

1 **Tuning the speed-accuracy trade-off to maximize reward rate**
2 **in multisensory decision-making**

3

4 Running title: Maximizing reward rate in multisensory decision-making

5

6 Jan Drugowitsch^{1,2,3}, Gregory C. DeAngelis^{1,5}, Dora E. Angelaki^{4,5}, Alexandre
7 Pouget^{1,3,5}

8

9 ¹ Department of Brain and Cognitive Sciences, University of Rochester,
10 Rochester, New York 14627, USA

11 ² Institut National de la Santé et de la Recherche Médicale, École Normale
12 Supérieure, 75005 Paris, France

13 ³ Département des Neurosciences Fondamentales, Université de Genève, CH-
14 1211 Geneva 4, Switzerland

15 ⁴ Department of Neuroscience, Baylor College of Medicine, Houston, Texas
16 77030, USA

17

18 ⁵ These authors contributed equally to this work

19

20 Corresponding author:

21 Jan Drugowitsch

22 University Medical Center (CMU)

23 Dept of Neuroscience, room 8012

24 1 rue Michel-Servet

25 CH-1211 Geneva 4

26 Switzerland

27 jdrugo@gmail.com

28

29

30

31

32

33 **Acknowledgements**

34 Experiments and D.E.A. were supported by NIH grant R01 DC007620. G.C.D. was
35 supported by NIH grant R01 EY016178. A.P. was supported by grants from the
36 National Science Foundation (BCS0446730), a Multidisciplinary University
37 Research Initiative (N00014-07-1-0937), the Air Force Office of Scientific
38 Research (FA9550-10-1-0336), and the James McDonnell Foundation.

39

40 **Abstract**

41 For decisions made under time pressure, effective decision making based on
42 uncertain or ambiguous evidence requires efficient accumulation of evidence
43 over time, as well as appropriately balancing speed and accuracy, known as the
44 speed/accuracy trade-off. For simple unimodal stimuli, previous studies have
45 shown that human subjects set their speed/accuracy trade-off to maximize
46 reward rate. We extend this analysis to situations in which information is
47 provided by multiple sensory modalities. Analyzing previously collected data (J.
48 Drugowitsch, DeAngelis, Klier, Angelaki, & Pouget, 2014), we show that human
49 subjects adjust their speed/accuracy trade-off to produce near-optimal reward
50 rates. This trade-off can change rapidly across trials according to the sensory
51 modalities involved, suggesting that it is represented by neural population codes
52 rather than implemented by slow neuronal mechanisms such as gradual changes
53 in synaptic weights. Furthermore, we show that deviations from the optimal
54 speed/accuracy trade-off can be explained by assuming an incomplete gradient-
55 based learning of these trade-offs.

56

57 **Introduction**

58 In the uncertain and ambiguous world we inhabit, effective decision making not
59 only requires efficient processing of sensory information, but also evaluating
60 when enough information has been accumulated to commit to a decision. One
61 can make fast, but uninformed and thus inaccurate, decisions or one can elect to
62 make slower, but well-informed, choices. Choosing this so-called speed-accuracy
63 trade-off (SAT) becomes even more complex if several sensory modalities
64 provide decision-related information. For example, the strategy for crossing a
65 busy street will be very different in bright daylight, when one can rely on both
66 eyes and ears to detect oncoming vehicles, as compared to complete darkness, in
67 which case the ears will prove to be the more reliable source of information.

68 The SAT has been extensively studied for perceptual decisions based on
69 information provided by a single sensory modality. For the most commonly
70 studied visual modality, it has been shown that animals accumulate evidence
71 near-optimally over time (Kiani & Shadlen, 2009). In this context, the efficiency

72 of the chosen SAT is assessed in comparison to *diffusion models*, a family of
73 models that trigger decisions as soon as a drifting and diffusing particle reaches
74 one of two bounds (Ratcliff, 1978). In these models, which describe the SAT
75 surprisingly well despite their simplicity (Palmer, Huk, & Shadlen, 2005; Ratcliff,
76 1978; Ratcliff & McKoon, 2008), the drift represents the available sensory
77 information, and the diffusion causes variability in decision times and choices.
78 The level of the bound controls the SAT, with a higher bound leading to slower,
79 more accurate choices. Instructed changes to the SAT have been shown to be
80 well captured by changes to only the bound in a diffusion model (Palmer et al.,
81 2005; Reddi, Asrress, & Carpenter, 2003; Reddi & Carpenter, 2000). Without
82 being explicitly instructed to make either fast or accurate decisions, well-trained
83 human subjects are known to adjust their SAT to maximize their reward rate
84 (Balci et al., 2011; Simen et al., 2009), or a combination of reward rate and choice
85 accuracy (Bogacz, Hu, Holmes, & Cohen, 2010). These SAT adjustments are also
86 well captured by tuning the corresponding diffusion model bounds. Thus, we can
87 define the SAT directly in terms of these bounds: a constant SAT refers to
88 behavior predicted by diffusion models with constant bounds, and a SAT that
89 changes across trials requires a diffusion model with bounds that vary on the
90 same time-scale.

