# Temporal integration is a robust feature of perceptual decisions

Alexandre Hyafil[1,2], Jaime de la Rocha[3], Cristina Pericas[3], Leor N. Katz[4], Alexander C. Huk[5], Jonathan W. Pillow[2]

[1]Centre de Recerca Matemàtica, Bellaterra (Spain). contact: alexandre.hyafil (at) gmail.com

[2]Princeton Neuroscience Institute, Princeton University, Princeton, New Jersey (USA)

[3]Institut d'Investigacions Biomèdiques August Pi i Sunyer (IDIBAPS), Barcelona, Spain

[4]Laboratory of Sensorimotor Research, National Eye Institute, National Institutes of Health

[5]Fuster Laboratory for Cognitive Neuroscience, Departments of Psychiatry & Biobehavioral Sciences and Ophthalmology, UCLA

Making informed decisions in noisy environments requires integrating sensory information over time. However, recent work has suggested that it may be difficult to determine whether an animal's decision-making strategy relies on evidence integration or not. In particular, strategies based on extrema-detection or random snapshots of the evidence stream may be difficult or even impossible to distinguish from classic evidence integration. Moreover, such non-integration strategies might be surprisingly common in experiments that aimed to study decisions based on integration. To determine whether temporal integration is central to perceptual decision making, we developed a new model-based approach for comparing temporal integration against alternative "non-integration" strategies for tasks in which the sensory signal is composed of discrete stimulus samples. We applied these methods to behavioral data from monkeys, rats, and humans performing a variety of sensory decision-making tasks. In all species and tasks, we found converging evidence in favor of temporal integration. First, in all observers across studies, the integration model better accounted for standard behavioral statistics such as psychometric curves and psychophysical kernels. Second, we found that sensory samples with large evidence do not contribute disproportionately to subject choices, as predicted by an extrema-detection strategy. Finally, we provide a direct confirmation of temporal integration by showing that the sum of both early and late evidence contributed to observer decisions. Overall, our results provide experimental evidence suggesting that temporal integration is an ubiquitous feature in mammalian perceptual decision-making. Our study also highlights the benefits of using experimental paradigms where the temporal stream of sensory evidence is controlled explicitly by the experimenter, and known precisely by the analyst, to characterize the temporal properties of the decision process.

**INTRODUCTION**

Perceptual decision-making is thought to rely on the temporal integration of noisy sensory information on a timescale of hundreds of milliseconds to seconds. Temporal integration corresponds to summing over time the evidence provided by each new sensory stimulus, and optimizes perceptual judgments in face of noise (Bogacz et al. 2006; Gold and Shadlen 2007). A perceptual decision can then be made on the basis of this accumulated evidence, either as some threshold on accumulated evidence is reached, or if some internal or external cue signals the need to initiate a response.

Although many behavioral and neural results are consistent with this integration framework, temporal integration is a feature that has often been taken for granted rather than explicitly tested. Recently, the claim that standard perceptual decision-making tasks rely on (or even frequently elicit) temporal integration has been challenged by theoretical results showing that non-integration strategies can produce behavior that carries superficial signatures of temporal integration (Stine et al. 2020). These signatures include the relationship between stimulus difficulty, stimulus duration and behavioral accuracy, the precise temporal weighting of sensory information on the decisions, and the patterns of reaction times.

Here, we propose new analytical tools for directly assessing integration and non-integration strategies from fixed-duration or variable-duration paradigms where, critically, the experimenter controls the fluctuations in perceptual evidence over time within each trial (discrete-sample stimulus, or DSS). By leveraging these controlled fluctuations, our methods allow us to make direct comparisons between integration and non-integration strategies. We apply these tools to assess temporal integration in data from monkeys, humans and rats that performed a variety of perceptual decision-making tasks with DSS. Applying these analyses to these behavioral datasets yields strong evidence that perceptual decision-making tasks in all three species rely on temporal integration. Temporal integration, a critical element of many major theories of perception at both the neural and behavioral levels, is indeed a robust and pervasive aspect of mammalian behavior. Our results also illuminate the power of targeted stimulus design and statistical analysis to test specific features of behavior.

**RESULTS**

**Integration and non-integration models**

In a typical perceptual evidence-integration experiment (Figure 1A), an observer is presented in each trial with a time-varying stimulus and must report which of two possible stimulus categories it belongs to. Typical examples include judging whether a dynamic visual stimulus is moving leftwards or rightwards (Yates et al. 2017; Katz et al. 2015); whether the orientation of a set of gratings is more aligned with cardinal or diagonal directions (Wyart et al. 2012); whether a combination of tones is dominated by high or low frequencies (Morillon, Schroeder, and Wyart 2014; Hermoso-Mendizabal et al. 2020; Znamenskiy and Zador 2013); or which of two acoustic streams is more intense or dense (Brunton, Botvinick, and Brody 2013; Pardo-Vazquez et al. 2019). Such paradigms have been used extensively in humans, nonhuman primates and rodents. Here we focus on experiments in which observers report their choice at

81     the end of a period whose duration is controlled by the experimenter (Kiani and Shadlen 2009;
82     Wyart et al. 2012; Brunton, Botvinick, and Brody 2013; Raposo et al. 2012), in contrast to so-
83     called "reaction time" tasks, in which the observer can respond after viewing as brief a portion
84     of the stimulus as they wish (Roitman and Shadlen 2002; Znamenskiy and Zador 2013; Pardo-
85     Vazquez et al. 2019; Hermoso-Mendizabal et al. 2020).

86     Moreover, we focus on experimental paradigms in which the sensory evidence in favor of each
87     category arrives in a sequence of discrete *samples*. Samples can correspond to motion pulses
88     (Yates et al. 2017), individual gratings (Wyart et al. 2012), acoustic tones (Morillon, Schroeder,
89     and Wyart 2014; Hermoso-Mendizabal et al. 2020; Znamenskiy and Zador 2013) numbers
90     (Bronfman et al. 2015) or symbols representing category probabilities (Yang and Shadlen
91     2007). We refer to this configuration as the discrete-sample stimulus (DSS) paradigm. In this
92     paradigm, the perceptual evidence provided by each sample can be controlled independently,
93     allowing for detailed analyses of how different samples contribute to the behavioral response.
94     The DSS framework can be contrasted with experiments in which the experimenter specifies
95     only the mean stimulus strength on each trial, and variations in sensory evidence over time
96     are not finely controlled or are not easily determined from the raw spatio-temporal stimulus.

97     Tasks using the DSS paradigm are classically thought to rely on sequential accumulation of
98     the stimulus evidence (Bogacz et al. 2006), which we refer to here as temporal integration.
99     Figure 1A shows an example stimulus sequence composed of $n$ samples that provide differing
100    amounts of evidence in favor of one alternative vs. another ("A" vs. "B"). The accumulated
101    evidence fluctuates as new samples are integrated and finishes at a positive value indicating
102    overall evidence for stimulus category A (Figure 1B). This integration process can be
103    formalized by defining the the decision variable or accumulated evidence $x_i$ and its updating
104    dynamics across stimulus samples: $x_i = x_{i-1} + m_i$ where $m_i = S_i + \varepsilon_i$ represents a noisy
105    version of the true stimulus evidence $S_i$ in the $i$-th sample corrupted by sensory noise $\varepsilon_i$. The
106    binary decision $r$ is simply based on the sign of the accumulated evidence $x_n$ at the end of the
107    sample sequence (composed of $n$ samples): $r = A$ if $x_n > 0$, and $r = B$ if $x_n < 0$. This
108    procedure corresponds to the normative strategy with uniform weighting that maximizes
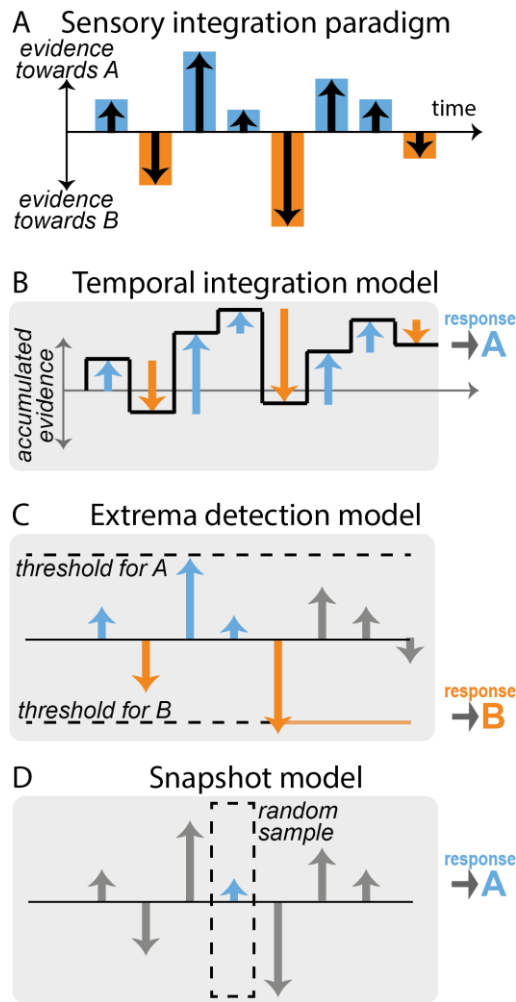
**Figure 1. A.** Schematic of a typical fixed-duration perceptual task with discrete-sample stimuli (DSS). A stimulus is composed of a discrete sequence of *n* samples (here, *n=6*). The subjects must report at the end of the sequence whether one specific quality of the stimulus was "overall" leaning more towards one of two possible categories A or B. Evidence in favor of category A or B varies across samples (blue and orange bars). **B.** Temporal integration model. The relative evidence in favor of each category is accumulated sequentially as each new sample is presented (black line), resulting in temporal integration of the sequence evidence. The choice is determined by the end point of the accumulation process: here, the overall evidence in favor of category A is positive, so response A is selected. **C.** Extrema detection model. A decision is made whenever the instantaneous evidence for a given sample (blue and orange arrows) reaches a certain fixed threshold (dotted lines). The selected choice corresponds to the sign of the evidence of the sample that reaches the threshold (here, response B). Subsequent samples are ignored (gray bars). **D.** Snapshot model. Here, only one sample is attended. Which sample is attended is determined in each trial by a stochastic policy. The response of the model simply depends on the evidence of the attended sample. Other samples are ignored (gray bars). Variants of the model include attending *K>1* sequential samples.

accuracy. For such perfect integration, $x_n = \Sigma_i S_i + \Sigma_i \varepsilon_i$, so that the probability of response A is $p(r = A) = \Phi(\Sigma_i \beta S_i)$ where $\Phi$ is the cumulative normal distribution function (the normative weight for the stimuli β depends on the noise variance $Var(\varepsilon)$ and the number of samples through $\beta = 1/\sqrt{n\,Var(\varepsilon)}$). Departures from optimality in the accumulation process such as accumulation leak, categorization dynamics, sensory adaptation or sticky boundaries may however yield unequal weighting of the different samples (Yates et al. 2017; Brunton, Botvinick, and Brody 2013; Prat-Ortega et al. 2021; Bronfman, Brezis, and Usher 2016). To accommodate for these, we allowed the model to take any arbitrary weighting of the samples: $p(r = A) = \Phi(\beta_0 + \Sigma_i \beta_i S_i)$ (see Methods for details). The mapping from final accumulated evidence to choice was probabilistic, to account for the effects of noise from different sources in the decision-making process (Drugowitsch et al. 2016).

Although it has been commonly assumed that observers use evidence integration strategies to perform these psychophysical tasks, recent work has suggested that observers may employ non-integration strategies instead (Stine et al. 2020). Here we consider two specific alternative models. The first non-integration model corresponds to an *extrema-detection* model (Waskom and Kiani 2018; Stine et al. 2020; Ditterich 2006). In this model, observers do not integrate evidence across samples but instead base their decision on extreme or salient bits of evidence. More specifically, the observer commits to a decision based on the first sample *i* in the stimulus sequence that exceeds one of the two symmetrical thresholds, i.e. such that $|m_i| \geq \theta$. In our example stimulus, the first sample that reaches this threshold in evidence

130　space is the fifth sample, which points towards stimulus category B, so response B is selected
131　(Figure 1C). This policy can be viewed as a memory-less decision process with sticky bounds.
132　If the stimulus sequence contains no extreme samples, so that neither threshold is reached,
133　the observer selects a response at random. (Following (Stine et al. 2020), we also explored
134　an alternative mechanism where in such cases the response is based on the last sample in
135　the sequence).

136　The second non-integration model corresponds to the *snapshot model* (Stine et al. 2020; Pinto
137　et al. 2018). In this model, the observer attends to only one sample *i* within the stimulus
138　sequence, and makes a decision based solely on the evidence from the attended sample: $r =$
139　$A$ if $m_i > 0$, and $r = B$ if $m_i < 0$. The position in the sequence of the attended sample is
140　randomly selected on each trial. In our example, the fourth sample is randomly selected, and
141　since it contains evidence towards stimulus category A, response A is selected (Figure 1D).
142　We considered variants of this model that gave it additional flexibility, including: allowing the
143　prior probability over the attended sample to depend on its position in the sequence using a
144　non-parametric probability mass function estimated from the data; allowing for deterministic
145　vs. probabilistic decision-making rule based on the attended evidence; including attentional
146　lapses that were either fixed to 0.02 (split equally between leftward and rightward responses)
147　or estimated from behavioral data. We finally considered a variant of the snapshot where the
148　decision was made based on a sub-sequence of *K* consecutive samples within the main
149　stimulus sequence ($1 \leq K < n$), rather than based on a single sample.

150

**Standard behavioral statistics favor integration accounts of pulse-based motion perception in primates**

153　To compare the three decision-making models defined above (i.e., temporal integration,
154　extrema-detection, snapshots), we first examined behavioral data from two monkeys
155　performing a fixed-duration motion integration task (Yates et al. 2017). In this experiment,
156　each stimulus was composed of a sequence of 7 motion samples of 150 ms each where the
157　motion strength towards left or right was manipulated independently for each sample. At the
158　end of the stimulus sequence, monkeys reported with a saccade whether the overall sequence
159　contained more motion towards the left or right direction. The animals performed 72137 and
160　33416 trials for monkey N and monkey P respectively, allowing for in-depth dissection of their
161　response patterns.

162　We fit the three models (and their variants) to the responses for each animal individually (see
163　Supplementary Figure 1 for estimated parameters for the different models). We then
164　simulated the fitted model and computed, for simulated and experimental data, the
165　psychophysical kernels capturing the weights of the different sensory samples based on their
166　position in the stimulus sequence (Figure 2B). Psychophysical kernels were non-monotonic
167　and differed in shape between the two animals, probably reflecting the complex contributions
168　of various dynamics and sub-optimalities along the sensory and decision pathways (Yates et
169　al. 2017).

