# Neural Correlates of Optimal Multisensory Decision Making

Han Hou[1,2], Qihao Zheng[1,2*], Yuchen Zhao[1,2*], Alexandre Pouget[3#], Yong Gu[1#§]

[1]Institute of Neuroscience, Key Laboratory of Primate Neurobiology, CAS Center for Excellence

in Brain Science and Intelligence Technology, Chinese Academy of Sciences, Shanghai, China

[2]University of Chinese Academy of Sciences, Beijing, China

[3]University of Geneva, Geneva, Switzerland

* These authors contributed equally to this work.

# These senior authors contributed equally to this work.

§ Correspondence should be addressed to: Yong Gu (guyong@ion.ac.cn)

14

# Abstract

16 Perceptual decisions are often based on multiple sensory inputs whose reliabilities rapidly vary

17 over time, yet little is known about how our brain integrates these inputs to optimize behavior. Here

18 we show multisensory evidence with time-varying reliability can be accumulated near optimally,

19 in a Bayesian sense, by simply taking time-invariant linear combinations of neural activity across

20 time and modalities, as long as the neural code for the sensory inputs is close to an invariant linear

21 probabilistic population code (ilPPC). Recordings in the lateral intraparietal area (LIP) while

22 macaques optimally performed a vestibular-visual multisensory decision-making task revealed that

23 LIP population activity reflects an integration process consistent with the ilPPC theory. Moreover,

24 LIP accumulates momentary evidence proportional to vestibular acceleration and visual velocity

25 which are encoded in sensory areas with a close approximation to ilPPCs. Together, these results

26 provide a remarkably simple and biologically plausible solution to optimal multisensory decision

27 making.

28

29 **Keywords:** Perceptual decision making, multisensory integration, LIP, probabilistic population

30 code, vestibular, optic flow, self-motion perception

## Introduction

Most perceptual decisions are based on multiple sensory inputs whose reliabilities vary over time. For instance, a predator can rely on both auditory and visual information to determine when and where to strike a prey, but these two sources of information are not generally equally reliable, nor are their reliabilities constant over time: as the prey gets closer, the quality of the image and sound typically improves, thus increasing their reliabilities. Although such multisensory decision making happens frequently in the real world, the underlying neural mechanisms remain largely unclear.

The so-called drift-diffusion model (DDM) (**[1]Ratcliff, 1978; [2]Ratcliff and McKoon, 2008; [3]Ratcliff and Rouder, 1998; [4]Ratcliff and Smith, 2004**), a widely used model of perceptual decision making, cannot deal with such decisions optimally in its most standard form. DDMs have been shown to implement the optimal policy for decisions involving just one source of sensory evidence whose reliability is constant over time (**[5]Laming, 1968; [6]Bogacz, et al., 2006**). Under such conditions, DDMs can implement the optimal strategy by simply summing evidence over time until an upper or lower bound, corresponding to the two possible choices, is hit (**[6]Bogacz, et al., 2006**). This type of models lends itself to a straightforward neural implementation in which neurons simply add their sensory inputs until they reach a preset threshold (**[2]Ratcliff and McKoon, 2008; [7]Gold and Shadlen, 2007**).

When multiple sensory inputs are involved, the standard DDMs can accumulate sensory evidence optimally as long as the reliabilities of the evidence stay constant during a single trial and across trials. Under this scenario, optimal integration of evidence over time can be achieved by first taking a weighted sum of the momentary evidence at each time step, with weights proportional to the reliability of each sensory stream, followed by temporal integration (**[8]Drugowitsch, et al., 2014**). However, this strategy no longer works when the reliabilities change over time within a single trial.

56  In this case, the momentary evidence must be linearly combined with weights proportional to the

57  time-varying reliabilities, which requires that the synaptic weights change on a very fast time scale

58  since, in the real life, reliability can change significantly over tens of milliseconds. Moreover, when

59  the reliabilities of the sensory inputs are not known in advance, which is typically the case in real-

60  world situations, neurons cannot determine how to appropriately modulate their synaptic weights

61  until after the sensory inputs have been observed. Therefore, even if it is possible to extend standard

62  DDMs to time-varying reliability **([8]Drugowitsch, et al., 2014)**, it is unclear how such a solution

63  could be implemented biologically.

64

65  In contrast, there exists another class of models which does not necessarily involve changes in

66  synaptic strength. As long as the sensory inputs are encoded with what is known as "invariant linear

67  probabilistic population codes" (ilPPC), the neural solution for optimal multisensory integration is

68  remarkably simple: it only requires that neurons compute linear combinations of their inputs across

69  time or modalities using fixed—reliability-independent—synaptic weights **([9]Beck, et al., 2008;**

70  **[10]Ma, et al., 2006)**. This solution relies on one specific property of ilPPC: the reliability of the

71  neural code is proportional to the amplitude of the neural responses. As a result, when summing

72  two sensory inputs with unequal reliability, the sensory input with the lowest reliability contribute

73  less to the sum because of its lower firing rate. This is formally equivalent to weighting Gaussian

74  samples with their reliability in an extended DDM, except that there is no need for actual weight

75  changes with ilPPC **([10]Ma, et al., 2006)**. Hence, the ilPPC framework is a promising solution to

76  multisensory decision-making tasks, but it lacks physiological supports.

77

78  To investigate whether the brain may implement this solution, we recorded the activity of single

79  neurons in the lateral intraparietal area (LIP) in macaques trained to discriminate their heading

80  direction of self-motion based on multiple sensory inputs: vestibular signals, visual optic flow, or

81  both. Importantly, the vestibular and visual stimuli followed a Gaussian-shape velocity temporal

82    profile, producing naturally varied cue reliability over time within each trial. This behavioral

83    paradigm has been well-established for studying multisensory heading discrimination in the past

84    decade ([11]Fetsch*, et al.*, 2012; [12]Gu*, et al.*, 2008; [13]Fetsch*, et al.*, 2013). Nevertheless, these

85    previous studies focused on areas that encode momentary heading inputs, leaving it unknown how

86    these sensory inputs are further accumulated by downstream neurons (e.g. LIP) during perceptual

87    decision making.

88

89    We focus first in LIP because it is the most extensively studied brain region where buildup choice-

90    related activity has been found during visuomotor decisions in macaques ([7]**Gold and Shadlen,**

91    **2007;** [14]**Shadlen and Newsome, 2001;** [15]**Shadlen and Newsome, 1996;** [16]**Huk*, et al.*, 2017;**

92    [17]**Roitman and Shadlen, 2002)**. In addition, LIP receives abundant anatomical inputs

93    ([18]**Boussaoud*, et al.*, 1990)** from areas encoding momentary vestibular and visual self-motion

94    information for heading discrimination, such as the dorsal medial superior temporal (MSTd) area

95    ([12]**Gu*, et al.*, 2008;** [19]**Gu*, et al.*, 2006)** and the ventral intraparietal area (VIP) ([20]**Chen*, et al.*, 2011c;**

96    [21]**Chen*, et al.*, 2013)**. It is therefore expected that the activity of LIP neurons should carry buildup

97    choice signals germane to the formation of multisensory decisions. Note that two recent rodent

98    studies ([22]**Nikbakht*, et al.*, 2018;** [23]**Raposo*, et al.*, 2014)** also have described multisensory decision

99    signals in rat posterior parietal cortex, a region analogous to its primate counterpart. However, these

100   studies did not characterize the computational solution implemented by these neural circuits, which

101   is precisely the question we investigate here. Specifically, we explored whether the response of LIP

102   neurons is consistent with the ilPPC theory in which neurons take fixed linear combinations of their

103   sensory inputs without any need for complex, time-dependent, modality-specific, reweighting of

104   the sensory inputs during multisensory decision making.

# Results

## Optimal multisensory decision‑making behavior on macaques

We trained two macaque monkeys to perform a vestibular-visual multisensory decision-making task ([12]Gu, *et al.*, 2008) (**Figure 1a**). On each trial, the monkeys experienced a 1.5s-fixed-duration forward motion with a small deviation either to the left or to the right of the dead ahead. At the end of the trial, the animals were required to report the perceived heading direction by making a saccade decision to one of the two choice targets (**Figure 1b**). We randomly interleaved three cue conditions over trials: a vestibular condition and a visual condition in which heading information was solely provided by inertial cues and optic flow, respectively, and a combined condition consisting of congruent vestibular and visual cues. Importantly, both the vestibular and visual stimuli followed a Gaussian-shape velocity temporal profile, peaking at the middle of the 1.5-s stimulus duration. This modulation of velocity over time has an important implication for the reliability of the sensory inputs provided to the animals. Indeed, previous psychophysical studies have established that a model in which the reliability of the visual flow field is proportional to velocity and the reliability of the vestibular signal is proportional to the acceleration, provides the best fits to the behavioral data ([8]Drugowitsch, *et al.*, 2014). Therefore, this stimulus allows us to test how neural circuits accumulate multisensory evidence whose reliability varies over time with distinct temporal profiles (see below).

To quantify the monkeys' behavioral performance, we plotted psychometric curves for each cue condition (**Figure 1c**). Consistent with the previous results ([12]Gu, *et al.*, 2008), the monkeys made more accurate decisions in the combined condition, as evidenced by a steeper psychometric function and a smaller psychophysical threshold (**Figure 1c**). Across all recording sessions and for both monkeys, the psychophysical threshold of the combined condition was significantly smaller than those of single cue conditions and close to the threshold predicted by optimal Bayesian
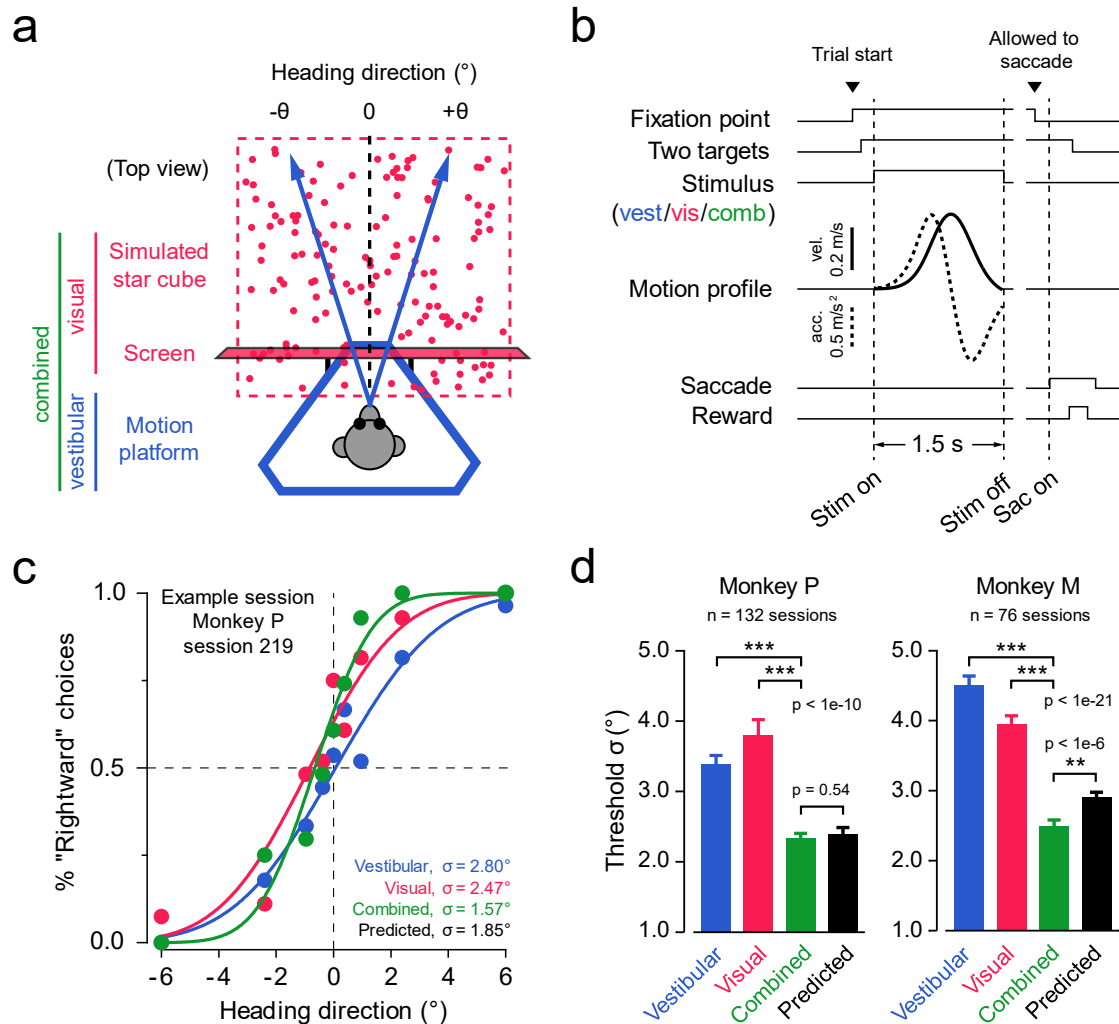
**Figure 1  Optimal cue integration in vestibular-visual multisensory decision-making task.**

**(a)** Schematic drawing of the experimental setup (top view). The vestibular (blue) and visual (red) stimuli of self-motion were provided by a motion platform and an LCD screen mounted on it, respectively. The monkey was seated on the platform and physically translated within the horizontal plane (blue arrows), whereas the screen rendered optic flow simulating what the monkey would see when moving through a three-dimensional star field (red dots). In a combined condition (green), both vestibular and visual stimuli were presented synchronously. The monkey's task was to discriminate whether the heading direction was to the left or the right of the straight ahead (black dashed line). **(b)** Task timeline. The monkey initiated a trial by fixating at a fixation point, and two choice targets appeared. The monkey then experienced a 1.5-s forward self-motion stimulus with a small leftward or rightward component, after which the monkey reported his perceived heading by making a saccadic eye movement to one of the two targets. The self-motion speed followed a Gaussian-shape profile. **(c)** Example psychometric functions from one session. The proportion of "rightward" choices is plotted against the headings for three cue conditions respectively. Smooth curves represent best-fitting cumulative Gaussian functions. **(d)** Average psychophysical thresholds from two monkeys for three conditions and predicted thresholds calculated from optimal cue integration theory (black bars). Error bars indicate s.e.m.; p values were from paired t-test.

131    multisensory integration (**[24]Knill and Richards, 1996**) (**Figure 1d**). Therefore, the monkeys can

132    integrate vestibular and visual cues near-optimally during our multisensory decision-making task.

133

## Heterogeneous multisensory choice signals in LIP

135    Next, we set out to explore how these optimal decisions were formed in the brain. We recorded

136    from 164 single, well-isolated neurons in LIP of two monkeys while they were performing the task

137    (**Supplementary Figure 1**). As expected, we found buildup choice-related signals in LIP neurons

138    under all cue conditions. As shown in PSTHs of the example cells (**Figure 2a** and **Supplementary**

139    **Figure 2**), there was generally an increasing divergence between the neuron's firing rate on trials

140    in which the monkey chose the target in the neuron's response field (IN choices, solid curves) and

141    trials in which the opposite target was chosen (OUT choices, dashed curves). Importantly, in all cue

142    conditions, the buildup choice signals tended to be stronger for heading directions more distant

143    away from straight ahead (**Supplementary Figure 3**), suggesting that the response of LIP neurons

144    reflects the accumulation of visual and vestibular sensory evidence for heading judgments.