91 Here, we extend the analysis of how human decision-makers adjust their
92 SAT to situations in which they receive information from multiple sensory
93 modalities. We have previously shown that, even in the case of multiple
94 modalities and time-varying evidence reliability, humans are able to accumulate
95 evidence across time and modalities in a statistically near-optimal fashion (J.
96 Drugowitsch et al., 2014). This analysis was based on a variant of diffusion
97 models that retains optimal evidence accumulation even for multiple sources of
98 evidence whose reliability varies differentially over time. As we focused on
99 evidence accumulation in that study, we were agnostic as to how the SAT varied
100 across stimulus conditions; thus, we left the model bounds, which controlled the
101 SAT, as free parameters that were adjusted to best explain the subjects' behavior.

102 In this follow-up study, we use the previously devised model to analyze
103 whether and how effectively human subjects adjust their SAT if they have
104 evidence from multiple modalities at their disposal. Specifically, we find that

105 subjects adjust their SAT on a trial by trial basis, depending on whether the
106 stimuli are unisensory or multisensory. Moreover, the changes in SAT result in
107 reward rates that are close to those achievable by the best-tuned model, a
108 finding that is robust to changes in assumptions about how the reward rate is
109 computed. Finally, we demonstrate that small deviations from the optimal SAT
110 seem to stem from an incomplete reward rate maximization process. Overall, our
111 findings hint at decision-making strategies that are more flexible than previously
112 assumed, with SATs that are efficiently changed on a trial-by-trial basis.

113

114 **Results and Discussion**

115 Our analysis is based on previously reported behavioral data from human
116 subjects performing a reaction-time version of a heading discrimination task
117 based on optic flow (visual condition), inertial motion (vestibular condition), or a
118 combination of both cues (combined condition) (J. Drugowitsch et al., 2014).
119 Reliability of the visual cue was varied randomly across trials by changing the
120 motion coherence of the optic flow. Subjects experienced forward translation
121 with a small leftward or rightward deviation, and were instructed to report as
122 quickly and as accurately as possible whether they moved leftward or rightward
123 relative to straight ahead.

124 First, we ask whether subjects can adjust their SAT from trial to trial.
125 Having related changes in the SAT to changes in diffusion model bounds, this is
126 akin to asking if their behavior could arise from a diffusion model with a bound
127 that changes on a trial-by-trial basis. Our diffusion model necessitates the use of
128 a *scaled* bound, which is the constant *actual* bound per modality divided by the
129 diffusion standard deviation that depends on optic flow coherence. The use of
130 such a scaled bound prohibits us from fitting actual bound levels, but rather
131 scaled versions thereof. For the same reason, we cannot unambiguously predict
132 the behavior that would emerge from a model with actual bounds matched
133 across modalities (i.e., a constant SAT). Therefore, we instead rely on a
134 qualitative argument about how such matched bounds would be reflected in the
135 relation between decision speed and accuracy across modalities.

136 As Figure. 1a illustrates for subject B2, increasing the coherence of the
137 optic flow caused subjects to make faster, more accurate choices. This pattern
138 was similar if only the visual modality (solid blue lines, Figure. 1a) or both
139 modalities were present (solid red lines, Figure. 1a). This result is qualitatively
140 compatible with the idea that subjects used a single SAT within conditions in
141 which the same modality (visual/vestibular) or modality combination
142 (combined) provided information about heading. Within the framework of
143 diffusion models with fixed actual bounds on the diffusing particle¹, such a single
144 SAT predicts that, once the amount of evidence per unit time (in our case
145 controlled by the coherence) increases, choices ought to be on average either
146 faster, more accurate, or both in combination, but never slower or less accurate.
147 However, our data violate this prediction, thus showing that the SAT changes
148 across conditions. Consider, for example, the choice accuracy and reaction times
149 of subject B2 in both the visual-only (top blue circle, Figure. 1a) and combined
150 condition (top red square, Figure. 1a) trials at 70% motion coherence. Although
151 the combined condition provides more evidence per unit time due to the
152 additional presence of the vestibular modality, responses in the combined
153 condition are less accurate than in the visual-only condition, violating the idea of
154 a single SAT (that is, a fixed diffusion model bound) across conditions. The same
155 pattern emerged across all subjects, whose choices in the combined condition
156 were on average significantly less accurate than in the visual condition (for 70%
157 coherence; one-tailed Wilcoxon signed-rank $W=54$, $p<0.002$). As these stimulus
158 conditions were interleaved across trials, our results clearly indicate that
159 subjects were able to change their SAT on a trial-by-trial basis.

160 Next, we explore whether these adjustments in the SAT serve to maximize
161 subjects' *reward rate*. Even though subjects did not receive an explicit reward for
162 correct trials, we assumed that correct decisions evoke an internal reward of
163 magnitude one. Therefore, we computed reward rate as the fraction of correct

¹ A less common alternative to bounding the diffusing particle in diffusion models is to bound the posterior belief. In case of the latter, changing the amount of evidence per unit time only affects the response time but not its accuracy, which remains unchanged. In rare cases, a bound on the diffusing particle equals a bound on the posterior belief (Jan Drugowitsch, Moreno-Bote, Churchland, Shadlen, & Pouget, 2012), but this is not the case in our context.