170

171　The temporal profile of the kernel was perfectly matched by the integration model, almost by
172　design, as we gave full flexibility to the model to adjust the sample weights. The snapshot

173 model was provided with similar flexibility, as the prior probability of attending each sample
174 could be fully adjusted to the monkey decisions. Surprisingly, however, the snapshot model
175 could not match the experimental psychophysical kernel as accurately. It consistently
176 underestimated the magnitude of weighting in monkey P (Figure 2B, bottom row). The
177 extrema-detection model was not endowed with such flexibility of sensory weighting. On the
178 contrary, since the decision was based on the first sample in the sequence reaching a certain
179 criterion, this inevitably generates a primacy effect in the psychophysical kernels - or at best
180 a flat weighting (Stine et al. 2020). The model thus failed to capture the non-monotonic
181 psychophysical kernels from animal data.

182 Next, we looked at the psychometric curves and choice accuracy predictions of each fitted
183 model (Figure 2C-D). Stine and colleagues have argued that integration and non-integration
184 models can capture the psychometric curves equally well (Stine et al. 2020). For both animals,
185 the accuracy and psychometric curves were accurately captured by the integration model.  In
186 line with Stine and colleagues, we also found that both non-integration models could
187 reproduce the shape of the psychometric curve in monkey N, although the quantitative fit was
188 always better for the integration than non-integration models. By contrast both non-integration
189 models failed to capture the psychometric curve for monkey P (Figure 2B, bottom row). More
190 systematically, the overall accuracy, which is an aggregate measure of the psychometric
191 curve, clearly differs between models, as the accuracy of the non-integration models
192 systematically deviated from animal data for both animals (Figure 2C). In other words, all
193 models produce the same type of psychometric curves up to a scaling factor, and this scaling
194 factor (directly linked to the model accuracy) is key to differentiate model fits. For the snapshot
195 model in monkey P, this discrepancy was explained because the model, limited to using one
196 stimulus sample, could not reach the performance of the model (compare the maximum
197 accuracy of the model indicated by the blue mark with the accuracy of the animal), as the
198 snapshot model is limited to making decisions based on one sensory sample only. (This also
199 explains why the psychophysical kernel of the snapshot model underestimated the true kernel
200 in monkey P). For the extrema-detection model in monkey P and for both non-integration
201 models in the other animal (monkey N) and for the extrema-detection model, the model
202 accuracy is not bounded below the subject's accuracy. In such cases, the model can produce
203 better-than-observed accuracy for certain parameter ranges, but these are not the parameters
204 found by the maximum likelihood procedure, probably because they produce a pattern of
205 errors that is inconsistent with the observed pattern of errors. This indicates an inability of the
206 models to match the pattern of errors of the animal (see Discussion).

207 Finally, we assessed quantitatively which model provided the best fit, while correcting for
208 model complexity using the Akaike Information Criterion (AIC, Figure 2A). In both monkeys,
209 AIC favored the integration model over the two non-integration models by a very large margin.
210 We also explored whether variants of the extrema-detection and snapshot models could
211 provide a better match to the behavioral metrics considered above (Supp Figure 2 & 3). We
212 found using the AIC metric that the integration model was preferred over all variants of both
213 non-integration models, for both monkeys. Moreover, these model variants could not replicate
214 the psychophysical kernels as well as the integration model did (Supp Figure 2 & 3). In
215 conclusion, while psychometric curves may not always discriminate between integration and
216 non-integration strategies, other metrics including psychophysical kernels, predicted accuracy
217 and quality of fit (AIC) support temporal integration in monkey perceptual decisions. For one
218 model in one monkey (the snapshot model in monkey P), even the simple metric of overall

accuracy compellingly supported temporal integration (Fig. 2C). For the other monkey and/or model, where the distinction was less clear, our model-based approach allowed us to leverage these other metrics to reveal strong support for the temporal integration model (Fig. 2A-C). Although these data relies only on two experimental subjects, we show below further evidence supporting the integration model in humans and rats.
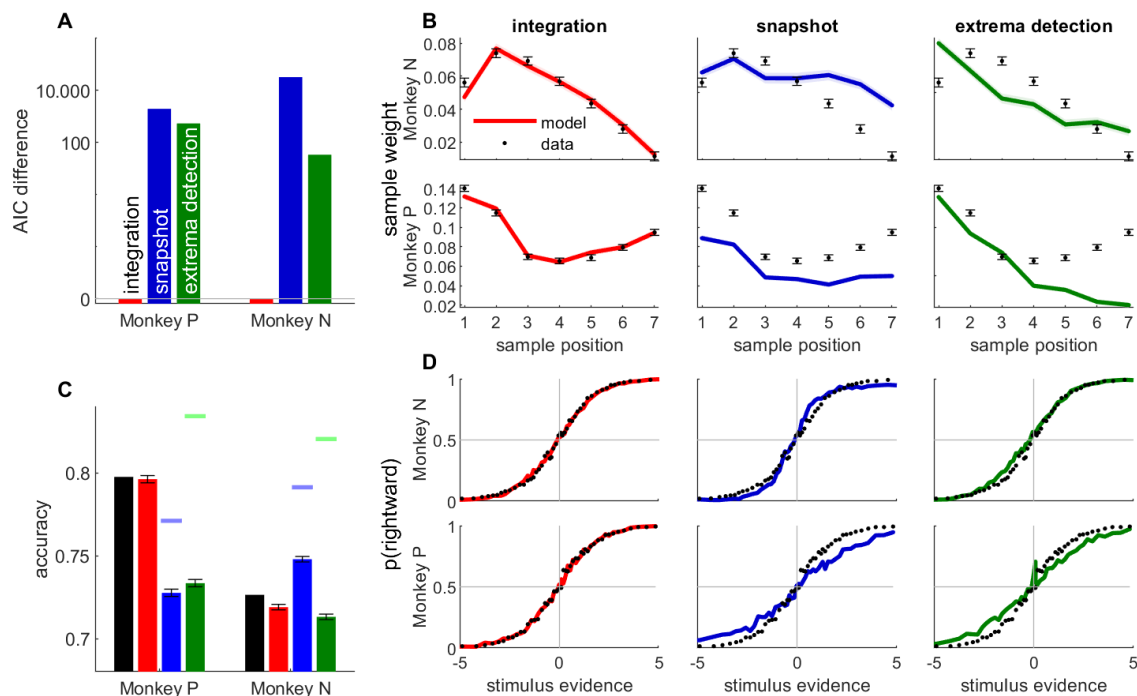


**Figure 2. The integration model better described monkey behavior than non-integration models.** **A.** Difference between AIC of models (temporal integration: red bar; snapshot model: blue; extrema-detection model: green) and temporal integration model for each monkey. Positive values indicate poorer fit to data. **B.** Psychophysical kernels for behavioral data (black dots) vs. simulated data from temporal integration model (left panel, red curve), snapshot model (middle panel, blue curve) and extrema-detection model (right panel, green curve) for the two animal (monkey N: top panels; monkey P: bottom panels). Each data point represents the weight of the motion pulse at the corresponding position on the animal/model response. Error bars and shadowed areas represent the standard error of the weights for animal and simulated data, respectively. **C.** Accuracy of animal responses (black bars) vs simulated data from fitted models (colour bars), for each monkey. Blue and green marks indicate the maximum performance for the snapshot and extrema-detection models, respectively. Error bars represent standard error of the mean. **D.** Psychometric curves for animal (black dots) and simulated data (colour lines) for monkey N, representing the proportion of rightward choices per quantile of weighted stimulus evidence.

**Monkey responses where the largest evidence sample is at odds with the overall stimulus sequence are inconsistent with the extrema-detection model**

7

244 While formal model comparison leads us to reject the non-integration models in favor of the
245 integration models, it is informative to examine qualitative features of the animal strategies
246 and identify how non-integration models failed to capture them. We started by designing two
247 analyses aimed at testing whether choices were consistent with the extrema-detection model,
248 namely by testing whether choices were strongly correlated with the largest-evidence
249 samples. In the first analysis, we looked at the subset of trials where the evidence provided
250 by the largest-evidence sample in the sequence was at odds with the total evidence in the
251 sequence: we show one example in Figure 3B, where the largest evidence sample points
252 towards response B, while the overall evidence points towards response A. These '*disagree
253 trials'* represent a substantial minority of the whole dataset: 1865 trials (2.6%) in monkey N,
254 1831 trials (5.5%) in monkey P. If integration is present, the response of the animal should in
255 general be aligned with the total evidence from the sequence (Figure 3A, red bars). By
256 contrast, if it followed the extrema-detection model (Figure 1C), it should in general follow the
257 largest evidence sample (Figure 3A, green bars). In both monkeys, animal choices were more
258 often than not aligned with the integrated evidence (Figure 3A, black bars), as predicted by
259 the integration model. The responses generated from the extrema-detection model tended to
260 align more with the largest evidence sample, although that behaviour was somehow erratic
261 (for monkey N) due to the large estimated decision noise in the model. This rules out that
262 monkey decisions rely on a memoryless strategy of simply detecting large evidence samples,
263 discarding all information provided by lower evidence samples. Our results complement a
264 previous analysis on disagree trials in this task (Levi et al. 2018), by explicitly comparing
265 monkey behavior to model predictions.

266

267 We reasoned that the extrema-detection would also leave a clear signature in the "subjective
268 weight" of the samples, defined as the impact of each sample on the decision as a function of
269 absolute sample evidence (Yang and Shadlen 2007; Waskom and Kiani 2018; Nienborg and
270 Cumming 2007). The extrema-detection model predicts that, in principle, samples whose
271 evidence is below the threshold have little impact on the decision, while samples whose
272 evidence is above the threshold have full impact on the decision. By contrast, the integration
273 model predicts that subjective weight should grow linearly with sample evidence. We
274 estimated subjective weights from monkey choices using a regression method similar in spirit
275 to previous methods (Yang and Shadlen 2007; Waskom and Kiani 2018), taking the form
276 $p(r_t = A) = \sigma(\beta_0 \Sigma_{i \in [1..n]} \beta_i f(S_{ti}))$. Here $f$ is a function that captures the subjective weight of
277 the sample as a function of its associated evidence. Whereas previous methods estimated
278 subjective weights assuming a uniform psychophysical kernel, our method estimated
279 simultaneously subjective weights $f(S)$ and the psychophysical kernel $\beta$, thus removing
280 potential estimation biases due to unequal weighting of sample evidence (see Methods). In
281 both monkeys, we indeed found that the subjective weight depends linearly on sample
282 evidence for low to median values of sample evidence (motion pulse lower than 6), in
283 agreement with the integration model (Supp. Figure 4). Surprisingly however, simulated data
284 of the extrema-detection model displayed the same linear pattern for low to median values of
285 sample evidence. We realized this was due to the very high estimated sensory noise (Supp
286 Fig 1), such that, according to the model, even samples with minimal sample evidence were
287 likely to reach the extrema-detection threshold. In other words, unlike the previous analyses,
288 inferring the subjective weights used by animals was inconclusive as to whether animals
289 deployed the extrema-detection strategy. This somewhat surprising dependency reinforces

290    the importance of validating intuitions by fitting and simulating models (Wilson and Collins
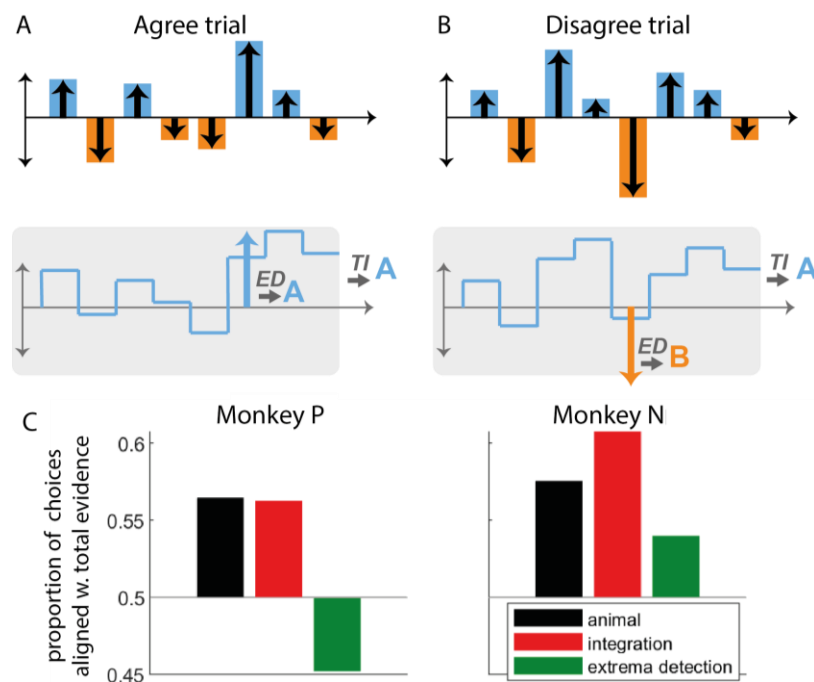291    2019).



**Figure 3. The pattern of animal choices is incompatible with extrema-value based decisions. A.** Example of an 'agree trial' where the total stimulus evidence (accumulated over samples) and the evidence from the largest-evidence sample point towards the same response (here, response A). In this case, we expect that temporal integration (TI) and extrema-detection (ED) will produce similar responses (here, A). **B.** Example of a 'disagree trial', where the total stimulus evidence and evidence from the largest-evidence sample point towards opposite responses (here A for the former; B for the latter). In this case, we expect that integration and extrema-detection models will produce opposite responses. **C.** Proportion of choices out of all *disagree trials* aligned with total evidence, for animal (black bars), integration (red) and extrema-detection model (green).


## Choice dependence on early and late stimulus evidence show direct evidence for temporal integration

305    Following model comparisons favoring integration over both snapshot and extrema-detection
306    models, the  immediately previous analysis  relied on a special subset of trials to provide an
307    additional, and perhaps more intuitive, signature of integration, which ruled out extrema-
308    detection as a possible strategy of either monkey. We next employed another novel analysis
309    specifically designed to tease apart unique signatures of the integration and snapshot models.
310    More specifically, we tested whether decisions were based on the information from only one
311    part of the sequence, as predicted by the snapshot model, or from the full sequence, as
312    predicted by the integration model. To facilitate the analysis, we defined *early evidence $E_t$* by
313    grouping evidence from the first three samples in the sequence, and *late evidence $L_t$,* as the
314    grouped evidence from the last four samples. We then displayed the proportion of rightward
315    responses as a function of both early and late evidence in a graphical representation that we
316    call *integration map* (Figure 4A). A pure integration strategy corresponds to summing early
317    and late evidence equally, which can be formalized as $p(r) = \sigma(E_t + L_t)$, where $\sigma$ is a

9

318 sigmoidal function. Because this only depends on the sum $E_t + L_t$, the probability of response
319 is invariant to changes in the $(E_t, L_t)$ space along the diagonal, which leaves the sum
320 unchanged. These diagonals correspond to isolines of the integration map (Figure 4A, left;
321 Supp Figure 5A). In other words, straight diagonal isolines in the integration map reflect the
322 fact that the decision only depends on the sum of evidence $E_t + L_t$. Straight isolines thus
323 constitute a specific signature of evidence integration.