145

146    To better quantify the choice-related signals, we used a ROC analysis to generate an index of choice

147    divergence (CD) (**[23]Raposo*, et al.*, 2014**) that measures the strength of the choice signals (**Figure**

148    **2b**). The four cells illustrated in Fig. 2 exhibited canonical ramping choice signals, but their CDs

149    varied greatly across cue conditions. For example, for Cell 1, the CD was largest in the combined

150    condition, modest in the visual condition, and smallest in the vestibular condition. By contrast, for

151    Cell 4, the CD was largest in the vestibular condition. The heterogeneity of choice signals was also

152    manifest at the population level.  Approximately half of the LIP neurons exhibited statistically

153    significant CD ($p < 0.05$, two-sided permutation test) in each cue condition (vestibular: 52%, visual:

154    46%, combined: 59%; **Figure 2c**), but these three subpopulations did not fully overlap. While more

155    than   two   thirds   of   neurons   (76%)   had   significant   choice   signals   in   *any*   of   the   three
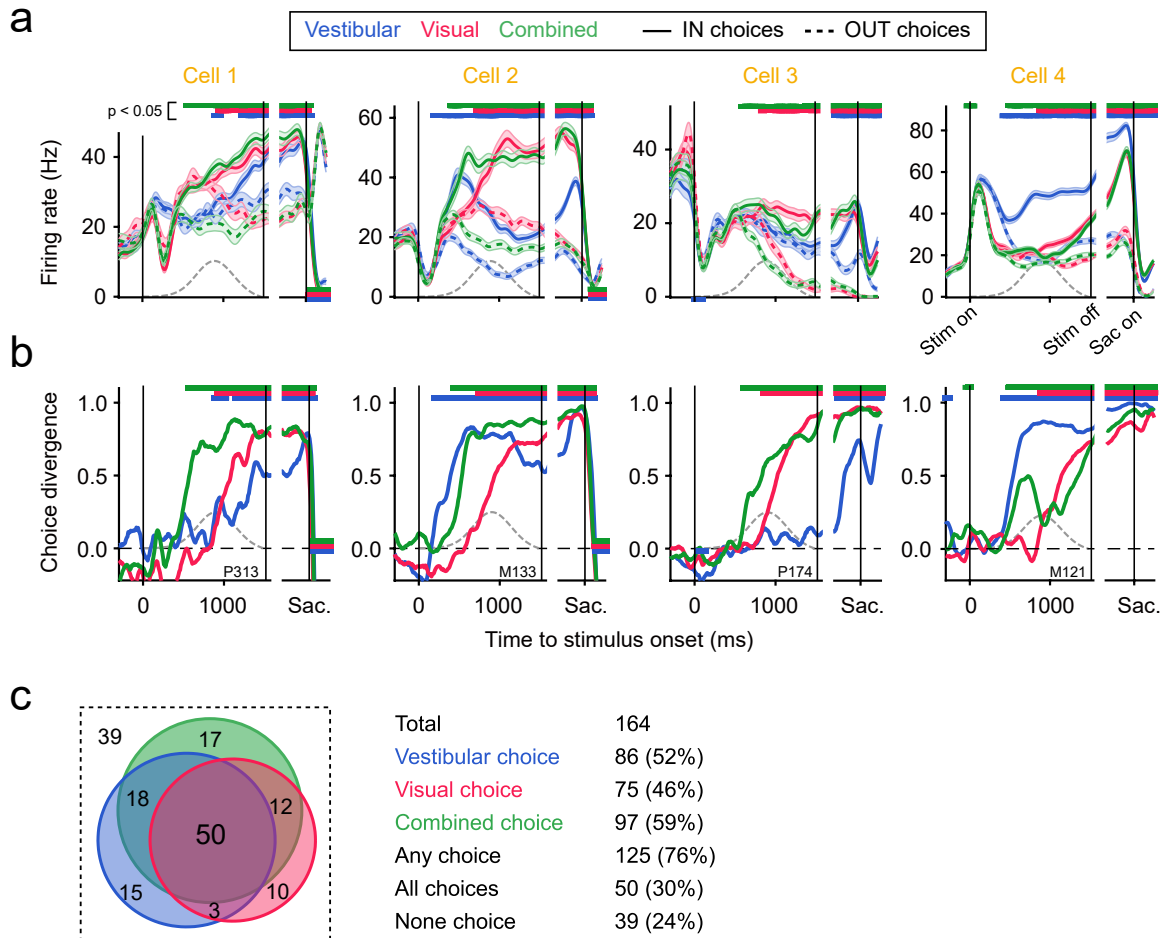
**Figure 2  Heterogeneous choice signals in LIP population.**

**(a)** Peri-stimulus time histograms (PSTHs) of four example cells. Spike trains were aligned to stimulus onset (left subpanels) and saccade onset (right subpanels), respectively, and grouped by cue condition and monkey's choice. Vestibular, blue; visual, red; combined, green. Toward the cell's response field (RF), or IN choices, solid curves; away from the cell's RF, or OUT choices, dashed curves. Mean firing rates were computed from 10-ms time windows and smoothed with a Gaussian ($\sigma$ = 50 ms); only correct trials or 0° heading trials were included. Shaded error bands, s.e.m. Horizontal color bars represent time epochs in which IN and OUT trials have significantly different firing rates (p < 0.05, t-test), with the color indicating cue condition and the position indicating the relationship between IN and OUT firings (IN > OUT, top; IN < OUT, bottom). Gray dashed curves represent the actual speed profile measured by an accelerometer attached to the motion platform. **(b)** Choice divergence (CD) of the same four cells. CD ranged from -1 to 1 and was derived from ROC analysis for PSTHs in each 10-ms window (see Methods). Horizontal color bars are the same as in **a** except that p-values were from permutation test (n = 1000). **(c)** Venn diagram showing the distribution of choice signals. Numbers within colored areas indicate the numbers of neurons that have significant grand CDs (CD computed from all spikes in 0–1500 ms) under the corresponding combinations of cue conditions.

156

157    conditions ("Any choice" cells in **Figure 2c**), only a third of neurons (30%) had significant choice

158    signals in *all* of the three conditions ("All choices" cells in **Figure 2c**).

159

160    Apart from the heterogeneous choice signals, LIP also encodes heterogeneous sensory modality

161    signals. For example, Cell #12 in **Supplementary Figure 2** exhibited differentiated firing rates

162    across cue conditions without much choice-related signal. In fact, as shown in **Supplementary**

163    **Figure 4a,** the majority of LIP neurons actually carried mixed choice and modality signals,

164    exhibiting a category-free like neural representation as previously seen in rat posterior parietal

165    cortex (**[23]Raposo*, et al.***, 2014)**. However, although randomly mixed at the single neuron level, the

166    choice and modality signals can still be linearly decoded from the LIP population (**Supplementary**

167    **Figure 5**). Therefore, we ignore the mixed modality signals thereafter, since they are irrelevant to

168    our heading discrimination task and orthogonal to the decision signals that we really care about.

169

170    Another potential difficulty in interpreting LIP activity arises from the fact that LIP neurons also

171    multiplex a combination of temporally overlapping decision- and non-decision- signals (**[25]Park*, et***

172    ***al.*, 2014; [26]Meister*, et al.*, 2013)**. In particular, the signal of saccade preparation may interfere with

173    the one reflecting evidence accumulation (**[14]Shadlen and Newsome, 2001)**. However, this was not

174    likely to be an issue in our study. In our fixed-duration task, we introduced a 300–600 ms delay

175    between the stimulus offset and the time at which the monkey was allowed to saccade (see

176    Methods). Moreover, the monkeys tended to stop integrating evidence around 500 ms prior to the

177    stimulus offset (see **Figure 3b** and below), further separating in time the processes of evidence

178    accumulation and saccade preparation. Therefore, the premotor activity of LIP should not play a

179    significant role in our analysis of multisensory evidence accumulation.

180

181    **LIP integrates vestibular acceleration and visual velocity**

182    Despite the high degree of heterogeneity, there was nonetheless a property shared amongst the LIP

183    neurons, namely, the temporal dynamics of the ramping activity was significantly faster in the

184    vestibular and combined conditions than in the visual condition (**Figure 3**). This was evident not

185    only in the averaged rate-based or ROC-based measures ("Any choice" cells, **Figure 3a, b**), but

186    also in the cell-by-cell analysis (**Figure 3c**). Notably, the averaged divergence time under the

187    vestibular and combined conditions aligned well to the acceleration peak of the Gaussian-shape

188    motion profile, whereas the divergence time under the visual condition better aligned to the velocity

189    peak (**Figure 3c**, dashed curves). This suggests that the physical quantities being integrated over

190    time are speed for the visual stimulus and acceleration for the vestibular stimulus.

191

192    An alternative explanation, however, might be that the apparent ~400 ms interval between the

193    vestibular and visual ramping was caused purely by a difference in their sensory latencies rather

194    than in their underlying physical quantities. For example, LIP activity could have been driven by

195    an ultrafast vestibular signal but a slow visual signal, both of which followed the velocity of the

196    motion. To test this, we designed an experiment in which we used two distinct velocity profiles, a

197    wide one and a narrow one (**Figure 3d**). These profiles were designed to have temporally aligned

198    velocity peaks but misaligned acceleration peaks. If our original physical-quantity hypothesis was

199    correct, we would expect the visual ramping to remain nearly the same under both profiles, while

200    the vestibular ramping should start earlier for the wide profile than for the narrow one, thus

201    reflecting the earlier acceleration peak under the wide profile. In contrast, if the sensory-latency

202    hypothesis was correct, there should be no shift in either the vestibular or visual ramping across the

203    two profiles. Our data matches the first prediction (**Figure 3e, f**). In other words, the temporal

204    discrepancy between the vestibular and visual ramping activities indeed resulted from different

205    physical quantities underlying the momentary evidence fed into LIP. This physiological finding

206    echoed a recent psychophysical study showing that, at the behavioral level, human subjects

207    optimally integrate vestibular and visual momentary evidence with reliability following the
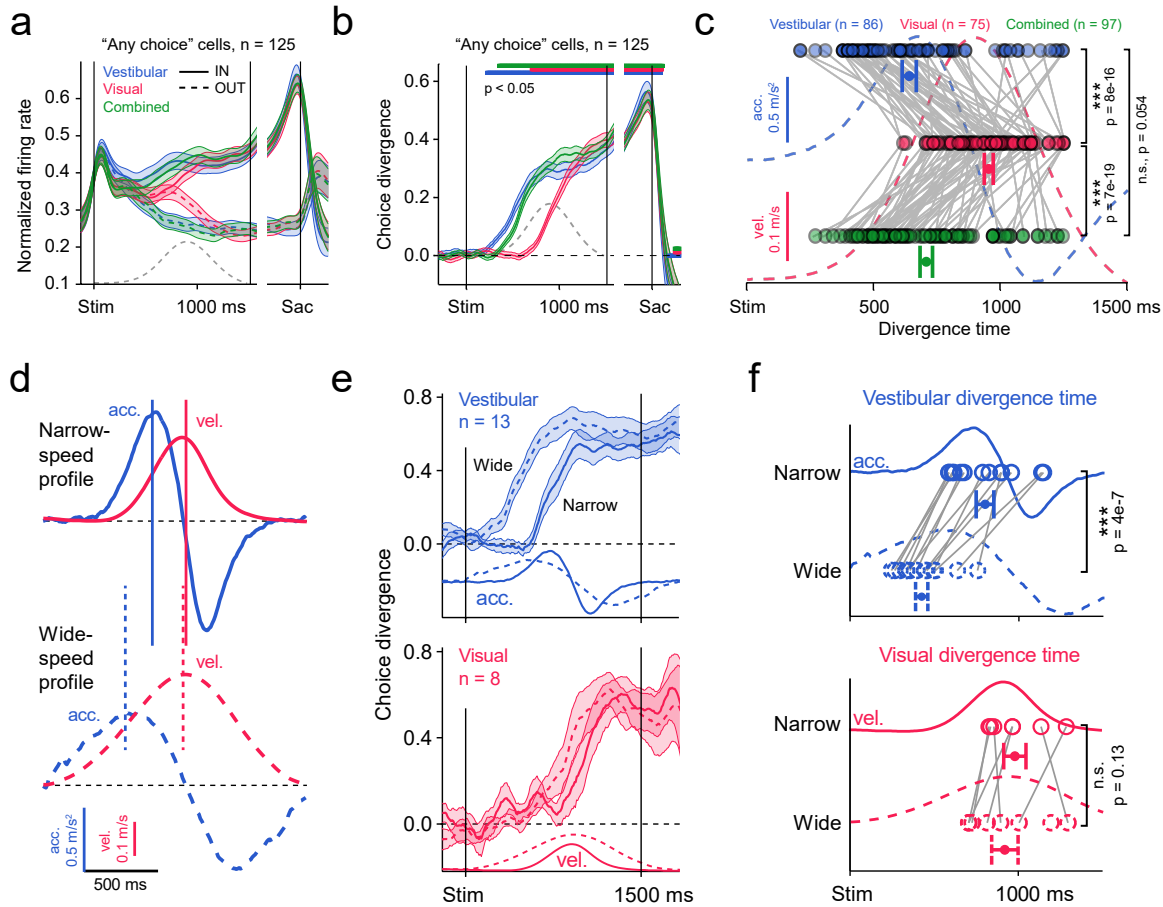
**Figure 3  LIP integrates vestibular acceleration but visual speed.**

**(a and b)** Population average of normalized PSTHs (**a**) and CD (**b**) from 125 "any choice" cells. The vestibular (blue) and combined (green) CDs ramped up much earlier than the visual one (red). Horizontal color bars indicate the time epochs in which population CDs are significantly larger than zero (p < 0.05, t-test). Gray dashed curve, the actual Gaussian speed profile; shaded error bands, s.e.m. **(c)** Divergence time of cells with significant grand CD for each condition. Divergence time was defined as the first occurrence of a 250-ms window in which CD was consistently larger than zero (p < 0.05, permutation test). Gray lines connect data from the same cells; acceleration and speed profiles shown in the background. Data points with horizontal error bars, mean ± s.e.m. of population divergence time; p values, t-test. **(d)** Two motion profiles used to isolate contributions of acceleration and speed to LIP ramping. Top and solid, the narrow-speed profile; bottom and dashed, the wide-speed profile; blue, acceleration; red, speed. Note that by widening the speed profile, we shifted the time of acceleration peak forward (blue vertical lines) while keeping the speed peak unchanged (red vertical lines). **(e)** Vestibular and visual CDs under the two motion profiles. **(f)** Comparison of divergence time between narrow and wide profiles. Note that the vestibular divergence time was significantly shifted, whereas the visual one was not, indicating that LIP integrates sensory evidence from vestibular acceleration and visual speed.

208

209     amplitude of acceleration and velocity, respectively (**[8]Drugowitsch*, et al.***, 2014)**.

210

## Network model implementing ilPPC for multisensory decision making

212     Next, we developed a neural model of multisensory decision making (refer to as M1 thereafter)

213     which takes as input vestibular neurons tuned to acceleration and visual neurons tuned to velocity

214     as observed *in vivo* (equation (2) and (3) in **Methods**; **Figure 4a**). These inputs converge onto an

215     integration layer which takes the sum of the visual and vestibular inputs, as well as integrates this

216     summed input over time. This layer projects in turn to an output layer, labeled LIP, which sums the

217     integrated visuo-vestibular inputs with the activity from another input layer encoding the two

218     possible targets to which the animal can eventually saccade (**Figure 4b, c**). As long as the input

219     layers encode the sensory inputs with ilPPC, this simple network can be shown analytically to

220     implement the Bayes optimal solution even when the reliability of the sensory inputs vary over

221     time as is the case in our experiment (**[9]Beck*, et al.***, 2008; [10]Ma*, et al.***, 2006)**. Note that separating

222     the integration layer from the LIP layer is not critical to our results. We did so to reflect the fact

223     that current experimental data suggest that LIP may not be the layer performing the integration *per*

224     *se*, but may only reflect the results of this integration (**[27]Katz*, et al.***, 2016)**.

