially

637    nontrivial statistical relationship between decision confidence and resource investments (as

638    proxied using, e.g., response times), *across trials with identical decision parameters*. This was

639    already exemplified in Figure 4 above (cf. panel D).

640    To make this distinction clearer, we performed another set of simulations aiming at evaluating

641    the impact of decision difficulty. Note that difficult decisions can be defined as those decisions

642    where the reliability of value representations improve very slowly. Within the MCD framework,

643    increasing decision difficulty can thus be modelled by decreasing type #1 effort efficacy. We

644    systematically varied $\beta$ from 2 to 8 (having set all the other decision parameters to 4), simulated

645    500 sample path trajectories of value mode dynamics for each difficulty level, and evaluated

646    the ensuing effort investments and achieved confidence levels. Figure 6 below summarizes

647    the simulation results.

648

649    One can see that the net effect of increasing decision difficulty (or equivalently, decreasing

650    type #1 effort efficacy) is to increase resource investment and decrease confidence. This holds

651    for both oMCD and its prospective variant. This means that, on average, reported confidence

652    levels will tend to correlate negatively with resource investments, *across difficulty levels* (at

653    least for this setting of decision parameters). However, for oMCD, this negative relationship

654    between resource investments and reported confidence levels is also true *within each difficulty*

655    *level* (across trials). This has no equivalent under prospective MCD. In addition, the shape of

656    this relationship is preserved across difficulty levels. This is because type #1 effort efficacy

657    induces rather small distortions on oMCD's confidence thresholds (cf. Figure 3 above).

658

659

660    Figure 6 also reveals how oMCD's optimal control policy prospectively anticipates the impact

661    of decision difficulty. In brief, the decay rate of oMCD's confidence threshold increases with

662    decision difficulty, because expected confidence gains become more costly. However, this is

663    overcompensated by the corresponding decrease in the ascent rate of expected confidence,

664    which will delay the time at which confidence eventually reaches the optimal threshold. This

665 eventually determines the way oMCD trades effort against confidence: difficult decisions are

666 given more deliberation time than easy decisions (this is also true for prospective MCD).

667 Note that the effect of difficulty on resource investment, as well as the shape of the

668 effort/confidence relationship, depends on the setting of decision parameters. In other words,

669 these effects do not generalize to all decision parameter settings. For example, increasing

670 decision difficulty will eventually decrease resource investments. Also, the sign of the

671 correlation between confidence and resource investments across difficulty levels may not

672 always align with the sign of this correlation within each difficulty level.

673

674


675     2. How optimal is oMCD's policy?

676 One of oMCD's main claims is that it is possible to derive a quasi-optimal decision control

677 policy, without detailed knowledge of the underlying value computations. But how well does

678 oMCD perform, when compared to ideal policies that rely on such detailed knowledge? To

679 address this question, we compare both resource investments and achieved confidence levels

680 under either oMCD or the ideal control policy, for both decision scenarios (see Supplementary

681 Methods 1 and 2 in the Supplementary Information for mathematical details regarding the

682 derivation of the corresponding ideal policies).

683 We thus conducted the two following sets of Monte-Carlo simulations series. For each decision

684 scenario, we simulate sample path trajectories of moments of value representations, under the

685 corresponding type of value computations. Each trajectory effectively corresponds to a dummy

686 decision trial, given some setting of the relevant decision parameters. Note that only a subset

687 of these parameters is common to all decision scenarii (cost/benefit parameters, i.e.: $R$, $\alpha$ and

688 $\nu$), whereas other parameters are typically decision-specific (*bayesian value denoising*: signal

689 reliability $\Sigma$ and prior variance $\sigma_0$, *progressive attribute integration*: attribute moments $\eta$ and $\zeta$

690 as well as attribute weights $w$). For each decision parameter setting, we derive both the ideal

691    control policy and oMCD's control policy (by approximating the effort efficacy parameters that

692    correspond to the decision-specific parameters). We then collect the resource investments and

693    achieved confidence that are induced by these policies, when applied on sample path

694    trajectories of value representation moments. Now, how do ideal and oMCD policies compare

695    across different settings of decision parameters?

696    Figure 7 below summarizes the comparison of ideal and oMCD policies under the *Bayesian*

697    *value denoising* scenario. This comparison is made across 200 sets of randomly drawn

698    decision parameters $\alpha$, $\nu$, $\Sigma$ and $\sigma_0$. For parameter setting, we derive the average effort

699    investment and achieved confidence level across 500 sample path trajectories of moments of

700    value representations.

701

702    One can see that variations in decision-relevant parameter settings induce very similar

703    variations in average resource investments, achieved confidence and net benefits under both

704    decision control policies. Also, although oMCD's policy yields both more effort costs (in terms

705    of resource investments) and more benefits (in terms of achieved confidence), these effects

706    compensate each other and oMCD's ensuing net benefits are comparable to those of the ideal

707    control policy. Moreover, despite oMCD's approximation of type #2 effort efficacy, it does not

708    seem to have a systematic impact on the similarity between the two policies. These results

709    imply that oMCD provides a tight approximation to the ideal policy for *Bayesian value*

710    *denoising*.

711    Now Figure 8 below summarizes the comparison of ideal and oMCD control policies under the

712    *progressive attribute integration* scenario (200 sets of randomly drawn decision parameters $\alpha$,

713    $\nu$, $\eta$, $\zeta$ and $w$, with $k = 10$).

714

715    As before, one can see that variations in decision-relevant parameter settings induce very

716    similar variations in average resource investments, achieved confidence and net benefits

717    under both control policies. Moreover, despite oMCD's approximation of type #1 effort efficacy,

718    it does not seem to have a systematic impact on the similarity between the two policies. These

719    results imply that oMCD provides an accurate approximation to the ideal control policy for

720    *progressive attribute integration*.

721    Taken together, these results mean that the MCD architecture operates a quasi-optimal

722    decision control that generalizes across decision processes without requiring detailed

723    knowledge about underlying value computations.

724

725        3.  How critical is the definition of MCD's benefit term?

726    The working assumption of MCD is that decision confidence serves as the main benefit term

727    of the resource allocation problem (cf. Equations 1-2). The advantage of this assumption is

728    that it applies to any kind of decision process, irrespective of the underlying computations.

