1     **Prior expectations induce pre-stimulus sensory templates**

2     Peter Kok[1,2]*, Pim Mostert[1] and Floris P. de Lange[1]

3     1.   Radboud University Nijmegen, Donders Institute for Brain, Cognition and Behaviour, Kapittelweg 29, 6525 EN

4          Nijmegen, The Netherlands

5     2.   Princeton University, Princeton Neuroscience Institute, 301 Peretsman-Scully Hall, Princeton, NJ 08544

6

7     **\*Corresponding author**

8     Princeton Neuroscience Institute

9     Princeton University

10    301 Peretsman-Scully Hall, Princeton, NJ 08544

11    Phone: 1 609 258 8729

12    E-mail: pkok@princeton.edu

13

19 **Abstract**

20

21 Perception can be described as a process of inference, integrating bottom-up sensory inputs and top-

22 down expectations. However, it is unclear how this process is neurally implemented. It has been

23 proposed that expectations lead to pre-stimulus baseline increases in sensory neurons tuned to the

24 expected stimulus, which in turn affects the processing of subsequent stimuli. Recent fMRI studies have

25 revealed stimulus-specific patterns of activation in sensory cortex as a result of expectation, but this

26 method lacks the temporal resolution necessary to distinguish pre- from post-stimulus processes. Here,

27 we combined human MEG with multivariate decoding techniques to probe the representational content

28 of neural signals in a time-resolved manner. We observed a representation of expected stimuli in the

29 neural signal well before they were presented, demonstrating that expectations indeed induce a pre-

30 activation of stimulus templates. These results suggest a mechanism for how predictive perception can

31 be neurally implemented.

## Introduction

Perception is heavily influenced by prior knowledge[1–3]. Accordingly, many theories cast perception as a process of inference, integrating bottom-up sensory inputs and top-down expectations[4–6]. However, it is unclear how this integration is neurally implemented. It has been proposed that prior expectations lead to baseline increases in sensory neurons tuned to the expected stimulus[7–9], which in turn leads to improved neural processing of matching stimuli[10,11]. In other words, expectations may induce stimulus templates in sensory cortex, prior to the actual presentation of the stimulus. Alternatively, top-down influences in sensory cortex may exert their influence only after the bottom-up stimulus has been initially processed, and the integration of the two sources of information may become apparent only during later stages of sensory processing[12].

The evidence necessary to distinguish between these hypotheses has been lacking. fMRI studies have revealed stimulus-specific patterns of activation in sensory cortex as a result of expectation[9,13], but this method lacks the temporal resolution necessary to distinguish pre- from post-stimulus periods. Here, we combined MEG with multivariate decoding techniques to probe the representational content of neural signals in a time-resolved manner[14–17]. We trained a forward model to decode the orientation of task-irrelevant gratings from the MEG signal[18,19], and applied this decoder to trials in which participants expected a grating of a particular orientation to be presented. This analysis revealed a neural representation of the expected grating that resembled the neural signal evoked by an actually presented grating. This representation was present already before stimulus presentation, demonstrating that expectations can indeed induce the pre-activation of stimulus templates.

3

## Results

Participants were exposed to auditory cues that predicted the likely orientation (45° or 135°) of an upcoming grating stimulus (Fig. 1a-b). This grating was followed by a second grating that differed slightly from the first in terms of orientation and contrast. In separate runs of the MEG session, participants performed either an orientation or contrast discrimination task on the two gratings (see Methods for details).

**Behavioural results.** Participants were able to discriminate small differences in orientation (3.9° ± 0.5°, accuracy = 74.0% ± 1.6%, mean ± sem) and contrast (4.6% ± 0.3%, accuracy = 76.6% ± 1.5%) of the cued gratings. There was no significant difference between the two tasks in terms of either accuracy ($F_{1,22}$ = 3.38, $p$ = 0.080) or reaction time (mean RT = 621 ms vs. 603 ms, $F_{1,22}$ = 1.46, $p$ = 0.24). Overall, accuracy and reaction times were not influenced by whether the cued grating had the expected or the unexpected orientation (accuracy: $F_{1,22}$ = 0.21, $p$ = 0.65; RT: $F_{1,22}$ = 0.03, $p$ = 0.87), nor was there an interaction between task and expectation (accuracy: $F_{1,22}$ = 0.96, $p$ = 0.34; RT: $F_{1,22}$ = 0.42, $p$ = 0.52). Note that these discrimination tasks were orthogonal to the expectation manipulation, in the sense that the expectation cue provided no information about the likely correct choice.

During the grating localiser (Fig. 1c, see Methods for details), participants correctly detected 91.2% ± 1.6% (mean ± sem) of fixation flickers, and incorrectly pressed the button on 0.2% ± 0.1% of trials, suggesting that participants were successfully engaged by the fixation task.

**MEG results – Localiser orientation decoding.** As mentioned, participants were exposed to auditory cues that predicted the likely orientation of an upcoming grating stimulus. The question we wanted to answer was whether the expectations induced by these auditory cues would evoke templates of the

4

77    visual stimuli prior to the presentation of the gratings. To be able to uncover such sensory templates, we

78    trained a decoding model to reconstruct the orientation of (task-irrelevant) visual gratings (Fig. 1c) from

79    the MEG signal, in a time-resolved manner. First, we found that this model was highly accurate at

80    reconstructing the orientation of such gratings from the MEG signal (Fig. 2). Grating orientation could be

81    decoded across an extended period of time (from 40 to 655 ms post-stimulus, $p < 0.001$, and from 685

82    to 730 ms, $p = 0.018$), peaking around 120-160 ms post-stimulus (Fig. 2c). Furthermore, in the period

83    around 100 to 330 ms post-stimulus, orientation decoding generalised across time, meaning that a

84    decoder trained on the evoked response at, for example, 120 ms post-stimulus could reconstruct the

85    grating orientation represented in the evoked response around 300 ms, and vice versa (Fig. 2d). In other

86    words, certain aspects of the representation of grating orientation were sustained over time.

87

88    **MEG results – Expectation induces stimulus templates.** Our main question pertained to the presence of

89    visual grating templates induced by the auditory expectation cues during the main experiment.

90    Therefore, we applied our model trained on task-irrelevant gratings to trials containing gratings that

91    were either validly or invalidly predicted, respectively (Fig. 3a). In both conditions, the decoding model

92    trained on task-irrelevant gratings succeeded in accurately reconstructing the orientation of the gratings

93    presented in the main experiment (valid expectation: cluster from training time 60 to 410 ms and

94    decoding time 60 to 400 ms, $p < 0.001$, and from training time 205 to 325 ms and decoding time 400 to

95    495 ms, $p = 0.045$; invalid expectation: cluster from training time 75 to 225 ms and decoding time 75 to

96    330 ms, $p = 0.0012$, and from training time 250 to 360 ms and decoding time 195 to 355 ms, $p = 0.027$).