225

226     In an ilPPC, the gain, or amplitude, of the tuning curves of the neurons should be proportional to

227     the reliability of the encoded variable. For instance, in the case of vestibular neurons, the amplitude

228     of the tuning curves to heading should scale with acceleration. In vivo, however, the responses of

229     vestibular and visual neurons are not fully consistent with the assumption of ilPPC because while

230     the amplitude does increase with reliability, in some neurons, the baseline activity decreases with

231     reliability (equation (2) and (3) in **Methods** and **Supplementary Figure 6a**). This violation of the

232     ilPPC assumption implies that a simple sum of activity could incur an information loss. Fortunately,

233     this information loss is small for a population of neurons with tuning properties similar to what has
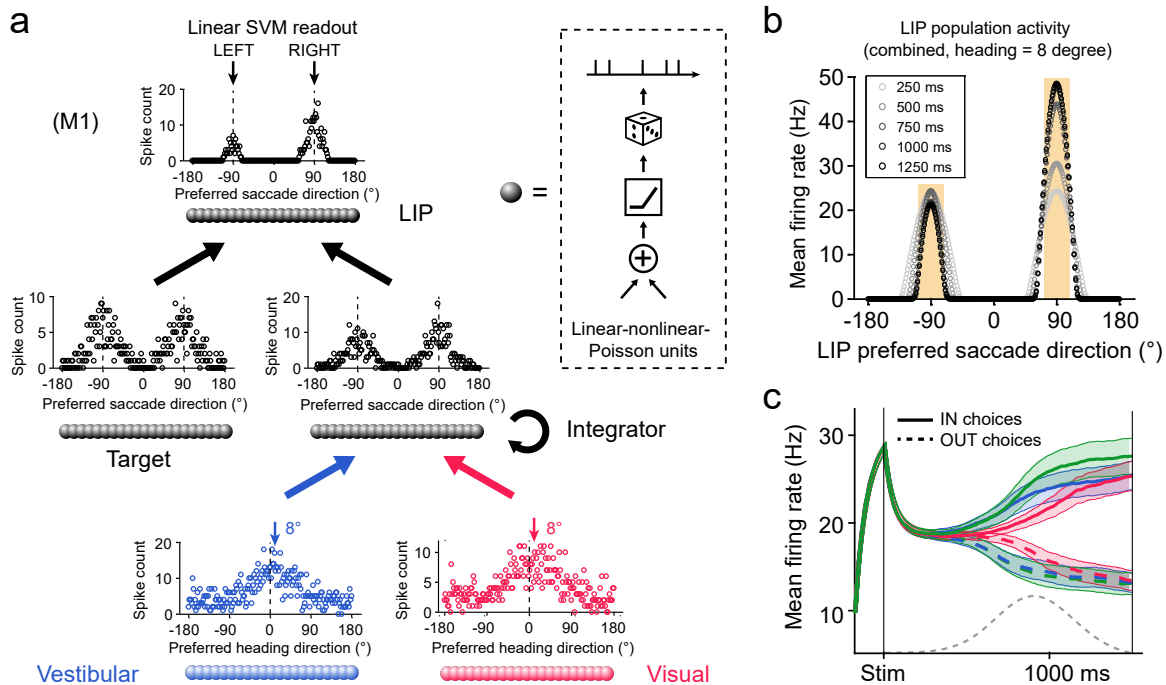
**Figure 4  Neural network model with invariant linear probabilistic population codes (ilPPC).**

**(a)** Network architecture of model M1. The model consists of three interconnected layers of linear-nonlinear-Poisson units (inset). Units in Vestibular and Visual layers have bell-shape ilPPC-compatible tuning curves for heading direction and receive heading stimuli with temporal dynamics following acceleration and speed, respectively. The intermediate Integrator layer simply sums the incoming spikes from the two sensory layers over time and transforms the tuning curves for heading direction to that for saccade direction (-90°, leftward choice; +90°, rightward choice). The LIP layer receives the integrated heading inputs from the Integrator layer, together with visual responses triggered by the two saccade targets. LIP units also have lateral connections implementing short-range excitation and long-range inhibition. Once a decision boundary is hit, or when the end of the trial is reached (1.5 s), LIP activity is decoded by a linear support vector machine for action selection (see **Methods**). Circles indicate representative patterns of activity for each layer; spike counts from 800–1000 ms; combined condition, 8° heading. **(b)** Population firing rate in the LIP layer at five different time points (the same stimulus as in **a**, averaged over 100 repetitions). **(c)** Average PSTHs across LIP population. Trials included three cue conditions and nine heading directions (±8°, ±4°, ±2°, ±1°, 0°). To mimic the experimental procedure, only units with preferred saccade direction close to ±90° were used (with deviation less than 20°; yellow shaded area in **b**). Notations are the same as in **Figure 2a and Figure 3a**.

234

235    been reported experimentally and information limiting correlations ([28]**Moreno-Bote, et al., 2014**).

236    Indeed, we found numerically that the information loss was around 5% over a wide range of

237    parameters values (Fano factor, mean correlation, baseline changes, and so on) (**Supplementary**

238    **Modeling** and **Supplementary Figure 6**).

239

240    Importantly, we also endowed the network with a stopping mechanism which terminates sensory

241    integration whenever a function of the LIP population activity reaches a preset threshold (see

242    Methods). Our experiment is not a reaction time experiment and may not require, in principle, such

243    a stopping bound. However, as can be seen in **Figure 3b**, LIP population response saturates around

244    1s, suggesting that evidence integration stops prematurely. This is indeed consistent with the

245    previous results suggesting that animals and humans use a stopping bound even in fixed duration

246    experiments ([29]**Kiani, et al., 2008**).

247

### LIP data are compatible with the iIPPC framework

249    In the first set of simulations on M1, we adjusted the height of the stopping bounds and found that

250    the model can replicate the near optimal animals' performance (**Figure 5a**). We then plotted the

251    activity of a typical output neuron (in the LIP layer) in all three conditions. As expected, the activity

252    in the combined condition is roughly equal to the sum of the vestibular-only activity and visual-

253    only activities (**Figure 5b**), at least in the first half of the trial. In the second half of the trial, the

254    activity in the combined condition deviates strongly from the sum because the traces correspond to

255    averages across trials that terminated at different times on different trials due to the stopping bound.

256

257    Neurons in M1 are homogeneous in the sense that they all take a perfect sum of their vestibular and

258    visual inputs. Importantly, however, optimal integration does not require such a perfect sum; it can

259    also be achieved with random linear combinations of vestibular and visual inputs
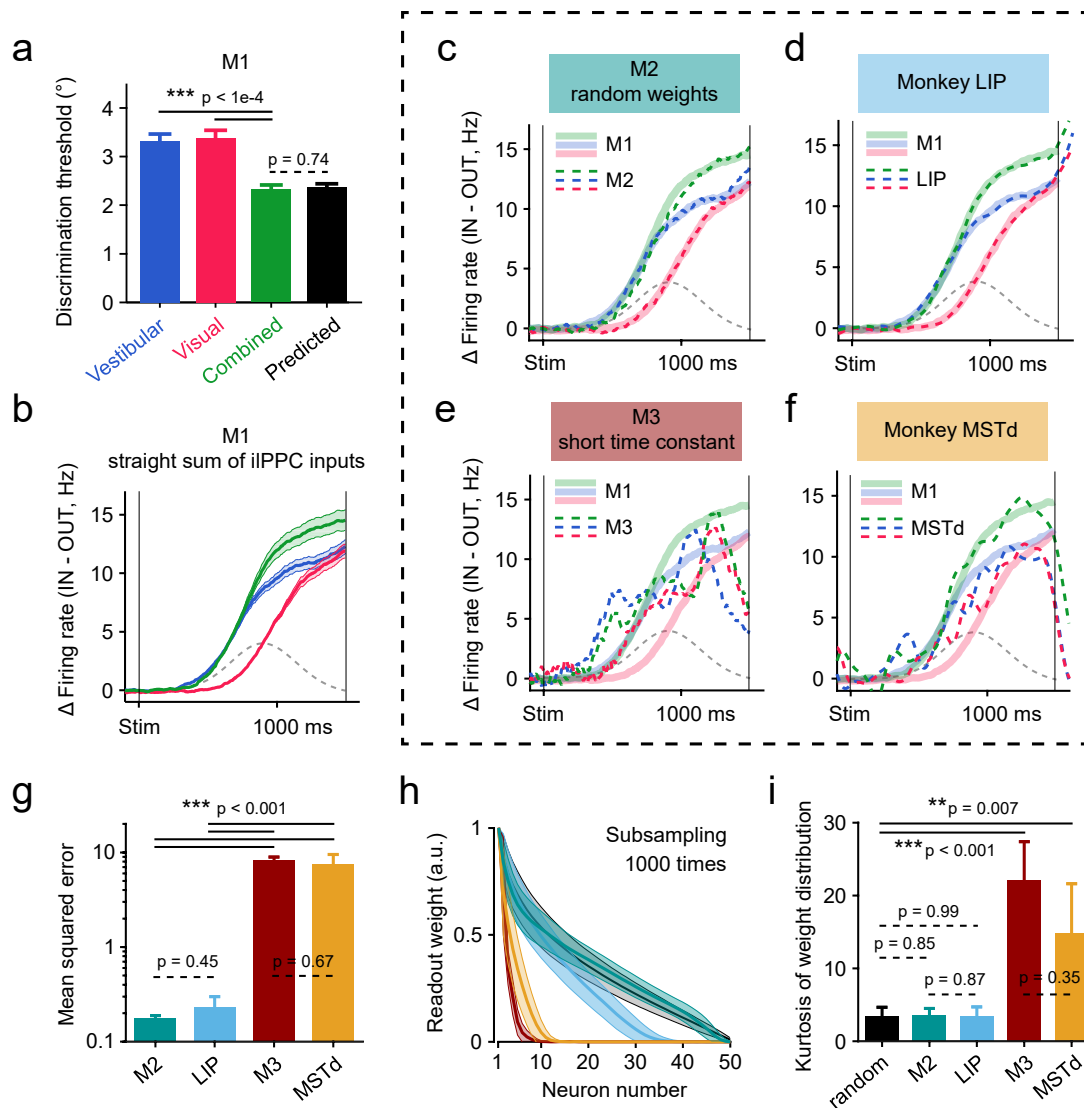
**Figure 5  Optimal ilPPC model M1 can be linearly approximated by M2 and LIP but not by M3 and MSTd.**

**(a)** Model M1 exhibited near-optimal behavior as the monkey. The psychophysical threshold under the combined condition (green) was indistinguishable from the Bayesian optimal prediction (black). **(b)** Ramping activity of M1 computed as the difference of PSTHs for IN and OUT trials. Activities from hypothetical units in the LIP layer with preferred direction close to ±90° were averaged together (see **Figure 4c** and **Methods**). Since M1 is optimal and homogeneous, we refer to M1's activities as "optimal traces" (see the main text). Notations are the same as before. **(c)** Optimal traces from M1 (thick shaded bands) can be linearly reconstructed by population activities obtained from a heterogenous model M2 (dashed curves). Model M2 had the same network architecture as M1 except that it relies on random combinations of ilPPC inputs in the integration layer (see **Methods**). **(d)** Optimal traces can also be linearly reconstructed by heterogenous single neuron activities from the LIP data. The similarity between **c** and **d** suggests that both model M2 and monkey LIP are heterogeneous variations of to the optimal ilPPC model M1. **(e and f)** In contrast, the optimal traces cannot be reconstructed from activities of a suboptimal model M3 (**e**) or from the MSTd data (**f**), presumably because the time constants in M3 and MSTd were too short. **(g)** Mean squared error of the fits in panels **c–f**. Error bars and p values were from subsampling test (n = 50 neurons, 1000 times). **(h)** Normalized readout weights ordered by magnitude. Shaded error bands indicate standard deviations of the subsampling distributions. **(i)** The kurtosis of the distributions of weights. The black curve in (**h**) and black bar in (**e**) were from random readout weights (see **Methods**).

260

261    (**[10]Ma, et al., 2006**). Accordingly, we simulated a second model, refer to as M2, in which the visual

262    and vestibular weights of each neuron were drawn from lognormal distributions (**Figure 5c** and see

263    **Methods**). Like M1, model M2 can be tuned to reproduce the Bayes optimal discrimination

264    thresholds (**Supplementary Figure 7a, b**). However, in contrast to model M1, the neurons showed

265    a wide range of response profiles closer to what we observed in vivo (**Supplementary Figure 7c**).

266    In particular, we found that the distribution of visual and vestibular weights was similar in the

267    model and in LIP data (**Supplementary Figure 7d**).

268

269    Since model M2 is a linear combination away from model M1, we tested whether the response of

270    M1 neurons could be estimated by linearly combining the response of M2 neurons. Multivariate

271    linear regression confirmed that M1 response profiles could indeed be perfectly reproduced by

272    linearly combining M2 responses (**Figure 5c**). Since LIP neurons also appear to be computing

273    random linear combinations of visual and vestibular inputs, the same result should hold for LIP

274    responses. This is indeed what we found: the response of M1 neurons can be closely approximated

275    by linearly combining the response of LIP neurons (**Figure 5d, g and Supplementary Figure 9**).

276

277    This last result is key: it suggests that LIP neurons behave quite similarly to the neurons in M2. The

278    two sets of neurons, however, differ quite significantly in how they integrate their inputs over time.

279    LIP neurons display a wide variety of temporal profiles (see **Supplementary Figure 2**), suggesting

280    that very few neurons act like perfect temporal integrators, in contrast to M2 neurons. Nonetheless,

281    the fact that linear combinations of LIP neurons could reproduce the response of M1 neurons

282    indicates that LIP responses provide a basis set sufficiently varied to allow perfect integration at

283    the population level, a result consistent with what has been recently reported in the posterior parietal

284    cortex of rats engaged in a perceptual decision making task (**[30]Scott, et al., 2017**).

285

286    In addition to this second model, we simulated a third model (M3) in which the time constant of

287     the integration layer was reduced to 100 ms. Interestingly, we found that it was not possible to

288     linearly combine the responses of M3 output neurons to reproduce the traces of the optimal model

289     M1 (**Figure 5e, g**), thus emphasizing the importance of long integration time constant for fitting

290     the optimal model. We also wondered whether M1 could be fitted by the response of MSTd neurons,

291     which are known to combine visual and vestibular responses and whose time constant are believed

292     to be of the same order as model M3. We found that the fit to M1 from MSTd neurons was markedly

293     worse than those obtained from M2 and LIP but was close to that from M3 (**Figure 5f, g**). Moreover,

294     only a small fraction of cells contributed significantly to this fit, in sharp contrast to what we

295     observed in M2 and LIP (**Figure 5h, i**). In fact, the late phase of M1 responses was captured mostly

296     by MSTd cells with short time constants who seemed sensitive to deceleration, rather than

297     integrating cells (**Supplementary Figure 8**).

298

299     Finally, we computed the shuffled Fisher information over time for the models and the experimental

300     data (**Figure 6**). The Fisher information in a neuronal population is a measure inversely

301     proportional to the square of the discrimination threshold of an ideal observer ([31]**Beck**, *et al.*, **2011;**

302     [32]**Seung and Sompolinsky, 1993**). The shuffled Fisher information is a related measure

303     corresponding to the information in a data set in which neurons are recorded one at a time, as

304     opposed to simultaneously, which is the case for our data set ([33]**Series**, *et al.*, **2004)** (see **Methods**).

305     Our network simulations revealed that the shuffled Fisher information should increase over time in

306     all conditions, reflecting the temporal accumulation of evidence (**Figure 6a**). In addition, we

307     observed that this rise in information starts earlier in the vestibular condition than in the visual one

308     because of the temporal offset between acceleration and velocity. In the combined condition, the

309     Fisher information follows at first the vestibular condition before exceeding the vestibular trace

310     once the visual information becomes available. Remarkably, the shuffled Fisher information

311     estimated from the LIP responses follows qualitatively the same trend as the ones observed in the

312     model (**Figure 6b**). In contrast to M2 and LIP neurons, shuffled Fisher information in M3 and
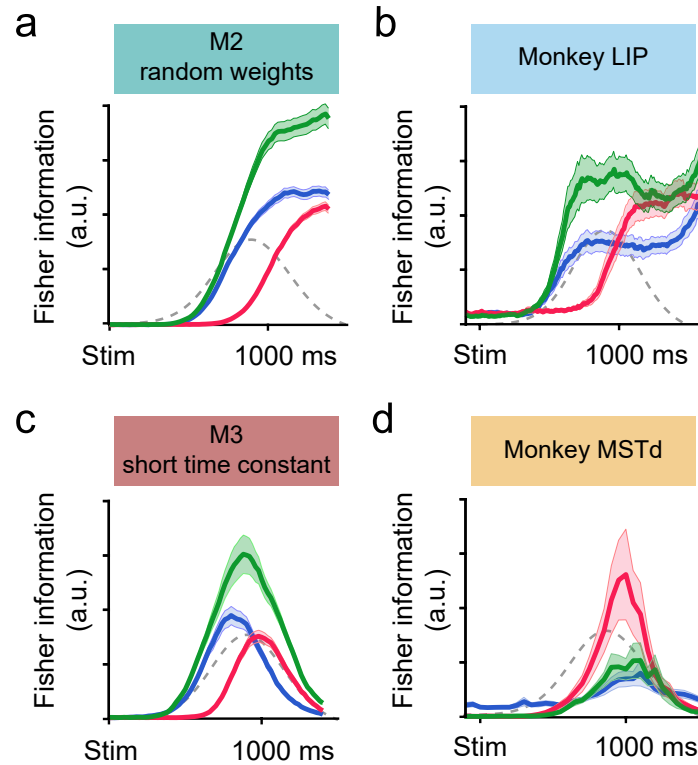
**Figure 6  Shuffled Fisher information for the model and the experimental data.**

**(a)** Shuffled Fisher information of M2 calculated by $I_{shuffled} = \sum_i f_i'^2/\sigma_i^2$, where $f_i'$ denotes the derivative of the local tuning curve of the $i$th neuron and $\sigma_i^2$ denotes the averaged variance of its responses around 0° (see **Methods**). Both correct and wrong trials were included. Shaded error bands, s.e.m. estimated from bootstrap. Note that the absolute value of shuffled Fisher information is arbitrary. **(b-d)** Same as in **a** but for the monkey LIP data, the M3 responses, and the monkey MSTd data, respectively. Note that LIP is similar to M2, and MSTd to M3.