164 decisions across all trials, divided by the average time between the onset of
165 consecutive trials. We proceed in two steps: first, we ask whether subjects have a
166 higher reward rate across trials of the multisensory condition compared to both
167 unimodal conditions. This is an important question because we have found
168 previously that subjects accumulate evidence optimally across modalities (J.
169 Drugowitsch et al., 2014), which implies that, with proper setting of the SAT,
170 they should be able to obtain higher reward rates in the multisensory condition
171 compared to the unimodal conditions. As shown in Figure 1b, reward rate is
172 indeed greater, for all subjects, when both sensory modalities are presented than
173 for either modality alone (both unimodal vs. combined: Wilcoxon signed-ranks
174 $W=0, p<0.002$). This confirms that subjects combined evidence across modalities
175 to improve their choices.

176 We now turn to the question of whether subjects tune their SATs to
177 maximize the reward rate. For this purpose, we focus on the reward rate across
178 all trials rather than for specific stimulus conditions, as subjects might, for
179 example, trade off decision accuracy in unimodal conditions with decision speed
180 in the combined condition. To determine how close subjects were to maximizing
181 their reward rate, we needed to compute the best achievable reward rate. To do
182 this, we tuned the bounds of our modified diffusion model to maximize its
183 reward rate, while keeping all other model parameters, including the non-
184 decision times and choice biases, fixed to those resulting from fits to the
185 behavior of individual subjects. As a starting point, we allowed bounds to vary
186 freely for each stimulus modality and each motion coherence, to provide the
187 greatest degrees of freedom for the maximization. As described further below,
188 we also performed the same analysis with more restrictive assumptions. We call
189 the reward rate resulting from this procedure the *optimal* reward rate. This
190 reward rate was subject-dependent, and was used as a baseline against which
191 the empirical reward rates were compared.

192 Figure 2a shows the outcome of this comparison. As can be seen, all but
193 one subject featured a reward rate that was greater than 90% of the optimum,
194 with two subjects over 95%. As a comparison, the best performance when
195 completely ignoring the stimulus and randomly choosing one option at trial
196 onset (i.e. all actual bounds set to zero) causes a significant 25% to 30% drop in

197 reward rate (subjects vs. random: Wilcoxon signed-rank $W=55$, $p<0.002$). Thus,
198 subjects featured near-optimal reward rates that were significantly better than
199 those resulting from rapid, uninformed choices.

200 Our analysis of the subjects' reward rate relative to the optimum is fairly
201 robust to assumptions we make about how this reward rate and its optimum are
202 defined. Thus far, we have assumed implicit, constant rewards for correct
203 decisions and the absence of any losses for the passage of time or incorrect
204 choices. However, accumulating evidence is effortful, and this effort might offset
205 the eventual gains resulting from correct choices. In fact, previous work suggests
206 that human decision makers incur such a cost, possibly related to mental effort,
207 in the range of 0.1-0.2 units of reward per second for accumulating evidence (Jan
208 Drugowitsch et al., 2012). Importantly, this cost modulates both the subjects' and
209 the optimal reward rate, causing the median reward rate across subjects to
210 actually rise slightly to 95.4% and 95.1% (costs of 0.1 and 0.2) of the optimum
211 value (Figure 2b, second and third columns), compared to the cost-free median
212 of 93.7%.

213 The optimal reward rates so far were obtained from a model in which we
214 allowed independent bounds for each stimulus modality and each motion
215 coherence, which implies that subjects can rapidly and accurately estimate
216 coherence. Using instead the more realistic assumption (J. Drugowitsch et al.,
217 2014) that bounds only vary across modalities while coherence modulates
218 diffusion variance but not bound height, we reduce the number of parameters
219 and thus degrees of freedom for reward rate maximization. As a result, subjects'
220 reward rates relative to the optimum rise slightly (median 94.5%), where the
221 optimal model is now restricted to use the same bound across all coherences
222 (Figure 2b, fourth column). Furthermore, we have assumed the model to feature
223 the same choice biases as the subjects. These biases reduce the probability of
224 performing correct choices, and thus the reward rate, such that removing them
225 from our model boosts the model's optimal reward rate. As a consequence,
226 removing these biases causes a consistent drop in subjects' relative reward rate
227 (Figure 2b, last two columns). Even then, reward rates are still around 90% of
228 the optimum (median 87.8% and 88.7% for free and parametric bounds,
229 respectively). If instead of featuring the observed behavior, subjects were to

230 ignore the stimulus and randomly choose one option at trial onset, they would
231 incur a significant drop in reward rate for all of the different assumptions about
232 how we define this optimum (e.g. with/without accumulation cost, ...) as
233 outlined above (subject vs. random, blue vs. red in Figure. 2b: Wilcoxon signed-
234 rank $W=55$, $p<0.002$, except cost 0.2: $W=54$, $p<0.004$).