324 We contrasted this integration map with the one obtained from a non-integration strategy
325 (Figure 4A middle panel; Supp Figure 5A). There we assumed that the decision depends either
326 on the early evidence or on the late evidence, as in the snapshot model, with equal probability.
327 This can be formalized as $p(r) = 0.5\sigma(E_t) + 0.5\sigma(L_t)$. In this case, if late evidence is null
328 $(\sigma(L_t) = 0.5)$ and early evidence is very strong toward the right $(\sigma(E_t) \simeq 1)$ the overall
329 probability for rightward response is $p(r) = 0.75$. This probability contrasts with that obtained
330 in the integration case where the early evidence would dominate and lead to an overwhelming
331 proportion of rightward responses, i.e. $p(r) \simeq 1$. The 25% of leftwards responses yielded by
332 the non-integration model correspond to trials where only the late (uninformative) part of the
333 stimulus is attended and a random response to the left is drawn. More generally, in regions of
334 the space in which either early or late evidence take large absolute values, their corresponding
335 probability of choice saturates to 0 or 1, when that evidence is attended, so the overall
336 response probability becomes only sensitive to the other evidence. As a result, the
337 equiprobable lines bend towards the horizontal and vertical axes (Figure 4A middle). Finally,
338 to compare predictions from both integration and non-integration models to monkey behavior,
339 we plotted the integration maps for both monkeys (Figure 4A, right; Supp Figure 5A). The
340 isolines were almost straight diagonal lines and showed no consistent curvature towards the
341 horizontal and vertical axes. This provides direct evidence that monkey responses depend
342 directly on the sum of early and late evidence— a clear signature of temporal integration.

343 We derived subsequent tests based on the integration map. We computed conditional
344 psychometric curves as the probability for rightward responses as a function of early evidence
345 $E_t$, conditioned on late evidence value $L_t$ (Figure 4B; Supp Figure 5B). From the integration
346 formula $p(r) = \sigma(E_t + L_t)$, we see that a change in late evidence value corresponds to a
347 horizontal shift of the conditional psychometric curves. By contrast, according to the non-
348 integration formula $p(r) = 0.5\sigma(E_t) + 0.5\sigma(L_t)$, conditioning on different values of late
349 evidence adds a fixed value to the response probability irrespective of early evidence, a
350 vertical shift akin to that introduced by lapse responses (Figure 4B middle panel). The
351 conditional psychometric curves for monkeys (Figure 4B right panel; Supp Fig 5 & 6) displayed
352 horizontal shifts as late evidence was changed, consistently with the integration hypothesis.
353 We sought to quantify these shifts in better detail. To this purpose, we fitted each conditional
354 psychometric curve with the formula $p(r) = (1 - \pi_L - \pi_R)\,\sigma(\alpha E_t + \beta) + \pi_R$, where $\pi_L, \pi_R, \alpha$
355 and $\beta$ correspond to the left lapse, right lapse, sensitivity and lateral bias parameters,
356 respectively (Figure 4C, Supp Fig 5 & 6). The integration model predicts that the bias
357 parameter $\beta$ should vary linearly with $L_t$, while lapse parameters should remain null (Figure
358 4D, left panel). By contrast, the non-integration model predicts that the horizontal shift
359 parameter $\beta$ should remain constant while left and right lapse parameters $(\pi_L, \pi_R)$ should vary
360 (middle panel), as these lapse parameters correspond to the trials where early evidence is not
361 attended and the response depends simply on late evidence. Both monkeys showed a very
362 strong linear dependence between late evidence and the horizontal shift $\beta$ (Figure 4D, right
363 panel; see also Supp Fig 5), further supporting that late evidence is summed to early evidence.

364 By contrast, the lapse parameters showed no consistent relationship with late evidence $L_t$
365 (Figure 4E, right panel). Finally, we directly assessed the similarities between the integration
366 maps from monkey responses and from simulated responses for the three models (integration,
367 snapshot, extrema-detection). The model-data correlation was larger in the integration model
368 than in the non-integration strategies for both monkeys (Figure 4E; unpaired t-test on
369 bootstrapped $r$ values: $p<0.001$ for each animal and comparison against extrema-detection
370 and against snapshot model). Overall, integration maps allow to dissect how early and late
371 parts of the stimulus sequence are combined to produce a behavioral response. In both
372 monkeys, these maps carried signatures of temporal integration. For monkey P, the integration
373 model and the data look very similar. For monkey N, there is still a qualitative dependency that
374 deviates from non-integration, but which is not as uniquely matched to the integration strategy
375 (although the imperfect coverage of the two-dimensional space impedes further investigations).
376 Thus, complementing the statistical model tests favoring integration, this richer visualization
377 allows the data to show us that some degree of integration is occurring, albeit not perfect.
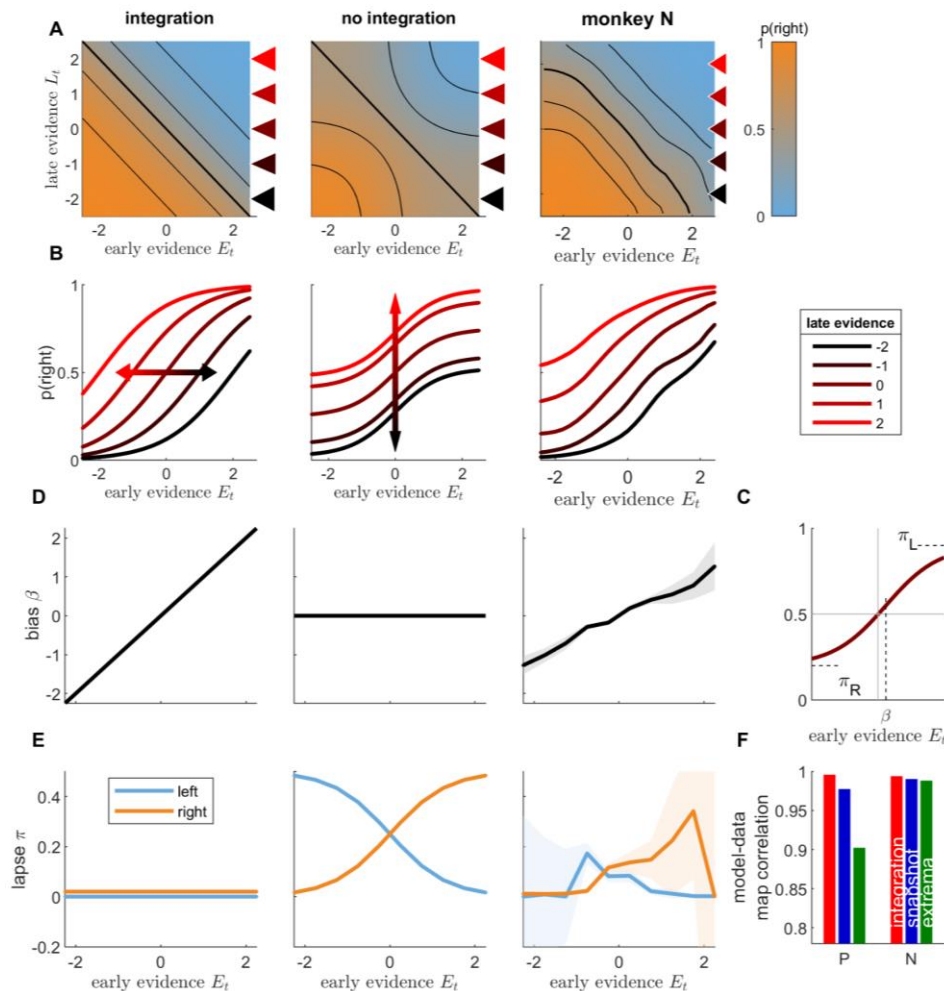
378



379

380 **Figure 4. Integration of early and late evidence into animal responses is incompatible with the**
381 **snapshot model. A.** Integration map representing the probability of rightward responses (orange: high

11

382    probability; blue: low probability) as a function of early stimulus evidence $E_t$ and late stimulus evidence
383    $L_t$, illustrated for a toy integration model (where $p(right) = \sigma(E_t + L_t)$; left panel) and a toy non-
384    integration model ($p(right) = 0.5\sigma(E_t) + 0.5\sigma(L_t)$; middle panel), and computed for monkey N
385    responses (right panel). Black lines represent the isolines for p(rightwards)=0.15, 0.3, 0.5, 0.7 and 0.85.
386    **B.** Conditional psychometric curves representing the probability for rightward response as a function of
387    early evidence $E_t$, for different values of late evidence $L_t$ (see inset for $L_t$ values), for toy models and
388    monkey N. The curves correspond to horizontal cuts in the integration maps at $L_t$ values marked by
389    colour triangles in panel A. **C.** Illustration of the fits to conditional psychometric curves. The value of the
390    bias $\beta$, left lapse $\pi_L$ and right lapse $\pi_R$ are estimated from the conditional psychometric curves for each
391    value of late evidence. **D.** Lateral bias as a function of late evidence for toy models and monkey N.
392    Shaded areas represent standard error of weights for animal data. **E.** Lapse parameters (blue: left lapse;
393    orange: right lapse) as a function of late evidence for toy models and monkey N. **F.** Pearson correlation
394    between integration maps for animal data and integration maps for simulated data, for each animal.
395    Red: integration model; blue: snapshot model; green: extrema-detection model.

396

### Temporal integration in human visual orientation judgments

398    Overall, all our analyses converged to support the idea that monkey decisions in a fixed-
399    duration motion discrimination task relied on temporal integration. We explored whether the
400    same results would hold for two other species and perceptual paradigms. We first analyzed
401    the behavioral responses from 9 human subjects performing a variable-duration orientation
402    discrimination task (Cheadle et al. 2014). In each trial, a sequence of 5 to 10 gratings with a
403    certain orientation were shown to the subject, and the subject had to report whether they
404    thought the gratings were overall mostly aligned to the left or to the right diagonal. In this task,
405    the experimenter can control the evidence provided by each sample by adjusting the
406    orientation of the grating. We performed the same analyses on the participant responses than
407    on monkey data. As for monkeys, we found that the integration model nicely captured
408    psychometric curves, participant accuracy and psychophysical kernels (Figure 5A-C, red
409    curves and symbols). By contrast, both non-integration models failed to capture these patterns
410    (Figure 5A-C, blue and green curves and symbols). The accuracy from both models
411    consistently underestimated participant performance: 8 and 6 out of 9 subjects outperformed
412    the maximum performance for the snapshot and extrema-detection models, respectively
413    (Supp. Figure 7). This suggests that human participants achieved such accuracy by integrating
414    sensory evidence over successive samples. Moreover, subjects overall weighted more later
415    samples (Figure 5C), which is inconsistent with the extrema-detection mechanism. A formal
416    model comparison confirmed that in each participant, the integration model provided a far
417    better account of subject responses than either of the non-integration models did (Figure 5D).
418    We then assessed how subjects combined information from weak and strong evidence
419    samples into their decisions, using the same analyses as for monkeys. As predicted by the
420    integration model, but not by the extrema-detection model, humans choices consistently
421    aligned with the total stimulus evidence and not simply with the strongest evidence sample
422    (Figure 5E). Finally, the average integration map for early and late evidence within the stimulus
423    sequence displayed nearly linear diagonal isolines, showing that both were integrated into the
424    response (Figure 5F). Integration maps from participants correlated better with maps predicted
425    by the integration model than with maps predicted by either of the alternative non-integration
426    strategies (Figure 5G; two-tailed t-test on bootstrapped *r* values: p<0.001 for 7 out 9
427    participants in the integration vs snapshot comparison; in all 9 participants for the integration

428    vs extrema-detection comparison). Overall, these analyses show converging evidence that
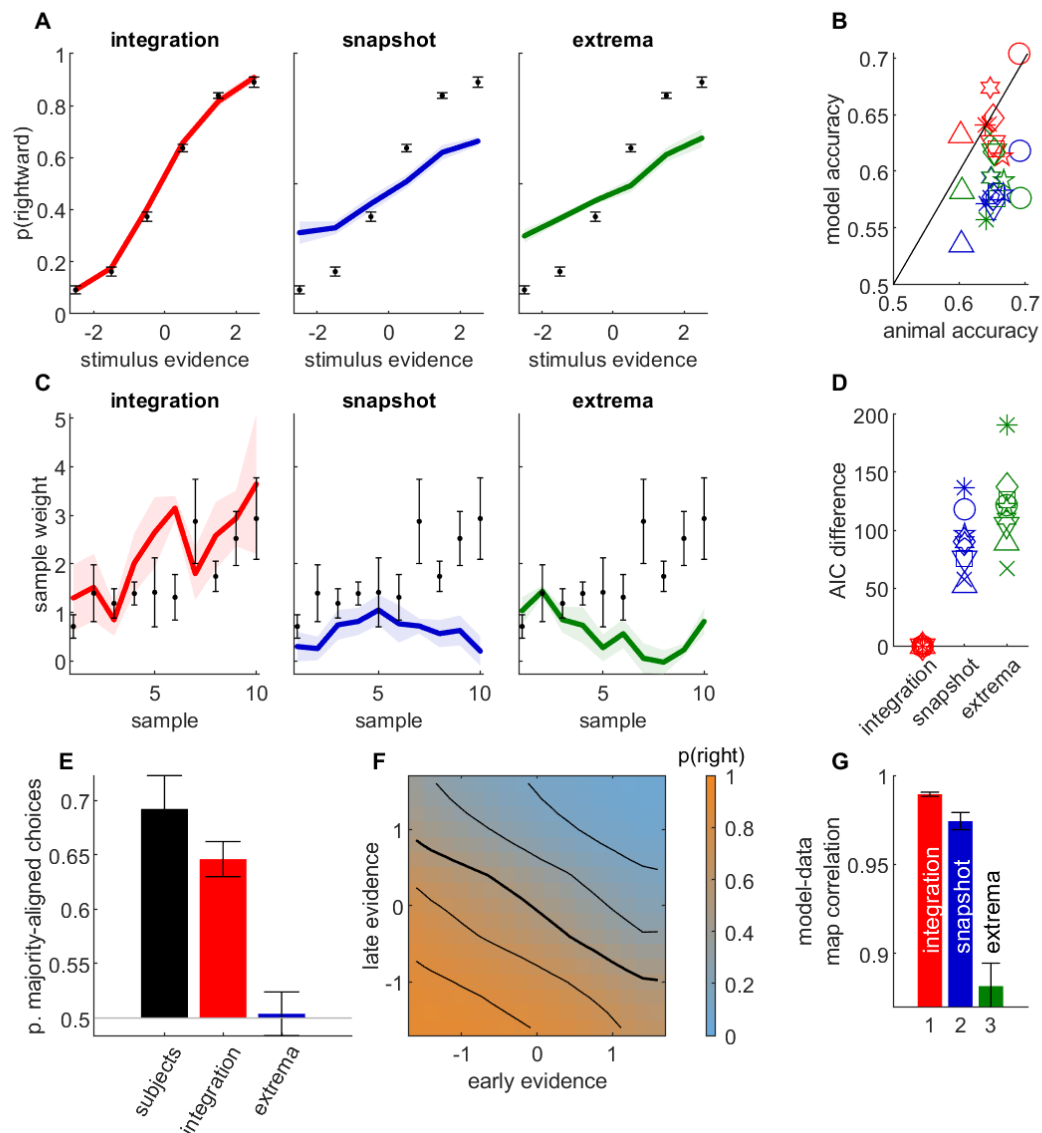429    human decisions in an orientation discrimination task rely on temporal integration.