97         If the cues induced sensory templates of the expected grating, one would expect these to be

98    revealed in the difference in decoding between valid and invalidly predicted gratings (see Material and

99    Methods for details of the subtraction logic). Indeed, this subtraction/analyses demonstrates that the

100   auditory expectation cues induce orientation-specific neural signals (Fig. 3a, bottom panel). These

101    signals were present already 40 ms before grating presentation, and extended into the post-stimulus

102    period (from decoding time -40 to 230 ms, $p$ = 0.0092, and from 300 to 530 ms, $p$ = 0.016). Furthermore,

103    these signals were uncovered when the decoder was trained on around 120 to 160 ms post-stimulus

104    during the grating localiser (Fig. 3b), suggesting that these cue-induced signals were similar to those

105    evoked by task-irrelevant gratings. In other words, the auditory expectation cues evoked orientation-

106    specific signals that were similar to sensory signals evoked by the corresponding actual grating stimuli.

107          In sum, expectations induced pre-stimulus sensory templates that influenced post-stimulus

108    representations as well; invalidly expected gratings had to 'overcome' a pre-stimulus activation of the

109    opposite orientation, while validly expected gratings were facilitated by a compatible pre-stimulus

110    activation (Supplementary Fig. 1a). The post-stimulus carryover of these expectation signals lasted

111    throughout the trial (Supplementary Fig. 1b).

112          As in previous studies using a similar paradigm[11,20], there was no interaction between the effects

113    of the expectation cue and the task (orientation vs. contrast discrimination) participants performed (no

114    clusters with $p$ < 0.4).

115          In the current study, there was no difference in the overall amplitude of the neural response

116    evoked between validly and invalidly expected gratings (no clusters with $p$ < 0.4, Supplementary Fig. 2).

117

## Discussion

119

120    Here, we show that expectations can induce sensory templates of the expected stimulus already before

121    the stimulus appears. These results extend previous fMRI studies demonstrating stimulus-specific

122    patterns of activation in sensory cortex induced by expectations, which could not resolve whether these

123    templates indeed reflected pre-stimulus expectations, or instead stimulus specific error signals induced

124    by the unexpected omission of a stimulus[9,13].

125    The fact that expectation signals were revealed by a decoder trained on physically presented

126    (but task-irrelevant) gratings suggests that these expectation signals resemble activity patterns induced

127    by actual stimuli. The expectation signal remained present throughout the trial, extending into the post-

128    stimulus period, suggesting the tonic activation of a stimulus template. These results are in line with a

129    recent monkey electrophysiology study[10], which showed that neurons in the face patch of IT cortex

130    encode the prior expectation of a face appearing, both prior to and following actual stimulus

131    presentation. When the subsequently presented stimulus is noisy or ambiguous, such a pre-stimulus

132    template could conceivably bias perception towards the expected stimulus[21–24].

133    What is the source of these cue-induced expectation signals? One candidate region is the

134    hippocampus, which is known to be involved in encoding associations between previously unrelated,

135    discontiguous stimuli[25], such as the auditory tones and visual gratings used in the present study.

136    Furthermore, fMRI studies have revealed predictive signals in the hippocampus[13,26,27], and Reddy and

137    colleagues[28] reported anticipatory firing to expected stimuli in the medial temporal lobe, including the

138    hippocampus. One intriguing possibility is that predictive signals from the hippocampus are fed back to

139    sensory cortex[13,29,30].

140    In addition to expectation, several other cognitive phenomena have been shown to induce

141    stimulus templates in sensory cortex, such as preparatory attention[17,31], mental imagery[32–34], and

7

142    working memory[35,36]. In fact, explicit task preparation can also induce pre-stimulus sensory templates

143    that last into the post-stimulus period[17]. Note that in the current study the task did not require explicit

144    use of the expectation cues, the task response was in fact orthogonal to the expectation. Furthermore,

145    there was no difference in the expectation signal between runs in which grating orientation was task-

146    relevant (orientation discrimination task) and when it was irrelevant (contrast discrimination task),

147    suggestion expectation may be a relatively automatic phenomenon[11,37]. In fact, neural modulations by

148    expectation have even been observed during states of inattention[38], sleep[39] and in patients experiencing

149    disorders of consciousness[40]. One important question for future research will be to establish whether

150    the same neural mechanism underlies the different cognitive phenomena that are capable of inducing

151    stimulus templates in sensory cortex, or whether different top-down mechanisms are at work. Indeed, it

152    has been suggested that expectation and attention, or task preparation, may have different underlying

153    neural mechanisms[20,41,42]. For instance, predictive coding theories suggest that attention may modulate

154    sensory signals in the superficial layers of sensory cortex, while predictions modulate the response in

155    deep layers[5,43].

156          One may wonder why the current study does not report a modulation of the overall neural

157    response by expectation, while previous studies have found an increased neural response to unexpected

158    stimuli[37,44–48], including some using an almost identical paradigm as the current study[11,20]. Of course, the

159    current study reports a null effect, from which it is hard to draw firm conclusions. However, it is possible

160    that the type of measurement of neural activity plays a role in the absence of the effect. Most previous

161    studies reporting expectation suppression in visual cortex used fMRI, while the current study used MEG.

162    It is possible that the BOLD signal, a mass-action signal that integrates synaptic and neural activity, as

163    well as integrating over time, is sensitive to certain neural effects that MEG, which is predominantly

164    sensitive to synchronised activity in pyramidal neurons oriented perpendicular to the cortical surface, is

8

165  not. It is even possible that within MEG, different types of sensors (i.e. magnetometers, planar and axial

166  gradiometers) differ in their sensitivity to expectation suppression[49].

167  Recent theories of sensory processing state that perception reflects the integration of bottom-

168  up inputs and top-down expectations, but ideas diverge on whether the brain continuously generates

169  stimulus templates in sensory cortex to pre-empt expected inputs[10,23,50,51], or rather engages in

170  perceptual inference only after receiving sensory inputs[52,53]. Our results are in line with the brain being

171  proactive, constantly forming predictions about future sensory inputs. These findings bring us closer to

172  uncovering the neural mechanisms by which we integrate prior knowledge with sensory inputs to

173  optimise perception.

## Methods

174

175

176 **Participants.** Twenty-three (15 female, age 26 ± 9, mean ± SD) healthy individuals participated in the

177 experiment. All participants were right-handed and had normal or corrected-to-normal vision. The study

178 was approved by the local ethics committee (CMO Arnhem-Nijmegen, The Netherlands) under the

179 general ethics approval ("Imaging Human Cognition", CMO 2014/288), and the experiment was

180 conducted in accordance with these guidelines. All participants gave written informed consent

181 according to the declaration of Helsinki.