235 Despite exhibiting near-optimal reward rates, all subjects feature small
236 deviations from optimality. These deviations may result from incomplete
237 learning of the optimal SAT. We only provided feedback about the correctness of
238 choices in early stages of the experiment, until performance stabilized, and
239 subjects did not receive feedback during the main experiment. Nevertheless,
240 subjects' speed/accuracy trade-off remained rather stable after removing
241 feedback, which includes all trials we analyzed. Thus, incomplete learning in the
242 initial training period should be reflected equally in all of these trials. To test the
243 incomplete-learning hypothesis, we assumed that subjects adjusted their
244 strategy in small steps by using gradient-based information about how the
245 reward rate changed in the local neighborhood of the currently chosen bounds.
246 For our argument, it does not matter if the gradient-based strategy was realized
247 through stochastic trial-and-error or more refined approaches involving analytic
248 estimates of the gradient, as long as it involved an unbiased estimate of the
249 gradient. What is important, however, is that such an approach would lead to
250 faster learning along directions of steeper gradients (Figure. 3a). As a result,
251 incomplete learning should lead to near-optimal bounds along directions having
252 a steep gradient, but large deviations from the optimal bound settings along
253 directions having shallow gradients.

254 To measure the steepness of the gradient for different near-optimal
255 bounds, we used the reward rate's curvature (that is, its second derivative) with
256 respect to each of these bounds. If these bounds were set by incomplete gradient
257 ascent, we would expect bounds associated with a strong curvature to be near-
258 optimal (red dimension in Figure. 3a; large curvature, close to optimal bound in
259 inset) and bounds in directions of shallow curvature to be far away from their
260 optimum (blue dimension in Figure. 3a; small curvature, distant from optimal
261 bound in inset). In contrast, strongly mis-tuned bounds associated with a large
262 curvature (points far away from either axis in Figure. 3b) would violate this

263 hypothesis. If we plot reward rate curvature against the distance between
264 estimated and optimal bounds, the data clearly show the predicted relationship
265 (Figure. 3b). Specifically, reward rate curvature is generally moderate to strong
266 in the vestibular-only and combined conditions, and most of these bounds are
267 found to be near-optimal. In contrast, curvature is rather low for the visual
268 condition, and many of the associated bounds are far from their optimal settings.
269 This is exactly the pattern one would expect to observe if deviations from
270 optimality result from a prematurely terminated gradient-based learning
271 strategy. This analysis rests on the assumption that the manner in which reward
272 rate varies with changes in the bounds is well approximated by a quadratic
273 function. If this were the case, then the estimated loss in reward rate featured by
274 the subjects when compared to the tuned model should also be well
275 approximated by this quadratic function. These two losses are indeed close to
276 each other for most subjects (Figure. 3c), thus validating the assumption.

277 Previous studies have suggested that deviations from optimal bound
278 settings may arise if subjects are uncertain about the inter-trial interval (Bogacz
279 et al., 2010; Zacksenhouse, Bogacz, & Holmes, 2010). With such uncertainty,
280 subjects should set their bound above that deemed to be optimal when the inter-
281 trial interval is perfectly known. A similar above-optimal bound would arise if
282 subjects are either uncertain about the optimal bound, or have difficulty in
283 maintaining their bounds at the same level across trials. This is because the
284 reward rate drops off more quickly below than above the optimal bounds
285 (Figure. 4a). Thus, if the subject's bounds fluctuate across trials, or the subjects
286 are uncertain about the optimal bounds, they should aim at setting their bounds
287 above rather than below this optimum. Indeed, this would minimize the
288 probability that the bound would fluctuate well below the optimal value, which
289 would result in a very sharp drop in reward rate. However, our data indicate
290 that, in contrast to previous findings from single-modality tasks (Bogacz et al.,
291 2010; Simen et al., 2009), subjects consistently set their bounds below the
292 optimum level (Figure 4b). In other words, they make faster and less accurate
293 decisions than predicted by either of the above considerations. Figure 1a (data
294 vs. tuned) illustrates an extreme case for subject B2, in which the best reward
295 rate is achieved in some conditions by waiting until stimulus offset. While not

296 always as extreme as shown for this subject, a distinct discrepancy between
297 observed and reward rate-maximizing behavior exists for all subjects, and is a
298 reflection of the fact that near-optimal reward rates can be achieved with
299 remarkably different joint tunings of reaction times and choice accuracy.

300 What are the potential neural correlates of the highly flexible decision
301 bounds and associated SATs that are reflected in the subjects' behavior? One
302 possibility is the observed bound on neural activity (Churchland, Kiani, &
303 Shadlen, 2008; Kiani, Hanks, & Shadlen, 2008; Roitman & Shadlen, 2002; Schall,
304 2003) in the lateral intraparietal cortex in monkeys, an area that seems to reflect
305 the accumulation of noisy and ambiguous evidence (Yang & Shadlen, 2007). It
306 still needs to be clarified if similar mechanisms are involved in our experimental
307 setup, in which we observed modality-dependent trial-by-trial changes in the
308 SAT. In contrast to suggestions from neuroimaging studies (Green, Biele, &
309 Heekeren, 2012), such trial-by-trial changes are unlikely to emerge from slow
310 changes in connectivity. A more likely alternative, that is compatible with
311 neurophysiological findings, is a neuronal "urgency signal" that modulates this
312 trade-off by how quickly it drives decision-related neuronal activity to a common
313 decision threshold (Hanks, Kiani, & Shadlen, 2014). Although only observed for
314 blocked designs, a similar modality-dependent urgency signal could account for
315 the trial-by-trial SAT changes of our experiment, and qualitatively mimic a
316 change in diffusion model bounds. Currently, our model can only predict changes
317 in scaled decision boundaries, which conflate actual boundary levels with the
318 diffusion standard deviation. It does not predict how the actual bound level
319 changes, which is the quantity that relates to the magnitude of such an urgency
320 signal. In general, quantitatively relating diffusion model parameters to neural
321 activity strongly depends on how specific neural populations encode
322 accumulated evidence, which has only been investigated for cases that are
323 substantially simpler (e.g., Kira, Yang, & Shadlen, 2015) than the ones we
324 consider here.