430
431    **Figure 5. Behavioral data from orientation discrimination task in humans provides further**
432    **evidence for temporal integration.** **A.** Psychometric curves for human data and simulated data,
433    averaged across participants (*n*=9). Legend as in figure 2C. **B.** Simulated model accuracy (y-axis) vs
434    participant accuracy (x-axis) for integration model (red), snapshot model (blue) and extrema-detection
435    model (green). Each symbol corresponds to a participant. **C.** Psychophysical kernel for human data and
436    simulated data, averaged across participants. Legend as in A. **D.** Difference in AIC between each model
437    and the integration model. Legend as in B. **E.** Proportion of choices aligned with total stimulus evidence
438    in disagree trials, for participant data (black bars) and simulated models, averaged over participants. **F.**
439    Integration map for early and late stimulus evidence, computed as in Figure 4A, averaged across
440    participants. **G.** Correlation between integration map of participants and simulated data for integration,
441    snapshot and extrema-detection models, averaged across participants. Colour code as in B. Error bars
442    represent the standard error of the mean across participants in all panels.

443

**Temporal integration in rat acoustic intensity judgments**

Finally, we analyzed data from 5 rats performing a fixed-duration auditory task where the animals had to discriminate the side with larger acoustic intensity (Pardo-Vazquez et al. 2019). The relative intensity of the left and right acoustic signals was modulated in sensory samples of 50 ms, so that the stimulus sequence provided time-varying evidence for the rewarded port. The stimulus sequence was composed of either 10 or 20 acoustic samples of 50 ms each, for a total duration of 500 or 1000 ms. We applied the same analysis pipeline as for monkey and human data. The integration model provided a much better account of rat choices than non-integration strategies, based on psychometric curves (Fig. 6A), predicted accuracy (Fig. 6B), psychophysical kernel (Fig. 6C) and model comparison using AIC (Fig. 6D). Similar to humans and monkeys, rats tended to select the side corresponding to the total stimulus evidence and not the largest sample evidence in "disagree" trials, as predicted by the integration model (Fig. 6E). Finally, the integration map was largely consistent with an integration strategy (Fig. 6F), and correlated more strongly with simulated maps from the integration model (unpaired t-test on bootstrapped $r$ values: $p<0.001$ for each animal and comparison against extrema-detection and against snapshot model).
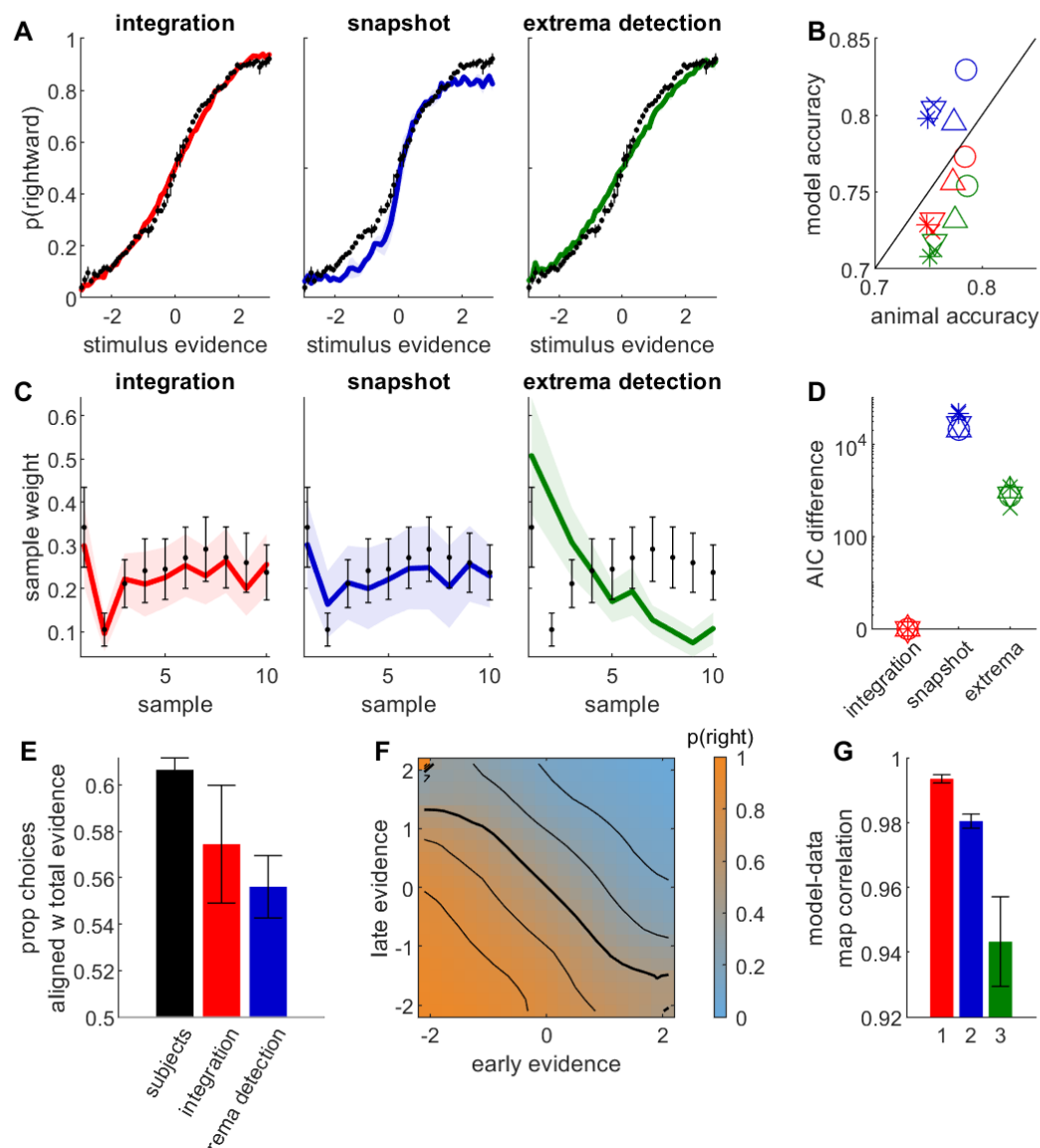
**Figure 6. Behavioral data from auditory discrimination task in 5 rats provides further evidence for temporal integration**. Rats were rewarded for correctly identifying the auditory sequence of larger intensity (number of samples: 10 or 20; stimulus duration: 500 or 1000 ms). Legend as in Figure 5. Psychophysical kernels are computed only for 10-sample stimuli (in 4 animals). See Supp Figure 8 for psychophysical kernels with 20-sample stimuli.

## DISCUSSION

We investigated the presence of temporal integration in perceptual decisions in monkeys, humans and rats through a series of standard and innovative analyses of response patterns. In all analyses we contrasted predictions from one integration and two non-integration computational models of behavioral responses (Figure 1). For each non-integration model, we considered multiple variants to explore the maximal flexibility offered by each framework to capture animal behavior. For our datasets, evidence in favor of integration was easy to achieve

15

476 using standard model comparison technique as well as comparing simulated psychometric
477 curves and psychophysical kernels to their experimental counterparts (Figure 2). Our results
478 are in line with previous evidence for temporal integration in perceptual decisions of humans
479 and mice (Pinto et al. 2018; Stine et al. 2020; Waskom and Kiani 2018). Importantly, we also
480 suggest new analyses targeted at revealing specific signatures of temporal integration.

481

482 In some cases, we could link the failure of the non-integration model to a fundamental limitation
483 of the model. For example, the extrema-detection model cannot explain the non-monotonic
484 psychophysical kernels of monkeys or the increasing psychophysical kernels in humans. This
485 is because the decision in that mode is based on the first sample to reach a certain fixed
486 criterion, so it will always produce a primacy effect, i.e., a decreasing psychophysical kernel.
487 Although this effect can be small, and in practice yields approximately flat kernels (Stine et al.
488 2020), it cannot produce increasing or non-monotonic kernels.

489

490 Another strong limitation of non-integration models (both the extrema detection and the
491 snapshot model) is that accuracy is limited by the fact that decisions depend on a single
492 sample. We found that that boundary performance (i.e. the maximum performance that a
493 model can reach) was actually lower than subject accuracy for most human participants, *de*
494 *facto* ruling out these non-integration strategies for these participants. This is consistent to
495 what was observed in a constrast discrimination DSS task where human subjects had to make
496 judgments about image sequences spanning up to tens of seconds each (Waskom and Kiani
497 2018). It clearly contrasts however with results from (Stine et al. 2020) where the non-
498 integration strategies matched the accuracy of human subjects performing the classical
499 random-dot-motion task. This discrepancy may be related to the different sources of noise in
500 the two paradigms. In DSS tasks, because the sensory evidence provided by the stimulus at
501 each moment is controlled by the experimenter, the unpredictability of human responses
502 essentially stems from internal noise at the level of sensory processing and temporal
503 integration (Waskom and Kiani 2018; Drugowitsch et al. 2016). By contrast, in the random dot
504 motion task (Kiani, Hanks, and Shadlen 2008), which is a non-DSS task because the
505 experimenter does not typically specify differing amounts of motion in each time epoch within
506 a single trial, typically elicits more variable responses due to the presence of stimulus noise.
507 This overall increased noise level leads to a looser relationship between the stimulus condition
508 and the behavioral responses, which can thus be accounted for by a larger spectrum of
509 computational mechanisms. These issues have been addressed by forcing "pulses" of a
510 certain stimulus strength and/or by performing post hoc analyses to estimate signal and noise
511 (Kiani, Hanks, and Shadlen 2008) but these are partial solutions that DSS paradigms solve by
512 design. This illustrates the benefits of using experimental designs where variability in stimulus
513 information can be fully controlled and parametrized by the experimenter, as these paradigms
514 discriminate more precisely between different models of perceptual decisions.

515 In at least one monkey, although quantitative metrics such as penalized log-likelihood and fits
516 to psychometric curves clearly pointed to the integration model as the best account to
517 behavior, the qualitative failure modes of the non-integration strategies (especially the
518 snapshot model) was not immediately clear. Although we tried variants for each non-
519 integration model, there remained a possibility that our precise implementation failed to

520  account for monkey behavior but that other possible implementations would. Note that the
521  extrema-detection and snapshot are two of the many possible non-integration strategies. A
522  generic form for non-integration strategies corresponds to a policy that implements position-
523  dependent thresholds on the instantaneous sensory evidence. In this framework, the extrema-
524  dependent model corresponds to the case with a position-independent threshold, while the
525  snapshot model corresponds to a null bound for one sample and infinite bounds for all other
526  samples. To rule out these more complex strategies, we conducted additional analyses that
527  specifically targeted core assumptions of the integration and non-integration strategies.

528  First, the extrema-detection model fails to account for the data because it predicts that largest-
529  evidence samples should have a disproportionate impact on choices. However, this does not
530  occur, as monkeys and humans tend to respond according to the total evidence and not the
531  single large-evidence sample (Figure 3C and 5E) - see (Levi et al. 2018) for a similar analysis.
532  All non-integration strategies share the property that on each trial the decision should only rely
533  either on the early or the late part of the trial. We thus directly examined the assumptions of
534  integration and non-integration models by assessing how the evidence from the early and late
535  parts of each stimulus sequence is combined to produce a decision. We introduced *integration
536  maps* (Figure 4) to inspect such integration: isolines of the integration maps will be rectilinear
537  if and only if early and late evidence are summed, in other words if and only if temporal
538  integration takes place. Unequal weighting of evidence would still produce rectilinear isolines,
539  albeit with a different angle. By contrast, a non-integration scenario when on each trial only a
540  single piece of evidence contributes to the decision predicts isolines that bend towards the
541  axes. Integration maps from monkey, human and rat subjects nicely matched the predictions
542  of the integration models, proving that their decisions do rely on temporal integration. Note
543  that this innovative analysis technique could be used to probe integration of evidence not only
544  at temporal level but also between different sources of evidence. Indeed, there has been an
545  intense debate about whether sensory information from different spatial locations or different
546  modalities are integrated prior to reaching a decision, or whether decisions are taken
547  separately for each source before being merged, which can be viewed as extensions to the
548  snapshot model (Pannunzi et al. 2015; Otto and Mamassian 2012; Lorteije et al. 2015; Hyafil
549  and Moreno-Bote 2017). Our integration analysis could provide new answers to this old
550  debate.

551

552  Integration maps can be computed not only for choice patterns but for any type of behavioral
553  or neural marker of cognition. We computed a neural integration map (Supp. Figure 9) by
554  looking at the average spike activity of Lateral Intra Parietal (LIP) neurons as a function of
555  early and late evidence, for neurons recorded while the monkeys performed the motion
556  discrimination experiment (Yates et al. 2017). The neural integration map clearly showed
557  rectilinear isolines, as predicted by an integration model of neural spiking. By contrast, neural
558  implementations of the snapshot and extrema-detection predicted strongly curved isolines.
559  The activity of LIP neurons correlates with the evidence accumulated over the presentation of
560  the stimulus in favor of either possible choices (Gold and Shadlen 2007). This result shows
561  that the activity of individual LIP neurons indeed reflects the temporal integration of sensory
562  information that drives animal behavior.

563  We have focused in this study on paradigms where the stimulus duration is fixed by the
564  experimenter, and subjects could only respond after stimulus extinction. Stine et al proposed

565 a method for distinguishing integration from non-integration strategies mixing experiments
566 where stimulus duration is controlled by the experimenter and experiments where the stimulus
567 plays until the subject responds ("reaction time paradigms"). Our study shows an alternate way
568 to differentiate integration and non-integration strategies that does not require these conditions,
569 and may therefore be applied to existing datasets.