182

183 **Stimuli.** Grayscale luminance-defined sinusoidal grating stimuli (spatial frequency: 1.0 cycles/°) were

184 generated using MATLAB (MathWorks, Natick, MA) in conjunction with the Psychophysics Toolbox[54].

185 Gratings were displayed in an annulus (outer diameter: 15° of visual angle, inner diameter: 1°),

186 surrounding a black fixation bull's eye (4 cd/m$^2$), on a gray (580 cd/m$^2$) background. The visual stimuli

187 were presented with an LCD projector (1024 × 768 resolution, 60 Hz refresh rate) positioned outside the

188 magnetically shielded room, and projected on a translucent screen via two front-silvered mirrors. The

189 projector lag was measured at 36 ms, which was corrected for by shifting the time axis of the data

190 accordingly. The auditory cue consisted of a pure tone (500 or 1000 Hz, 250 ms duration, including 10

191 ms on and off-ramp time), presented over MEG-compatible earphones.

192

193 **Experimental design.** Each trial consisted of an auditory cue, followed by two consecutive grating stimuli

194 (750 ms SOA between auditory and first visual stimulus) (Fig. 1a). The two grating stimuli were

195 presented for 250 ms each, separated by a blank screen (500 ms). A central fixation bull's eye (0.7°) was

196 presented throughout the trial, as well as during the intertrial interval (ITI, 2250 ms). The auditory cue

197 consisted of either a low- (500 Hz) or high-frequency (1000 Hz) tone, which predicted the orientation of

10

198     the first grating stimulus (45° or 135°) with 75% validity (Fig. 1b). In the other 25% of trials, the first

199     grating had the orthogonal orientation. Thus, the first grating had an orientation of either exactly 45° or

200     135°, and a luminance contrast of 80%. The second grating differed slightly from the first in terms of

201     both orientation and contrast (see below), as well as being in antiphase to the first grating (which had a

202     random spatial phase). The contingencies between the auditory cues and grating orientations were

203     flipped halfway through the experiment (i.e., after four runs), and the order was counterbalanced over

204     subjects.

205          In separate runs (64 trials each, ~4.5 minutes), subjects performed either an orientation or a

206     contrast discrimination task on the two gratings. When performing the orientation task, subjects had to

207     judge whether the second grating was rotated clockwise or anticlockwise with respect to the first

208     grating. In the contrast task, a judgment had to be made on whether the second grating had lower or

209     higher contrast than the first one. These tasks were explicitly designed to avoid a direct relationship

210     between the perceptual expectation and the task response. Subjects indicated their response (response

211     deadline: 750 ms after offset of the second grating) using an MEG-compatible button box. The

212     orientation and contrast differences between the two gratings were determined by an adaptive

213     staircase procedure[55], being updated after each trial. This was done to yield comparable task difficulty

214     and performance (~ 75% correct) for the different tasks. Staircase thresholds obtained during one task

215     were used to set the stimulus differences during the other task, in order to make the stimuli as similar as

216     possible in both contexts. As in previous studies using a similar paradigm[11,20], there was no interaction

217     between the effects of the expectation cue and the task participants performed, and therefore we

218     collapsed over the two tasks in all MEG analyses.

219          All subjects completed eight runs (four of each task, alternating every two runs, order was

220     counterbalanced over subjects) of the experiment, yielding a total of 512 trials. The staircases were kept

221     running throughout the experiment. Before the first run, as well as in between runs four and five, when

222    the contingencies between cue and stimuli were flipped, subjects performed a short practice run

223    containing 32 trials of both tasks (~4.5 minutes).

224        Interleaved with the main task runs, subjects performed eight runs of a grating localiser task (Fig.

225    1c). Each run (~2 min) consisted of 80 grating presentations (ITI uniformly jittered between 1000 and

226    1200 ms). The grating annuli were identical to those presented during the main task (80% contrast, 250

227    ms duration, 1.0 cycles/°, random spatial phase). Each grating had one of eight orientations (spanning

228    the 180° space, starting at 0°, in steps of 22.5°), each of which was presented ten times per run in

229    pseudorandom order. A black fixation bull's eye (4 cd/m$^2$, 0.7° diameter, identical to the one presented

230    during the main task runs) was presented throughout the run. On 10% of trials (counterbalanced across

231    orientations), the black fixation point in the centre of the bull's eye (0.2°, 4 cd/m$^2$) briefly turned gray

232    (324 cd/m$^2$) during the first 50 ms of grating presentation. Participants task was to press a button

233    (response deadline: 500 ms) when they perceived this fixation flicker. This simple task was meant to

234    ensure central fixation, while rendering the gratings task-irrelevant. Trials containing fixation flickers

235    were excluded from further analyses.

236        Finally, participants were exposed to a tone localiser (~1.5 min), presented at the start, end, and

237    halfway through the MEG session. These runs consisted of 81 presentations of the two tones used in the

238    main experiment. Data from these runs were not analysed further.

239        Prior to the MEG session (1–3 days), all participants completed a behavioural session.  The aim

240    of this session was to familiarise participants with the tasks and to initialise the staircase values for both

241    the orientation and the contrast discrimination task (see above). The behavioural session consisted of

242    written instructions and 32 practice trials of each task, followed by four runs (~4.5 min each) of the main

243    experiment (each task twice, alternating between runs, cue contingencies switching between the

244    second and third run). Finally, participants were exposed to one run each of the grating and tone

245    localiser, to familiarise them with the procedure.

12

246

247    **MEG recording and preprocessing.** Whole-head neural recordings were obtained using a 275-channel

248    MEG system with axial gradiometers (CTF Systems, Coquitlam, BC, Canada) located in a magnetically

249    shielded room.  Throughout the experiment, head position was monitored online, and corrected if

250    necessary, using three fiducial coils that were placed on the nasion and on earplugs in both ears [56]. If

251    subjects had moved their head more than 5 mm from the starting position they were repositioned

252    during block breaks. Furthermore, both horizontal and vertical electrooculograms (EOGs), as well as an

253    electrocardiogram (ECG) were recorded to facilitate removal of eye- and heart-related artifacts. The

254    ground electrode was placed at the left mastoid. All signals were sampled at a rate of 1200 Hz.