313

314     MSTd followed the profile expected for neurons with short time constant: it simply reflected the

315     velocity profile of the stimulus and did not exhibit the plateau expected from a decision area

316     (**Figure 6c, d**).

317

318     Taken together, our results are consistent with the notion that MSTd neurons provide the visual

319     momentary evidence for decision making, while LIP circuits, or circuits upstream from LIP,

320     implement the near optimal solution of model M1, in the sense that the LIP population activity is a

321     mere linear transformation away from that solution.

322 **Discussion**

323    Integrating ever-changing sensory inputs from different sources across time is crucial for animals

324    to optimize their decisions in a complex environment, yet little is known about the underlying

325    mechanisms, either experimentally or theoretically. In the current study, we present, to the best of

326    our knowledge, the first electrophysiological data on multisensory decision making from non-

327    human primates. We found that LIP neurons in the macaque posterior parietal cortex encode

328    ramping decision signals not only for the visual condition, as widely shown in the literature, but

329    also for the vestibular and combined conditions, except with distinct temporal dynamics.

330    Importantly, these data are compatible with an ilPPC framework where optimal multisensory

331    evidence accumulation is achieved by simply summing sensory inputs across both modalities and

332    time, even with mismatched temporal profiles of cue reliabilities and with heterogeneous sensory-

333    motor representation. Therefore, our results provide the first neural correlate of optimal

334    multisensory decision making.

335

336    **Distinct visual and vestibular temporal dynamics in LIP**

337    By comparing the temporal dynamics of LIP population under different modalities, we found that

338    LIP neurons accumulate vestibular acceleration and visual speed, which serve as momentary

339    evidence for their respective modalities. These findings may seem confusing at first glance, since

340    it is more intuitive to assume that neural circuits would combine evidence with the same temporal

341    dynamics across cues, namely, either visual and vestibular speed or visual and vestibular

342    acceleration (**[19]Gu, *et al.*, 2006; [34]Chen, *et al.*, 2011a; [35]Fetsch, *et al.*, 2010; [36]Smith, *et al.*, 2017**).

343    In support of this idea, recent studies have found a remarkable transformation from acceleration-

344    dominated to speed-dominated vestibular signal along the vestibular pathway, i.e. from peripheral

345    otolith organs to the central nervous system (**[37]Laurens, *et al.*, 2017**), as well as a moderate but

346    noticeable further transformation along several sensory cortices (**[19]Gu, *et al.*, 2006; [34]Chen, *et al.*,**

347  **2011a;** [35]**Fetsch,** *et al.***, 2010;** [37]**Laurens,** *et al.***, 2017)**. Given that visual motion responses are

348  typically dominated by speed **(**[19]**Gu,** *et al.***, 2006;** [38]**Lisberger and Movshon, 1999)**, one would

349  think that the brain may deliberately turn the vestibular signal from acceleration- to speed-sensitive

350  to facilitate the combination with the visual signal.

351

352  However, if the vestibular momentary evidence is proportional to acceleration corrupted by white

353  noise across time, integrating this evidence to obtain a velocity signal would not simplify decision

354  making. On the contrary, this step would introduce temporal correlations **(**[39]**Churchland,** *et al.***,**

355  **2011)**, in which case, even with ilPPC, a simple sum of the momentary evidence would no longer

356  be optimal **(**[6]**Bogacz,** *et al.***, 2006)**. Instead, downstream circuits would have to compute a weighted

357  sum of the sensory evidence, which would effectively differentiate the momentary evidence before

358  summing them. In other words, optimal integration would effectively recover the original

359  acceleration signals. Our results, along with previous psychophysical results **(**[8]**Drugowitsch,** *et al.***,**

360  **2014)**, strongly suggest that the brain does not go through this extra step and uses the acceleration

361  signals as momentary evidence instead.

362

363  ## Multisensory convergence in the brain for heading decision

364  One of the long-standing questions about multisensory integration is whether integration takes

365  place early or late along the sensory streams **(**[40]**Bizley,** *et al.***, 2016)**. There are clear signs of

366  multisensory responses in relatively early- or mid- stage of sensory areas, thus supporting the early

367  theory **(**[41]**Gu, 2018)**. Our results are more consistent with the late-convergence theory in which

368  multisensory momentary evidence are combined across modalities and time in decision areas such

369  as LIP. However, this dichotomy between early and late theories does not necessarily make sense

370  given the recurrent nature of the cortical circuitry. In a highly recurrent network, it is notoriously

371  difficult to identify a node as a primary site of integration. Thus, integration might take place

372   simultaneously across multiple sites but in such a way that the output of the computation is

373   consistent across sites. For example, **Deneve, *et al.* [42]** demonstrated how this could take place in a

374   large recurrent network performing optimal multisensory integration, though their work did not

375   consider the problem of temporal integration.

376

377   It might be possible to gain further insight into the distributed nature of multisensory decision

378   making by combining the previous models with the one we have presented here. Such an extended

379   model might explain why vestibular momentary evidence is tuned to velocity by the time they

380   appear in MSTd (**[37]Laurens, *et al.*, 2017; [41]Gu, 2018)**, and why this velocity tuned vestibular input

381   does not appear to be integrated in LIP. It could also shed light on recent physiological experiments

382   in which electrical microstimulation and chemical inactivation in MSTd could dramatically affect

383   heading discrimination based on optic flow while this effect was largely negligible in the vestibular

384   condition (**[43]Gu, *et al.*, 2012)**. By contrast, and in accord with our finding that LIP integrates

385   vestibular acceleration, inactivating the vestibular cortex PIVC, where vestibular momentary

386   evidence is dominated by acceleration (**[34]Chen, *et al.*, 2011a; [37]Laurens, *et al.*, 2017)**, substantially

387   diminished the macaque's heading ability based on vestibular cue (**[44]Chen, *et al.*, 2016)**. Note,

388   however, a detailed construction of such a model lies beyond the scope of the present study but will

389   eventually be required for a multi-area theory of multisensory decision making.

390

## Computational models for multisensory decision making

392   Our results indicate that, at the population level, LIP implements an optimal solution for

393   multisensory decision making under the assumption that the sensory inputs are encoded with ilPPC.

394   This assumption is not perfectly satisfied in our experiment since the visual and vestibular inputs

395   deviate from pure ilPPCs, but we saw that this deviation introduces only a minor information loss.

396   While these results provide the first experimental support for the ilPPC theory of multisensory

397    decision making, it will be important to test in future experiments other predictions of this

398    framework. In particular, the ilPPC theory predicts that LIP activity encodes a full probability

399    distribution over choices given the evidence so far (**[9]Beck***, et al.***, 2008)**. Testing this prediction

400    thoroughly requires simultaneous recording of LIP ensemble, manipulating the cue reliability

401    (motion profile or visual coherence) on a trial-by-trial basis, and preferably engaging the animals

402    in a reaction-time task, all of which should be addressed in future studies.

403

404    There are of course other models of decision making which could potentially account for the

405    responses we have observed in LIP (**[45]Chandrasekaran, 2017)**. In particular, it has been argued

406    that LIP is part of a network of areas implementing point attractor networks (**[46]Wong and Wang,**

407    **2006; [47]Wang, 2002)**. However, it is not immediately clear how this approach can be generalized

408    to the type of decision we have considered here. Indeed, as we have seen, the optimal solution

409    depends critically on the code that is used to encode the momentary evidence. To the extent that

410    this code is close to an ilPPC, the optimal solution is to sum the inputs spikes, in which case one

411    needs a line attractor network, which is effectively what our network approximates. Therefore, as

412    long as these previous models of decision making are fine-tuned to approximate line attractor

413    networks, and as long as they are fed ilPPCs as inputs, the two classes of models would be

414    equivalent.

415

416    Training recurrent neural network (RNNs) on our task (**[48]Mante***, et al.***, 2013; [49]Song***, et al.***, 2017)**

417    provides a third alternative for modeling multisensory decision making. We also tried this approach

418    and found that the resulting network was capable of reproducing the behavioral thresholds of the

419    animal while exhibiting a wide variety of single neuron responses similar to what we saw in LIP (

420    **Supplementary Figure 10**). Nonetheless, this approach has one major drawback: it makes it very

421    difficult to understand how the network solves the task. We could try to reverse engineer the

422    network, but given that an analytical solution can be derived from first principles for our task, and

423    given that this solution is close to what we observed in LIP, it is unclear what insight could be

424    gained from the recurrent network. In contrast, our ilPPC model provides a close approximation to

425    the optimal solution, consistent with the experimental results, along with a clear understanding as

426    to why this approach is optimal.

427

428

## Methods

### Subjects and Apparatus

All animal procedures were approved by the Animal Care Committee of Shanghai Institutes for Biological Sciences, Chinese Academy of Sciences and have been described previously in detail ([12]Gu, *et al.*, 2008; [19]Gu, *et al.*, 2006). Briefly, two male adult rhesus monkeys, Monkey P and Monkey M, weighing ~8 kg, were chronically implanted with a lightweight plastic ring for head restraint and a scleral coil for monitoring eye movements (Riverbend Instruments). During experiments, the monkey sat comfortably in a primate chair mounted on top of a custom-built virtual reality system, which consisted of a motion platform (MOOG MB-E-6DOF/12/1000KG) and an LCD screen (~30 cm of view distance and ~90° × 90° of visual angle; HP LD4201), presenting vestibular and visual motion stimuli to the monkey, respectively. The stimuli were controlled by customized C++ software and synchronized with the electrophysiological recording system by TEMPO (Reflective Computing, U.S.A).

To tune the synchronization between vestibular and visual stimuli, we rendered a virtual world-fixed crosshair on the screen while projected a second crosshair at the same place on the screen using a real world-fixed laser pen. When the platform was moving, we carefully adjusted a delay parameter in the C++ software (with 1 ms resolution) until the two crosshairs aligned precisely together all the time, as verified by a high-speed camera (Meizu Pro 5) and/or a pair of back-to-back mounted photodiodes. This synchronization procedure was repeated occasionally over the whole period of data collection.

### Behavioral Tasks

#### *Memory-guided Saccade Task*

We used the standard memory-guided saccade task ([50]Barash, *et al.*, 1991) to characterize and

454    select LIP cells for recording in the main decision-making experiments. The monkey fixated at a

455    central fixation point for 100 ms and then a saccade target flashed briefly (500 ms) in the periphery.

456    The monkey was required to maintain fixation during the delay period (1000 ms) until the fixation

457    point extinguished and then saccade to the remembered target location within 1000 ms for a liquid

458    reward. For all tasks in the present study, at any time when there existed a fixation point, trials were

459    aborted immediately if the monkey's gaze deviated from a $2° \times 2°$ electronic window around the

460    fixation point.

461

462    ***Multisensory Heading Discrimination Task***

463    In the main experiments, we trained the monkeys to report their direction of self-motion in a two-

464    alternative forced-choice heading discrimination task **([12]Gu, *et al.*, 2008)** (**Figure 1**). The monkey

465    initiated a trial by fixating on a central, head-fixed fixation point, and two choice targets then

466    appeared. The locations of the two targets were determined case-by-case for each recording session

467    (see below). After fixating for a short delay (100 ms), the monkey then began to experience a fixed-

468    duration (1.5 s) forward motion in the horizontal plane with a small leftward or rightward

469    component relative to straight ahead. The animals were required to maintain fixation during the

470    presentation of the motion stimuli. At the end of the trial, the motion ended, and the monkey was

471    required to maintain fixation for another 300–600 ms random delay (uniformly distributed) until

472    the fixation point disappeared, at which point the monkey was allowed to make a saccade choice

473    toward one of the two targets to report his perceived heading direction (left or right).

474

475    Across trials, nine heading angles ($\pm8°$, $\pm4°$, $\pm2°$, $\pm1°$, and $0°$) and three cue conditions (vestibular,

476    visual, and combined) were jointly interleaved, resulting in 27 unique stimulus conditions, each of

477    which was repeated $15 \pm 3$ (median $\pm$ m.a.d.) times per one session. In a vestibular or a visual trial,

478    heading information was solely provided by inertial motion (real movement of the motion platform)

479    or optic flow (simulated movement through a star field on the display), respectively, whereas in a

480    combined trial, congruent vestibular and visual cues were provided synchronously. To maximize

481    the behavioral benefit of cue integration, we balanced the monkey's performance under the

482    vestibular and the visual conditions by manipulating the motion coherence of the optic flow (the

483    percentage of dots that moved coherently). The visual coherence was 12% and 8% for monkey P

484    and M, respectively.

485

486    To ensure that the reliabilities of sensory cues varied throughout each trial, we used Gaussian-shape,

487    rather than constant, velocity profiles for all motion stimuli. In the main experiments, the Gaussian

488    profile had a displacement $d = 0.2\ m$ and a standard deviation $\sigma = 210\ ms$ (half duration at about

489    60% of the peak velocity), resulting in a peak velocity $v_{max} = 0.37\ m/s$ and a peak acceleration

490    $a_{max} = 1.1\ m/s^2$. In the experiment where we sought to independently vary the peak times of

491    velocity and acceleration (**Figure 3**), two additional sets of motion parameters were used. For the

492    narrow-speed profile, $d = 0.10\ m$, $\sigma = 150\ ms$, $v_{max} = 0.37\ m/s$, and $a_{max} = 1.1\ m/s^2$; for

493    the wide-speed profile, $d = 0.25\ m$, $\sigma = 330\ ms$, $v_{max} = 0.31\ m/s$, and $a_{max} = 0.6\ m/s^2$.

494

495    ## Electrophysiology

496    We carried out extracellular single-unit recordings as described previously (**[12]Gu, *et al.*, 2008**) from

497    four hemispheres in two monkeys. For each hemisphere, reliable area mapping was first achieved

498    through cross-validation between structural MRI data and electrophysiological properties,

499    including transition patterns of gray/white matter along each penetration, sizes of visual

500    receptive/response field, strengths of spatial tuning to visual and vestibular heading stimuli, and

501    activities in the memory-guided saccade task. Based on the mapping results, Area LIP was

502    registered by its spatial relationships with other adjacent areas (VIP, Area 5, MSTd, etc.), its weak

503    sensory encoding of heading information, and its overall strong saccade-related activity

504    (**Supplementary Figure 1**). Our recording sites located in the ventral division of LIP, extending

505     from 7–13 mm lateral to the midline and -5 mm (posterior) to +3 mm (anterior) relative to the

506     interaural plane.