729    However, as we hinted in the introduction, for the specific case of value-based decisions, there

730    exists another natural candidate definition of the benefit term, i.e.: the value of the chosen

731    option. One may argue that changing the definition of the benefit term effectively changes the

732    nature of the resource allocation problem. So how critical is MCD's working assumption? Is

733    oMCD robust to such alternative setting of the resource allocation problem?

734    On the computational side of things, the derivation of the ensuing optimal control policy is very

735    similar to that of oMCD. Since the value of the chosen option is, by definition, the maximum

736    value over the choice set, we refer to this policy as *max(value)*. It is relatively easy to show

737    that oMCD and *max(value)* share one common important feature, i.e.: the critical quantity that

738    triggers decisions is the absolute difference $\left|\Delta\mu\left(t\right)\right|$ in value modes. However, in contrast to

739    oMCD, *max(value)* is insensitive to the variance of value representations (and hence to type

740     #1 effort efficacy). We refer the interested reader to the Supplementary Methods 3 in the

741     Supplementary Information for mathematical details regarding the derivation of *max(value)*'s

742     policy.

743     So do *max(value)* and oMCD policies respond similarly to variations in MCD parameters? To

744     address this question, we performed the following series of Monte-Carlo simulations. First, we

745     sample a set of MCD parameters ($\alpha$, $\beta$, $\gamma$, $\nu$ and $\kappa$) randomly. Second, we derive the optimal

746     control threshold dynamics under both *max(value)* and oMCD policies. Third, we extract the

747     mean response time, confidence, and net benefits over 500 random simulations of moments

748     of value representations sample paths (according to Equation 1). We then repeat the three

749     steps above 200 times. The results of this analysis are summarized in Figure 9 below.

750

751     Although oMCD tends to invest fewer resources than *max(value)* on average, it also achieves

752     smaller confidence levels. This is essentially because the confidence mapping (cf. Equation 8)

753     enforces an upper bound on oMCD's benefit term. Comparatively, *max(value)* thus tolerates

754     stronger effort costs. Nevertheless, both effects compensate each other and both control

755     policies eventually yield very similar outcomes in terms of net benefits. Unsurprisingly, each

756     policy is (slightly) better than the other at maximizing its own benefit on average. More

757     importantly, variations in decision parameter settings induce very similar variations in average

758     resource investments, achieved confidence levels and net benefits. This result suggests that

759     both frameworks are much less different than intuitively thought of, at least in terms of

760     empirically observable decision features (choice, deliberation time, confidence). Moreover,

761     type #1 effort efficacy, which induces variations in oMCD's policy that have no equivalent in

762     *max(value)*, does not seem to have a systematic impact on the similarity between the two

763     policies. In conclusion, oMCD can be thought of as providing a quasi-optimal policy for

764     maximizing the value of the chosen option. In other terms, oMCD is robust to violations of its

765     working assumptions.

766

767      4. Does MCD reproduce established empirical results?

768 As we highlighted before, MCD is agnostic about the underlying decision process. However,

769 what eventually determines the choice that is made is the inner workings of value

770 representation updates. This is important, since some of the decision features may depend

771 upon, e.g., whether the system eventually arrives at a choice that is consistent with the

772 comparison of options' values or not. Inspecting these kinds of effects thus requires performing

773 Monte-Carlo simulations under distinct decision processes (here: *Bayesian value denoising*

774 and *progressive attribute integration*).

775

776 Let us first consider the *Bayesian value denoising* scenario. First, we simulated $10^4$ stochastic

777 dynamics of Bayesian value belief updates according to Equations 20-21, having set the

778 decision parameters as follows: *R=1, α=0.1, v=2, σ₀=10, μ₀=0, Σ=100*, and randomly sampling

779 trial-specific hidden value signals *V* under the ideal observer's prior belief. Note that we chose

780 this parameter setting because it reproduces the empirically observed rate of value-

781 consistent/value-inconsistent decisions (see Figure 12 below). Second, we identified the

782 oMCD-optimal confidence threshold dynamics, having set the effort efficacy parameters to

783 their analytical approximation (cf. Equation 23 and related derivations). We then store the

784 ensuing resource investments and achieved confidence levels, as well as the choices of the

785 decision system (as given by the comparison of value modes at decision time). Figure 10 below

786 summarizes the results of this Monte-Carlo simulations series.

787

788 First, one can see that the MCD approximation of within-trial choice confidence dynamics is

789 reasonably accurate (panel A), and smoothly trades errors at early and late decision times.

790 Second, on average, resource investment decreases with the absolute difference in hidden

791 option values (cf. black line in panel B). Third, above and beyond the effect of option value

792 difference, resource investment decreases when choice confidence increases (cf. blue and

793    red lines in panel B). This derives from the shape of the oMCD confidence threshold dynamics

794    (cf. Figure 3). Fourth, the consistency of choice with value is higher for high-confidence choices

795    than for low-confidence choices (panel C). This observation derives from performing a logistic

796    regression of choice against hidden value, when splitting trials according to whether they yield

797    a high or a low level of confidence [41]. Fifth, on average, choice confidence decreases with the

798    absolute difference in hidden option values (cf. black line in panel D). Note that the oMCD

799    framework also predicts that confidence is higher for choices that are consistent with the

800    comparison of hidden values than for inconsistent choices (cf. red and blue lines in panel D).

801    This suggests that MCD possesses some level of metacognitive sensitivity [42], i.e., it reports

802    lower confidence when making a decision that is at odds with the hidden (unknown) value.

803    Under the assumption that decision time proxies resource investment, these are standard

804    results in empirical studies of value-based decision making [7,13,41,43]. Interestingly, when

805    focusing on choices that are inconsistent with the comparison of hidden values, the impact of

806    value difference on confidence reverses, i.e., choice confidence *decreases* with the absolute

807    difference in hidden values. This relates to known results in the context of perceptual decision

808    making [44]. We note that these results depend upon effort cost parameters. In particular,

809    metacognitive sensitivity tends to decrease in parameter regimes where the dynamics of

810    oMCD confidence thresholds stop the decisions very early (e.g. low cost power and/or high

811    unitary effort cost). This may explain the loss of metacognitive sensitivity that concurs with

812    mental fatigue, which effectively increases one's sensitivity to cognitive effort [45].