255       The data were preprocessed offline using FieldTrip[57] (www.fieldtriptoolbox.org). In order to

256    identify artifacts, the variance (collapsed over channels and time) was calculated for each trial. Trials

257    with large variances were subsequently selected for manual inspection and removed if they contained

258    excessive and irregular artifacts. Independent component analysis was subsequently used to remove

259    regular artifacts, such as heartbeats and eye blinks. Specifically, for each subject, the independent

260    components were correlated to both EOGs and the ECG to identify potentially contaminating

261    components, and these were subsequently inspected manually before removal. For the main analyses,

262    data were low-pass filtered using a two-pass Butterworth filter with a filter order of 6 and a frequency

263    cutoff of 40 Hz. To rule out that the temporal smoothing caused by low-pass filtering may have

264    artificially decreased the onset latency of neural signals, we repeated the decoding analyses (see below)

265    on data that were not low-pass filtered (Supplementary Fig. 3). Here, only notch filters were applied at

266    50, 100 and 150 Hz to remove line noise and its harmonics. Finally, main task data were baseline

267    corrected on the interval of −250 to 0 ms relative to auditory cue onset, and grating localiser data were

268    baseline corrected on the interval of -200 to 0 ms relative to visual grating onset.

269

270     **Event-related field analysis.** Event-related fields (ERFs) were calculated per participant, and subjected

271     to a planar gradient transformation[58] before averaging across participants. The planar transformation

272     simplifies the interpretation of the sensor-level data because it typically places the maximal signal above

273     the source. To avoid differences in the amount of noise when comparing conditions with different

274     numbers of trials, we matched the trial count by randomly selecting a subsample of trials from the

275     conditions with more trials (i.e., valid expectations).

276

277     **Orientation decoding analysis.** To probe sensory representations in the visual cortex, we used a forward

278     modelling approach to reconstruct the orientation of the grating stimuli from the MEG signal[17–19,59]. The

279     forward modelling approach was two-fold. First, a theoretical forward model was postulated that

280     described the measured activity in the MEG sensors, given the orientation of the presented grating.

281     Second, this forward model was used to obtain an inverse model that specified the transformation from

282     MEG sensor space to orientation space. The forward and inverse models were estimated on the basis of

283     the grating localiser data. The inverse model was then applied to the data from the main experiment, in

284     order to generalise from sensory signals evoked by task-irrelevant gratings to the gratings and

285     expectation signals evoked in the main task. To test the performance of the model we also applied it to

286     the localiser data itself, using a cross-validation approach in which in each iteration one trial of each

287     orientation was used at the test set, and the remaining data were used as the training set.

288          The forward model was based on work by Brouwer and Heeger[18,19] and involved 32 hypothetical

289     channels, each with an idealised orientation tuning curve. Each channel consisted of a half-wave-

290     rectified sinusoid raised to the fifth power, and the 32 channels were spaced evenly within the 180°

291     orientation space, such that a tuning curve with any possible orientation preference could be expressed

292     exactly as a weighted sum of the channels. Arranging the hypothesised channel activities for each trial

14

293    along the columns of a matrix **C** (32 channels × $n$ trials), the observed data could be described by the

294    following linear model:

295                                    $$\mathbf{B} = \mathbf{WC} + \mathbf{N}$$

296    where **B** are the ($m$ sensors × $n$ trials) MEG data, **W** is a weight matrix ($m$ sensors × 32 channels) that

297    specifies how channel activity is transformed into sensory activity, and **N** are the residuals (i.e., noise).

298        In order to obtain the inverse model, we estimated an array of spatial filters that, when applied

299    to the data, aimed to reconstruct the underlying channel activities as accurately as possible. In doing so,

300    we extended Brouwer and Heeger's[18,19] approach in three respects. First, since the MEG signal in

301    (nearby) sensors is correlated, we took into account the correlational structure of the noise. Second, we

302    estimated a spatial filter for each orientation channel independently. As a result, the number of

303    channels used in our model was not constrained, whereas the maximum number of channels would

304    otherwise be dependent on the number of presented orientations. In practice, this resulted in

305    smoothing in orientation space, because the channels were not truly independent. Third, each filter was

306    normalised such that the magnitude of its output matched the magnitude of the underlying channel

307    activity it was designed to recover. Prior to estimating the inverse model, **B** and **C** were demeaned such

308    that their average over trials equalled zero, for each sensor and channel, respectively.

309        As stated above, the inverse model was estimated on the basis of the grating localiser data. On

310    each localiser trial, one of eight orientations was presented (see above), and the hypothetical responses

311    of each of the channels could thus be calculated for each trial, resulting in the response row vector

312    $\mathbf{c}_{train,i}$ , of length $n_{train}$ trials, for each channel $i$. The weights on the sensors $\mathbf{w}_i$ could now be obtained

313    through least squares estimation, for each channel:

314                    $$\mathbf{w}_i = \mathbf{B}_{train}\mathbf{c}_{train,i}^{\mathsf{T}}\left(\mathbf{c}_{train,i}\mathbf{c}_{train,i}^{\mathsf{T}}\right)^{-1}$$

15

315     where $\mathbf{B}_{train}$ are the ($m$ sensors $\times$ $n_{train}$ trials) localiser MEG data. Subsequently, the optimal spatial

316     filter $\mathbf{v}_i$ to recover the activity of the $i$-th channel was obtained as follows[16]:

317
$$\mathbf{v}_i = \frac{\tilde{\Sigma}_i^{-1}\mathbf{w}_i}{\mathbf{w}_i^{\mathsf{T}}\tilde{\Sigma}_i^{-1}\mathbf{w}_i}$$

318     where $\tilde{\Sigma}_i$ is the regularised covariance matrix for channel $i$. Incorporating the noise covariance in the

319     filter estimation leads to the suppression of noise that arises from correlations between sensors. The

320     noise covariance was estimated as follows:

321
$$\hat{\Sigma}_i = \frac{1}{n_{train}-1}\boldsymbol{\varepsilon}_i\boldsymbol{\varepsilon}_i^{\mathsf{T}}$$

322
$$\boldsymbol{\varepsilon}_i = \mathbf{B}_{train} - \mathbf{w}_i\mathbf{c}_{train,i}$$

323     where $n_{train}$ is the number of training trials. For optimal noise suppression, we improved this estimation

324     by means of regularization by shrinkage, using the analytically determined optimal shrinkage parameter

325     (for details, see[60]), yielding the regularised covariance matrix $\tilde{\Sigma}_i$ .

326          Such a spatial filter was estimated for each hypothetical channel, yielding an $m$ sensors $\times$ 32

327     channel filter matrix $\mathbf{V}$. Given that we performed our decoding analysis in a time-resolved manner, $\mathbf{V}$

328     was estimated at each time point of the training data, in steps of 5 ms, resulting in array of filter

329     matrices, or decoders. To improve the signal-to-noise ratio, the data were first averaged within a

330     window of 29.2 ms centred on the time point of interest.  The window length of 29.2 ms was based on

331     an a priori chosen length of 30 ms, but minus one sample such that the window contained an odd

332     number of samples for symmetric centring[16]. These filter matrices could now be applied to estimate the

333     orientation channel responses in independent data – in this case, the trials from the main experiment:

334
$$\mathbf{C}_{test} = \mathbf{V}^{\mathsf{T}}\mathbf{B}_{test}$$

335    where $\mathbf{B}_{test}$ are the ($m$ sensors $\times$ $n_{test}$ trials) main experiment data. These channel responses were

336    estimated at each time point of the test data, in steps of 5 ms, with the data being averaged within a

337    window of 29.2 ms at each step. This procedure resulted in a four-dimensional (training time $\times$ testing

338    time $\times$ 32 channel $\times$ $n_{test}$ ) matrix of estimated channel responses for each trial in the main experiment.