507

508     Once we encountered a well-isolated single unit in LIP, we first explored its response field (RF) by

509     hand (using a flashing patch) and then examined its electrophysiological properties using the

510     memory-guided saccade task. The saccade target in each trial was randomly positioned at one of

511     the 8 locations 45° apart on a circle centered on the fixation point (5°–25° radius, optimized

512     according to the cell's RF location). We calculated online the memory-saccade spatial tuning for

513     three response epochs: (1) visual response period, 75–400 ms from target onset; (2) delay period,

514     25–900 ms from target offset; and (3) presaccadic period, 200–50 ms before the saccade onset

515     (**Supplementary Figure 2**). The cell's spatio-temporal tunings were used to refine its RF location

516     (via vector sum) and to determine its inclusion in the subsequent decision-making task. Since the

517     decision-related activity of LIP neurons cannot be strongly predicted by the persistent activity

518     during the delay period alone (**[26]Meister, *et al.*, 2013)** (**Supplementary Figure 4b**), we adopted a

519     wider cell selection criterion than conventionally used, in which we included cells that have

520     significant spatial selectivity for *any* of the three response epochs (**[26]Meister, *et al.*, 2013)** (one-

521     way ANOVA, $p < 0.05$, 3–5 repetitions). If the cell met this criterion, then we recorded its decision-

522     related activity while engaging the monkey in the main multisensory decision-making task, with

523     the two choice targets being positioned in its RF and 180° opposite to its RF, respectively.

524

525     Although we collected data from a relatively broad sample of LIP neurons, we nevertheless had

526     two sampling biases during this process. First, we were biased toward cells with strong persistent

527     activity so that our multisensory data could be better compared with previous unisensory data in

528     the decision-making literature, where in most cases only these cells were recorded. Second, we

529     were biased toward cells with RF close to the horizontal line through the fixation point. Unlike the

530     classical random dot stimuli whose motion direction on the fronto-parallel plane can be aligned

531    with the cell's RF (and the choice targets) session by session, our self-motion stimuli were always

532    on the horizontal plane and thus were not adjustable according to the cell's RF on the fronto-parallel

533    plane. As a result, the subjects had to make an additional mapping from their perceived heading

534    directions (always left or right) to the choice targets (often inclined, and in extreme cases, up or

535    down). Therefore, to make the task more intuitive to the monkeys and to minimize the potential

536    influence of this mapping step on neural activity, we discarded a cell if the angle between the

537    horizontal line and the line connecting the fixation point to its RF exceeded 60°, although we

538    observed little change in monkeys' behavior even when this angle approached 80°.

539

540    ## Data Analysis

541    *Psychophysics*

542    To quantify the behavioral performance for both the monkeys and the model in the multisensory

543    decision-making task, we constructed psychometric curves by plotting the proportion of "rightward"

544    choices as a function of heading (**Figure 1c**) and fitted them with cumulative Gaussian functions

545    (**[12]Gu, et al., 2008**). The psychophysical threshold for each cue condition was defined as the

546    standard deviation of their respective Gaussian fit. The Bayesian optimal prediction of

547    psychophysical threshold under the combined condition $\sigma_{prediction}$ was solved from the inverse

548    variance rule (**[24]Knill and Richards, 1996**)

549    $$\frac{1}{\sigma^2_{prediction}} = \frac{1}{\sigma^2_{vestibular}} + \frac{1}{\sigma^2_{visual}}$$

550    where $\sigma_{vestibular}$ and $\sigma_{visual}$ represent psychophysical thresholds under the vestibular and visual

551    conditions, respectively.

552

553    *Choice-related neural activities*

554    We constructed peri-stimulus time histograms (PSTHs) for two epochs of interest in a trial, the

555    decision formation epoch and the saccade epoch, by aligning raw spike trains to the stimulus onset

556    and the saccade onset, respectively. Firing rates were computed in non-overlapping 10-ms bins and

557    smoothed over time by convolving with a Gaussian kernel ($\sigma = 50\ ms$). Unless otherwise noted,

558    only correct trials were used in the following analyses, except for the ambiguous 0° heading where

559    we included all complete trials.

560

561    To illustrate the choice-related activity of a cell, we grouped the trials according to the monkey's

562    choice, i.e., trials ending up with a saccade toward the cell's RF (IN trials) versus trials ending up

563    with a saccade away from the cell's RF (OUT trials), and computed the averaged PSTHs of these

564    two groups of trials for each cue condition (**Figure 2a**). When averaged across cells, each cell's

565    PSTHs were normalized such that the cell's overall firing rate had a dynamic range of [0, 1] (**Figure**

566    **3**). To quantify the strength of choice signals and better visualize ramping activities, we calculated

567    choice divergence (**[23]Raposo, et al., 2014**) for each 10-ms time bin and for each cue condition using

568    receiver operating curve (ROC) analysis (**Figure 2b**). Choice divergence ranged from -1 to 1 and

569    was defined as $2 \times (\mathrm{AUC} - 0.5)$, where AUC represents the area under the ROC curve derived

570    from PSTHs of IN and OUT trials. To capture the onset of choice signals, we computed a

571    divergence time defined as the time of the first occurrence of a 250-ms window (25 successive 10-

572    ms bins) in which choice divergence was consistently and significantly larger than 0 (**Figure 3c,**

573    **f**). The statistical significance of choice divergence (p < 0.05, relative to the chance level of 0) was

574    assessed by two-sided permutation test (1000 permutations). We also calculated a grand choice

575    divergence which ignored temporal information and used all the spikes in the decision formation

576    epoch (0–1500 ms from the stimulus onset). The same permutation test was performed on the grand

577    choice divergence to determine whether a cell had overall significant choice signals for a certain

578    cue condition (for example, in **Figure 2c**).

579

580     ***Linear Fitting of Mean Firing Rates***

581     We fitted a linear weighted summation model to predict neural responses under the combined

582     condition with those under the single cue conditions, using (**[12]Gu, *et al.*, 2008**)

583 $$r_{combined} = w_{vestibular} r_{vestibular} + w_{visual} r_{visual} + C$$

584     where $C$ is a constant, and $r_{combined}$, $r_{vestibular}$, and $r_{visaul}$ are mean firing rates across a trial (0–

585     1500 ms from stimulus onset) for the three cue conditions, respectively. The weights for single cue

586     conditions, $w_{vestibular}$ and $w_{visual}$, were determined by the least-squares method and plotted

587     against each other to evaluate the heterogeneity of choice signals in the population for both LIP

588     data and the model (**Supplementary Figure 7d**).

589

590     ***Fisher Information Analysis***

591     To compute Fisher information (**[32]Seung and Sompolinsky, 1993**), the full covariance matrix of

592     the population responses is needed, but this requires simultaneously recording from hundreds of

593     neurons, which is not accessible to us yet. Instead, we calculated the shuffled Fisher information,

594     which corresponds to the information in a population of neurons in which correlations have been

595     removed (typically via shuffling across trials, hence the name). Shuffled Fisher information is given

596     by (**[33]Series, *et al.*, 2004; [51]Gu, *et al.*, 2010**):

597 $$I_{shuffled} = \sum_{i=1}^{N} \frac{f_i'^2}{\sigma_i^2} \tag{1}$$

598     where $N$ is the number of neurons in the population; for the $i$th neuron, $f_i'$ denotes the derivative

599     of its local tuning curve, and $\sigma_i^2$ denotes the averaged variance of its responses around 0° heading.

600     The tuning curve $f_i$ was constructed from both correct and wrong trials grouped by heading angles,

601     using spike counts in 250-ms sliding windows (advancing in 10-ms steps), and its derivative $f_i'$

602     was obtained from the slope of a linear fit of $f_i$ against headings. The variance $\sigma_i^2$ was computed

603     for each heading angle and then averaged. To estimate the standard errors of $I_{shuffled}$, we used a

604    bootstrap procedure in which random samples of neurons were drawn from the population by

605    resampling with replacement (1000 iterations). To compare the experimental data with the model,

606    we repeated all the above steps on artificial LIP neurons in the model M2 and M3 (see below), with

607    the inter-neuronal noise correlation being ignored as well (**Figure 6**).

608

609    Two caveats are noteworthy when interpreting the Fisher information results. First, since the slope

610    of tuning curve $f'$ is squared in the right-hand side of equation (1), the Fisher information will

611    always be non-negative regardless of the sign of $f'$. As a result, even when the motion speed was

612    zero at the beginning of a trial, the population Fisher information already had a positive value

613    because of the noisy tuning curves during that period. Second, since we ignored inter-neuronal

614    noise correlations, $I_{shuffled}$ is most likely very different from the true Fisher information and thus

615    its value is arbitrary (**[33]Series, et al., 2004**). Nonetheless, if we assume the noise correlation

616    structure of LIP population is similar across cue conditions, we can still rely on the qualitative

617    temporal evolution of $I_{shuffled}$ to appreciate how multisensory signals are accumulated across time

618    and cues in LIP.

619

## Network Simulation of ilPPC Framework

621    ***The responses of visual and vestibular neurons closely approximate ilPPC***

622    As mentioned previously, an important assumption of ilPPC is that the amplitude of the sensory

623    tuning curves be proportional to the nuisance parameters (in our case visual speed and vestibular

624    acceleration) (**[9]Beck, et al., 2008**). To check whether this is the case for the visual neurons, we

625    analyzed the spatio-temporal tuning curves of neurons in area MSTd (data from (**[19]Gu, et al., 2006**)).

626    We noticed that, for some neurons, the average tuning curves are not fully consistent with the ilPPC

627    assumption (**Supplementary Figure 6a**). Briefly, the mean firing rate of an MSTd neuron at time

628    $t$ in response to a visual stimulus with heading $\theta$ can be well captured by

629
$$f(\theta, t) = v(t)(A \exp[K(\cos(\theta - \theta_i) - 1)] - C) + B \qquad (2)$$

630 where $\theta_i$ denotes the preferred heading of the neuron $i$ and $v(t)$ is the velocity profile; $A$, $K$, $C$,

631 and $B$ correspond to the amplitude, the width, the null inhibition, and the baseline of its tuning

632 curve, respectively. The ilPPC framework requires the $v(t)$ term to be separable, namely, $f(\theta, t) =$

633 $h(\theta)g(v(t))$, where $h(\theta)$ is a pure spatial component and $g(v(t))$ is a multiplicative gain function

634 ([9]**Beck, *et al.*, 2008**; [10]**Ma, *et al.*, 2006**). In equation (2), this requirement is equivalent to $C = 0$

635 and $B = 0$, however, we found that some MSTd neurons often had non-zero baselines ($C > 0$ and

636 $B > 0$). This will be harmful to the optimality of the ilPPC framework because, for example, when

637 $v(t) = 0$ (and thus the sensory reliability is zero), MSTd neurons still tend to generate background

638 spikes, which will bring nothing but noise into the simply summed population activity of

639 downstream areas in an ilPPC network.

640

641 To estimate the information loss due to this deviation, we simulated a population of MSTd neurons

642 with heterogeneous spatio-temporal tuning curves similar to what has been found experimentally

643 ([19]**Gu, *et al.*, 2006**). We calculated the information that can be decoded from the population by a

644 series of optimal decoders $I_{optimal}$ and that can be recovered by the ilPPC solution $I_{ilPPC}$. We

645 assumed that the population responses in MSTd contains differential correlations ([28]**Moreno-Bote,**

646 ***et al.*, 2014)** such that the discrimination threshold of an ideal observer of MSTd activity was of

647 the same order as the animal's performance. Under such conditions, we found that the information

648 loss $(I_{optimal} - I_{ilPPC})/I_{optimal}$ was around 5%. Detailed calculations of information loss are

649 provided in the **Supplementary Materials**. Therefore, the population response of MSTd neurons

650 provide a close approximation to an ilPPC, in the sense that simply summing the activity of MSTd

651 neurons over time preserve 95% of the information conveyed by these neurons.

652

653 We also checked whether the ilPPC assumption holds in the case of vestibular neurons. Equation

654 (2) above still provides a good approximation to vestibular tuning curves, except that $C$ is close to

655 zero for most neurons ([37]**Laurens, *et al.*, 2017**), in which case the information less is even less

656 pronounced.

657

658 ***Network Model Implementing the ilPPC solution (Model M1)***

659 We extended a previous ilPPC network model for unisensory decision making ([9]**Beck, *et al.*, 2008**)

660 to our multisensory decision-making task. Two sensory layers, the vestibular layer and the visual

661 layer, contained 100 linear-nonlinear-Poisson (LNP) neurons with bell-shape tuning curves to the

662 heading direction (equation (2)). For the $i$th neuron in the vestibular or visual layer, the probability

663 of firing a spike at time step $[t_n - \delta t, t_n]$ was given by

$$p(r_i^\bullet(t_n) = 1) = [\delta t(g_\bullet(t)(A\exp[K(\cos(\theta - \theta_i) - 1)] - C) + B) + n_i]^+ \qquad (3)$$

665 where $\bullet \in \{\text{VEST}, \text{VIS}\}$, $A, K, C, B, \theta$, and $\theta_i$ have the same meanings as in equation (2), $n_i$ is a

666 correlated noise term, and $[\cdot]^+$ is the threshold-linear operator: $[x]^+ = \max(x, 0)$. The spatial

667 tuning was gain-modulated by a time-dependent function $g_\bullet(t)$, which modeled the reliability of

668 the sensory evidence and took the form

$$g_{\text{VEST}}(t) = c_{\text{VEST}}|\hat{a}(t)|, \quad g_{\text{VIS}}(t) = c_{\text{VIS}}\hat{v}(t)$$

670 in which $\hat{a}(t)$ and $\hat{v}(t)$ are the same acceleration and velocity profiles as the experiments but with

671 the maximum values normalized to 1, respectively, whereas $c_{\text{VEST}}$ and $c_{\text{VIS}}$ are scaling parameters

672 used to control the signal-to-noise ratio of sensory inputs and to balance the behavior performance

673 between the two cue conditions like in the experiments. The noise $n_i$ in equation (3) was generated

674 by convolving independent Gaussian noise with a circular Gaussian kernel,

$$n_i = \sum_j A_\eta \exp\left(K_\eta(\cos(\theta_i - \theta_j) - 1)\right)\eta_j$$

676 where $\eta_j \sim i.i.d. N(0,1)$, and $A_\eta$ and $K_\eta$ were set to $10^{-5}$ and 2, respectively. Other parameters

677 we used were: $A = 60$ Hz, $K = 1.5, C = 10$ Hz, $B = 20$ Hz, $c_{\text{VEST}} = c_{\text{VIS}} = 2.4, \delta t = 1$ ms.

678    Note that in equation (3), the gain $g_\bullet(t)$ cannot be factored out because $B > 0$, which is the same

679    case as in MSTd (equation (2)). Accordingly, the neural code of M1's sensory layers is not exact

680    ilPPC ([9]**Beck, *et al.*, 2008**). However, it is still a close approximation to ilPPC, since we have shown

681    in the previous section that MSTd is 95% ilPPC-compatible.

682

683    The two sensory layers then projected to 100 LNP neurons in the integrator layer. We distinguished

684    the integrator layer from the LIP layer because there are reasons to believe that LIP reflects the

685    integration of the evidence but may not implement the integration *per se* ([27]**Katz, *et al.*, 2016**). The

686    integrator layer summed the sensory responses across both cues and time,

687
$$m_i^{\text{INT}}(t_{n+1}) = m_i^{\text{INT}}(t_n) + g_{stim}(t_n)\left(\sum_j W_{ij}^{\text{INTVEST}} r_j^{\text{VEST}}(t_n) + \sum_j W_{ij}^{\text{INTVIS}} r_j^{\text{VIS}}(t_n)\right) \quad (4)$$

688    where $m_i^{\text{INT}}$ denotes the membrane potential proxy of neuron $i$, $W_{ij}^{\text{INTVEST}}$ and $W_{ij}^{\text{INTVIS}}$ are

689    matrices for the feedforward weights from the vestibular and visual layer to the integrator layer,

690    respectively, and $g_{stim}(t_n)$ is an attentional gain factor (see below). Note that we ignored the issue

691    of how neural circuits perform perfect integration and just assumed that they do. We could have

692    simulated one of the known circuit solutions to this problem ([52]**Goldman, 2009**), but this would

693    not have affected our results, while making the simulation considerably more complicated.

