

A retinotopic reference frame for space throughout human visual cortex

**Martin Szinte^{1,2,*}, Gilles de Hollander^{2,3*}, Marco Aqil^{2,4,5*},
Inês Veríssimo^{2,5}, Serge Dumoulin^{2,4,5,6}, Tomas Knapen^{2,4,5}**

1. Institut de Neurosciences de la Timone, CNRS & Aix-Marseille Université, UMR 7289, Marseille, France.

2. Spinoza Centre for Neuroimaging, Royal Dutch Academy of Sciences, Amsterdam, Netherlands.

3. Zurich Center for Neuroeconomics, Department of Economics, University of Zurich, Switzerland

4. Computational Cognitive Neuroscience and Neuroimaging, Netherlands Institute for Neuroscience, Amsterdam, Netherlands

5. Department of Cognitive Psychology, Vrije Universiteit Amsterdam, Amsterdam, Netherlands.

6. Department of Experimental Psychology, Utrecht Universiteit, Utrecht, Netherlands.

* First-author equal contribution

Summary

We perceive a stable visual world across eye movements, despite the drastic retinal transients these movements produce. To explain vision's spatial stability, it has been suggested that the brain encodes the location of attended visual stimuli in an external, or spatiotopic, reference frame. However, spatiotopy is seemingly at odds with the fundamental retinotopic organization of visual inputs. Here, we probe the spatial reference frame of vision using ultra-high-field (7T) fMRI and single-voxel population receptive field mapping, while independently manipulating both gaze direction and spatial attention. To manipulate spatial attention, participants performed an equally demanding visual task on either a bar stimulus that traversed the visual field, or a small foveated stimulus. To dissociate retinal stimulus position from its real-world position the entire stimulus array was placed at one of three distinct horizontal screen positions in each run. We found that population receptive fields in all cortical visual field maps shift with the gaze, irrespective of how spatial attention is deployed. This pattern of results is consistent with a fully retinotopic reference frame for visual-spatial processing. Reasoning that a spatiotopic reference frame could conceivably be computed at the level of entire visual areas rather than at the level of individual voxels, we also used Bayesian decoding of stimulus location from the BOLD response patterns in visual areas. We found that decoded stimulus locations also adhere to the retinotopic frame of reference, by shifting with gaze position. Again, this result holds for all visual areas and irrespective of the deployment of spatial attention. We conclude that visual locations are encoded in a retinotopic reference frame throughout the visual hierarchy.

Keywords

Retinotopy, population receptive field, Bayesian decoding, ultra-high-field fMRI

1 Introduction

2 Eye movements rapidly alter the projection of visual objects on the retina. Nevertheless, our subjective
3 visual perception is stable and invariant to these eye movements. This “*space constancy*” phenomenon
4 presents researchers with a conundrum: how do our brains build a stable impression of the world based
5 on such fleeting, jumbled sensory inputs? One major step to understanding space constancy would be to
6 answer the question of how the information of visual objects’ location is represented by the visual system.

7 Visual location is first encoded by the retina¹, whose topography is maintained as spatial information
8 is transmitted to subsequent processing stage^{2,3}. Foundationally, visual locations are encoded in a common
9 “retinotopic” reference frame^{4,5}, keeping the retinal photoreceptors organization. Neural responses in this
10 retinotopic reference frame are modulated by gaze direction^{6,7}. This combination of action-related and
11 visual sources of information is thought to allow the brain to localize objects in the outside world⁸⁻¹⁰. But
12 the precise mechanisms underpinning world-centered localization, including the role of spatial attention,
13 remain a matter of debate¹¹⁻¹⁴. Specifically, an open fundamental question is whether gaze modulation of
14 visual information leads to the encoding of locations of visual stimuli in a spatiotopic frame of reference.