339    Each trials' channel responses were shifted such that the channel with its hypothetical peak response at

340    the orientation presented on that trial (i.e. 45° or 135°) ended up in the position of the 0° channel,

341    before averaging over trials within each condition (i.e., valid vs. invalid expectation). Thus, the presented

342    orientation was defined as 0°, by convention. Note that for 3D surface plots that show the evolution of

343    channel responses over time (e.g., Fig. 2b), the response of the 90° channel (i.e., orthogonal to the

344    presented orientation) was used as a baseline, to avoid negative numbers for visualisation purposes.

345        To quantify decoding performance, the channel responses for a given condition were converted

346    into polar form and projected onto a vector with angle 0° (the presented orientation, see above).

347
$$r = |z|\cos\left(\arg\left(z\right)\right), \qquad z = \mathbf{c}\ e^{2i\varphi}$$

348    where $\mathbf{c}$ is a vector of estimated channel responses, and $\varphi$ is the vector of angles at which the channels

349    peak (multiplied by 2 to project the 180° orientation space onto the full 360° space). The scalar

350    projection $r$ indicates the strength of the decoder signal for the orientation presented on screen. (Note

351    that subtracting the estimated response of the 90° channel from that of the 0° channel yielded virtually

352    identical results, data not shown.) This quantification yielded (training time $\times$ testing time) temporal

353    generalisation matrices of orientation decoding performance.

354        In order to isolate any orientation-specific neural signals evoked by the expectation cues, we

355    applied the following subtraction logic. On valid expectation trials, the expected and presented

356    orientations are identical, and thus the orientation signal induced by both the cue and stimulus be

357    expected to be positive, by convention. On invalid expectation trials on the other hand, the expected

17

358    and presented orientations are orthogonal, and thus the orientation signal induced by the stimulus

359    would be positive and the signal induced by cue would be expected to be negative. Thus, subtracting the

360    orientation decoding signal on invalid trials from that on valid trials would subtract out the stimulus-

361    evoked signal while revealing any cue-induced orientation signal.

362

363    **Statistical testing.** Neural signals evoked by the different conditions were statistically tested using

364    nonparametric cluster-based permutation tests[61]. For ERF analyses, we averaged over the spatial (sensor)

365    dimension, on the basis of independent localisation of the 10 sensors that showed the strongest visual-

366    evoked activity during the grating localiser between 50 and 150 ms post-stimulus. Therefore, our

367    statistical analysis considered one-dimensional (temporal) clusters. For orientation decoding analyses,

368    the data consisted of two-dimensional (training time $\times$ testing time) decoding performance matrices,

369    and the statistical analysis thus considered two-dimensional clusters. For both one- and two-

370    dimensional data, univariate $t$-statistics were calculated for the entire matrix and neighbouring elements

371    that passed a threshold value corresponding to a $p$-value of 0.01 (two-tailed) were collected into

372    separate negative and positive clusters. Elements were considered neighbours if they were directly

373    adjacent, either cardinally or diagonally. Cluster-level test statistics consisted of the sum of $t$-values

374    within each cluster, and these were compared to a null distribution of test statistics created by drawing

375    10,000 random permutations of the observed data. A cluster was considered significant when its $p$-value

376    was below 0.05 (two-tailed).

18

## References

377 **References**

378

379   1. Helmholtz, H. von. *Treatise on physiological optics*. **III,** (The Optical Society of America, 1866).

380   2. Gregory, R. L. Knowledge in perception and illusion. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **352,** 1121–

381      1127 (1997).

382   3. Kersten, D., Mamassian, P. & Yuille, A. Object Perception as Bayesian Inference. *Annu. Rev. Psychol.*

383      **55,** 271–304 (2004).

384   4. Lee, T. S. & Mumford, D. Hierarchical Bayesian inference in the visual cortex. *J. Opt. Soc. Am. A* **20,**

385      1434 (2003).

386   5. Friston, K. A theory of cortical responses. *Philos. Trans. R. Soc. B Biol. Sci.* **360,** 815–836 (2005).

387   6. Summerfield, C. & de Lange, F. P. Expectation in perceptual decision making: neural and

388      computational mechanisms. *Nat. Rev. Neurosci.* **15,** 745–756 (2014).

389   7. Wyart, V., Nobre, A. C. & Summerfield, C. Dissociable prior influences of signal probability and

390      relevance on visual contrast sensitivity. *Proc. Natl. Acad. Sci.* **109,** 3593–3598 (2012).

391   8. SanMiguel, I., Widmann, A., Bendixen, A., Trujillo-Barreto, N. & Schroger, E. Hearing Silences: Human

392      Auditory Processing Relies on Preactivation of Sound-Specific Brain Activity Patterns. *J. Neurosci.* **33,**

393      8633–8639 (2013).

394   9. Kok, P., Failing, M. F. & de Lange, F. P. Prior Expectations Evoke Stimulus Templates in the Primary

395      Visual Cortex. *J. Cogn. Neurosci.* **26,** 1546–1554 (2014).

396   10. Bell, A. H., Summerfield, C., Morin, E. L., Malecek, N. J. & Ungerleider, L. G. Encoding of Stimulus

397      Probability in Macaque Inferior Temporal Cortex. *Curr. Biol.* **26,** 2280–2290 (2016).

398   11. Kok, P., Jehee, J. F. M. & de Lange, F. P. Less Is More: Expectation Sharpens Representations in the

399      Primary Visual Cortex. *Neuron* **75,** 265–270 (2012).

400    12. Rao, V., DeAngelis, G. C. & Snyder, L. H. Neural Correlates of Prior Expectations of Motion in the

401         Lateral Intraparietal and Middle Temporal Areas. *J. Neurosci.* **32,** 10063–10074 (2012).

402    13. Hindy, N. C., Ng, F. Y. & Turk-Browne, N. B. Linking pattern completion in the hippocampus to

403         predictive coding in visual cortex. *Nat. Neurosci.* **19,** 665–667 (2016).

404    14. Cichy, R. M., Pantazis, D. & Oliva, A. Resolving human object recognition in space and time. *Nat.*

405         *Neurosci.* **17,** 455–462 (2014).