694

695    The feedforward connections $W_{ij}^{\text{INT}\bullet}$ map the negative and positive heading directions onto the

696    two saccade targets, i.e., neurons preferring $-90°$ and $+90°$ in the integrator layer, respectively, by

697
$$W_{ij}^{\text{INT}\bullet} = a\exp\left(k\left(\cos\left(\theta_i^{\text{INT}} - \hat{\theta}\right) - 1\right)\right)\left|\sin(\theta_j^{\bullet})\right|$$

698    in which a step function $\hat{\theta}$ controls the mapping,

699
$$\hat{\theta} = \begin{cases} -\pi/2, & \text{if } \theta_j^{\bullet} \leq 0 \\ \pi/2, & \text{if } \theta_j^{\bullet} > 0 \end{cases}.$$

700    We used $a = 20$ and $k = 4$ in our simulations. After the linear step, the membrane potential proxy

701    was used to determine the probability of the $i$th integrator neuron firing a spike between times $t_n$

702    and $t_n + \delta t$,

703
$$p(r_i^{\mathrm{INT}}(t_n) = 1) = [m_i^{\mathrm{INT}}(t_n)]^+.$$

704

705    Finally, the LIP layer received excitatory inputs from the integrator layer, together with visual

706    inputs triggered by the two saccade targets (sent from the target layer). In addition, there were also

707    lateral connections in LIP to prevent saturation. In the linear step, the membrane potential proxy of

708    the $i$th LIP neuron followed

709    $$m_i^{\mathrm{LIP}}(t_{n+1}) = \left(1 - \frac{\delta t}{\tau}\right) m_i^{\mathrm{LIP}}(t_n) + \frac{1}{\tau}\left(\sum_j W_{ij}^{\mathrm{LIPINT}} r_j^{\mathrm{INT}}(t_n) + \sum_j W_{ij}^{\mathrm{LIPTARG}} r_j^{\mathrm{TARG}}(t_n) + \sum_j W_{ij}^{\mathrm{LIP}} r_j^{\mathrm{LIP}}(t_n)\right) \quad (5)$$

710    where the time constant, $\tau$, was set to 100 ms; $W_{ij}^{\mathrm{LIPINT}}$ and $W_{ij}^{\mathrm{LIPTAR}}$ are weight matrices for the

711    feedforward connections from the integrator layer and the target layer to the LIP layer, respectively,

712    and $W_{ij}^{\mathrm{LIP}}$ is the matrix for the recurrent connections within LIP. We used translation-invariant

713    weights for all these connections,

714    $$W_{ij} = a \exp\left(k(\cos(\theta_i - \theta_j) - 1)\right) + b.$$

715    For $W_{ij}^{\mathrm{LIPINT}}$, we used $a = 15, k = 10, b = -3.6$; for $W_{ij}^{\mathrm{LIPTARG}}$, we used $a = 8, k = 5, b = 0$;

716    and for $W_{ij}^{\mathrm{LIP}}$, we used $a = 5, k = 10, b = -3$. The term $r_j^{\mathrm{TARG}}(t_n)$ in equation (5) denotes the

717    visual response of the $j$th neuron in the target layer induced by the two saccade targets,

718    $$p(r_j^{\mathrm{TAR}}(t_n) = 1) = s_{targ}(t_n) \sum_{m=1}^{2} p_{targ} \exp\left(k_{targ}(\cos(\theta_j^{\mathrm{TAR}} - \theta_m) - 1)\right)$$

719    where $\theta_1 = -\pi/2$ and $\theta_2 = \pi/2$, $p_{targ} = 0.050$, and $k_{targ} = 4$. The term $s_{targ}(t_n)$ modeled the

720    saliency of the targets: $s_{targ}(t_n) = 1$ before stimulus onset and $s_{targ}(t_n) = 0.6$ afterwards.

721

722    After the linear step done in equation (5), the probability of observing a spike from the $i$th LIP

723    neuron for the next time step was given by, again,

724
$$p\big(r_j^{\mathrm{LIP}}(t_{n+1}) = 1\big) = [m_i^{\mathrm{LIP}}(t_{n+1})]^+. \tag{6}$$

725

726 ***Decision Bound and Action Selection***

727 To let the model make decisions, we endowed it with a stopping bound such that the evidence

728 integration terminated when the peak activity in the LIP layer reached a threshold value. This

729 mechanism generates premature decisions in our fixed duration task, which have been observed in

730 the previous experiments ([29]**Kiani, *et al.*, 2008**) as well as ours (see the main text). Specifically,

731 once the firing rate of any neuron in the LIP layer (determined from equation (6)) exceeded $\Theta^\bullet =$

732 37 Hz for a vestibular or a visual trial and $\Theta^{\mathrm{COMB}} = 42$ Hz for a combined trial, we blocked the

733 sensory inputs to the integrator layer by setting the gain factor in equation (4) to zero:

734
$$g_{stim}(t_n) = \begin{cases} 1, & \text{if } t_n < t_\Theta \\ 0, & \text{if } t_n \geq t_\Theta \end{cases}$$

735 where $t_\Theta$ denotes the time of bound crossing. The instantaneous population activity at this time

736 point $\boldsymbol{r}^{\mathrm{LIP}}(t_\Theta)$ was then used to determine the model's choice, while the network dynamics

737 continued to evolve until the end of the 1.5-s trial.

738

739 To read out the model's choice, we trained a linear support vector machine (SVM) to classify the

740 heading direction from $\boldsymbol{r}^{\mathrm{LIP}}(t_\Theta)$. We ran the network for 100 trials, used $\boldsymbol{r}^{\mathrm{LIP}}(t_\Theta)$ in 30 trials to

741 train the SVM, and then applied the SVM on the remaining 70 trials to make decisions and generate

742 psychometric functions of the model (with bootstrap 1000 times, **Figure 5a** and **Supplementary**

743 **Figure 7a**). The SVM acts like (or even outperforms) a local optimal linear estimator (LOLE)

744 trained by gradient descent ([33]**Series, *et al.*, 2004**). Importantly, such decoders could be

745 implemented with population codes in a biologically realistic point attractor network tuned for

746 optimal action selection in a discrimination task ([53]**Deneve, *et al.*, 1999**), which could correspond

747 to downstream areas such as the motor layer of the superior colliculus ([9]**Beck, *et al.*, 2008**).

748

749     *Heterogeneous ilPPC Network (M2)*

750     In model M2, we generalized the homogeneous ilPPC network described above (model M1) to a

751     heterogeneous one. Instead of taking perfect sums like in model M1, neurons in the integration

752     layer of the model computed random linear combinations of vestibular and visual inputs. It is

753     indeed been widely shown that integration weights *in vivo* are heterogeneous and are well-captured

754     by "long-tailed" lognormal distributions (see for example (**[54]Song, et al., 2005**)). To simulate this

755     in M2, we drew each synaptic weight $w_{\mathrm{M2}}$ in M2 from a lognormal distribution

$$p(w_{\mathrm{M2}} = x) = \frac{1}{\sqrt{2\pi}\sigma x} \exp\left(-\frac{(\log x - \mu)^2}{2\sigma^2}\right) \tag{7}$$

757     where $\mu$ and $\sigma$ were chosen such that the expectation $e(w_{\mathrm{M2}})$ and the standard deviation $s(w_{\mathrm{M2}})$

758     of $w_{\mathrm{M2}}$ were both equal to its counterpart synaptic weight $w_{\mathrm{M1}}$ in model M1:

$$e(w_{\mathrm{M2}}) = s(w_{\mathrm{M2}}) = w_{\mathrm{M1}}.$$

760     The parameters $\mu$ and $\sigma$ in equation (7) were related to $e$ and $s$ through

$$\mu = \log\left(e^2/\sqrt{e^2 + s^2}\right)$$
$$\sigma = \sqrt{\log(s^2/e^2 + 1)} \quad .$$

762     If $w_{\mathrm{M1}} < 0$, a negative sign was added to the resulting $w_{\mathrm{M2}}$, since lognormal distributions are

763     always non-negative.

764

765     *Network with Short Integration Time Constant (M3)*

766     We also simulated a sub-optimal model M3 in which the network does not integrate evidence over

767     time. This was done by replacing equation (4) with

768     $m_i^{\mathrm{INT}}(t_{n+1}) = \left(1 - \frac{\delta t}{\tau}\right) m_i^{\mathrm{INT}}(t_n) + \frac{1}{\tau} g_{stim}(t_n) \left(\sum_j W_{ij}^{\mathrm{INTVEST}} r_j^{\mathrm{VEST}}(t_n) + \sum_j W_{ij}^{\mathrm{INTVIS}} r_j^{\mathrm{VIS}}(t_n)\right)$

769     where $\tau = 100\ ms$ and other terms are the same as in equation (4).

770

771     *Linear Reproduction of M1 Response*

772 To test whether the responses of the optimal and homogeneous model M1 can be linearly

773 reproduced from responses of M2, M3, and the experimental data, we first calculated the "optimal

774 traces" from M1 (**Figure 5b**), using

$$\Delta PSTH_{\mathrm{M1}}^{\bullet} = <PSTH_{M1,i}^{\bullet,+}> - <PSTH_{M1,i}^{\bullet,-}>$$

776 Where $\bullet$ denotes three cue conditions (vestibular, visual, and combined), $PSTH_{M1,i}^{\bullet,+}$ and

777 $PSTH_{M1,i}^{\bullet,-}$ denote averaged PSTH for the $i$th LIP unit in the network M1 when the network makes

778 correct choices towards the neuron's preferred direction and null direction, respectively, and $<\cdot>$

779 denotes averaging across cells. To mimic the experimental procedure, only cells whose preferred

780 directions were close to $\pm 90°$ (with deviations less than 20°) were used. Similarly, we extracted

781 single cell activities from M2, M3, the LIP data, and the MSTd data (**[19]Gu, _et al._, 2006**)

$$\Delta PSTH_{*,i}^{\bullet} = PSTH_{*,i}^{\bullet,+} - PSTH_{*,i}^{\bullet,-}$$

783 where $* \in \{\mathrm{M2, M3, LIP\ data, MSTd\ data}\}$. Then we optimized sets of linear weights $\boldsymbol{w}_*$ to

784 minimize the cost function

$$E_* = \sum_{\bullet} \sum_{n} \left( \Delta PSTH_{\mathrm{M1}}^{\bullet}(t_n) - \sum_{i} w_{*,i} \Delta PSTH_{*,i}^{\bullet}(t_n) \right)^2 \qquad (8)$$

786 where, for example, $w_{\mathrm{LIP},i}$ represents the weight of the neuron $i$ in the LIP data when a

787 downstream area reads out LIP dynamics linearly to reproduce the optimal traces. To reduce

788 overfitting, we partitioned the data into two subsets along time by randomly assigning the time bins

789 into two sets, one for fitting ($T_{\mathrm{fit}}$) and the other for validating ($T_{\mathrm{valid}}$). During fitting, when the

790 validating error $E_{*,t_n \in T_{\mathrm{valid}}}$ started increasing, we stopped the iteration, a procedure known as early

791 stopping. The fitting results are shown in **Figure 5c–f**. Note that the $\Delta PSTH$s in the cost function

792 (equation (8)) grouped all the heading angles together. The results were qualitatively similar when

793 the cost function included error terms calculated from each heading angle separately, i.e.,

794
$$E_* = \sum_{\bullet} \sum_{n} \sum_{|h|} \left( \Delta PSTH_{\mathrm{M1}}^{\bullet,|h|}(t_n) - \sum_{i} w_{*,i} \Delta PSTH_{*,i}^{\bullet,|h|}(t_n) \right)^2 \qquad (9)$$

795    where $|h|$ denotes the absolute value of heading angle ($0°$, $1°$, $2°$, $4°$, $8°$). The reconstructions of

796    M1 traces with LIP activities using equation (9) are shown in **Supplementary Figure 9**.

797

798    To assess the robustness of the linear reconstruction, we randomly subsampled the same number of

799    neurons (n = 50, without replacement) from the four data sets, performed the linear fitting, and

800    repeated this procedure for 1000 times. The mean squared error and the distribution of readout

801    weights of the fittings are shown in **Figure 5g, h**. To examine whether only a small fraction of cells

802    contributed heavily to the fittings or whether the majority of cells did, we compared the

803    distributions of weights from the four data sets with the distribution of weights from a random

804    linear decoder. To do so, for each subsampling, we also generated a set of random readout weights

805    from a rectified Gaussian distribution (**Figure 5h**, black curve) and computed the kurtosis of the

806    distribution of weights from the random decoder as well as those from the four data sets (**Figure**

807    **5i**). The p-values were derived from the empirical subsampling distributions (two-tailed).

808

809    ## Data and Code Availability

810    MATLAB code for the network model and the information loss calculation is available at the

811    following public repository: https://github.com/hanhou/Multisensory-PPC. Experimental data and

812    code for data analysis are available upon request to the authors.

813

814    # Acknowledgments

818    the Strategic Priority Research Program of CAS (XDBS01070201), the Shanghai Municipal

819    Science and Technology Major Project (2018SHZDZX05) to Y.G and by grants from the Simons

820    Collaboration for the Global Brain and the Swiss National Science Foundation (#31003A_165831)

821    to A.P.

822

## Author Contributions

824    H.H. and Y.G. conceived the project and designed the experiments. H.H., Q.Z., and Y.Z. performed

825    the experiments. H.H. analyzed the data. H.H. and A.P. developed the models and implemented the

826    simulations. H.H., A.P., and Y.G. wrote the manuscript.

827

## Competing financial interests

829    The authors declare no competing financial interests.

830

## References

832    1. Ratcliff, R. A theory of memory retrieval. *Psychological review* **85**, 59 (1978).
833    2. Ratcliff, R. & McKoon, G. The diffusion decision model: theory and data for two-choice decision
834    tasks. *Neural computation* **20**, 873-922 (2008).
835    3. Ratcliff, R. & Rouder, J.N. Modeling response times for two-choice decisions. *Psychological*
836    *science* **9**, 347-356 (1998).
837    4. Ratcliff, R. & Smith, P.L. A comparison of sequential sampling models for two-choice reaction
838    time. *Psychol Rev* **111**, 333-367 (2004).
839    5. Laming, D.R.J. Information theory of choice-reaction times. (1968).
840    6. Bogacz, R., Brown, E., Moehlis, J., Holmes, P. & Cohen, J.D. The physics of optimal decision
841    making: a formal analysis of models of performance in two-alternative forced-choice tasks. *Psychol*
842    *Rev* **113**, 700-765 (2006).
843    7. Gold, J.I. & Shadlen, M.N. The neural basis of decision making. *Annual review of neuroscience*
844    **30**, 535-574 (2007).
845    8. Drugowitsch, J., DeAngelis, G.C., Klier, E.M., Angelaki, D.E. & Pouget, A. Optimal
846    multisensory decision-making in a reaction-time task. *Elife* **3** (2014).
847    9. Beck, J.M.*, et al.* Probabilistic population codes for Bayesian decision making. *Neuron* **60**, 1142-
848    1152 (2008).
849    10. Ma, W.J., Beck, J.M., Latham, P.E. & Pouget, A. Bayesian inference with probabilistic
850    population codes. *Nat Neurosci* **9**, 1432-1438 (2006).
851    11. Fetsch, C.R., Pouget, A., DeAngelis, G.C. & Angelaki, D.E. Neural correlates of reliability-
852    based cue weighting during multisensory integration. *Nat Neurosci* **15**, 146-154 (2012).