570 Other studies have shown how integration and non-integration strategies can be disentangled
571 in free reaction-time task paradigms. Specifically, different models make different predictions
572 regarding how the total sample evidence presented before response time should vary with
573 response time (Glickman and Usher 2019; Zuo and Diamond 2019). Glickman and Usher used
574 these predictions to rule out non-integration strategies in a counting task in humans, and Zuo
575 and Diamond found evidence for evidence integration to bound when rats discriminate
576 textures using whisker touches (Zuo and Diamond 2019). Furthermore, decisions in self-paced
577 paradigms are influenced by the sensory evidence from the early part of the stimulus (Winkel
578 et al. 2014), ruling out the proposal that they would only depend on the sensory evidence at
579 the time of decisions (Thura et al. 2012). Of note, the absence of integration seems a more
580 viable strategy when the duration of the stimulus is controlled externally and the benefits of
581 integrating in terms of accuracy might not compensate for its cognitive cost. In free-reaction
582 time paradigms, waiting for a long sequence of samples and selecting its response based on
583 a single sample does not seem a particularly efficient strategy. If the cognitive cost of
584 integration is high, it is more beneficial to interrupt the stimulus sequence early with a rapid
585 response. Such rapid responses are commonly seen and can be attributed either to urgency
586 signals modulating the integration of stimulus evidence (Drugowitsch et al. 2012) or to action
587 initiation mechanisms that time the response after a specific time (e.g. one or two samples)
588 following stimulus onset (Hernández-Navarro et al. 2021). Here, we have shown that even in
589 paradigms where the stimulus duration is controlled by the experimenter, mammals often
590 integrate sensory evidence over the entire stimulus.

591 In conclusion, we have found strong evidence for temporal integration in perceptual tasks
592 across species (monkeys, humans and rats) and perceptual domain (visual motion, visual
593 orientation and auditory discrimination). Thus, although the time scale of integration can be
594 adapted to the statistics of the environment (Ossmy et al. 2013; Glaze, Kable, and Gold 2015;
595 Kilpatrick et al. 2019), the principle that stimulus evidence is integrated over time appears as
596 a hallmark of perception. This evidence was gathered by leveraging experimentally-controlled
597 sensory evidence at each sensory sample composing a stimulus, and novel model-based
598 statistical analysis. We speculate that temporal integration is a ubiquitous feature of perceptual
599 decisions due to hard-wired neural integrating circuits, such as recurrent stabilizing
600 connectivity in sensory and perceptual areas (Wang 2008; Wimmer et al. 2015).

601

602

603 **METHODS**

604

605 **Monkey experiment.**

606 We present here the most relevant features of the behavioral protocol - see (Yates et al. 2017)
607 for further experimental details. Two adult rhesus macaques (subject N, a 10-year old female;

and subject P, a 14-year old male) performed a motion discrimination task. On each trial, a stimulus consisting of a hexagonal grid (5-7 degrees, scaled by eccentricity) of Gabor patches (0.9 cycle per degree; temporal frequency 5 Hz for Monkey P; 7 Hz for Monkey N) was presented. Monkeys were trained to report the net direction of motion in a field of drifting and flickering Gabor elements with an eye movement to one of two targets. Each trial motion stimulus consisted of seven consecutive motion pulses, each lasting 9 or 10 video samples (150 ms or 166 ms; pulse duration did not vary within a session), with no interruptions or gaps between the pulses. The strength and direction of each pulse $S_{ti}$ for trial $t$ and sample $i$ was set by a draw from a Gaussian rounded to the nearest integer value. The difficulty of each trial was modulated by manipulating the mean and variance of the Gaussian distribution. Monkeys were rewarded based on the empirical stimulus and not on the stimulus distribution. We analyzed a total of 112 sessions for monkey N and 60 sessions for monkey P, with a total of 72137 and 33416 valid trials, respectively. These sessions correspond to sessions with electrophysiological recordings reported in (Yates et al. 2017) and purely behavioral sessions. All experimental protocols were approved by The University of Texas Institutional Animal Care and Use Committee (AUP-2012-00085, AUP-2015-00068) and in accordance with National Institute of Health standards for care and use of laboratory animals

**Human experiment.**

9 adult subjects (5 males, 4 females; aged 19-30) performed an orientation discrimination task whereby on each trial they reported in each trial whether a series of gratings were perceived to be mostly tilted clockwise or counterclockwise (Drugowitsch et al. 2016). Each discrete-sample stimulus consisted of five to ten gratings. Each grating was a high-contrast Gabor patch (colour: blue or purple; spatial frequency = 2 cycles per degree; SD of Gaussian envelope = 1 degree) presented within a circular aperture (4 degrees) against a uniform gray background. Each grating was presented during 100 ms, and the interval between gratings was fixed to 300 ms. The angles of the gratings were sampled from a von Mises distribution centered on the reference angle ($\alpha_0 = 45$ degrees for clockwise sequences, 135 degrees for anticlockwise sequences) and with a concentration coefficient $\kappa = 0.3$. The normative evidence provided by sample $i$ in trial $t$ in favor of the clockwise category corresponds to how well the grating orientation $\alpha_{ti}$ aligns with the reference orientation, i.e. $S_{ti} = 2\kappa \, cos(2(\alpha_{ti} - \alpha_0))$ .

Each sequence was preceded by a rectangle flashed twice during 100 ms (the interval between the flashes and between the second flash and the first grating varied between 300 and 400 ms). The participants indicated their choice with a button press after the onset of a centrally occurring dot that succeeded the rectangle mask and were made with a button press with the right hand. Failure to provide a response within 1000 ms after central dot onset was classified as invalid trial. Auditory feedback was provided 250 ms after participant response (at latest 1100 ms after end of stimulus sequence). It consisted of an ascending tone (400 Hz/800 Hz; 83 ms/167 ms) for correct responses; descending tone (400 Hz/ 400 Hz; 83 ms/167 ms) for incorrect responses; a low tone (400 Hz; 250 ms) for invalid trials.

Trials were separated by a blank interstimulus interval of 1,200-1,600 ms (truncated exponential distribution of mean 1,333 ms). Experiments consisted of 480 trials in 10 blocks of 48. It was preceded with two blocks of initiation with 36 trials each. In the first initiation block, there was only one grating in the sequence, and it was perfectly aligned with one of the

653 reference angles. In the second initiation block, sequences of gratings were introduced, and
654 the difficulty was gradually increased (the distribution concentration linearly decreased from
655 $\kappa = 1.2$ to $\kappa = 0.3$). Invalid trials (mean 6.9 per participant, std 9.4) were excluded from all
656 regression analyses. The study was approved by the local ethics committee (approval
657 2013/5435/I from CEIm- Parc de Salut MAR).

658

659 **Rat experiment.**

660 Rat experiments were approved by the local ethics committee of the University of Barcelona
661 (Comité d'Experimentació Animal, Barcelona, Spain, protocol number Ref 390/14). 5 male
662 Long-Evans rats (no genetic modifications; 350-650g; 8-10 weeks-old at the beginning of the
663 experiment), pair-housed and kept on stable conditions of temperature (23ºC) and humidity
664 (60%) with a constant light-dark cycle (12h:12h, experiments were conducted during the light
665 phase). Rats had free access to food, but water was restricted to behavioral sessions. Free
666 water during a limited period was provided on days with no experimental sessions.

667 Rats performed a fixed-duration auditory discrimination task where they had to classify noisy
668 stimuli based on the intensity difference between the two lateral speakers (Pardo-Vazquez et
669 al. 2019; Hermoso-Mendizabal et al. 2020). A LED on the center port indicated that the rat
670 could start the trial by poking in that center port. After this poke, rats had to hold their snouts
671 in the central port during 300 ms (i.e. fixation). Following this period, an acoustic DSS was
672 played. Rats had to remain in the central port during the entire presentation of the stimulus.
673 At stimulus offset, the center LED went off and rats could then come out of the center port and
674 head towards one of the two lateral ports. Entering the lateral port associated with the speaker
675 that generated the larger sound intensity led to a reward of 24 µl of water (correct responses),
676 while entering the opposite port lead to a 5 s timeout accompanied with a bright light during
677 the entire period (incorrect responses). If rats broke fixation during the pre-stimulus fixation
678 period or during the stimulus presentation, the sound was interrupted, the center LED
679 remained on, and the rat had to initiate a new trial starting by center fixation followed by a new
680 stimulus. Fixation breaks were not included in any of the analyses. Stimulus duration was 0.5
681 s (10 samples) or 1 s (20 samples). Two rats performed 0.5-second stimuli only (77810 and
682 54803 valid trials, respectively); one rat performed 1 s stimuli only (42474 valid trials); the
683 remaining two rats performed a mixture of 0.5 and 1 s stimuli trials randomly interleaved (5016
684 trials and 65212 valid trials, respectively for one animal; 7374 and 38829 trials for the other
685 animal). In each trial $k$ one stimulus $S^X_k(t)$ was played in each speaker ($X$=R for the Right
686 speaker and $X$=L for the Left speaker). Each stimulus was an amplitude modulated (AM)
687 broadband noise defined by $S^X_k(t) = [1 + sin(f_{AM}t + \varphi)]a^X_k(t)\,\xi_X(t)$ where $f_{AM}$=20 Hz
688 (sensory samples lasted 50 ms), the phase delay $\varphi = 3\pi/2$ and $\xi_X(t)$ were broadband noise
689 bursts. The amplitudes of each sound in each frame were $a^L_k(t) = (1 + S_{k,f})/2$ and $a^R_k(t) =$
690 $(1 - S_{k,f})/2$ with $S_{k,f}(t)$ being the instantaneous evidence that was drawn independently in
691 each frame $f$ from a transformed Beta distribution with support [-1,1]. With this parametrization
692 of the two sounds the sum of the two envelopes was constant in all frames $a^L_k(t) + a^R_k(t)$=1.
693 There were 7 x 5 stimulus conditions, each defined by a Beta distribution, spanning 7 mean
694 values (-1, -0.5, -0.15, 0, 0.15, 0.5 and 1) and 5 different standard deviations (0, 0.11, 0.25,
695 0.57 and 0.8). In around the first half of the sessions, only sample sequences in which the
696 total stimulus evidence matched the targeted nominal evidence were used. This effectively

697 introduced weak correlations between samples. In the second half of the sessions, this
698 condition was removed and samples in each stimulus were drawn independently from the
699 corresponding Beta distribution.

700
701 **Integration model**

702 The integration model for human participants corresponds to a logistic regression model,
703 where the probability of selecting the right choice $p(r_t)$ at trial $t$ depends on the weighted sum
704 of the sample evidence: $p(r_t) = \sigma(\beta_0 + \Sigma_{i \in [1..n]} \beta_i S_{ti})$, where $\beta_0$ is a lateral bias, $S_{ti}$ is the
705 signed sample evidence at sample $i$ ; $\beta_i$ is the sensory weight associated with the $i$th sample in
706 the stimulus sequence; and $\sigma(x) = (1 + e^{-x})^{-1}$ is the logistic function. The vector $\beta_i$'s allowed
707 to capture different shapes of psychophysical kernels (e.g. primacy effects, recency effects)
708 which can emerge due to a variety of suboptimalities in the integration process (leak, attractor
709 dynamics, sticky bounds, sensory after-effects, etc.) (Brunton, Botvinick, and Brody 2013;
710 Yates et al. 2017; Prat-Ortega et al. 2021; Bronfman, Brezis, and Usher 2016).

711 For the monkey and rat data, we included a session-dependent modulation gain $\gamma_t$ to capture
712 the large variations in performance in monkeys across the course of sessions (see Supp
713 Figure 1A):

714
$$p(r_t) = \sigma(\beta_0 + \gamma_t \Sigma_{i \in [1..n]} \beta_i S_{ti})$$

715 This model corresponds to a bilinear logistic regression model which pertains to the larger
716 family of Generalized Unrestricted Models (GUMs) (Adam and Hyafil 2020). Parameters $(\beta, \gamma)$
717 were fitted using the Laplace approximation as described in (Adam and Hyafil 2020). The
718 modulation gain was omitted when applied to human data, yielding a classical logistic
719 regression model.

720
721 **Snapshot model**

722 In the snapshot model, decisions are based on each trial based upon a single sample. The
723 model also includes the possibility for left and right lapses. In each trial, the attended sample
724 is drawn from a multinomial distribution of parameters $(\pi_1, .. \pi_n, \pi_L, \pi_R)$, where the first terms
725 $\pi_i$ $(1 \leq i \leq n)$ correspond to the probability of attending sample $i,$ and $\pi_L$ and $\pi_R$ correspond to
726 the probability of left and right lapses, respectively. Upon selecting sample $i,$ the probability for
727 selecting the right choice is given by the function $H_i(S_t)$. In the deterministic version of the
728 model, $H_i$ is simply determined by the sign of the $i$-th sample evidence: $H_i(S_t) = 1$ if $S_{ti} > 0$,
729 $H_i(S_t) = 0$ if $S_{ti} < 0$, and $H_i(S_t) = 0.5$ if $S_{ti} = 0$ (i.e. random guess if the sample has null
730 evidence). We also define similar functions for lapse responses: $H_R(S) = 1$ and $H_L(S) = 0$,
731 irrespective of the stimulus. In the non-deterministic version of the model, the probability
732 $H_i(S_{ti})$ is determined by a logistic function of the attended sample evidence $H_i(S_t) = \sigma(\beta_i S_{ti})$
733 where $\beta_i$ describes a sensitivity parameter. The deterministic case can be viewed as the limit
734 of the non-deterministic case when all sensitivity parameters $\beta_i$ diverge to $+\infty$, i.e. when
735 sensory and decision noise are negligible.

736 The overall probability for selecting right choice (marginalizing over the attended sample,
737 which is a hidden variable) can be captured by a mixture model :

738 $$p(r_t) = \pi_R + \Sigma_{i \in [1..n]} \pi_i H_i(S_t) = \Sigma_{i \in [1..n,L,R]} \pi_i H_i(S_t)$$

739 The mixture coefficients $\pi_i$ $(i = 1,..n,L,R)$ are constrained to be non-negative and sum up to
740 1. In the non-deterministic model, the parameters also include sensitivity parameters $\beta_i$. The
741 model is fitted using Expectation-Maximization (Bishop 2006). In the Expectation step, we
742 compute the responsibility $z_{ti}$, i.e. the posterior probability that the sample *i* was attended at
743 trial *t* (for *i=L, R,* the probability that the trial corresponded to a lapse trial):

744

745 $z_{ti} = \pi_i \theta(S_{ti})/\Sigma_j \pi_j H(S_{tj})$    for rightward responses ($R_t = 1$)
746 $z_{ti} = \pi_i(1 - H(S_{ti}))/\Sigma_j \pi_j(1 - H(S_{tj}))$    for leftward responses ($R_t = 0$)

747

748 In the Maximization step, we update the value of the parameters by maximizing the Expected
749 Complete Log-Likelihood (ECLL): $Q(\pi, \beta) = \Sigma_{ti} z_{ti} \log p(r_t; \pi, \beta)$. Maximizing over the mixture
750 coefficients with the unity-sum constraint provides the classical update: $\pi_i = \Sigma_{ti} z_{ti}/N$, where
751 *N* is the total number of trials. In the non-deterministic model, maximizing the ECLL over
752 sensitivity parameters is equivalent to fitting a logistic regression model with weighted
753 coefficients $z_{ti}$, which is a convex problem. Best fitting parameters can be found using Newton-
754 Raphson updates on the parameters:

755 $\beta_i^{(new)} = \beta_i - \frac{\partial Q/\partial \beta_i}{\partial^2 Q/\partial \beta_i^2}$                                     with
756 $\partial Q/\partial \beta_i = \Sigma_t z_{ti}(p(r_t) - R_t)$ and $\partial^2 Q/\partial \beta_i^2 = \Sigma_t z_{ti} S_{ti}^2 p(r_t)(1 - p(r_t))$

757 To speed up the computations, in each M step, we only performed one Newton-Raphson
758 update for each sensitivity parameter, rather than iterating the updates fully until convergence.
759 The EM procedure was run until convergence, assessed by an increment in the log-likelihood
760 $L(\pi, \beta)$ of less than $10^{-9}$ after one EM iteration. The log-likelihood for a given set of parameters
761 is given by $L(\pi, \beta) = \Sigma_t \log p(r_t)$. The EM iterative procedure was repeated with 10 different
762 initializations of the parameters to avoid local minima.