813

814    Let us now consider the *progressive attribute integration* scenario. We essentially reproduced

815    the same analysis as above, while simulating stochastic dynamics of value computations by

816    attribute integration according to Equations 24-25, and setting the model parameters to yield

817    a similar rate of value-consistent choices ($R=1$, $\alpha=3$, $v=4$, $k=20$, $\eta_k=1$, $\varsigma_k=1$). Figure 11 below

818    summarizes the results of this Monte-Carlo simulations series.

819

820 In brief, one can see that we qualitatively reproduce the above relationships between effort

821 investment, confidence and choice consistency. This is important, since this means that these

822 relationships tend to generalize across different decision processes. However, this

823 equivalence is only qualitative, and does not always hold. For example, reducing the unitary

824 effort cost eventually renders the oMCD confidence threshold dynamics concave. For

825 *progressive attribute integration*, this reverses the impact of the difference in option values

826 onto confidence for value-inconsistent choices back again. This does not seem to happen

827 under *Bayesian value denoising*.

828

829 For completeness, we re-analyzed the data reported in our previous investigation of (the

830 prospective variant of) the metacognitive control of decisions [7]. In brief, participants were native

831 French speakers, with no reported history of psychiatric or neurological illness. A total of 41

832 people (28 women; age: mean = 28, SD = 5, min = 20, max = 40) participated in this study (no

833 participant was excluded). All participants rated the pleasantness of a series of food items, and

834 performed two-alternative forced choices between pairs of (pseudo-randomly selected) items.

835 In addition to participants' value ratings and choice, we also collected choice confidence,

836 decision time, and subjective effort rating. We note that in this context, within-decision value

837 computations may rely either on retrieving previously experienced food samples from episodic

838 memory [46,47], or on integrating value-relevant attributes (e.g., tastiness and healthiness)

839 derived from cognitive decompositions of choice options [30,48]. Both cognitive scenarios map

840 onto *Bayesian value denoising* (which would average over memory samples) and *progressive*

841 *attribute integration* processes, respectively.

842 We already verified the main predictions of the prospective MCD model, in terms of the

843 relationship between pre-choice (default) value ratings and decision time/effort, as well as the

844 ensuing decision-related variables (i.e. change-of-mind, confidence, choice-induced

845 preference change, etc). As we already discussed, prospective and online variants of MCD

846    make very similar predictions for these kinds of relationships. We now reproduce the above

847    analyses (cf. Figures 10 and 11), which disclose predictions that are specific to the oMCD

848    framework. Figure 12 below summarizes the results of these analyses.

849

850    Note that subjective effort ratings are commensurate with response times, which suggests that

851    effort intensity shows little variations when compared to effort durations. We will comment on

852    this in the Discussion section below. In any case, one can see that the overall pattern of

853    relationships between resource investments (as proxied by either decision time or reported

854    mental effort), choice confidence and item values is qualitatively similar to that predicted from

855    the online MCD model (cf. Figures 10 and 11 above). Note that all the oMCD predictions

856    discussed above are statistically significant in our empirical data:

857    • Effect of DV on reported effort (all trials): $t(40)=-7.6$, mean $r=-0.25 \pm 0.07$ (95% CI),
858       $p<10^{-4}$

859    • Effect of DV on reported effort (high confidence): $t(40)=-5.7$, mean $r=-0.18 \pm 0.07$ (95%
860       CI), $p<10^{-4}$

861    • Effect of DV on reported effort (low confidence): $t(40)=-5.0$, mean $r=-0.14 \pm 0.05$ (95%
862       CI), $p<10^{-4}$

863    • Effort difference (high versus low confidence): $t(40)=-7.3$, mean effort difference$=-0.19$
864       $\pm 0.05$ (95% CI), $p<10^{-4}$

865    • Effect of DV on decision time (all trials): $t(40)=-7.78$, mean $r=-0.19 \pm 0.05$ (95% CI),
866       $p<10^{-4}$

867    • Effect of DV on decision time (high confidence): $t(40)=-5.9$, mean $r=-0.15 \pm 0.05$ (95%
868       CI), $p<10^{-4}$

869    • Effect of DV on decision time (low confidence): $t(40)=-3.9$, mean $r=-0.10 \pm 0.05$ (95%
870       CI), $p=0.0002$

- Response time difference (high versus low confidence): $t(40)=-7.0$, mean RT difference=$-0.62 \pm 0.17$ (95% CI), $p<10^{-4}$

- Effect of DV on choice (all trials): $t(40)=25.2$, mean effect size=$1.56 \pm 0.12$ (logistic regression, 95% CI), $p<10^{-4}$

- Effect of DV on choice (high confidence): $t(40)=32.6$, mean effect size=$2.02 \pm 0.12$ (logistic regression, 95% CI), $p<10^{-4}$

- Effect of DV on choice (low confidence): $t(40)=10.4$, mean effect size=$0.84 \pm 0.16$ (logistic regression, 95% CI), $p<10^{-4}$

- Effect of DV on choice (high versus low confidence): $t(40)=13.8$, mean effect size difference =$1.17 \pm 0.16$ (logistic regression, 95% CI), $p<10^{-4}$

- Effect of DV on confidence (all trials): $t(40)=8.5$, mean r=$0.27 \pm 0.06$ (95% CI), $p<10^{-4}$

- Effect of DV on confidence (value-consistent): $t(40)=10.6$, mean r=$0.27 \pm 0.05$ (95% CI), $p<10^{-4}$

- Effect of DV on confidence (value-inconsistent): $t(40)=-4.22$, mean r=$-0.18 \pm 0.09$ (95% CI), $p<10^{-4}$

- Confidence difference (value-consistent versus value-inconsistent): $t(40)=10.8$, mean confidence difference =$0.10 \pm 0.02$ (95% CI), $p<10^{-4}$

where DV stands for difference in option values, all statistical significance tests are one-sided and derive from standard random effect analyses (sample size: n=41). We note that these analyses were not part of a preregistration protocol.