325 Further qualitative evidence for neural mechanisms that support trial-by-
326 trial changes in the SAT comes from monkeys performing a visual search task
327 with different, visually cued, response deadlines (Heitz & Schall, 2012). Even
328 though the different deadline conditions were blocked, analysis of FEF neural

activity revealed a change in baseline activity that emerged already in the first trial of each consecutive block, hinting at flexible mechanisms that pre-emptively govern changes in SAT. In general, such changes in SAT are likely to emerge through orchestrated changes in multiple neural mechanisms, such as changes in baseline, visual gain, duration of perceptual processing, and the other effects observed by Heitz and Schall (2012), or through combined changes to perceptual processing and motor preparation, as suggested by Salinas, Scerra, Hauser, Costello, and Stanford (2014).

The observed SATs support the hypothesis that gradient-based information is used by subjects during the initial training trials to try to learn the optimal bound settings. We do not make strong assumptions about exactly how this training information is used, and even a very simple strategy of occasional bound adjustments in the light of positive or negative feedback is, in fact, gradient-based (albeit not very efficient) (e.g., Myung & Busemeyer, 1989). The clearest example of a strategy that is not gradient-based is one that does not at all adjust the SAT, or one that does so randomly, without regard to the error feedback that was given to subjects during the initial training period. Such strategies are not guaranteed to lead to the consistent curvature/bound distance relationship observed in Figure. 3b. For a single speed/accuracy trade-off, adjusting this trade-off has already been thoroughly investigated, albeit with conflicting results (Balci et al., 2011; Myung & Busemeyer, 1989; Simen, Cohen, & Holmes, 2006; Simen et al., 2009). Greater insight into the dynamics of learning this trade-off will require further experiments that keep the task stable throughout acquisition of the strategy, and reduce the number of conditions and potential confounds to explain the observed changes in behavior.

In summary, we have shown that subjects performing a multisensory reaction-time task tune their SAT to achieve reward rates close to those achievable by the best-tuned model. This near-optimal performance is invariant under various assumptions about how the reward rate is computed, and is, even under the most conservative assumptions, in the range of 90% of the optimal reward rate. Deviations from optimality are unlikely to have emerged from a strategy of setting bounds to make them robust to perturbations. Instead, our data support the idea that decision bounds have been tuned by a gradient-based

strategy. Such tuning is also in line with the observation of near-optimal reward rates, which are unlikely to result from a random bound-setting strategy. Overall, our study provides novel insights into the flexibility with which human decision makers choose between speed and accuracy of their choices.

Materials and Methods

Seven subjects (3 male) aged 23-38 years participated in a reaction-time version of a heading discrimination task with three different coherence levels of the visual stimulus. Of these subjects, three (subjects B, D, F; 1 male) participated in a follow-up experiment with six coherence levels. The six-coherence version of their data is referred to as B2, D2, and F2. More details about the subjects and the task can be found in J. Drugowitsch et al. (2014). Not discussed in this reference is the inter-trial interval, which is the time from decision to stimulus onset in the next trial. This interval is required to compute the reward rate, and was 6s on average across trials.

Unless otherwise noted, we used a variant of the modified diffusion model described in J. Drugowitsch et al. (2014) to fit the subjects' behavior, and we tuned its parameters to maximize reward rates. Rather than using a constant decision bound for each modality and parameterizing how the diffusion variance depends on the coherence of visual motion (as in J. Drugowitsch et al., 2014), the model variant used here allowed for a separate bound/variance combination per modality and coherence. Thus, it featured 7 bound parameters for the 3-coherence experiments, and 13 bound parameters for the 6-coherence experiments. This variant was chosen to increase the model's flexibility when maximizing its reward rate. The original model variant with constant bounds and a changing variance led to qualitatively comparable results (Figure. 1a, "tuned", and Figure. 2b, "parametric bounds").

For each subject, we adjusted the model's parameters to fit the subject's behavior as in J. Drugowitsch et al. (2014), through a combination of posterior sampling and gradient ascent. Based on these maximum-likelihood parameters, we then found the model parameters that maximized reward rate by adjusting the bound/variance parameters using gradient ascent on the reward rate, while

394 keeping all other model parameters fixed. To avoid getting trapped in local
395 maxima, we performed this maximization 50 times with random re-starts, and
396 chose the parameters that led to the overall highest reward rate. When
397 performing the maximization, we only modified the parameters controlling the
398 bounds, while keeping all other parameters fixed to the maximum-likelihood
399 values. The latter differed across subjects, such that this maximization led to
400 different maximum reward rates for different subjects. For the “no bias” variant
401 in Figure. 2b, we set the choice biases to zero before performing the reward rate
402 maximization.