763

764 Note that for monkey and rat data, since we observed large variations in performance across
765 sessions, the model based its choices on session-gain modulated evidence $\underline{S_{ti}} = \gamma_t S_{ti}$ instead
766 raw evidence $S_{ti}$ (this had no impact for the deterministic variant since $\underline{S_{ti}}$ and $S_{ti}$ always have
767 the same sign). We fitted the model from individual subject responses either with lapses $\pi_L$ and
768 $\pi_R$ as free parameters, or fixed to $\pi_L = \pi_R = 0.01$. Figures in the main manuscript correspond
769 to the deterministic snapshot model with fixed lapses. We also studied variants of the snapshot
770 model where decisions in each trial are based on *K* attended samples, i.e depends on
771 $(S_{ti},..S_{t,i+K-1})$ with $1 \leq K \leq n - 1$ and $1 \leq i \leq n - K + 1$ is the first attended sample. In the
772 deterministic case, the choice is directly determined by the sign of the sum of the signed
773 evidence for the attended samples. In the non-deterministic case, the evidence for the
774 attended samples are weighted and passed through a sigmoid: $H_i(S_t) =$
775 $\sigma(\Sigma_{k \in [1..K]} \beta_{ki} S_{t,i+k-1})$. The model with a single attended sample presented above is equivalent
776 to this extended model when using $K = 1$. At the other end, using $K = n$ corresponds to the
777 temporal integration model (without the lateral bias).

778

779 **Extrema-detection model**

780 In the extrema-detection model, a choice is selected according to the first sample in the
781 sequence whose absolute evidence value reaches a certain threshold $\theta$, i.e $p(r_t|\theta) =$
782 $H(m_{ti}), |m_{ti}| \geq \theta, |m_{tj}| < \theta$ for all $j < i$. Here $m_{ti}$ is the sample evidence corrupted by
783 sensory noise $\varepsilon_{ti}$ which is distributed normally with variance $\sigma^2$: $m_{ti} = S_{ti} + \varepsilon_{ti}$ with $\varepsilon_{ti} \sim$
784 $N(0, \sigma^2)$. $H$ is the step function. If the stimulus sequence ends and no sample has reached the
785 threshold, then the decision is taken at chance. As described in (Waskom and Kiani 2018),
786 the probability for a rightward choice at trial $t$ can be expressed as:

787
$$p(r_t) = \Sigma_i \Phi\left(\frac{S_{ti} - \theta}{\sigma}\right)\Pi_{j<i}\left(1 - \Phi\left(\frac{S_{tj} - \theta}{\sigma}\right) - \Phi\left(\frac{-S_{tj} - \theta}{\sigma}\right)\right) + \frac{1}{2}\Pi_{j\leq n}\left(1 - \Phi\left(\frac{S_{tj} - \theta}{\sigma}\right)\right.$$
788
$$\left. - \Phi\left(\frac{-S_{tj} - \theta}{\sigma}\right)\right)$$

789

790 where $\Phi$ is the cumulative normal distribution. We also included the possibility for left and right
791 lapses with probability $\pi_L$ and $\pi_R$. Following Stine and colleagues (Stine et al. 2020), we
792 explored an alternative default rule called 'last sample' rule: if the stimulus extinguishes and
793 the threshold has not been reached, then the decision is based on the (noisy) last sample
794 rather than simply by chance. This changes the equation describing the probability for
795 rightward choices to:

796
$$p(r_t) = \Sigma_{i<n} \Phi\left(\frac{S_{ti} - \theta}{\sigma}\right)\Pi_{j<i}\left(1 - \Phi\left(\frac{S_{tj} - \theta}{\sigma}\right) - \Phi\left(\frac{-S_{tj} - \theta}{\sigma}\right)\right) + \Phi\left(\frac{S_{tn}}{\sigma}\right)\Pi_{j<n}\left(1 - \Phi\left(\frac{S_{tj} - \theta}{\sigma}\right)\right.$$
797
$$\left. - \Phi\left(\frac{-S_{tj} - \theta}{\sigma}\right)\right)$$

798 As for the snapshot model, we used the session-gain modulated evidence $\underline{S_{ti}}$ instead of raw
799 evidence $S_{ti}$ for fitting the model to monkey and rat data. The four parameters of the model
800 $(\theta, \sigma, \pi_L, \pi_R)$ were estimated from each subject data by maximizing the log-likelihood with
801 interior-point algorithm (function *fmincon* in Matlab) and 10 different initializations of the
802 parameters.

803

804

805 **Model validation and model comparison.**

806 Psychophysical kernels were obtained from subject data and simulated data by running a
807 logistic regression model: $p(r_t) = \sigma(\beta_0 + \Sigma_i \beta_i S_{ti})$. Standard errors of the weights $\beta_i$ were
808 obtained from the Laplace approximation. For psychometric curves, we first defined the
809 weighted stimulus evidence $T_t$ at trial $t$ as the session-modulated weighted sum of signed
810 sample evidence; with the weights obtained from the logistic regression model above
811 $T_t = \gamma_t \Sigma_i \beta_i S_{ti}$. We then divided the total stimulus evidence into 50 quantiles (10 for human
812 subjects) and computed the psychometric curve as the proportion of rightward choices for
813 each quantile.

23

814  The boundary performance for the snapshot and extrema-detection models corresponds to
815  the best choice accuracy out of all the parameterizations for each model. In the snapshot
816  model, the boundary performance corresponds to the deterministic version with no-lapse,
817  where the attended sample is always the sample $i^*$ whose sign better predicts the stimulus
818  category over all animal trials, i.e. $\pi_{i^*} = 1$ and $\pi_i = 0$ if $i \neq i^*$. For the extrema-detection model,
819  the boundary performance corresponds to the lapse-free model with no sensory noise ($\sigma = 0$)
820  and a certain value for threshold $\theta$ that is identified for each subject by simple parameter
821  search.

822  Finally, model selection was performed using the Akaike Information Criterion $AIC = 2p -$
823  $2L_{ML}$, where $p$ is the number of model parameters and $L_{ML}$ is the likelihood evaluated at
824  maximum                              likelihood                              parameters.
825

826  **Analysis of majority-driven choices**

827  We selected for each animal the subset of trials corresponding to when the largest evidence
828  sample was at odds with the total stimulus evidence, i.e. where $sign(S_{tj}, |S_{tj}| \geq |S_{ti}| \; \vee \; i) \neq$
829  $sign(\Sigma_i S_{ti})$. For this subset of trials, we computed the proportion of animal choices that were
830  aligned with the overall stimulus evidence. We repeated the analysis for simulated data from
831  the integration and extrema-detection models.

832

833  **Subjective weighting analysis**
834  In order to estimate the impact of each sample on the animal choice as a function of sample
835  evidence, we built and estimated the following statistical model

836  $$p(r_t = A) \; = \; \sigma(\beta_0 + \gamma_t \Sigma_{i \in [1..n]} \beta_i \, f(S_{ti}))$$

837  As can be seen, this model is equivalent to the temporal integration model under the
838  assumption that $f$ is a linear function. Rather, here we wanted to estimate the function $f$ (as
839  well as the session gain $\gamma_t$, lateral bias $\beta_0$ and sensory weight $\beta_i$). Including the session gain
840  was necessary for estimating $f$ accurately from the monkey and rat behavioral data, since the
841  distribution of pulse strength $S_{ti}$ was varied across sessions and could otherwise induce a
842  confound. We assumed that $f$ is an odd function, i.e. $f(-S_{ti}) = -f(S_{ti})$. This equation takes
843  the form of a Generalized Unrestricted Model and was fitted using the Laplace approximation
844  method as described in (Adam and Hyafil 2020). In the monkey experiment, sample evidence
845  could take only a finite number of values, so $f$ was simply estimated over these values. In the
846  human experiment, sample evidence could take continuous values. In this case, we defined a
847  Gaussian Process prior over $f$ with squared exponential kernel with length scale 0.1 and
848  variance 1.

849

850  **Integration of early and late evidence**

851  We designed a new analysis tool to characterize the statistical mapping from the
852  multidimensional stimulus space $\boldsymbol{S}_t = (S_{t1}, \ldots S_{tn}) \in \mathfrak{R}^n$ onto binary choices $r_t \in [0,1]$. We first
853  collapsed the stimulus sequence $\boldsymbol{S}_t$ onto the two-dimensional space defined by early evidence
854  $\boldsymbol{E}_t$ and late evidence $\boldsymbol{L}_t$ defined by $E_t = \gamma_t \Sigma_{1 \leq i \leq [n/2]} \beta_i S_{ti}$ and $L_t = \gamma_t \Sigma_{[n/2]+1 \leq i \leq n} \beta_i S_{ti}$, where the
855  weights $\beta_i$ and session gains $\gamma_t$ correspond to parameters estimated from the temporal

856   integration model (session gains were omitted for human participants). Next we plotted the
857   integration map which represents the probability for rightward choices as a function of $(E_t, L_t)$.
858   The map was obtained by smoothing data points with a two-dimensional gaussian kernel.
859   More specifically, for each pair value *(E,L)*, we selected the trials whose early and late
860   evidence values $E_t$ and $L_t$ fell within a certain distance to *(E,L)*, i.e. $d_t = dist((E,L)(E_t, L_t)) <$
861   2. We then computed the proportion of rightward choices for the selected trials, with a weight
862   for each trial depending on the distance to the pair value $w_t = N((E_t, L_t); (E,L), 0.1^2 I)$.
863   Because the space *(E,L)* was not sampled uniformly during the experiment, we represent the
864   density of trials by brightness. For each subject we obtained integration maps both from
865   subject data as well as from model simulations. For each model, we computed the Pearson
866   correlation between the maps obtained from the corresponding simulation and from the
867   subject data. We tested the significance of correlation measures between models by using a
868   bootstrapping procedure: we calculated the correlation measure *r* from 100 bootstraps for
869   each model and participant, and then performed an unpaired t-test between bootstrapped *r*.

870

871   Next, we analyzed the conditional psychometric curves, i.e. the psychometric curves for the
872   early evidence conditioned on the value of late evidence, which correspond to vertical cuts in
873   the integration map. To do so, we first binned late evidence $L_t$ by bins of width 0.5. Conditional
874   psychometric curve represent the probability of rightward choices as a function of early
875   evidence $E_t$, separately for each late evidence bin. For each late evidence bin, we also
876   estimated the corresponding bias $\beta$, left lapse $\pi_L$ and right lapse $\pi_R$ by fitting the following
877   function on the corresponding subset of trials:

878   $$p(r_t) = \pi_R + (1 - \pi_R - \pi_L)\sigma(\beta E_t)$$

879

### Analysis of LIP neuron activity

881   We analyzed the activity of 82 LIP neurons recorded over 43 sessions of the motion
882   discrimination tasks (Yates et al. 2017). We applied the following procedure to extract the
883   integration map for LIP neurons. For each neuron *n*, we computed the spike count $s_t^{(n)}$ in a
884   window of 500 ms width following each stimulus offset, which is where LIP neurons were found
885   to have maximal selectivity to motion evidence from the entire pulse sequence (Yates et al.
886   2017). We then applied a Poisson GLM $E(s_t^{(n)}) = exp(w_0^{(n)} + \Sigma_i w_i^{(n)} S_{ti})$ for each neuron *n*
887   to extract the impact of each sample *i* on the individual neural spike count $w_i^{(n)}$. For each trial
888   *t*, we used these weights to compute the neuron-weighted early and late evidence defined by
889   and $E_t^{(n)} = \Sigma_{1 \leq i \leq 3} w_i^{(n)} S_{ti}$  $L_t^{(n)} = \Sigma_{4 \leq i \leq 7} w_i^{(n)} S_{ti}$. Note that this weighting converts the
890   evidence onto the space defined by the preferred direction of the neuron, such that positive
891   evidence signals evidence towards the preferred direction and negative evidence signals
892   evidence towards the anti-preferred direction. We then merged the vectors for normalized
893   spike counts $\underline{s}_t^{(n)} = s_t^{(n)}/exp(w_0^{(n)})$, early evidence $E_t^{(n)}$ and late evidence $L_t^{(n)}$ across all
894   neurons. The normalized spike counts were binned by values of early and late evidence (bin
895   width: 0.02), and the average over each bin was computed after convolving with a two-
896   dimensional gaussian kernel of width 0.1. The neural integration map represents the average
897   normalized activity per bin.

898 Simulations of spiking data for the integration and non-integration models were proceeded as
899 follows. First, the neural integration model corresponds to linear summing with neuron-specific
900 weights which are then passed through an exponential nonlinearity; the spike counts for each
901 trial are generated using a Poisson distribution whose rate is equal to the nonlinear output
902 (Supp Figure 9a, top). This corresponds exactly to the generative process of the Poisson GLM
903 described above. For the extrema detection model (Supp Figure 9a middle), we hypothesized
904 that LIP activity would only be driven by the sample that reaches the threshold (and dictates
905 the animal response). To this end, we first simulated the behavioral extrema detection model
906 for all trials, using parameters $(\theta, \sigma, \pi_L, \pi_R)$ fitted from the corresponding animal, to identify
907 which sample $i$ reaches the subject-specific threshold. We then assumed that the spiking
908 activity of the neuron would follow the stimulus value at sample $i$ $S_{ti}$ (signed by the preferred
909 direction of the neuron $p^{(n)}$ through:

$$E_{ED}(s_t^{(n)}) = exp(w_0^{(n)} + p^{(n)} S_{ti} \Sigma_j w_j^{(n)}/2)$$

911 Again the spike count were generated from a Poisson distribution with rate $E_{ED}(s_t^{(n)})$.

912 Finally, for the snapshot model (Supp. Figure 9a bottom), we assumed that the neuron activity
913 would merely reflect the sensory value of the only sample it would attend. We assumed that
914 the probability mass function to attend each of the 7 samples would be neuron-specific, so we
915 used the normalized weights of the Poisson GLM for that specific neuron as defining such
916 probability (weights were signed by the neuron preferred direction so that the vast majority of
917 weights were positive; negative weights were ignored). For each trial, we thus randomly
918 sampled the attended sample $i$ using this probability mass function and then simulated the
919 spike count $s_t^{(n)}$ from a Poisson distribution with rate $E_{Snapshot}(s_t^{(n)}) = exp(w_0^{(n)} +$
920 $p^{(n)} S_{ti} \Sigma_j w_j^{(n)})$.