894 **Discussion**

895

896 In this work, we have presented the online/reactive metacognitive control of decisions or oMCD

897 framework.

898

899     1.  Limitations

900 To begin with, recall that we have framed oMCD as a solution to a resource allocation problem.

901 More precisely, we think of decision deliberation as involving the investment of costly cognitive

902 resources, which are necessary to process decision-relevant information. The outcome of such

903 resource allocation is to override default behavioral responses, which would otherwise be

904 triggered by automatic (e.g., reflexive, habitual or intuitive) brain processes. Under this view,

905 the brain faces the problem of adjusting *the amount* of resources to invest, which we equate

906 with the issue of effort regulation. This perspective is not novel: the notion of mental effort was

907 central to the early definition of automatic versus controlled processing, with the former

908 described as quick and effortless, and the latter as slower and effortful [49]. Since controlled

909 processes are slow, it is reasonable to assume that the brain may regulate effort simply by

910 adjusting its duration. This is the premise of our computational framework, which relies on the

911 theory of optimal stopping [21]. However, effort actually unfolds along two dimensions: duration

912 and intensity. This means that, in principle, both decision speed and confidence may be

913 increased at the cost of increasing effort intensity. Accordingly, investing cognitive control is

914 known to speed up responses in the context of, e.g., behavioral conflict tasks [50,51]. This raises

915 the question: what determines the brain's policy for trading effort intensity against effort

916 duration? A possibility is that this depends upon the nature of the cognitive resource that is

917 required for processing decision-relevant information. The issues of *how* to control resource

918 investment and *which* resource to invest are thus intertwined [2]. For example, one may think of

919 resources as being composed of cognitive modules, such as working memory or attention,

920 whose neurobiological underpinnings may induce distinct costs and/or limitations on effort

921   intensity and duration [52–54]. More generally, the effort intensity/duration tradeoff may be

922   eventually determined by the neurobiological constraints that are imposed on the neural

923   architecture that operates the processing of decision-relevant information [4,55]. For example,

924   value-based decision making may require the active maintenance of multiple value

925   representations that tend to interfere with each other, e.g., because they involve the same

926   neural population within the orbitofrontal cortex [32]. In this case, cognitive control may alter the

927   OFC neural code with the aim of temporarily dampening these interferences. In principle, the

928   associated neural mechanism may operate based on simple confidence monitoring (which

929   would proxy value conflict signals), without knowledge of the intricate architecture of value

930   coding in the OFC. We will test these ideas using artificial neural network models of MCD in

931   forthcoming publications.

932

933       2.   On the generality of oMCD control policy

934   One of the main assumptions behind MCD is that mental effort investment is regulated by a

935   unique controller that operates under agnostic assumptions about the inner workings of the

936   underlying decision system. This constraint somehow culminates in the simplicity of oMCD's

937   control architecture, which reduces to a monitoring of decision confidence. In this context, we

938   have shown that the optimal stopping policies of distinct decision processes (*Bayesian value*

939   *denoising* or *progressive attribute integration*) can be approximated using a simple calibration

940   of effort efficacy parameters. We have also highlighted the ensuing properties of oMCD : when

941   coupled with these different underlying decision systems, oMCD reproduces most established

942   empirical results in the field of value-based decision-making. In addition, we have shown that

943   oMCD is robust to alternative settings of the resource allocation problem. In particular, decision

944   confidence seems to be a reasonable proxy for the value of the chosen option, which is the

945   standard candidate titration for the benefit of value-value based decisions [14,24]. Taken together,

946   these results suggest that the architecture of oMCD control, which relies on the internal

947   monitoring of decision confidence, may generalize to most kinds of decision processes.

948     Preliminary investigations show that this holds for yet another important kind of value-based

949     decisions, whereby value computation is the output of a forward planning process on a decision

950     tree [56,57]. Arguably, this also holds for perceptual or evidence-based decisions. In this context,

951     decision confidence can be defined - somewhat more straightforwardly - as the subjective

952     probability of being correct [35]. As long as effort efficacy parameters can be simply identified,

953     the MCD architecture will provide an accurate approximation to the optimal resource allocation

954     policy. This is trivial when perceptual detection or discrimination processes can be described

955     as some form of *Bayesian denoising* of some perceptual variable of interest [23,40]. This would

956     also hold for perceptual categorization processes, which may rather resemble *attribute*

957     *integration* scenarios [19]. In fact, oMCD's potential generalizability derives from its agnostic

958     stance regarding the nature of information processing that takes place in the underlying

959     decision system. This is also why oMCD can in principle be extended to describe the

960     metacognitive control of other kinds of cognitive processes (e.g., reasoning or memory

961     encoding/retrieval). In this context, an interesting avenue of investigation would be to consider

962     the impact of metacognitive adaptation on the generalization of control policies across

963     cognitive domains. Note that, because we assume MCD's control architecture to be invariant

964     across contexts, it requires a systematic calibration (in terms of, e.g., effort costs and/or

965     efficacies) to guaranty the quasi-optimality of resource allocation. As we highlighted before,

966     we expect such calibration to converge very quickly (e.g., over a few training trials). This is

967     because effort efficacies can be learned from within-trial confidence dynamics. Nevertheless,

968     whether this specific kind of metacognitive adaptation is sufficient to recycle and adjust MCD's

969     control architecture to novel cognitive domains, as well as how it shapes cross-domain

970     metacognitive learning effects, is virtually unknown and would require specific empirical tests.