403 In all cases, the reward rate was computed as the fraction of correct
404 choices across trials, divided by the average trial time, which is the time between
405 the onsets of consecutive trials. Any non-zero evidence accumulation cost
406 (Figures. 2 and 4) was first multiplied with the average decision time (that is,
407 reaction time minus estimated non-decision time) across all trials, and then
408 subtracted from the numerator.

409 Our argument about the speed of convergence of steepest gradient ascent
410 is based on the assumption that bounds are updated according to $\theta^n = \theta^{n-1} +$
411 $\alpha \nabla f(\theta^{n-1})$, where θ^{n-1} and θ^n are the bound vectors before and during the n th
412 steepest gradient ascent step, $f(\theta)$ and $\nabla f(\theta)$ are the reward rate and its
413 gradient for bounds θ , and α is the step size. The speed of this procedure (i.e. the
414 bound change between consecutive steps) depends for each bound on the size of
415 the corresponding element in the reward rate gradient. For optimal bounds, this
416 gradient is zero, which makes the gradient itself unsuitable as a measure of
417 gradient ascent speed. Instead, we use the rate of change of this gradient close to
418 the bounds $\hat{\theta}$ estimated for individual subjects. This rate of change, called the
419 *curvature*, is proportional to the gradient close to $\hat{\theta}$, and therefore also
420 proportional to the speed at which $\hat{\theta}$ is approached. Close-to and at the optimal
421 reward rate, which is a maximum, this curvature is negative. As we were more
422 interested in its size than its sign, Figure. 3 shows the absolute value of this
423 curvature. We estimated this curvature at $\hat{\theta}$ by computing the Hessian of $f(\hat{\theta})$ by
424 finite differences [D’Errico, John (2006). Adaptive Robust Numerical
425 Differentiation. MATLAB Central File Exchange. Retrieved July 3, 2014], where
426 we used the model that allowed for a different bound level per modality and

427 coherence (7 and 13 bound parameters/dimensions for 3 and 6 coherence
428 experiment, respectively). Before computing the distance between estimated and
429 reward rate-maximizing bounds, we projected bound parameter vectors into the
430 eigenspace of this Hessian, corresponding to the orientations of decreasing
431 curvature strength. The absolute bound difference was then computed for each
432 dimension (i.e. modality and coherence) of this eigenspace separately, with the
433 corresponding curvature given by the associated eigenvalue (Figure. 3b).

434 In Figure. 3b, each bound dimension (i.e. modality and coherence, see
435 figure legend) is associated with a different color. As described in the previous
436 paragraph, this figure shows bound differences and curvatures not in the space
437 of original bound levels, but rather in a projected space. To illustrate this bound
438 coordinate transformation in the figure colors, we performed the same
439 coordinate transform on the RGB values associated with each dimension, to find
440 the colors associated with the dimensions of the projected space. The projected
441 colors (filled circles in Figure. 3b plot) closely match the original ones (Figure. 3b
442 legend), which reveals that the curvature eigenspace is well aligned to that of the
443 bound parameters. This indicates that the reward rate curvatures associated
444 with each of the bound parameters, that is, each modality/coherence
445 combination, are fairly independent. Due to the close match between projected
446 and original colors, we do not mention the color transformation in the legend of
447 Figure. 3.

448 Our analysis is also valid if subjects do not follow the reward rate gradient
449 explicitly. They could, for example, approximate this gradient stochastically on a
450 step-by-step basis. As long as the stochastic approximation is unbiased, our
451 argument still holds. One such stochastic approximation would be to test if a
452 change in a single bound (corresponding to a single trial) improves the noisy
453 estimate of the reward rate, that is, if $f(\theta^n) > f(\theta^{n-1}) + \varepsilon$, where only a single
454 element (i.e. bound) is changed between θ^{n-1} and θ^n , and ε is zero-mean
455 symmetric random noise. In this case, larger changes, which are more likely to
456 occur in directions of larger gradient, are more likely accepted. As a result, faster
457 progress is made along steeper directions, which is the basic premise upon
458 which our analysis is based.

To illustrate how the reward rate changed with bound height (Figure. 4), we assumed that all (7 or 13) bound parameters varied along a straight line drawn from the origin to the reward rate-maximizing parameter settings. To project the maximum-likelihood bound parameters from the subject fits onto this line (dots in Figure. 4a&b), we followed the iso-reward rate contour from these parameters until they intersected with the line. We also tried an alternative approach by projecting these parameters onto the line by vector projection, which resulted in a change of the reward rate, but otherwise led to qualitatively similar results as those shown in Figure. 4b. In both cases, the subjects' bound parameters were well below those found to maximize the reward rate.

Figure Legends

Figure 1. The SAT and reward rate for unimodal vs. combined conditions.