921 We simulated spiking activity for each neuron and for each integration and non-integration
922 model, and then used simulated data to compute neural integration maps exactly as described
923 above for the actual LIP neuron activity.

924

## Data and code availability

926 All experimental data (behavioral and neural data in monkeys, behavioral data in rats and
927 humans) and code to run the analysis will be made publicly available
928 at https://github.com/ahyafil prior to final publication
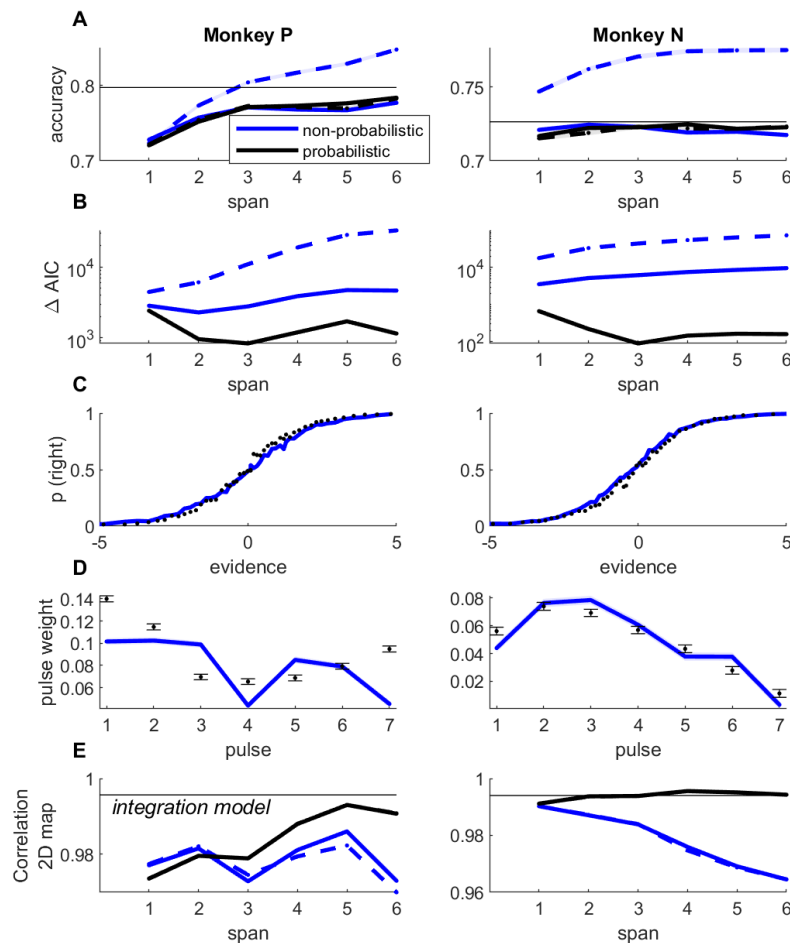
## ACKNOWLEDGMENTS

936

937

**REFERENCES**

Adam, Vincent, and Alexandre Hyafil. 2020. "Non-Linear Regression Models for Behavioral and Neural Data Analysis." *arXiv Preprint arXiv:2002.00920*. https://arxiv.org/abs/2002.00920.

Bishop, Christopher M. 2006. *Pattern Recognition and Machine Learning*.

Bogacz, Rafal, Eric Brown, Jeff Moehlis, Philip Holmes, and Jonathan D. Cohen. 2006. "The Physics of Optimal Decision Making: A Formal Analysis of Models of Performance in Two-Alternative Forced-Choice Tasks." *Psychological Review* 113 (4): 700–765.

Bronfman, Zohar Z., Noam Brezis, Rani Moran, Konstantinos Tsetsos, Tobias Donner, and Marius Usher. 2015. "Decisions Reduce Sensitivity to Subsequent Information." *Proceedings. Biological Sciences / The Royal Society* 282 (1810). https://doi.org/10.1098/rspb.2015.0228.

Bronfman, Zohar Z., Noam Brezis, and Marius Usher. 2016. "Non-Monotonic Temporal-Weighting Indicates a Dynamically Modulated Evidence-Integration Mechanism." *PLoS Computational Biology* 12 (2): e1004667.

Brunton, Bingni W., Matthew M. Botvinick, and Carlos D. Brody. 2013. "Rats and Humans Can Optimally Accumulate Evidence for Decision-Making." *Science* 340 (6128): 95–98.

Cheadle, Samuel, Valentin Wyart, Konstantinos Tsetsos, Nicholas Myers, Vincent de Gardelle, Santiago Herce Castañón, and Christopher Summerfield. 2014. "Adaptive Gain Control during Human Perceptual Choice." *Neuron* 81 (6): 1429–41.

Ditterich, Jochen. 2006. "Stochastic Models of Decisions about Motion Direction: Behavior and Physiology." *Neural Networks: The Official Journal of the International Neural Network Society* 19 (8): 981–1012.

Drugowitsch, Jan, Rubén Moreno-Bote, Anne K. Churchland, Michael N. Shadlen, and Alexandre Pouget. 2012. "The Cost of Accumulating Evidence in Perceptual Decision Making." *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience* 32 (11): 3612–28.

Drugowitsch, Jan, Valentin Wyart, Anne-Dominique Devauchelle, and Etienne Koechlin. 2016. "Computational Precision of Mental Inference as Critical Source of Human Choice Suboptimality." *Neuron* 92 (6): 1398–1411.

Glaze, Christopher M., Joseph W. Kable, and Joshua I. Gold. 2015. "Normative Evidence Accumulation in Unpredictable Environments." *eLife* 8 (1): 1.

Glickman, Moshe, and Marius Usher. 2019. "Integration to Boundary in Decisions between Numerical Sequences." *Cognition* 193 (December): 104022.

Gold, Joshua I., and Michael N. Shadlen. 2007. "The Neural Basis of Decision Making." *Annual Review of Neuroscience* 30: 535–74.

Hermoso-Mendizabal, Ainhoa, Alexandre Hyafil, Pavel E. Rueda-Orozco, Santiago Jaramillo, David Robbe, and Jaime de la Rocha. 2020. "Response Outcomes Gate the Impact of Expectations on Perceptual Decisions." *Nature Communications* 11 (1): 1057.

Hernández-Navarro, Lluís, Ainhoa Hermoso-Mendizabal, Daniel Duque, Jaime de la Rocha, and Alexandre Hyafil. 2021. "Proactive and Reactive Accumulation-to-Bound Processes Compete during Perceptual Decisions." *Nature Communications* 12 (1): 7148.

Hyafil, Alexandre, and Rubén Moreno-Bote. 2017. "Breaking down Hierarchies of Decision-Making in Primates." *eLife* 6 (June). https://doi.org/10.7554/eLife.16650.

Katz, Leor N., Jay A. Hennig, Lawrence K. Cormack, and Alexander C. Huk. 2015. "A Distinct Mechanism of Temporal Integration for Motion through Depth." *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience* 35 (28): 10212–16.

Kiani, Roozbeh, Timothy D. Hanks, and Michael N. Shadlen. 2008. "Bounded Integration in Parietal Cortex Underlies Decisions Even When Viewing Duration Is Dictated by the Environment." *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience* 28 (12): 3017–29.

Kiani, Roozbeh, and Michael N. Shadlen. 2009. "Representation of Confidence Associated with a Decision by Neurons in the Parietal Cortex." *Science* 324 (5928): 759–64.

Kilpatrick, Zachary P., William R. Holmes, Tahra L. Eissa, and Krešimir Josić. 2019. "Optimal Models of Decision-Making in Dynamic Environments." *Current Opinion in Neurobiology* 58 (October): 54–60.

Levi, Aaron J., Jacob L. Yates, Alexander C. Huk, and Leor N. Katz. 2018. "Strategic and Dynamic Temporal Weighting for Perceptual Decisions in Humans and Macaques." *eNeuro* 5 (5). https://doi.org/10.1523/ENEURO.0169-18.2018.

Lorteije, Jeannette A. M., Ariel Zylberberg, Brian G. Ouellette, Chris I. De Zeeuw, Mariano Sigman, and Pieter R. Roelfsema. 2015. "The Formation of Hierarchical Decisions in the Visual Cortex." *Neuron* 87 (6): 1344–56.

Morillon, Benjamin, Charles E. Schroeder, and Valentin Wyart. 2014. "Motor Contributions to the

Temporal Precision of Auditory Attention." *Nature Communications* 5 (October): 5255.

Nienborg, Hendrikje, and Bruce G. Cumming. 2007. "Psychophysically Measured Task Strategy for Disparity Discrimination Is Reflected in V2 Neurons." *Nature Neuroscience* 10 (12): 1608–14.

Ossmy, Ori, Rani Moran, Thomas Pfeffer, Konstantinos Tsetsos, Marius Usher, and Tobias H. Donner. 2013. "The Timescale of Perceptual Evidence Integration Can Be Adapted to the Environment." *Current Biology: CB* 23 (11): 981–86.

Otto, Thomas U., and Pascal Mamassian. 2012. "Noise and Correlations in Parallel Perceptual Decision Making." *Current Biology: CB* 22 (15): 1391–96.

Pannunzi, Mario, Alexis Pérez-Bellido, Alexandre Pereda-Baños, Joan López-Moliner, Gustavo Deco, and Salvador Soto-Faraco. 2015. "Deconstructing Multisensory Enhancement in Detection." *Journal of Neurophysiology* 113 (6): 1800–1818.

Pardo-Vazquez, Jose L., Juan R. Castiñeiras-de Saa, Mafalda Valente, Iris Damião, Tiago Costa, M. Inês Vicente, André G. Mendonça, Zachary F. Mainen, and Alfonso Renart. 2019. "The Mechanistic Foundation of Weber's Law." *Nature Neuroscience* 22 (9): 1493–1502.

Pinto, Lucas, Sue A. Koay, Ben Engelhard, Alice M. Yoon, Ben Deverett, Stephan Y. Thiberge, Ilana B. Witten, David W. Tank, and Carlos D. Brody. 2018. "An Accumulation-of-Evidence Task Using Visual Pulses for Mice Navigating in Virtual Reality." *Frontiers in Behavioral Neuroscience* 12. https://doi.org/10.3389/fnbeh.2018.00036.

Prat-Ortega, Genís, Klaus Wimmer, Alex Roxin, and Jaime de la Rocha. 2021. "Flexible Categorization in Perceptual Decision Making." *Nature Communications* 12 (1): 1283.

Raposo, D., J. P. Sheppard, P. R. Schrater, and A. K. Churchland. 2012. "Multisensory Decision-Making in Rats and Humans." *Journal of Neuroscience* 32 (11): 3726–35.

Roitman, Jamie D., and Michael N. Shadlen. 2002. "Response of Neurons in the Lateral Intraparietal Area during a Combined Visual Discrimination Reaction Time Task" 22 (21): 9475–89.

Stine, Gabriel M., Ariel Zylberberg, Jochen Ditterich, and Michael N. Shadlen. 2020. "Differentiating between Integration and Non-Integration Strategies in Perceptual Decision Making." *eLife* 9 (April). https://doi.org/10.7554/eLife.55365.

Thura, David, Julie Beauregard-Racine, Charles-William Fradet, and Paul Cisek. 2012. "Decision Making by Urgency Gating: Theory and Experimental Support." *Journal of Neurophysiology* 108 (11): 2912–30.

Wang, Xiao-Jing. 2008. "Decision Making in Recurrent Neuronal Circuits." *Neuron* 60 (2): 215–34.

Waskom, Michael L., and Roozbeh Kiani. 2018. "Decision Making through Integration of Sensory Evidence at Prolonged Timescales." *Current Biology: CB* 28 (23): 3850–56.e9.

Wilson, Robert C., and Anne Ge Collins. 2019. "Ten Simple Rules for the Computational Modeling of Behavioral Data." *eLife* 8 (November). https://doi.org/10.7554/eLife.49547.

Wimmer, Klaus, Albert Compte, Alex Roxin, Diogo Peixoto, Alfonso Renart, and Jaime de la Rocha. 2015. "Sensory Integration Dynamics in a Hierarchical Network Explains Choice Probabilities in Cortical Area MT." *Nature Communications* 6 (February): 6177.

Winkel, Jasper, Max C. Keuken, Leendert van Maanen, Eric-Jan Wagenmakers, and Birte U. Forstmann. 2014. "Early Evidence Affects Later Decisions: Why Evidence Accumulation Is Required to Explain Response Time Data." *Psychonomic Bulletin & Review* 21 (3): 777–84.

Wyart, Valentin, Vincent de Gardelle, Jacqueline Scholl, and Christopher Summerfield. 2012. "Rhythmic Fluctuations in Evidence Accumulation during Decision Making in the Human Brain." *Neuron* 76 (4): 847–58.

Yang, Tianming, and Michael N. Shadlen. 2007. "Probabilistic Reasoning by Neurons." *Nature* 447 (7148): 1075–80.

Yates, Jacob L., Il Memming Park, Leor N. Katz, Jonathan W. Pillow, and Alexander C. Huk. 2017. "Functional Dissection of Signal and Noise in MT and LIP during Decision-Making." *Nature Neuroscience* 20 (9): 1285–92.

Znamenskiy, Petr, and Anthony M. Zador. 2013. "Corticostriatal Neurons in Auditory Cortex Drive Decisions during Auditory Discrimination." *Nature* 497 (7450): 482–85.