971

972         3.   On the difference between prospective and online/reactive variants of MCD

973     Retrospectively, prospective and online/reactive variants of MCD solve the same

974     computational problem, i.e. maximizing the expected net benefit of resource allocation. We

975 have shown that their respective control policies share many common features. In particular,

976 they tend to respond similarly to changes in effort costs and/or efficacies. However, they differ

977 in at least two important aspects. First, although its algorithmic derivation is more sophisticated,

978 oMCD's control policy is computationally simpler than its prospective variant. This is because

979 it does not require an explicit comparison of all admissible resource investments prior to

980 decision deliberation. Rather, it relies on dynamical changes in decision confidence signals to

981 trigger a binary (yes/no) stopping decision. In other terms, the comparison between admissible

982 resource investments is performed implicitly, while the control system monitors the progress

983 of the underlying decision system. This renders the neurocomputational architecture of oMCD

984 very similar to basic Drift Diffusion Decision Models or DDMs, whose candidate neural

985 underpinnings have been partially identified [58–60]. Second, only oMCD predicts non trivial

986 second-order statistics on key decision features beyond those induced by changes in effort

987 costs and efficacies. For example, both prospective and online/reactive MCD typically predict

988 a negative correlation between reported confidence levels and response times *across difficulty*

989 *levels* (as induced by different type #1 effort efficacies), but only oMCD predicts such a

990 relationship *within each difficulty level* (across trials). The range and diversity of non trivial

991 second-order statistics that oMCD predicts is exemplified in Figures 10-11. We note that some

992 of these predicted statistical relationships are within the grasp of those existing variants of

993 DDMs that explicitly account for decision confidence. This holds, e.g., for the two-way

994 interaction between confidence and item values onto response time and choice [41]. Others may

995 be more specific to oMCD (and related ideal control policies), e.g., the inversion of the

996 value/confidence relationship for value-consistent and value-inconsistent choices. In any case,

997 these non trivial second-order statistics are the hallmark of online/reactive control policies. In

998 this context, what oMCD offers is a way to predict how these relationships should change,

999 would effort costs and/or efficacies be experimentally manipulated.

1000

1001     4. On extending MCD with goal hierarchies

1002    Whether MCD is operated online or not, it relies upon some prospective computation, which

1003    anticipates the costs and benefits of investing additional resources in the decision. In turn, the

1004    optimal cost-benefit tradeoff relies upon decision-specific features, such as decision

1005    importance and difficulty. The former is signalled by the weight parameter $R$ that scales

1006    confidence in the benefit term (cf. Equation 1). In our previous empirical work on MCD,

1007    participants were asked to decide between pairs of food items. In this context, we manipulated

1008    decision importance by instructing participants that they would have to eat the item they

1009    eventually chose (so-called "consequential decisions") or not. As predicted by the MCD

1010    framework, increasing decision importance systematically increases decision time, above and

1011    beyond the effect of option values [7]. In other terms, increasing decision importance may

1012    overcompensate the cost of mental effort by increasing the demand for confidence. More

1013    generally, we think of $R$ as the expected reward attached to the attainment of the

1014    superordinate goal, within which the decision is framed. Importantly, although $R$ is analogous

1015    to a reward, it is distinct from the values that are attached to the choice options. This does not

1016    mean that the values that decision systems attach to choice options are independent from the

1017    goal: recent research has demonstrated that option values are strongly influenced by how

1018    useful choice options are for achieving one's goal [12,61]. However, at least in principle,

1019    alternative choice options that would be instrumental for attaining an important goal may still

1020    have low value. For example, while starving, one may only have access to low

1021    quality/palatability food items. A possibility is to conceive of goals as being organized

1022    hierarchically, whereby superordinate goals are broken down into candidate subordinate goals

1023    [62,63]. According to MCD, the selection of subordinate goals would be under higher scrutiny

1024    when superordinate stakes increase (everything else being equal). Having said this, the

1025    urgency of attaining superordinate goals may also incur additional temporal costs for

1026    subordinate goal selection, which may overcompensate the increased demand for confidence

1027    (as would be the case for, e.g., starvation). We intend to investigate these kinds of issues in

1028    forthcoming publications.

1029

1030

1031

1032

1033

1034

1035

1036

**Data Availability Statement**

All empirical data (as well as analysis code) is available here: https://owncloud.icm-institute.org/index.php/s/wAsSPNndwZVlBlR.


**Code Availability Statement**

The matlab code that was used to generate all Figures in this manuscript is available here: https://owncloud.icm-institute.org/index.php/s/nXnbv2b3gtNz0Jj. This code is also available as part of the VBA academic freeware (https://mbb-team.github.io/VBA-toolbox/), which is versioned and regularly updated.


**Author Contributions**

Douglas Lee collected the empirical data which we re-analyze in this work.

Juliette Benon, Douglas Lee and Jean Daunizeau derived the mathematical model and analyzed the empirical data.

Juliette Benon, Douglas Lee, William Hopper, Morgan Verdeil, Mathias Pessiglione, Fabien Vinckier, Sebastien Bouret, Marion Rouault, Raphael Lebouc, Giovanni Pezzulo, Christiane Schreiweis, Eric Burguière and Jean Daunizeau contributed to elaborating the oMCD theoretical framework and wrote the paper.


**Competing interests**

We declare no competing interests.


**Acknowledgements**

## References

1. Kahneman, D. *Thinking, Fast and Slow*. (Macmillan, 2011).

2. Shenhav, A., Botvinick, M. M. & Cohen, J. D. The Expected Value of Control: An Integrative Theory of Anterior Cingulate Cortex Function. *Neuron* **79**, 217–240 (2013).

3. Shenhav, A. *et al.* Toward a Rational and Mechanistic Account of Mental Effort. *Annu. Rev. Neurosci.* **40**, 99–124 (2017).

4. Musslick, S., Shenhav, A., Botvinick, M. & D Cohen, J. A Computational Model of Control Allocation based on the Expected Value of Control. in (2015).

5. Lieder, F., Shenhav, A., Musslick, S. & Griffiths, T. L. Rational metareasoning and the plasticity of cognitive control. *PLOS Comput. Biol.* **14**, e1006043 (2018).

6. Griffiths, T. L., Lieder, F. & Goodman, N. D. Rational Use of Cognitive Resources: Levels of Analysis Between the Computational and the Algorithmic. *Top. Cogn. Sci.* **7**, 217–229 (2015).

7. Lee, D. G. & Daunizeau, J. Trading mental effort for confidence in the metacognitive control of value-based decision-making. *eLife* **10**, e63282 (2021).

8. Lee, D. G., Daunizeau, J. & Pezzulo, G. Evidence or Confidence: What Is Really Monitored during a Decision? *Psychon. Bull. Rev.* (2023) doi:10.3758/s13423-023-02255-9.