(a) Fraction of correct choices is plotted as a function of mean reaction time for subject B2. Blue/cyan: visual condition; green/lime: vestibular condition; red/orange: combined condition. Solid: data; dashed: model with parametric bound tuned to maximize reward rate. Motion coherence varies across data points in the red/orange and blue/cyan curves. The tuned model generally predicts slower and more accurate choices in the visual condition, leading to the longest-possible reaction time (2s stimulus time + non-decision time) for all but the highest stimulus coherence. (b) Reward rate for trials of the combined condition is plotted against reward rate for trials of the visual condition (open blue symbols), the vestibular condition (open green symbols) and both unimodal conditions in combination (gray filled symbols). Reward rates are computed as number of correct decisions per unit time for the respective trial subgroups, and are shown for each subject separately, with bootstrapped 95% confidence intervals.

Figure 2. Reward rates of subjects relative to the optimal reward rate. The optimal reward rate is the best reward rate achievable by a model with tuned decision bounds. (a) Each subject's reward rate is shown as a fraction of the

optimal reward rate (blue bars). In addition, the expected reward rate is shown for immediate random decisions (red bars). (b) Box-plots show relative reward rates for different assumptions regarding how reward rate is computed. ‘no cost’ corresponds to the case shown in panel *a*. ‘cost 0.1’ and ‘cost 0.2’ assume a cost per second for accumulating evidence over time. ‘parametric bounds’ uses the original bounds from J. Drugowitsch et al. (2014), rather than a separate bound parameter for each modality and coherence. The last two bars (‘unbiased’) remove the subjects’ decision biases before computing the optimal reward rate. All box-plots show the maximum/minimum relative reward rates (whiskers), the 25% and 75% percentiles (central bar), and the median (central line) value across subjects. Data are shown for the subjects’ reward rates (blue) and for immediate random choices (red).

Figure 3. Evidence for bound mistuning due to incomplete gradient-based learning. (a) The effects of incomplete gradient ascent on the relation between projected bound distance and local curvature (that is, second derivative of the reward rate at estimated bound) are illustrated for a fictional maximization problem with only two bounds. The grey trajectory shows a sequence of gradient ascent steps on the reward rate function, whose shape is illustrated by two iso-reward rate contours (black) around its maximum (cross). Stopping this gradient ascent procedure (large grey filled circle) before it reaches the optimum causes this stopping point to be close to the optimal bound in directions of large curvature (red), and farther away from the optimum in directions of shallow curvature (blue). (b) Curvature at the estimated bound location is plotted against the distance between the estimated and optimal bound (see text for details). This plot includes 7 (3 coherence condition) or 13 (6 coherence condition) data points per subject, one for each modality/coherence combination. Data for the visual, vestibular and combined conditions are shown in shades of blue/cyan, green, and red/yellow, respectively, and motion coherence is indicated by color saturation. (c) The reward rate loss (i.e., optimal model reward rate minus subject’s reward rate) as estimated from the model (abscissa) is plotted against the loss predicted by the quadratic approximation used in the analysis in (a)-(b), for each subject (ordinate). If the reward rate has a quadratic dependence on the

524 bounds, then all the data points would lie along the diagonal. Small deviations
525 from the diagonal indicate that the reward rate is indeed close-to-quadratic in
526 these bounds.

527

528 **Figure 4. Subjects' bound settings relative to the optimal bound.** (a) The
529 curves show how the reward rate changes with a simultaneous, linear change of
530 all bounds. From left to right, bound levels increase from zero to the (reward-
531 rate maximizing) optimal bound levels (unity values on the abscissa), and
532 continue to bound levels well above this optimum. Different colors correspond to
533 different assumptions about the cost for accumulating evidence over time. The
534 optimal bound levels (unity values on the abscissa) that maximize the reward
535 rate depend on these costs, and thus differ between the three curves. The
536 empirical bound level estimates for individual subjects do not lie on the straight
537 line that is defined by the simultaneous, linear change of all optimal bounds. To
538 evaluate where these empirical bounds lie with respect to the optimal bounds,
539 we found the closest point (along contours of equal reward rate) on this line for
540 the empirical bounds. These points are shown for subject A for different costs by
541 the filled circles. (b) The closest points are illustrated for all subjects, for
542 different accumulation costs. As can be seen, for any assumption for this cost, the
543 subjects' bounds are well below the optimal settings.

544

545 **References**

- 546 Balci, F., Simen, P., Niyogi, R., Saxe, A., Hughes, J. A., Holmes, P., & Cohen, J. D.
547 (2011). Acquisition of decision making criteria: reward rate ultimately
548 beats accuracy. *Atten Percept Psychophys*, 73(2), 640-657. doi:
549 10.3758/s13414-010-0049-7
- 550 Bogacz, R., Hu, P. T., Holmes, P. J., & Cohen, J. D. (2010). Do humans produce the
551 speed-accuracy trade-off that maximizes reward rate? *Q J Exp Psychol*
552 (*Hove*), 63(5), 863-891. doi: 914714341 [pii]
553 10.1080/17470210903091643
- 554 Churchland, A. K., Kiani, R., & Shadlen, M. N. (2008). Decision-making with
555 multiple alternatives. *Nat Neurosci*, 11(6), 693-702. doi: nn.2123 [pii]
556 10.1038/nn.2123
- 557 Drugowitsch, J., DeAngelis, G. C., Klier, E. M., Angelaki, D. E., & Pouget, A. (2014).
558 Optimal multisensory decision-making in a reaction-time task. *Elife*, 3.
559 doi: 10.7554/eLife.03005