406    15. King, J.-R. & Dehaene, S. Characterizing the dynamics of mental representations: the temporal

407         generalization method. *Trends Cogn. Sci.* **18,** 203–210 (2014).

408    16. Mostert, P., Kok, P. & de Lange, F. P. Dissociating sensory from decision processes in human

409         perceptual decision making. *Sci. Rep.* **5,** 18253 (2015).

410    17. Myers, N. E. *et al.* Testing sensory evidence against mnemonic templates. *eLife* **4,** e09000 (2015).

411    18. Brouwer, G. J. & Heeger, D. J. Decoding and Reconstructing Color from Responses in Human Visual

412         Cortex. *J. Neurosci.* **29,** 13992–14003 (2009).

413    19. Brouwer, G. J. & Heeger, D. J. Cross-orientation suppression in human visual cortex. *J. Neurophysiol.*

414         **106,** 2108–2119 (2011).

415    20. Kok, P., Van Lieshout, L. L. F. & De Lange, F. P. Local expectation violations result in global activity

416         gain in primary visual cortex. *Sci. Rep.* **6,** 37706 (2016).

417    21. Chalk, M., Seitz, A. R. & Series, P. Rapidly learned stimulus expectations alter perception of motion. *J.*

418         *Vis.* **10,** 2–2 (2010).

419    22. Kok, P., Brouwer, G. J., van Gerven, M. A. J. & de Lange, F. P. Prior Expectations Bias Sensory

420         Representations in Visual Cortex. *J. Neurosci.* **33,** 16275–16284 (2013).

421    23. Pajani, A., Kok, P., Kouider, S. & de Lange, F. P. Spontaneous Activity Patterns in Primary Visual

422         Cortex Predispose to Visual Hallucinations. *J. Neurosci.* **35,** 12947–12953 (2015).

423   24. St. John-Saaltink, E., Kok, P., Lau, H. C. & Lange, F. P. de. Serial Dependence in Perceptual Decisions

424        Is Reflected in Activity Patterns in Primary Visual Cortex. *J. Neurosci.* **36,** 6186–6192 (2016).

425   25. Wallenstein, G. V., Hasselmo, M. E. & Eichenbaum, H. The hippocampus as an associator of

426        discontiguous events. *Trends Neurosci.* **21,** 317–323 (1998).

427   26. Schapiro, A. C., Kustner, L. V. & Turk-Browne, N. B. Shaping of Object Representations in the Human

428        Medial Temporal Lobe Based on Temporal Regularities. *Curr. Biol.* **22,** 1622–1627 (2012).

429   27. Davachi, L. & DuBrow, S. How the hippocampus preserves order: the role of prediction and context.

430        *Trends Cogn. Sci.* **19,** 92–99 (2015).

431   28. Reddy, L. *et al.* Learning of anticipatory responses in single neurons of the human medial temporal

432        lobe. *Nat. Commun.* **6,** 8556 (2015).

433   29. Lavenex, P. & Amaral, D. G. Hippocampal-neocortical interaction: a hierarchy of associativity.

434        *Hippocampus* **10,** 420–430 (2000).

435   30. Bosch, S. E., Jehee, J. F. M., Fernandez, G. & Doeller, C. F. Reinstatement of Associative Memories in

436        Early Visual Cortex Is Signaled by the Hippocampus. *J. Neurosci.* **34,** 7493–7500 (2014).

437   31. Stokes, M., Thompson, R., Nobre, A. C. & Duncan, J. Shape-specific preparatory activity mediates

438        attention to targets in human visual cortex. *Proc. Natl. Acad. Sci.* **106,** 19569–19574 (2009).

439   32. Stokes, M., Thompson, R., Cusack, R. & Duncan, J. Top-Down Activation of Shape-Specific Population

440        Codes in Visual Cortex during Mental Imagery. *J. Neurosci.* **29,** 1565–1572 (2009).

441   33. Lee, S.-H., Kravitz, D. J. & Baker, C. I. Disentangling visual imagery and perception of real-world

442        objects. *NeuroImage* **59,** 4064–4073 (2012).

443   34. Albers, A. M., Kok, P., Toni, I., Dijkerman, H. C. & de Lange, F. P. Shared Representations for Working

444        Memory and Mental Imagery in Early Visual Cortex. *Curr. Biol.* **23,** 1427–1431 (2013).

445   35. Harrison, S. A. & Tong, F. Decoding reveals the contents of visual working memory in early visual

446        areas. *Nature* **458,** 632–635 (2009).

447    36. Serences, J. T., Ester, E. F., Vogel, E. K. & Awh, E. Stimulus-Specific Delay Activity in Human Primary

448        Visual Cortex. *Psychol. Sci.* **20,** 207–214 (2009).

449    37. den Ouden, H. E. M., Friston, K. J., Daw, N. D., McIntosh, A. R. & Stephan, K. E. A Dual Role for

450        Prediction Error in Associative Learning. *Cereb. Cortex* **19,** 1175–1185 (2009).

451    38. Näätänen, R. The role of attention in auditory information processing as revealed by event-related

452        potentials and other brain measures of cognitive function. *Behav. Brain Sci.* **13,** 201–233 (1990).

453    39. Nakano, T., Homae, F., Watanabe, H. & Taga, G. Anticipatory Cortical Activation Precedes Auditory

454        Events in Sleeping Infants. *PLoS ONE* **3,** e3912 (2008).

455    40. Bekinschtein, T. A. *et al.* Neural signature of the conscious processing of auditory regularities. *Proc.*

456        *Natl. Acad. Sci.* **106,** 1672–1677 (2009).

457    41. Summerfield, C. & Egner, T. Expectation (and attention) in visual cognition. *Trends Cogn. Sci.* **13,**

458        403–409 (2009).

459    42. Summerfield, C. & Egner, T. Feature-Based Attention and Feature-Based Expectation. *Trends Cogn.*

460        *Sci.* **20,** 401–404 (2016).

461    43. Kok, P., Bains, L. J., van Mourik, T., Norris, D. G. & de Lange, F. P. Selective Activation of the Deep

462        Layers of the Human Primary Visual Cortex by Top-Down Feedback. *Curr. Biol.* **26,** 371–376 (2016).

463    44. Summerfield, C., Trittschuh, E. H., Monti, J. M., Mesulam, M.-M. & Egner, T. Neural repetition

464        suppression reflects fulfilled perceptual expectations. *Nat. Neurosci.* **11,** 1004–1006 (2008).

465    45. Alink, A., Schwiedrzik, C. M., Kohler, A., Singer, W. & Muckli, L. Stimulus Predictability Reduces

466        Responses in Primary Visual Cortex. *J. Neurosci.* **30,** 2960–2966 (2010).

467    46. Meyer, T. & Olson, C. R. Statistical learning of visual transitions in monkey inferotemporal cortex.