9. Yeung, N. & Summerfield, C. Metacognition in human decision-making: confidence and error monitoring. *Philos. Trans. R. Soc. B Biol. Sci.* **367**, 1310–1321 (2012).

10. Balsdon, T., Wyart, V. & Mamassian, P. Confidence controls perceptual evidence accumulation. *Nat. Commun.* **11**, 1753 (2020).

11. Balsdon, T., Mamassian, P. & Wyart, V. Separable neural signatures of confidence during perceptual decisions. *eLife* **10**, e68491 (2021).

12. De Martino, B. & Cortese, A. Goals, usefulness and abstraction in value-based choice. *Trends Cogn. Sci.* (2022).

13. Rangel, A., Camerer, C. & Montague, P. R. A framework for studying the neurobiology of value-based decision making. *Nat. Rev. Neurosci.* **9**, 545–556 (2008).

14. Tajima, S., Drugowitsch, J. & Pouget, A. Optimal policy for value-based decision-making. *Nat. Commun.* **7**, 12400 (2016).

15. Smith, S. M. & Krajbich, I. Mental representations distinguish value-based decisions from perceptual decisions. *Psychon. Bull. Rev.* **28**, 1413–1422 (2021).

16. Bennett, D., Bode, S., Brydevall, M., Warren, H. & Murawski, C. Intrinsic Valuation of Information in Decision Making under Uncertainty. *PLOS Comput. Biol.* **12**, e1005020 (2016).

17. Bromberg-Martin, E. S. *et al.* A neural mechanism for conserved value computations integrating information and rewards. *Nat. Neurosci.* **27**, 159–175 (2024).

18. Jezzini, A., Bromberg-Martin, E. S., Trambaiolli, L. R., Haber, S. N. & Monosov, I. E. A prefrontal network integrates preferences for advance information about uncertain rewards and punishments. *Neuron* **109**, 2339-2352.e5 (2021).

19. Summerfield, C. & Tsetsos, K. Building Bridges between Perceptual and Economic Decision-Making: Neural and Computational Mechanisms. *Front. Neurosci.* **6**, (2012).

20. Lopez-Persem, A., Domenech, P. & Pessiglione, M. How prior preferences determine decision-making frames and biases in the human brain. *eLife* **5**, e20317 (2016).

21. Shiryaev, A. N. *Optimal Stopping Rules*. (Springer Science & Business Media, 2007).

22. Feinberg, E. A. & Shwartz, A. *Handbook of Markov Decision Processes: Methods and Applications*. (Springer Science & Business Media, 2012).

23. Drugowitsch, J., Moreno-Bote, R., Churchland, A. K., Shadlen, M. N. & Pouget, A. The Cost of Accumulating Evidence in Perceptual Decision Making. *J. Neurosci.* **32**, 3612–3628 (2012).

24. Tajima, S., Drugowitsch, J., Patel, N. & Pouget, A. Optimal policy for multi-alternative decisions. *Nat. Neurosci.* **22**, 1503–1511 (2019).

25. Fudenberg, D., Newey, W., Strack, P. & Strzalecki, T. Testing the drift-diffusion model. *Proc. Natl. Acad. Sci.* **117**, 33141–33148 (2020).

26. Ratcliff, R., Smith, P. L., Brown, S. D. & McKoon, G. Diffusion Decision Model: Current Issues and History. *Trends Cogn. Sci.* **20**, 260–281 (2016).

27. Pleskac, T. J. & Busemeyer, J. R. Two-stage dynamic signal detection: a theory of choice, decision time, and confidence. *Psychol. Rev.* **117**, 864–901 (2010).

28. Fellows, L. K. Deciding how to decide: ventromedial frontal lobe damage affects information acquisition in multi-attribute decision making. *Brain* **129**, 944–952 (2006).

29. Hunt, L. T., Dolan, R. J. & Behrens, T. E. J. Hierarchical competitions subserving multi-attribute choice. *Nat. Neurosci.* **17**, 1613–1622 (2014).

30. Lim, S.-L., O'Doherty, J. P. & Rangel, A. Stimulus Value Signals in Ventromedial PFC Reflect the Integration of Attribute Value Signals Computed in Fusiform Gyrus and Posterior Superior Temporal Gyrus. *J. Neurosci.* **33**, 8729–8741 (2013).

31. O'Doherty, J. P., Rutishauser, U. & Iigaya, K. The hierarchical construction of value. *Curr. Opin. Behav. Sci.* **41**, 71–77 (2021).

32. Pessiglione, M. & Daunizeau, J. Bridging across functional models: the OFC as a value-making neural network. *Behavioral Neuroscience* (in press) (2021).

33. Suzuki, S., Cross, L. & O'Doherty, J. P. Elucidating the underlying components of food valuation in the human orbitofrontal cortex. *Nat. Neurosci.* **20**, 1780–1786 (2017).

34. Daunizeau, J. Semi-analytical approximations to statistical moments of sigmoid and softmax mappings of normal variables. *ArXiv170300091 Q-Bio Stat* (2017).

35. Pouget, A., Drugowitsch, J. & Kepecs, A. Confidence and certainty: distinct probabilistic quantities for different goals. *Nat. Neurosci.* **19**, 366–374 (2016).

36. Lebreton, M., Abitbol, R., Daunizeau, J. & Pessiglione, M. Automatic integration of confidence in the brain valuation signal. *Nat. Neurosci.* **18**, 1159–1167 (2015).

37. Lopez-Persem, A. *et al.* Four core properties of the human brain valuation system demonstrated in intracranial signals. *Nat. Neurosci.* **23**, 664–675 (2020).

38. Papadimitriou, C. H. & Tsitsiklis, J. N. The Complexity of Markov Decision Processes. *Math. Oper. Res.* **12**, 441–450 (1987).

39. Bellman, R. *Dynamic Programming*. (1957).

40. Daunizeau, J. *et al.* Observing the observer (I): meta-bayesian models of learning and decision-making. *PloS One* **5**, e15554 (2010).

41. De Martino, B., Fleming, S. M., Garrett, N. & Dolan, R. J. Confidence in value-based choice. *Nat. Neurosci.* **16**, 105–110 (2012).