853    12. Gu, Y., Angelaki, D.E. & Deangelis, G.C. Neural correlates of multisensory cue integration in
854    macaque MSTd. *Nat Neurosci* **11**, 1201-1210 (2008).
855    13. Fetsch, C.R., DeAngelis, G.C. & Angelaki, D.E. Bridging the gap between theories of sensory
856    cue integration and the physiology of multisensory neurons. *Nature reviews. Neuroscience* **14**, 429-
857    442 (2013).
858    14. Shadlen, M.N. & Newsome, W.T. Neural basis of a perceptual decision in the parietal cortex
859    (area LIP) of the rhesus monkey. *Journal of neurophysiology* **86**, 1916-1936 (2001).
860    15. Shadlen, M.N. & Newsome, W.T. Motion perception: seeing and deciding. *Proceedings of the
861    National Academy of Sciences of the United States of America* **93**, 628-633 (1996).
862    16. Huk, A.C., Katz, L.N. & Yates, J.L. The Role of the Lateral Intraparietal Area in (the Study of)
863    Decision Making. *Annual review of neuroscience* **40**, 349-372 (2017).
864    17. Roitman, J.D. & Shadlen, M.N. Response of neurons in the lateral intraparietal area during a
865    combined visual discrimination reaction time task. *The Journal of neuroscience : the official
866    journal of the Society for Neuroscience* **22**, 9475-9489 (2002).
867    18. Boussaoud, D., Ungerleider, L.G. & Desimone, R. Pathways for Motion Analysis - Cortical
868    Connections of the Medial Superior Temporal and Fundus of the Superior Temporal Visual Areas
869    in the Macaque. *J Comp Neurol* **296**, 462-495 (1990).
870    19. Gu, Y., Watkins, P.V., Angelaki, D.E. & DeAngelis, G.C. Visual and nonvisual contributions to
871    three-dimensional heading selectivity in the medial superior temporal area. *The Journal of
872    neuroscience : the official journal of the Society for Neuroscience* **26**, 73-85 (2006).
873    20. Chen, A., DeAngelis, G.C. & Angelaki, D.E. Representation of vestibular and visual cues to
874    self-motion in ventral intraparietal cortex. *The Journal of neuroscience : the official journal of the
875    Society for Neuroscience* **31**, 12036-12052 (2011c).
876    21. Chen, A., DeAngelis, G.C. & Angelaki, D.E. Functional Specializations of the Ventral
877    Intraparietal Area for Multisensory Heading Discrimination. *The Journal of Neuroscience* **33**,
878    3567-3581 (2013).
879    22. Nikbakht, N., Tafreshiha, A., Zoccolan, D. & Diamond, M.E. Supralinear and Supramodal
880    Integration of Visual and Tactile Signals in Rats: Psychophysics and Neuronal Mechanisms. *Neuron*
881    **97**, 626-639.e628 (2018).
882    23. Raposo, D., Kaufman, M.T. & Churchland, A.K. A category-free neural population supports
883    evolving demands during decision-making. *Nat Neurosci* **17**, 1784-1792 (2014).
884    24. Knill, D.C. & Richards, W. *Perception as Bayesian inference* (Cambridge University Press,
885    1996).
886    25. Park, I.M., Meister, M.L., Huk, A.C. & Pillow, J.W. Encoding and decoding in parietal cortex
887    during sensorimotor decision-making. *Nat Neurosci* **17**, 1395-1403 (2014).
888    26. Meister, M.L., Hennig, J.A. & Huk, A.C. Signal multiplexing and single-neuron computations
889    in lateral intraparietal area during decision-making. *The Journal of Neuroscience* **33**, 2254-2267
890    (2013).
891    27. Katz, L.N., Yates, J.L., Pillow, J.W. & Huk, A.C. Dissociated functional significance of
892    decision-related activity in the primate dorsal stream. *Nature* **advance online publication** (2016).
893    28. Moreno-Bote, R.*, et al.* Information-limiting correlations. *Nat Neurosci* **17**, 1410-1417 (2014).
894    29. Kiani, R., Hanks, T.D. & Shadlen, M.N. Bounded integration in parietal cortex underlies
895    decisions even when viewing duration is dictated by the environment. *The Journal of neuroscience :
896    the official journal of the Society for Neuroscience* **28**, 3017-3029 (2008).
897    30. Scott, B.B.*, et al.* Fronto-parietal Cortical Circuits Encode Accumulated Evidence with a
898    Diversity of Timescales. *Neuron* **95**, 385-398.e385 (2017).
899    31. Beck, J., Bejjanki, V.R. & Pouget, A. Insights from a simple expression for linear fisher
900    information in a recurrently connected population of spiking neurons. *Neural computation* **23**,
901    1484-1502 (2011).
902    32. Seung, H.S. & Sompolinsky, H. Simple models for reading neuronal population codes.
903    *Proceedings of the National Academy of Sciences of the United States of America* **90**, 10749-10753

904     (1993).
905     33. Series, P., Latham, P.E. & Pouget, A. Tuning curve sharpening for orientation selectivity: coding
906     efficiency and the impact of correlations. *Nat Neurosci* **7**, 1129-1135 (2004).
907     34. Chen, A., DeAngelis, G.C. & Angelaki, D.E. A comparison of vestibular spatiotemporal tuning
908     in macaque parietoinsular vestibular cortex, ventral intraparietal area, and medial superior temporal
909     area. *The Journal of neuroscience : the official journal of the Society for Neuroscience* **31**, 3082-
910     3094 (2011a).
911     35. Fetsch, C.R*., et al.* Spatiotemporal Properties of Vestibular Responses in Area MSTd. *Journal
912     of neurophysiology* **104**, 1506-1522 (2010).
913     36. Smith, A.T., Greenlee, M.W., DeAngelis, G.C. & Angelaki, D.E. Distributed Visual–Vestibular
914     Processing in the Cerebral Cortex of Man and Macaque. *Multisensory Research* **30**, 91-120 (2017).
915     37. Laurens, J*., et al.* Transformation of spatiotemporal dynamics in the macaque vestibular system
916     from otolith afferents to cortex. *Elife* **6**, e20787 (2017).
917     38. Lisberger, S.G. & Movshon, J.A. Visual motion analysis for pursuit eye movements in area MT
918     of macaque monkeys. *The Journal of neuroscience : the official journal of the Society for
919     Neuroscience* **19**, 2224-2246 (1999).
920     39. Churchland, A.K*., et al.* Variance as a signature of neural computations during decision making.
921     *Neuron* **69**, 818-831 (2011).
922     40. Bizley, J.K., Jones, G.P. & Town, S.M. Where are multisensory signals combined for perceptual
923     decision-making? *Current opinion in neurobiology* **40**, 31-37 (2016).
924     41. Gu, Y. Vestibular signals in primate cortex for self-motion perception. *Current opinion in
925     neurobiology* **52**, 10-17 (2018).
926     42. Deneve, S., Latham, P.E. & Pouget, A. Efficient computation and cue integration with noisy
927     population codes. *Nat Neurosci* **4**, 826-831 (2001).
928     43. Gu, Y., Deangelis, G.C. & Angelaki, D.E. Causal links between dorsal medial superior temporal
929     area neurons and multisensory heading perception. *The Journal of neuroscience : the official
930     journal of the Society for Neuroscience* **32**, 2299-2313 (2012).
931     44. Chen, A., Gu, Y., Liu, S., DeAngelis, G.C. & Angelaki, D.E. Evidence for a Causal Contribution
932     of Macaque Vestibular, But Not Intraparietal, Cortex to Heading Perception. *The Journal of
933     neuroscience : the official journal of the Society for Neuroscience* **36**, 3789-3798 (2016).
934     45. Chandrasekaran, C. Computational principles and models of multisensory integration. *Current
935     opinion in neurobiology* **43**, 25-34 (2017).
936     46. Wong, K.F. & Wang, X.J. A recurrent network mechanism of time integration in perceptual
937     decisions. *The Journal of neuroscience : the official journal of the Society for Neuroscience* **26**,
938     1314-1328 (2006).
939     47. Wang, X.J. Probabilistic decision making by slow reverberation in cortical circuits. *Neuron* **36**,
940     955-968 (2002).
941     48. Mante, V., Sussillo, D., Shenoy, K. & Newsome, W. Context-dependent computation by
942     recurrent dynamics in prefrontal cortex. *Nature* **503**, 78-84 (2013).
943     49. Song, H.F., Yang, G.R. & Wang, X.J. Reward-based training of recurrent neural networks for
944     cognitive and value-based tasks. *Elife* **6** (2017).
945     50. Barash, S., Bracewell, R.M., Fogassi, L., Gnadt, J.W. & Andersen, R.A. Saccade-related activity
946     in the lateral intraparietal area. I. Temporal properties; comparison with area 7a. *Journal of
947     neurophysiology* **66**, 1095-1108 (1991).
948     51. Gu, Y., Fetsch, C.R., Adeyemo, B., Deangelis, G.C. & Angelaki, D.E. Decoding of MSTd
949     population activity accounts for variations in the precision of heading perception. *Neuron* **66**, 596-
950     609 (2010).
951     52. Goldman, M.S. Memory without Feedback in a Neural Network. *Neuron* **61**, 621-634 (2009).
952     53. Deneve, S., Latham, P.E. & Pouget, A. Reading population codes: a neural implementation of
953     ideal observers. *Nat Neurosci* **2**, 740-745 (1999).
954     54. Song, S., Sjostrom, P.J., Reigl, M., Nelson, S. & Chklovskii, D.B. Highly nonrandom features

955     of synaptic connectivity in local cortical circuits. *PLoS biology* **3**, e68 (2005).
956

957

## Supplementary Materials

**Supplementary Figure 1. Recording sites and reliable area mapping**

**Supplementary Figure 2. More example LIP cells**

**Supplementary Figure 3. Task-difficulty dependence of choice signals**

**Supplementary Figure 4. Macaque LIP is category-free**

**Supplementary Figure 5. De-mixing of choice and modality signals**

**Supplementary Figure 6. Information loss of ilPPC solution with heterogeneous MSTd population**

**Supplementary Figure 7. Model M2 achieves optimal behavior with heterogeneous units**

**Supplementary Figure 8. Example units in linear reconstruction of M1**

**Supplementary Figure 9. Linear reconstruction of M1 with cost function calculated for separated heading angles**

**Supplementary Figure 10. A trained recurrent neural network (RNN) performing multisensory decision-making task**

**Figure 1  Optimal cue integration in vestibular-visual multisensory decision-making task.**

**(a)** Schematic drawing of the experimental setup (top view). The vestibular (blue) and visual (red) stimuli of self-motion were provided by a motion platform and an LCD screen mounted on it, respectively. The monkey was seated on the platform and physically translated within the horizontal plane (blue arrows), whereas the screen rendered optical flow simulating what the monkey would see when moving through a three-dimensional star field (red dots). In a combined condition (green), both vestibular and visual stimuli were presented synchronously. The monkey's task was to discriminate whether the heading direction was to the left or the right of the straight ahead (black dashed line). **(b)** Task timeline. The monkey initiated a trial by fixating at a fixation point, and two choice targets appeared. The monkey then experienced a 1.5-s forward self-motion stimulus with a small leftward or rightward component, after which the monkey reported his perceived heading by making a saccadic eye movement to one of the two targets. The self-motion speed followed a Gaussian-shape profile. **(c)** Example psychometric functions from one session. The proportion of "rightward" choices is plotted against the headings for three cue conditions respectively. Smooth curves represent best-fitting cumulative Gaussian functions. **(d)** Average psychophysical thresholds from two monkeys for three conditions and predicted thresholds calculated from optimal cue integration theory (black bars). Error bars indicate s.e.m.; p values were from paired t-test.
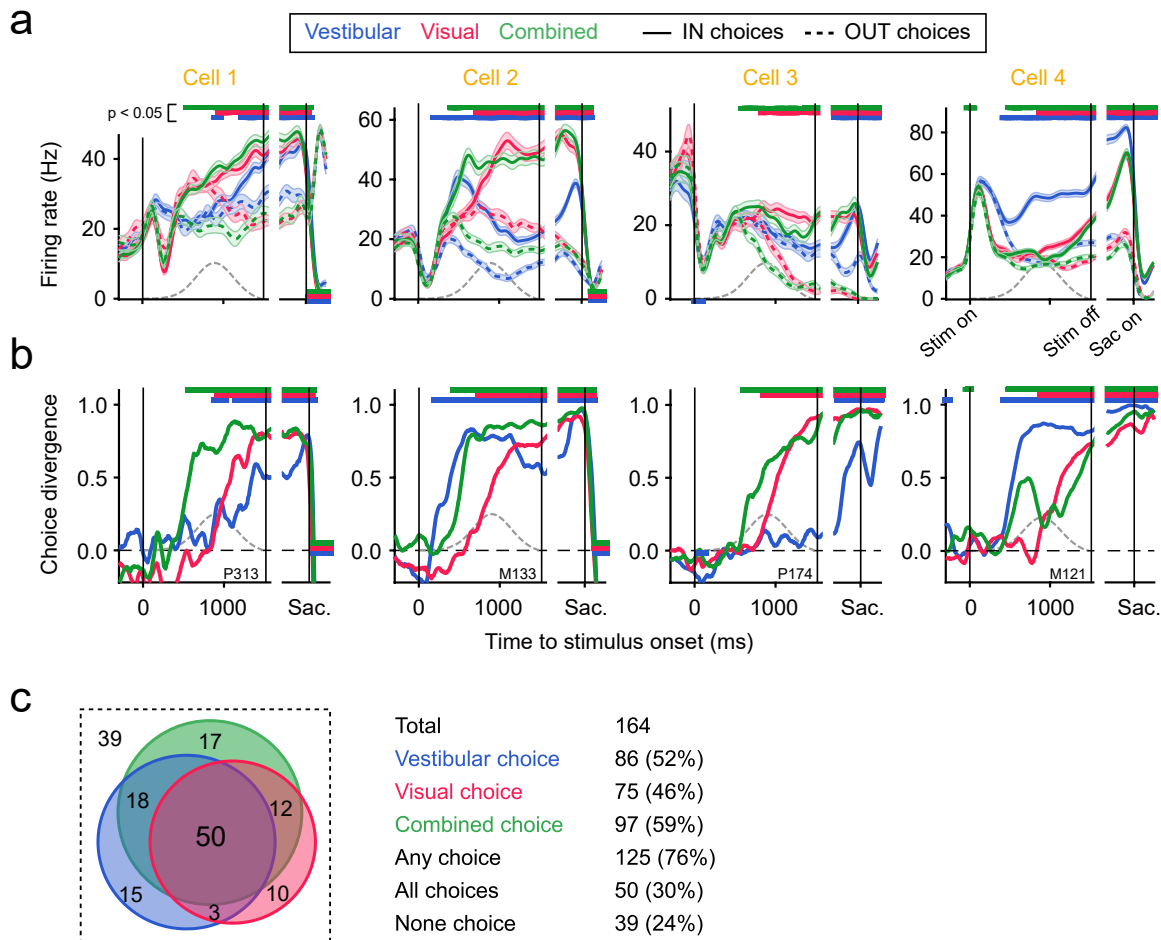
127