468        *Proc. Natl. Acad. Sci.* **108,** 19401–19406 (2011).

469    47. Todorovic, A., van Ede, F., Maris, E. & de Lange, F. P. Prior Expectation Mediates Neural Adaptation

470        to Repeated Sounds in the Auditory Cortex: An MEG Study. *J. Neurosci.* **31,** 9118–9123 (2011).

471  48. Wacongne, C. *et al.* Evidence for a hierarchy of predictions and prediction errors in human cortex.

472      *Proc. Natl. Acad. Sci.* **108,** 20754–20759 (2011).

473  49. Cashdollar, N., Ruhnau, P., Weisz, N. & Hasson, U. The Role of Working Memory in the Probabilistic

474      Inference of Future Sensory Events. *Cereb. Cortex* bhw138 (2016). doi:10.1093/cercor/bhw138

475  50. Berkes, P., Orban, G., Lengyel, M. & Fiser, J. Spontaneous Cortical Activity Reveals Hallmarks of an

476      Optimal Internal Model of the Environment. *Science* **331,** 83–87 (2011).

477  51. Fiser, A. *et al.* Experience-dependent spatial expectations in mouse visual cortex. *Nat. Neurosci.* **19,**

478      1658–1664 (2016).

479  52. Rao, R. P. & Ballard, D. H. Predictive coding in the visual cortex: a functional interpretation of some

480      extra-classical receptive-field effects. *Nat. Neurosci.* **2,** 79–87 (1999).

481  53. Bar, M. *et al.* Top-down facilitation of visual recognition. *Proc. Natl. Acad. Sci. U. S. A.* **103,** 449–454

482      (2006).

483  54. Brainard, D. H. The Psychophysics Toolbox. *Spat. Vis.* **10,** 433–436 (1997).

484  55. Watson, A. B. & Pelli, D. G. Quest: A Bayesian adaptive psychometric method. *Percept. Psychophys.*

485      **33,** 113–120 (1983).

486  56. Stolk, A., Todorovic, A., Schoffelen, J.-M. & Oostenveld, R. Online and offline tools for head

487      movement compensation in MEG. *NeuroImage* **68,** 39–48 (2013).

488  57. Oostenveld, R., Fries, P., Maris, E. & Schoffelen, J.-M. FieldTrip: Open Source Software for Advanced

489      Analysis of MEG, EEG, and Invasive Electrophysiological Data. *Comput Intell Neurosci.* **2011,** 1:1–1:9

490      (2011).

491  58. Bastiaansen, M. C. M. & Knösche, T. R. Tangential derivative mapping of axial MEG applied to event-

492      related desynchronization research. *Clin. Neurophysiol.* **111,** 1300–1305 (2000).

493  59. Garcia, J. O., Srinivasan, R. & Serences, J. T. Near-Real-Time Feature-Selective Modulations in

494      Human Cortex. *Curr. Biol.* **23,** 515–522 (2013).

495    60. Blankertz, B., Lemm, S., Treder, M., Haufe, S. & Müller, K.-R. Single-trial analysis and classification of

496       ERP components — A tutorial. *NeuroImage* **56,** 814–825 (2011).

497    61. Maris, E. & Oostenveld, R. Nonparametric statistical testing of EEG- and MEG-data. *J. Neurosci.*

498       *Methods* **164,** 177–190 (2007).

499

500

507     **Author contributions**

508     P.K. and F.d.L. designed the experiment, P.K. and P.M. conducted the experiment, P.K. and P.M.

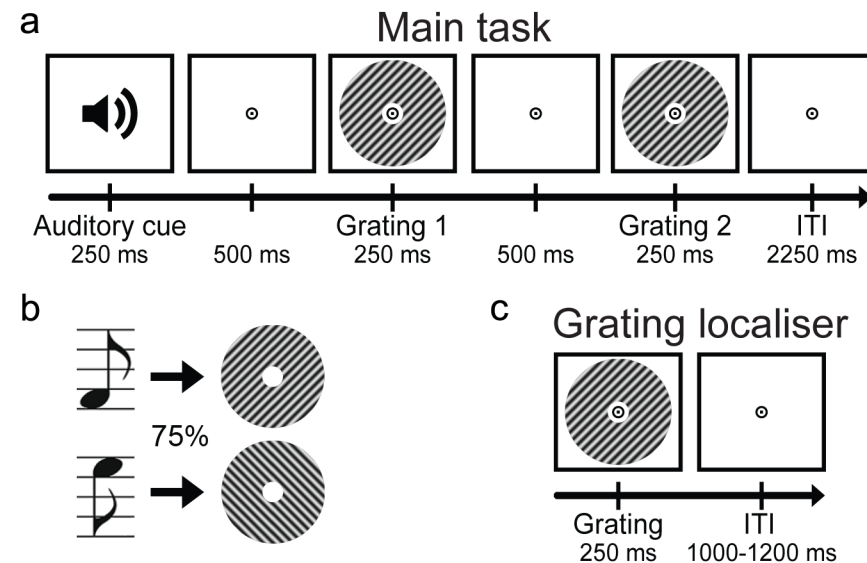509     analysed the data, P.K., P.M. and F.d.L. wrote the paper.

510

511     **Competing financial interests**

512     The authors declare no competing financial interests.

513

514 **Figures**

515



516
517 **Figure 1. Experimental paradigm.** (**a**) Each trial started with an auditory cue that predicted the

518 orientation of the subsequent grating stimulus. This first grating was followed by a second one, which