1146    42. Fleming, S. M. & Lau, H. C. How to measure metacognition. *Front. Hum. Neurosci.* **8**, (2014).

1147    43. Milosavljevic, M., Malmaud, J., Huth, A., Koch, C. & Rangel, A. The drift diffusion model can

1148        account for value-based choice response times under high and low time pressure. *Judgm. Decis.*

1149        *Mak.* **5**, 437–449 (2010).

1150    44. Kepecs, A., Uchida, N., Zariwala, H. A. & Mainen, Z. F. Neural correlates, computation and

1151        behavioural impact of decision confidence. *Nature* **455**, 227–231 (2008).

1152    45. Blain, B., Hollard, G. & Pessiglione, M. Neural mechanisms underlying the impact of daylong

1153        cognitive work on economic decisions. *Proc. Natl. Acad. Sci.* **113**, 6967–6972 (2016).

1154    46. Bakkour, A. *et al.* The hippocampus supports deliberation during value-based decisions. *eLife* **8**,

1155        e46080 (2019).

1156    47. Lebreton, M. *et al.* A critical role for the hippocampus in the valuation of imagined outcomes.

1157        *PLoS Biol.* **11**, e1001684 (2013).

1158    48. Hare, T. A., Camerer, C. F. & Rangel, A. Self-control in decision-making involves modulation of

1159        the vmPFC valuation system. *Science* **324**, 646–648 (2009).

1160    49. Schneider, W. & Shiffrin, R. M. Controlled and automatic human information processing: I.

1161        Detection, search, and attention. *Psychol. Rev.* **84**, 1–66 (1977).

1162    50. Lin, H., Ristic, J., Inzlicht, M. & Otto, A. R. The Average Reward Rate Modulates Behavioral and

1163        Neural Indices of Effortful Control Allocation. *J. Cogn. Neurosci.* **34**, 2113–2126 (2022).

1164    51. Otto, A. R., Braem, S., Silvetti, M. & Vassena, E. Is the juice worth the squeeze? Learning the

1165        marginal value of mental effort over time. *J. Exp. Psychol. Gen.* **151**, 2324–2341 (2022).

1166    52. Grujic, N., Brus, J., Burdakov, D. & Polania, R. Rational inattention in mice. *Sci. Adv.* **8**,

1167        eabj8935 (2022).

1168    53. Kool, W., Shenhav, A. & Botvinick, M. M. Cognitive Control as Cost-Benefit Decision Making.

1169        in *The Wiley Handbook of Cognitive Control* (ed. Egner, T.) 167–189 (John Wiley & Sons, Ltd,

1170        2017). doi:10.1002/9781118920497.ch10.

1171    54. Silvestrini, N., Musslick, S., Berry, A. S. & Vassena, E. An integrative effort: Bridging

1172        motivational intensity theory and recent neurocomputational and neuronal models of effort and

1173        control allocation. *Psychol. Rev.* **130**, 1081–1103 (2023).

1174   55. Petri, G. *et al.* Universal limits to parallel processing capability of network architectures.

1175        *ArXiv170803263 Q-Bio* (2017).

1176   56. Consul, S., Heindrich, L., Stojcheski, J. & Lieder, F. Improving Human Decision-making by

1177        Discovering Efficient Strategies for Hierarchical Planning. *Comput. Brain Behav.* **5**, 185–216

1178        (2022).

1179   57. Sezener, C. E. Computing the Value of Computation for Planning. *ArXiv* (2018).

1180   58. Gold, J. I. & Shadlen, M. N. The neural basis of decision making. *Annu. Rev. Neurosci.* **30**, 535–

1181        574 (2007).

1182   59. Lam, N. H. *et al.* Effects of Altered Excitation-Inhibition Balance on Decision Making in a

1183        Cortical Circuit Model. *J. Neurosci.* **42**, 1035–1053 (2022).

1184   60. Turner, B. M., van Maanen, L. & Forstmann, B. U. Informing cognitive abstractions through

1185        neuroimaging: the neural drift diffusion model. *Psychol. Rev.* **122**, 312–336 (2015).

1186   61. Castegnetti, G., Zurita, M. & De Martino, B. How usefulness shapes neural representations during

1187        goal-directed behavior. *Sci. Adv.* **7**, eabd5363 (2021).

1188   62. Bay, D. & Daniel, H. The theory of trying and goal-directed behavior: The effect of moving up

1189        the hierarchy of goals. *Psychol. Mark.* **20**, 669–684 (2003).

1190   63. Gozli, D. G. & Dolcini, N. Reaching Into the Unknown: Actions, Goal Hierarchies, and

1191        Explorative Agency. *Front. Psychol.* **9**, (2018).

1192   64. Lee, D. G. & Daunizeau, J. Trading mental effort for confidence in the metacognitive control of

1193        value-based decision-making. *eLife* **10**, e63282 (2021).

1194

1195

1196

1197 **Figure captions**

1198

1199 **Figure 1: derivation of oMCD's optimal control policy.** Net benefits (y-axis) are plotted

1200 against the value mode difference (x-axis). The red and green lines show the net benefit if the

1201 system were stopping at $t = T - 1$, and the expected net benefit at $t = T - 1$. Finally, the dotted

1202 black line shows the optimal net benefit at $t = T - 1$, and the dotted blue line shows its

1203 expectation at $t = T - 2$ (see main text).

1204

1205 **Figure 2: oMCD's optimal control policy. A**: The black dotted line shows the oMCD-optimal

1206 net benefit threshold. The blue line and shaded area depict the mean and standard deviation

1207 of net benefit dynamics (over the 1000 Monte-Carlo simulations), respectively. This reflects

1208 the possible variations of within-trial confidence dynamics. The vertical red line indicates the

1209 optimal resource allocation as obtained from the prospective variant of MCD, and the horizontal

1210 red line depicts the corresponding average net benefit level. **B**: The black dotted line shows

1211 the oMCD-optimal confidence threshold. The blue line and shaded area depict the mean and

1212 standard deviation of decision confidence (over the same Monte-Carlo simulations). The

1213 horizontal red line depicts the average confidence level that corresponds to the optimal

1214 resource allocation under prospective MCD.