**Figure 2  Heterogeneous choice signals in LIP population.**

**(a)** Peri-stimulus time histograms (PSTHs) of four examples cells. Spike trains were aligned to stimulus onset (left subpanels) and saccade onset (right subpanels), respectively, and grouped by cue condition and monkey's choice. Vestibular, blue; visual, red; combined, green. Toward the cell's response field (RF), or IN choices, solid curves; away from the cell's RF, or OUT choices, dashed curves. Mean firing rates were computed from 10-ms time windows and smoothed with a Gaussian ($\sigma$ = 50 ms); only correct trials or 0° heading trials were included. Shaded error bands, s.e.m. Horizontal color bars represent time epochs in which IN and OUT trials have significantly different firing rates ($p < 0.05$, t-test), with the color indicating cue condition and the position indicating the relationship between IN and OUT firings (IN > OUT, top; IN < OUT, bottom). Gray dashed curves represent the actual speed profile measured by an accelerometer attached to the motion platform. **(b)** Choice divergence (CD) of the same four cells. CD ranged from -1 to 1 and was derived from ROC analysis for PSTHs in each 10-ms window (see Methods). Horizontal color bars are the same as in **a** except that p-values were from permutation test (n = 1000). **(c)** Venn diagram showing the distribution of choice signals. Numbers within colored areas indicate the numbers of neurons that have significant grand CDs (CD computed from all spikes in 0–1500 ms) under the corresponding combinations of cue conditions.
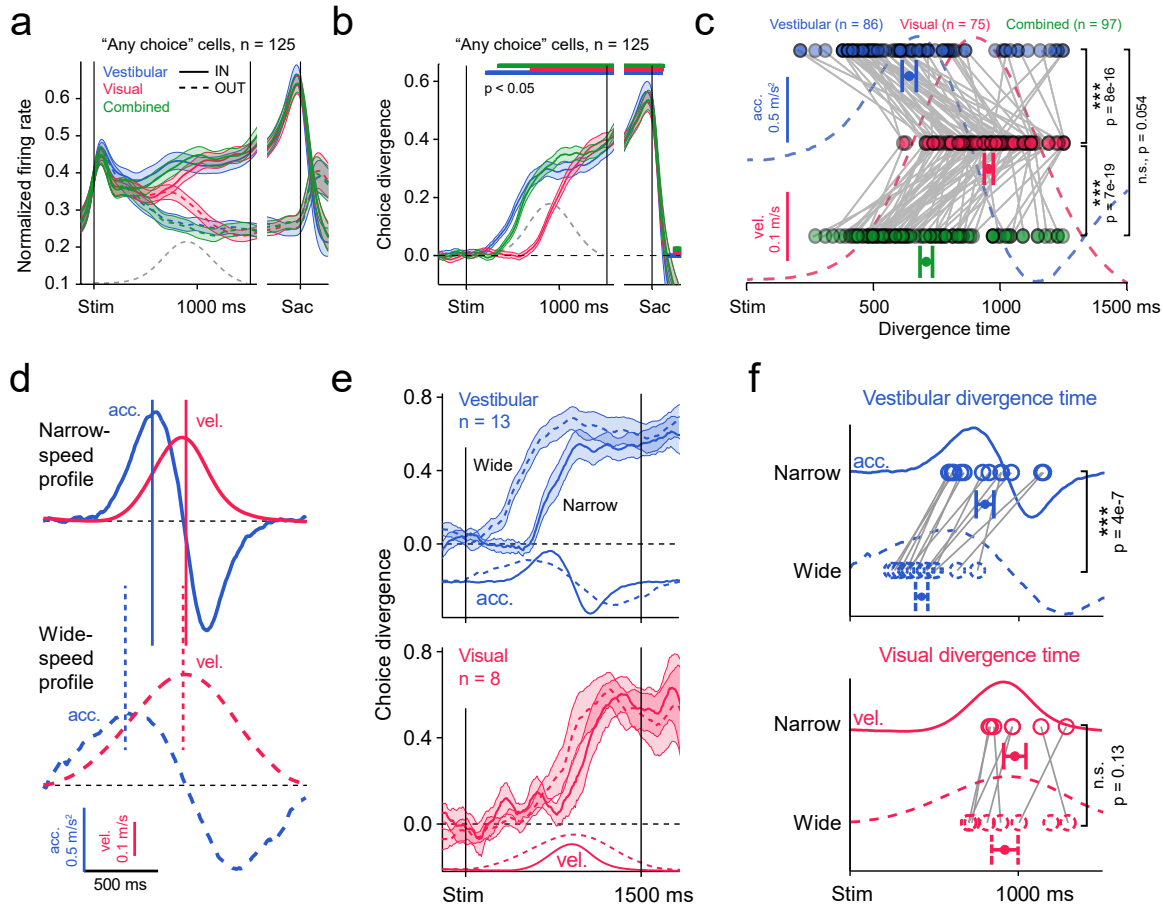
**Figure 3  LIP integrates vestibular acceleration but visual speed.**

**(a and b)** Population average of normalized PSTHs (**a**) and CD (**b**) from 125 "any choice" cells. The vestibular (blue) and combined (green) CDs ramped up much earlier than the visual one (red). Horizontal color bars indicate the time epochs in which population CDs are significantly larger than zero (p < 0.05, t-test). Gray dashed curve, the actual Gaussian speed profile; shaded error bands, s.e.m. **(c)** Divergence time of cells with significant grand CD for each condition. Divergence time was defined as the first occurrence of a 250-ms window in which CD was consistently larger than zero (p < 0.05, permutation test). Gray lines connect data from the same cells; acceleration and speed profiles shown in the background. Data points with horizontal error bars, mean ± s.e.m. of population divergence time; p values, t-test. **(d)** Two motion profiles used to isolate contributions of acceleration and speed to LIP ramping. Top and solid, the narrow-speed profile; bottom and dashed, the wide-speed profile; blue, acceleration; red, speed. Note that by widening the speed profile, we shifted the time of acceleration peak forward (blue vertical lines) while keeping the speed peak unchanged (red vertical lines). **(e)** Vestibular and visual CDs under the two motion profiles. **(f)** Comparison of divergence time between narrow and wide profiles. Note that the vestibular divergence time was significantly shifted, whereas the visual one was not, indicating that LIP integrates sensory evidence from vestibular acceleration and visual speed.
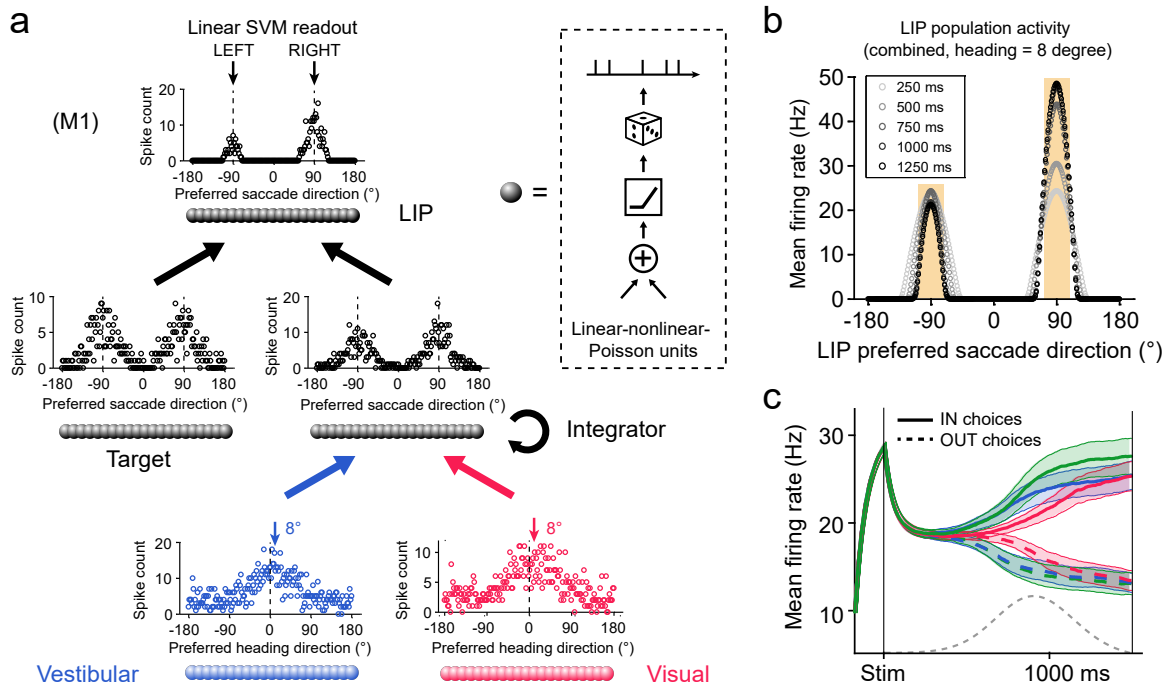
205

**Figure 4  Neural network model with invariant linear probabilistic population codes (ilPPC).**

**(a)** Network architecture of model M1. The model consists of three interconnected layers of linear-nonlinear-Poisson units (inset). Units in Vestibular and Visual layers have bell-shape ilPPC-compatible tuning curves for heading direction and receive heading stimuli with temporal dynamics following acceleration and speed, respectively. The intermediate Integrator layer simply sums the incoming spikes from the two sensory layers over time and transforms the tuning curves for heading direction to that for saccade direction (-90°, leftward choice; +90°, rightward choice). The LIP layer receives the integrated heading inputs from the Integrator layer, together with visual responses triggered by the two saccade targets. LIP units also have lateral connections implementing short-range excitation and long-range inhibition. Once a decision boundary is hit, or when the end of the trial is reached (1.5 s), LIP activity is decoded by a linear support vector machine for action selection (see **Methods**). Circles indicate representative patterns of activity for each layer; spike counts from 800–1000 ms; combined condition, 8° heading. **(b)** Population firing rate in the LIP layer at five different time points (the same stimulus as in **a**, averaged over 100 repetitions). **(c)** Average PSTHs across LIP population. Trials included three cue conditions and nine heading directions (±8°, ±4°, ±2°, ±1°, 0°). To mimic the experimental procedure, only units with preferred saccade direction close to ±90° were used (with deviation less than 20°; yellow shaded area in **b**). Notations are the same as in **Figure 2a and Figure 3a**.
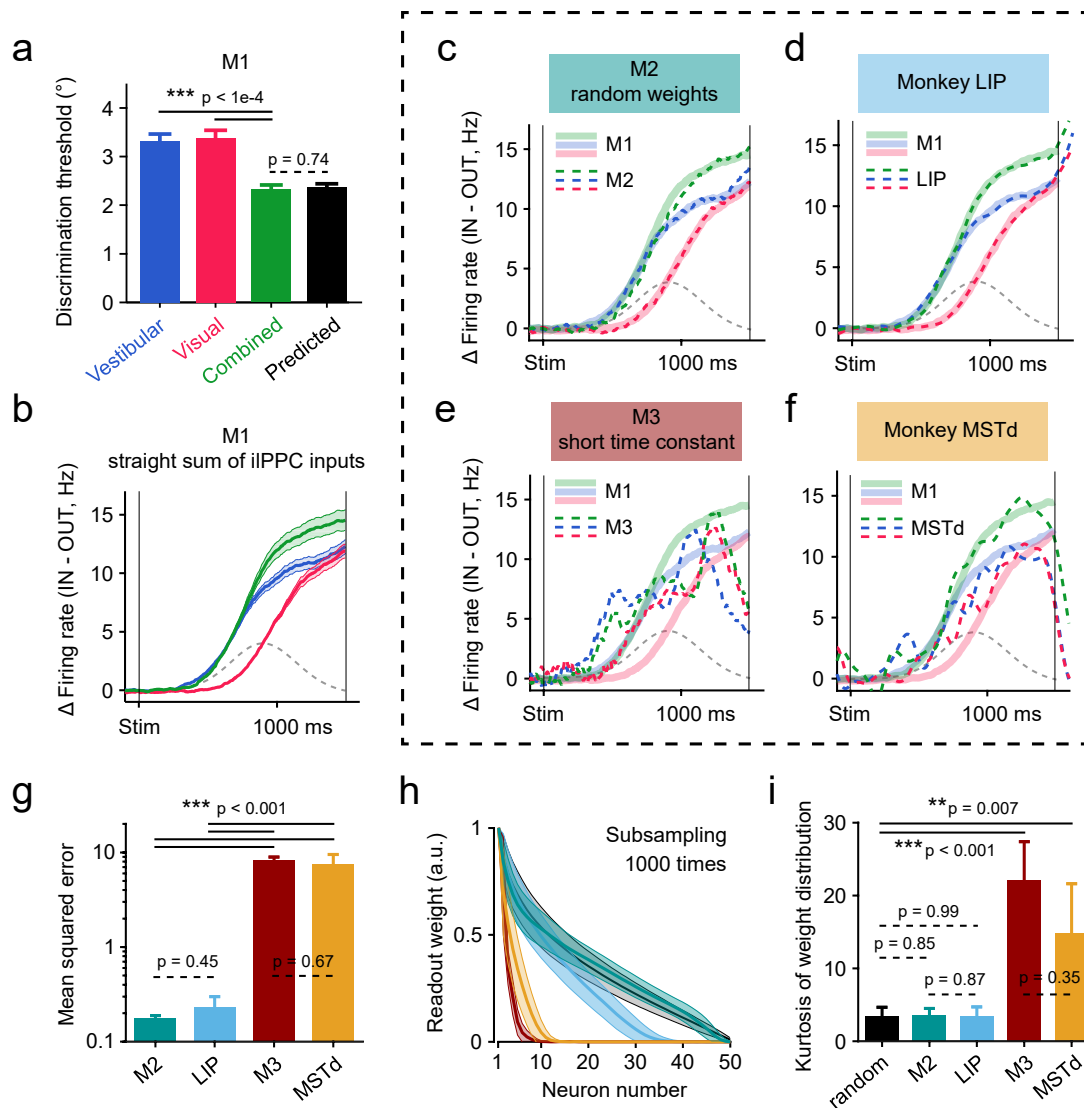
**Figure 5  Optimal ilPPC model M1 can be linearly approximated by M2 and LIP but not by M3 and MSTd.**

**(a)** Model M1 exhibited near-optimal behavior as the monkey. The psychophysical threshold under the combined condition (green) was indistinguishable from the Bayesian optimal prediction (black). **(b)** Ramping activity of M1 computed as the difference of PSTHs for IN and OUT trials. Activities from hypothetical units in the LIP layer with preferred direction close to ±90° were averaged together (see **Figure 4c** and **Methods**). Since M1 is optimal and homogeneous, we refer to M1's activities as "optimal traces" (see the main text). Notations are the same as before. **(c)** Optimal traces from M1 (thick shaded bands) can be linearly reconstructed by population activities obtained from a heterogenous model M2 (dashed curves). Model M2 had the same network architecture as M1 except that it relies on random combinations of ilPPC inputs in the integration layer (see **Methods**). **(d)** Optimal traces can also be linearly reconstructed by heterogenous single neuron activities from the LIP data. The similarity between **c** and **d** suggests that both model M2 and monkey LIP are heterogeneous variations of to the optimal ilPPC model M1. **(e and f)** In contrast, the optimal traces cannot be reconstructed from activities of a suboptimal model M3 (**e**) or from the MSTd data (**f**), presumably because the time constants in M3 and MSTd were too short. **(g)** Mean squared error of the fits in panels **c–f**. Error bars and p values were from subsampling test (n = 50 neurons, 1000 times). **(h)** Normalized readout weights ordered by magnitude. Shaded error bands indicate standard deviations of the subsampling distributions. **(i)** The kurtosis of the distributions of weights. The black curve in (**h**) and black bar in (**e**) were from random readout weights (see **Methods**).
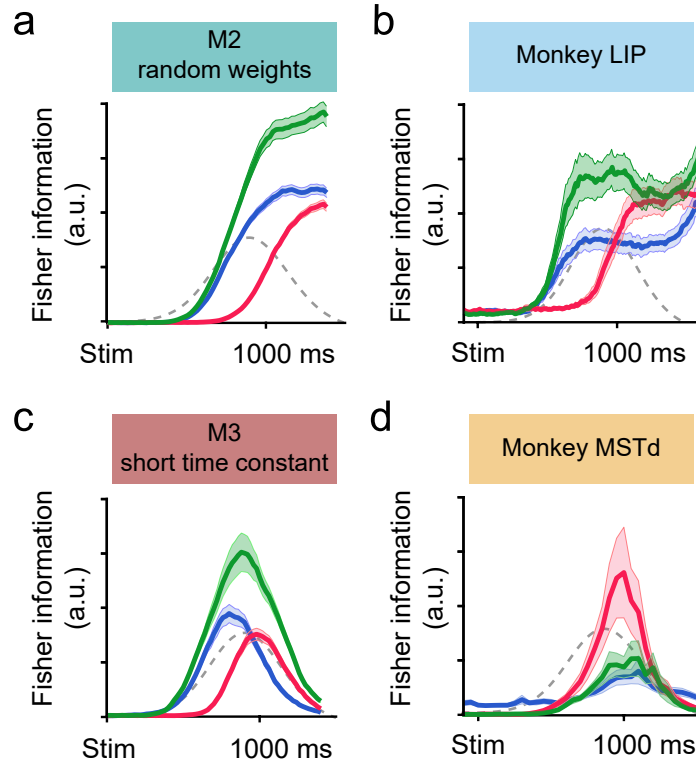
257

**Figure 6  Shuffled Fisher information for the model and the experimental data.**

**(a)** Shuffled Fisher information of M2 calculated by $I_{shuffled} = \sum_i f_i'^2 / \sigma_i^2$, where $f_i'$ denotes the derivative of the local tuning curve of the $i$th neuron and $\sigma_i^2$ denotes the averaged variance of its responses around 0° (see **Methods**). Both correct and wrong trials were included. Shaded error bands, s.e.m. estimated from bootstrap. Note that the absolute value of shuffled Fisher information is arbitrary. **(b-d)** Same as in **a** but for the monkey LIP data, the M3 responses, and the monkey MSTd data, respectively. Note that LIP is similar to M2, and MSTd to M3.

310