519 differed slightly from the first in terms of orientation and contrast. In separate runs, participants
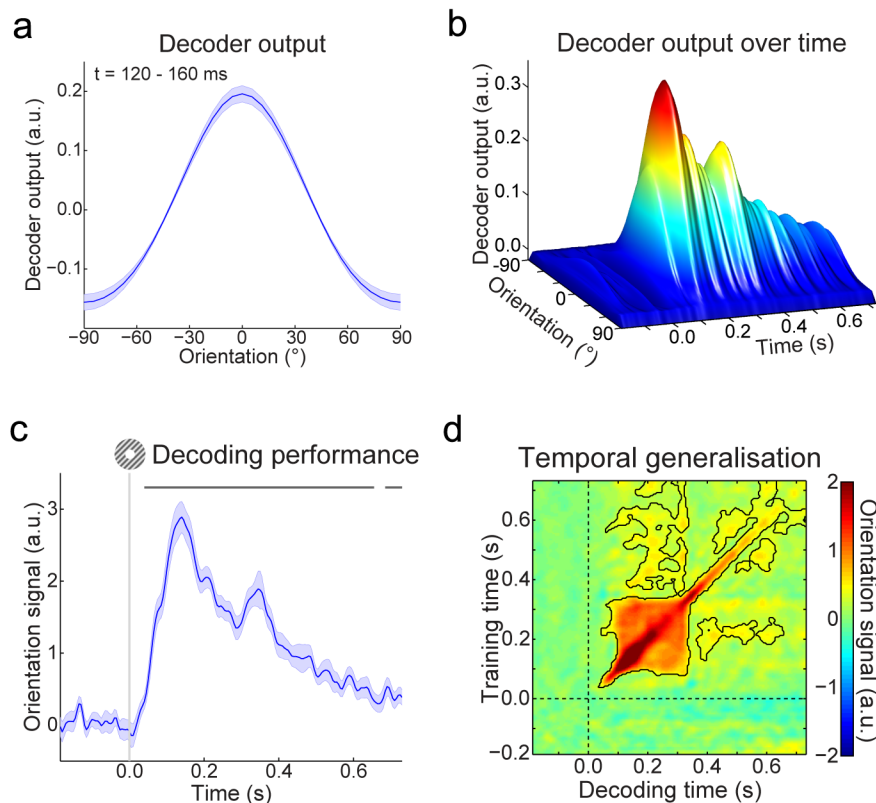
520 performed either an orientation or contrast discrimination task on the two gratings. (**b**) Throughout the

521 experiment, two different tones were used as cues, each one predicting one of the two possible
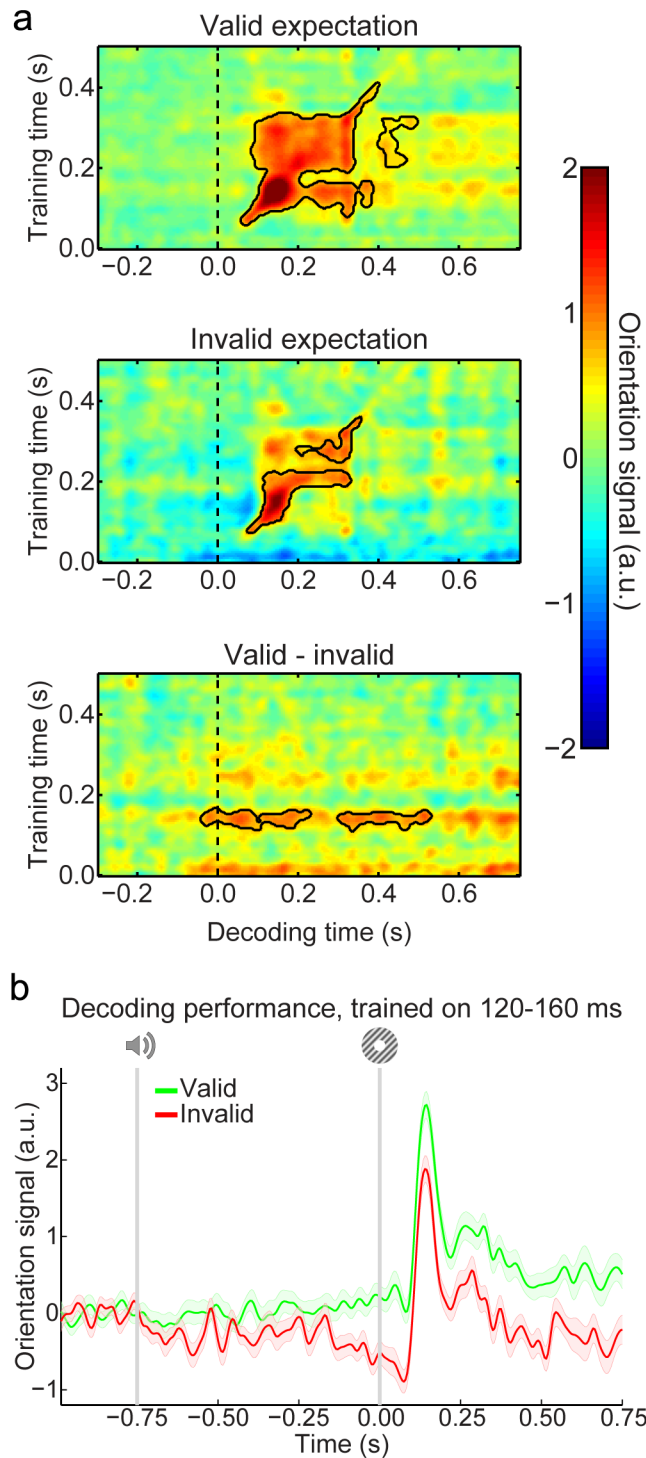
522 orientations (45° or 135°) with 75% validity. These contingencies were flipped halfway through the

523 experiment. (**c**) In separate grating localiser runs, participants were exposed to task-irrelevant gratings

524 while they performed a fixation dot dimming task.

525

526
527 **Figure 2. Localiser orientation decoding.** (**a**) The output of the decoder consisted of the responses of 32

528 hypothetical orientation channels, shown here decoders trained and tested on the MEG signal 120-160

529 ms post-stimulus during the grating localiser (cross-validated). Shaded region represent SEM. (**b**)

530 Decoder output over time, trained and tested in 5 ms steps (sliding window of 29.2 ms), showing the

531 temporal evolution of the orientation signal. (**c**) The response of the 32 orientation channels collapsed

532 into a single metric of decoding performance (see Methods), over time. Shaded region represent SEM,

533 horizontal lines indicate significant clusters ($p < 0.05$). (**d**) Temporal generalisation matrix of orientation

534 decoding performance, obtained by training decoders on each time point, and testing all decoders on all

535 time points (as above, steps of 5 ms and a sliding window of 29.2 ms). This method provides insight into

536 the sustained versus dynamical nature of orientation representations[15]. Solid black lines indicate

537 significant clusters ($p < 0.05$), dashed lines indicate grating onset (t = 0s).

538

**a**



**b**

**Figure 3. Expectation induces stimulus templates.** (**a**) Temporal generalisation matrices of orientation decoding during the main experiment. Decoders were trained on the grating localiser (training time on the y-axis) and tested on the main experiment (time on the x-axis; dashed vertical line indicates t = 0s, onset of the first grating). Decoding shown separately for gratings preceded by a valid expectation (top row), invalid expectation (middle row), and the subtraction of the two conditions (i.e., the expectation cue effect, bottom row). Solid black lines indicate significant clusters ($p < 0.05$). (**b**) Orientation decoding during the main task, averaged over training time $120 - 160$ ms post-stimulus during the grating localiser. That is, a horizontal slice through the temporal generalisation matrices above at the training time for which we see a significant cluster of expected orientation decoding, for visualisation. Shaded regions indicate SEM.