1215

1216 **Figure 3: Impact of decision parameters on oMCD's optimal confidence threshold**

1217 **dynamics. A**: Effect of type #1 effort efficacy. Optimal confidence threshold (y-axis, black dots)

1218 is plotted against decision time (x-axis), for different β levels (color code). **B**: Effect of type #2

1219 effort efficacy, same format. **C**: Effect of unitary effort cost, same format. **D**: Effect of cost

1220 power, same format.

1221

1222 **Figure 4: the performance of oMCD's optimal control policy. A**: the average amount of

1223 resources invested (y-axis) is shown under oMCD (black), prospective MCD (red), or oracle

1224    (green) policies. Errobars depict standard error around the mean (s.e.m.). **B**: Average

1225    confidence level at the time of decision, same format. **C**: The average net benefits, same

1226    format. **D**: Achieved confidence (y-axis) is plotted against resource investment deciles (x-axis)

1227    for all control policies (oMCD: black, MCD: red, oracle: green). The black dotted line shows

1228    oMCD's optimal confidence threshold.

1229

1230    **Figure 5: comparison between prospective MCD and oMCD.** A: the amount of resources

1231    invested under the prospective variant of MCD (x-axis) is plotted against the average amount

1232    of resources invested under oMCD (y-axis). Each dot corresponds to a specific set of decision

1233    parameters (200 samples). The color code indicates type #2 effort efficacy (blue: low $\gamma$, red:

1234    high $\gamma$). B: decision confidence, same format. C: net benefit, same format.

1235

1236    **Figure 6: Impact of difficulty level. A**: oMCD's mean resource investment (y-axis, black dots)

1237    is plotted as a function of type #1 effort efficacy (x-axis). Errorbars depict standard deviations

1238    across trials, and red diamonds show the resource investment under prospective MCD. **B**:

1239    Achieved confidence, same format. **C**: Achieved confidence (y-axis) is plotted against resource

1240    investments deciles (x-axis), for each difficulty level (color code: $\beta$ = type #1 effort efficacy),

1241    under oMCD's optimal policy. **D**: oMCD's confidence threshold (y-axis, plain lines) is plotted

1242    against decision time (x-axis), for each difficulty level (same color code as lower-left panel).

1243    Dashed lines show expected confidence, and dots show the corresponding resource

1244    investments under prospective MCD.

1245

1246    **Figure 7: *Bayesian value denoising*: comparison of oMCD and ideal control policies**. **A**:

1247    average resource investments under oMCD's policy (y-axis) are plotted against average

1248    resource investments under the ideal policy (x-axis), across parameter settings (dots). The

1249    color code indicates type #2 effort efficacy (blue: low $\gamma$, red: high $\gamma$). **B**: average achieved

1250    confidence, same format. **C**: average net benefit, same format.

1251

**Figure 8:** *Progressive attribute integration***: comparison of oMCD and ideal control policies**. Same format as Figure 7. The color code indicates type #1 effort efficacy (blue: low $\beta$, red: high $\beta$).

**Figure 9: Comparison of** *max(value)* **and oMCD control policies.** A: mean invested resources under oMCD's control policy (y-axis) and under *max(value)* policy (x-axis) are plotted against each other across random MCD parameter settings. The color code indicates type #1 effort efficacy (blue: low $\beta$, red: high $\beta$). **B**: mean confidence, same format. **C**: mean MCD's net benefit, same format. **D**: mean *max(value)* net benefit, same format.

**Figure 10: oMCD predictions under** *Bayesian value denoising.* **A**: The blue line and shaded area depict the mean and standard deviation of confidence trajectories (across the $10^4$ Monte-Carlo simulations), respectively. The blue dashed line shows the expected confidence under the corresponding MCD approximation, and the black dashed line shows the oMCD-optimal confidence threshold. **B**: Resource investment (y-axis) is plotted against the difference in hidden option values (x-axis), for all trials (black), high-confidence trials (blue) and low-confidence trials (red), respectively. **C**: The probability of choosing the first option (y-axis) is plotted against the difference in hidden option values (x-axis), for all trials (black), high-confidence trials (blue) and low-confidence trials (red), respectively. **D**: Achieved choice confidence (y-axis) is plotted against the difference in hidden option values (x-axis), for all trials (black), value-consistent trials (blue) and value-inconsistent trials (red), respectively.

**Figure 11: oMCD predictions under** *progressive attribute integration.* Same format as Figure 10.

**Figure 12: Re-analysis of behavioral data in a simple value-based decision making experiment** [7]**. A**: Reported mental effort (y-axis) is plotted against the difference in reported

1279   option values (x-axis), for all trials (black), high-confidence trials (blue) and low-confidence

1280   trials (red), respectively. **B:** Response time, same format. **C&D**: same format as Figure 10.

1281

1282

1283

**A: net benefit**

**B: confidence**

**A: type #1 effort efficacy**

**B: type #2 effort efficacy**

**C: unitary effort cost**

**D: effort cost power**

**A: resource investment**

**B: achieved confidence**

**C: net benefit**

**D: effort/confidence relationship**

**A: resource investment**

**B: confidence**

**C: net benefit**

**A: resource investment**

**B: confidence**

**C: effort/confidence relationship**

**D: optimal policy**

**A: resource investment**

**B: confidence**

**C: net benefit**

**A: resource investment**

**B: confidence**

**C: net benefit**

**A: resource investments**

**B: confidence**

**C: net benefit (oMCD)**

**D: net benefit (max(value))**

**A: oMCD policy**

confidence

resource

**B: resource investment**

resource investment

dv

- all trials
- high conf.
- low conf.

**C: value consistency**

P(choice)

dv

- all trials
- high conf.
- low conf.

**D: achieved confidence**

confidence

dv

- all trials
- value-consistent (74.4%)
- value-inconsistent (25.6%)

**A: oMCD policy**

confidence

resource

**B: resource investment**

resource investment

dv

- all trials
- high conf.
- low conf.

**C: value consistency**

P(choice)

dv

- all trials
- high conf.
- low conf.

**D: achieved confidence**

confidence

dv

- all trials
- value-consistent (74.9%)
- value-inconsistent (25.1%)