Zuo, Yanfang, and Mathew E. Diamond. 2019. "Rats Generate Vibrissal Sensory Evidence until Boundary Crossing Triggers a Decision." *Current Biology: CB* 29 (9): 1415–24.e5.
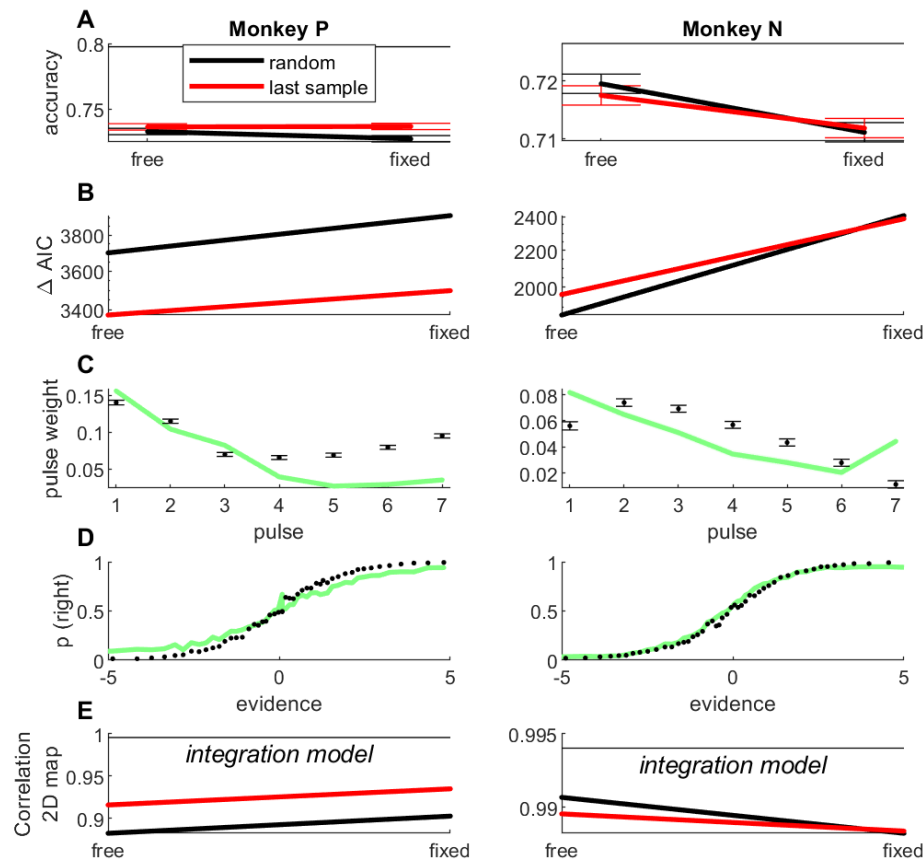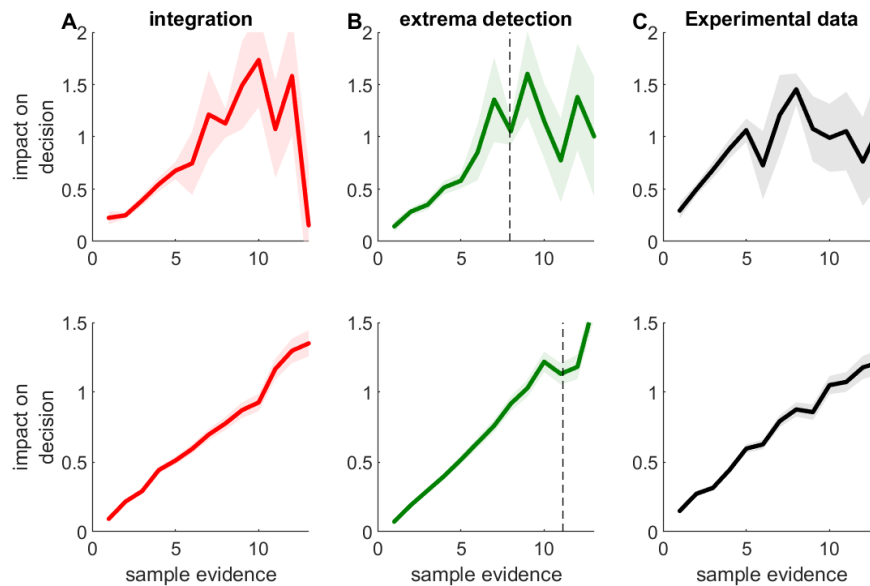
1053

## SUPPLEMENTARY FIGURES



**Supplementary Figure 1. Parameter fits for integration and non-integration models. A.** Modulation gain $\gamma$ per session for the integration model, for each animal (green: monkey P; purple: monkey N). **B.** Mixture coefficients $\pi_i$ of the snapshot model estimated for each monkey, representing the prior probability that each sample is attended on each trial. **C.** Parameters $T$ and $\sigma$ of the extrema-detection model, estimated for each monkey. Error bars correspond to the confidence interval obtained using the Laplace approximation.

**Supplementary Figure 2. Model fits for variants of the snapshot model. A.** Predicted accuracy for the snapshot model fitted to monkey data, as a function of memory span $K$, for fixed lapses (blue curve, $\pi_L = \pi_R = 0.01$) and lapses estimated from the data (black curve). Full lines represent the model with sensory noise ("probabilistic"), dotted lines represent the model without sensory noise ("non-probabilistic"). Memory span $K$ corresponds to the number of successive samples used to define the decision on each trial (see Methods). The horizontal bar corresponds to the average accuracy of the animal. **B.** AIC difference between snapshot and integration model. Legend as in A. Positive values indicate that the snapshot model provides a worse fit. **C.** Psychometric curve for the snapshot model with span $K=3$ samples, sensory noise and free lapse parameters (best snapshot model variant according to AIC). **D.** Psychophysical kernel for the same variant of the model. **E.** Correlation between data and model integration maps for variants of the snapshot model.
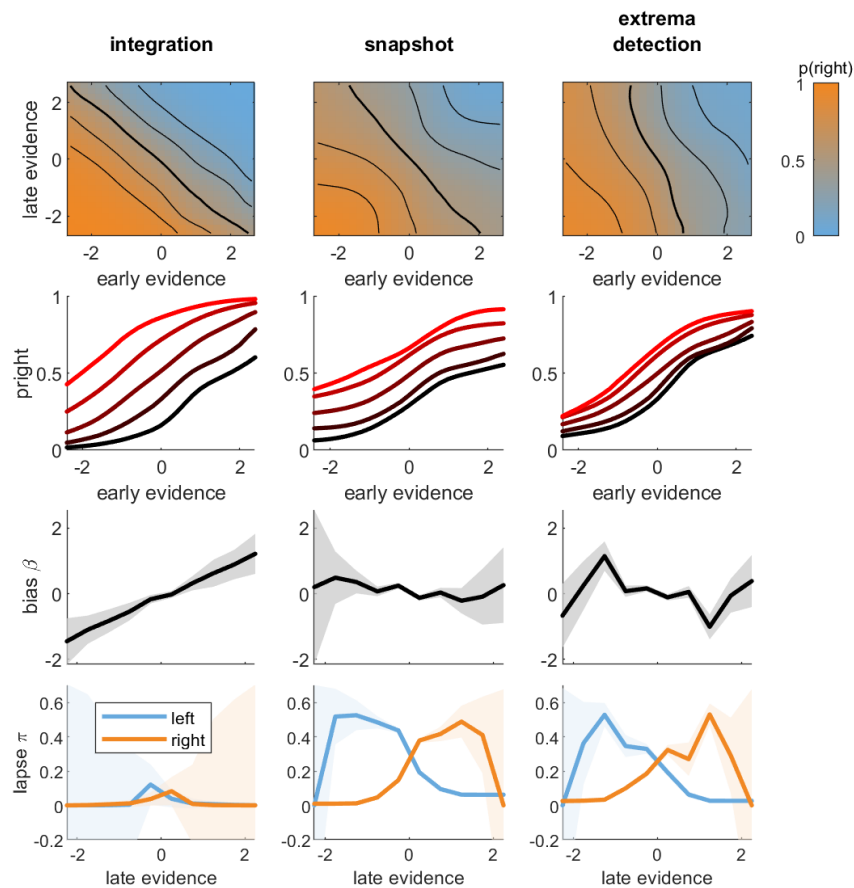
**Supplementary Figure 3. Model fits for variants of the extrema-detection model. A.** Predicted accuracy for the extrema-detection model fitted to the monkey data, for random (black curves) and last sample (red curve) default rule, and for fixed lapses ($\pi_L = \pi_R = 0.01$) or lapse parameters estimated from the data. The horizontal bar indicates animal accuracy. **B.** AIC difference between variants of the extrema-detection model and the integration model. Legend as in A. Positive values indicate that the extrema-detection model provides a worse fit. **C-D.** Psychometric curve (C) and psychophysical kernel (D) for the model variant that provided the best match to behavior in terms of predicted accuracy and AIC: free lapse parameters and last sample rule. **E.** Correlation between integration maps from animal data and simulated data (see Figure 4) for variants of the extrema-detection model. The horizontal bar marks the correlation between experimental data and the integration model.
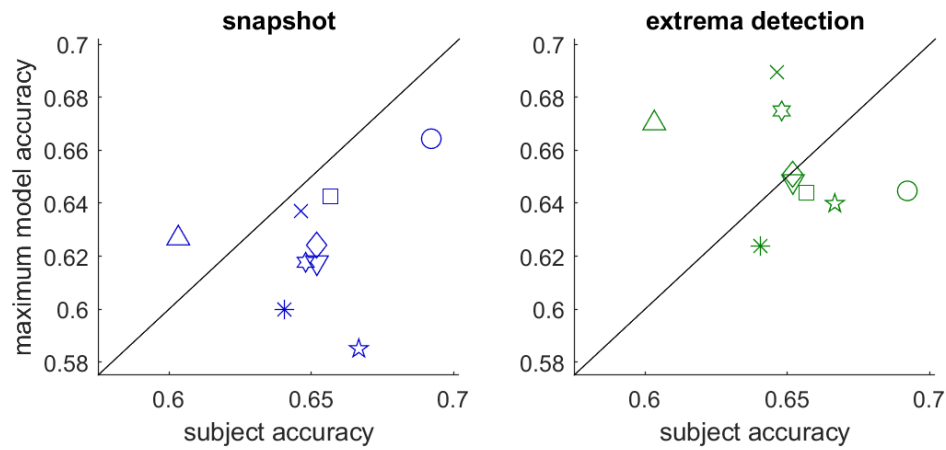
31

**Supplementary Figure 4. Subjective weights for animal data and simulated models.** Impact on decision of individual samples as a function of absolute sample evidence. Shaded area: standard error of the weight. Top row: monkey P; bottom row: monkey N. **A.** Integration model. **B.** extrema-detection model. The vertical dotted line marks the value of the threshold $T$ estimated from animal data. **C.** Impact on decision of individual pulses, estimated from each monkey.
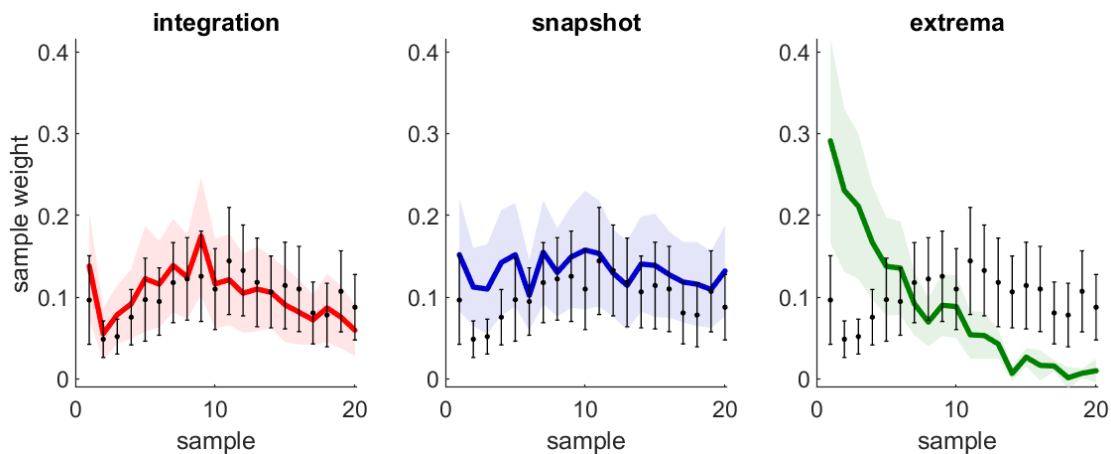


**Supplementary Figure 5. Integration of early and late evidence for monkey P. A.** Integration map. Legend as in Figure 4A. **B.** Conditional psychometric curves. Legend as in Figure 4B. **C.** Bias and lapse parameters from conditional psychometric curves, as a function of late evidence. Legend as in Figure 4D-E.
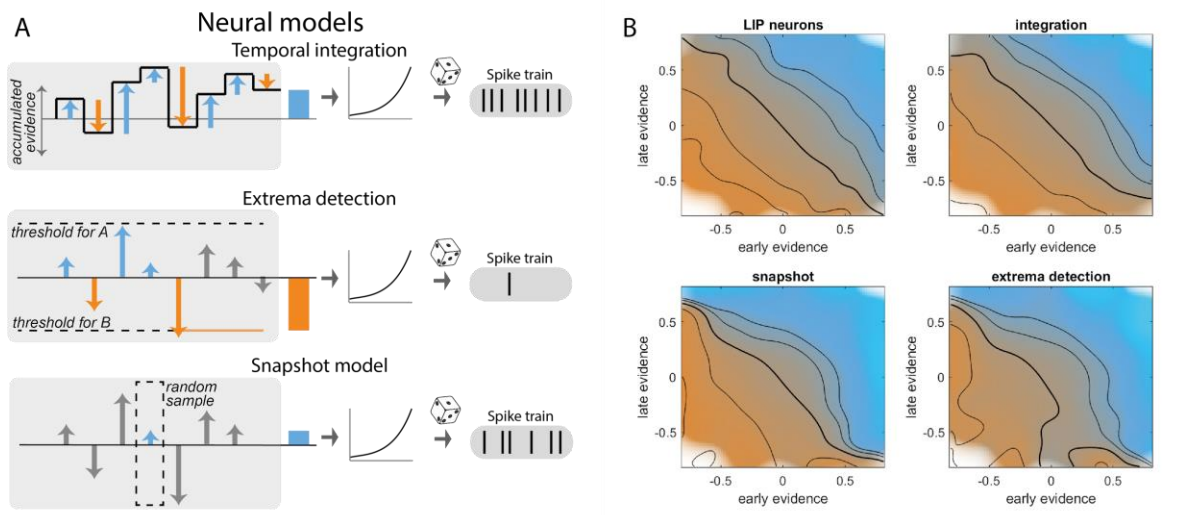
1097
1098 **Supplementary Figure 6. Integration between early and late evidence for simulated data from**
1099 **integration and non-integration models.** Data was simulated for each model from parameters
1100 estimated from monkey N. Left panels: integration model. Middle panels: snapshot models. Right
1101 panels: extrema-detection models. **A.** Integration maps. **B.** Conditional psychometric curves. **C.** Lateral
1102 bias and **D.** lapse parameters estimated from conditional psychometric curves, as a function late
1103 evidence. Legend as in Figure 4.

1104

**Supplementary Figure 7. Maximum accuracy of the non-integration models vs. human subject accuracy in the orientation discrimination task.** Left panel: snapshot model (with span $K$=1). Right panel: extrema-detection. Each symbol represents a subject.

1105
1106
1107

1108



1109

**Supplementary Figure 8. Psychophysical kernels for animals and models in rats ($n$=3) performing the DSS task with 20-sample stimuli.**

1110
1111

1112

1113

**Supplementary Figure 9. Individual LIP neurons integrate sensory information over stimulus sequence. A.** Neural models for temporal integration, extrema-detection and snapshot model. **B.** Integration map for LIP neurons, and simulated neurons following either integration, extrema-detection or snapshot model. Color represents the average normalized spike count per bins of neuron-weighted early and late evidence (see Methods). Isolines represent values of 0.4, 0.6, 1, 1.4 and 1.8.