560 Drugowitsch, J., Moreno-Bote, R., Churchland, A. K., Shadlen, M. N., & Pouget, A.
 561 (2012). The Cost of Accumulating Evidence in Perceptual Decision
 562 Making. *J Neurosci*, 32(11), 3612-3628. doi: 32/11/3612 [pii]
 563 10.1523/JNEUROSCI.4010-11.2012
 564 Green, N., Biele, G. P., & Heekeren, H. R. (2012). Changes in neural connectivity
 565 underlie decision threshold modulation for reward maximization. *J*
 566 *Neurosci*, 32(43), 14942-14950. doi: 10.1523/JNEUROSCI.0573-12.2012
 567 Hanks, T., Kiani, R., & Shadlen, M. N. (2014). A neural mechanism of speed-
 568 accuracy tradeoff in macaque area LIP. *Elife*, 3. doi: 10.7554/eLife.02260
 569 Heitz, R. P., & Schall, J. D. (2012). Neural mechanisms of speed-accuracy tradeoff.
 570 *Neuron*, 76(3), 616-628. doi: 10.1016/j.neuron.2012.08.030
 571 Kiani, R., Hanks, T. D., & Shadlen, M. N. (2008). Bounded integration in parietal
 572 cortex underlies decisions even when viewing duration is dictated by the
 573 environment. *J Neurosci*, 28(12), 3017-3029. doi: 28/12/3017 [pii]
 574 10.1523/JNEUROSCI.4761-07.2008
 575 Kiani, R., & Shadlen, M. N. (2009). Representation of confidence associated with a
 576 decision by neurons in the parietal cortex. *Science*, 324(5928), 759-764.
 577 doi: 324/5928/759 [pii]
 578 10.1126/science.1169405
 579 Kira, S., Yang, T., & Shadlen, M. N. (2015). A neural implementation of Wald's
 580 sequential probability ratio test. *Neuron*, 85(4), 861-873. doi:
 581 10.1016/j.neuron.2015.01.007
 582 Myung, I. J., & Busemeyer, J. R. (1989). Criterion Learning in a Deferred Decision-
 583 Making Task. *American Journal of Psychology*, 102(1), 1-16. doi: Doi
 584 10.2307/1423113
 585 Palmer, J., Huk, A. C., & Shadlen, M. N. (2005). The effect of stimulus strength on
 586 the speed and accuracy of a perceptual decision. *J Vis*, 5(5), 376-404. doi:
 587 10.1167/5.5.1
 588 /5/5/1/ [pii]
 589 Ratcliff, R. (1978). Theory of Memory Retrieval. *Psychological Review*, 85(2), 59-
 590 108.
 591 Ratcliff, R., & McKoon, G. (2008). The diffusion decision model: theory and data
 592 for two-choice decision tasks. *Neural Comput*, 20(4), 873-922. doi:
 593 10.1162/neco.2008.12-06-420
 594 Reddi, B. A., Asrress, K. N., & Carpenter, R. H. (2003). Accuracy, information, and
 595 response time in a saccadic decision task. *J Neurophysiol*, 90(5), 3538-
 596 3546. doi: 10.1152/jn.00689.2002
 597 Reddi, B. A., & Carpenter, R. H. (2000). The influence of urgency on decision time.
 598 *Nat Neurosci*, 3(8), 827-830. doi: 10.1038/77739
 599 Roitman, J. D., & Shadlen, M. N. (2002). Response of neurons in the lateral
 600 intraparietal area during a combined visual discrimination reaction time
 601 task. *J Neurosci*, 22(21), 9475-9489. doi: 22/21/9475 [pii]
 602 Salinas, E., Scerra, V. E., Hauser, C. K., Costello, M. G., & Stanford, T. R. (2014).
 603 Decoupling speed and accuracy in an urgent decision-making task reveals
 604 multiple contributions to their trade-off. *Front Neurosci*, 8, 85. doi:
 605 10.3389/fnins.2014.00085
 606 Schall, J. D. (2003). Neural correlates of decision processes: neural and mental
 607 chronometry. *Curr Opin Neurobiol*, 13(2), 182-186. doi:
 608 S0959438803000394 [pii]

609 Simen, P., Cohen, J. D., & Holmes, P. (2006). Rapid decision threshold modulation
610 by reward rate in a neural network. *Neural Netw*, 19(8), 1013-1026. doi:
611 10.1016/j.neunet.2006.05.038
612 Simen, P., Contreras, D., Buck, C., Hu, P., Holmes, P., & Cohen, J. D. (2009). Reward
613 rate optimization in two-alternative decision making: empirical tests of
614 theoretical predictions. *J Exp Psychol Hum Percept Perform*, 35(6), 1865-
615 1897. doi: 2009-22869-017 [pii]
616 10.1037/a0016926
617 Yang, T., & Shadlen, M. N. (2007). Probabilistic reasoning by neurons. *Nature*,
618 447(7148), 1075-1080. doi: 10.1038/nature05852
619 Zacksenhouse, M., Bogacz, R., & Holmes, P. (2010). Robust versus optimal
620 strategies for two-alternative forced choice tasks. *J Math Psychol*, 54(2),
621 230-246. doi: 10.1016/j.jmp.2009.12.004
622