15 Here, we leveraged population receptive field (pRF) estimation¹⁵ and ultra-high-field (7T) whole brain
16 fMRI to probe the spatial reference frame of vision at the level of local neural populations sampled by
17 individual voxels. We determined voxels’ pRF position while participants directed their gaze to three
18 distinct locations, allowing us to test whether pRFs are pinioned to the outside world (the *spatiotopic*
19 hypothesis), or fixed to the retina (the *retinotopic* hypothesis). To study the role of spatial attention,
20 participants performed an equally demanding orientation-discrimination task either at fixation or on a
21 moving visual-mapping stimulus. We found that throughout the cortical visual hierarchy, pRFs
22 systematically shift with gaze, as being fixed to the retina, irrespective of where spatial attention is
23 deployed. Next, to probe the visual reference frame at the level of entire visual field maps, we adapted a
24 Bayesian decoding method¹⁶⁻¹⁸, to infer the spatial properties of our visual stimuli from the pattern of
25 voxel responses in entire visual areas. In line with our single-voxel results, we found that decoded stimulus
26 locations are fixed to the retina, again favoring the retinotopic hypothesis. Altogether our results suggest
27 that location is represented throughout the visual hierarchy by populations of receptive fields operating
28 in a retinotopic reference frame.

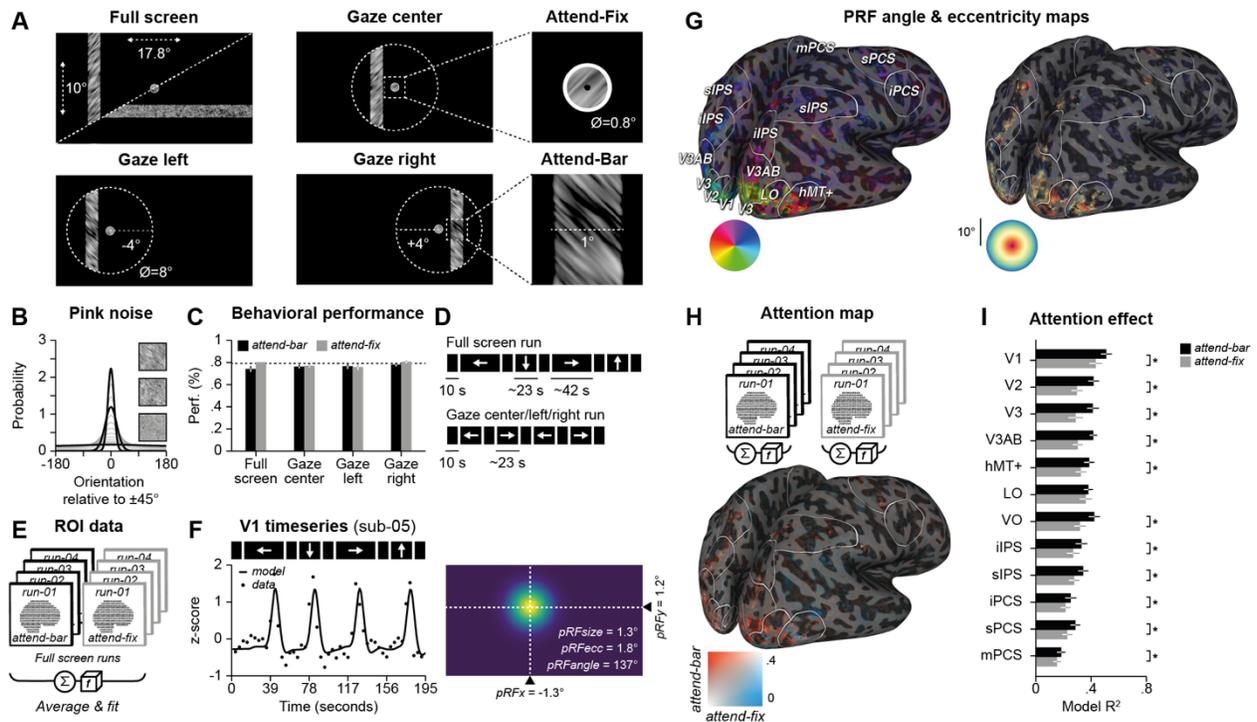


Figure 1. Methods, pRF modeling and attention effect. **A.** Participants fixated a bull's-eye in either *full screen* (top left panel), *gaze center* (top center), *gaze left* (bottom left) and *gaze right* runs (center left). They reported the orientation of pink noise patterns ($+45^\circ$ or -45°) presented either within the bar (*attend-bar*) or within the fixation bull's-eye (*attend-fix*). **B.** We titrated the difficulty of the orientation discrimination task by varying the dispersion coefficient of the orientation filter applied to the pink noise pattern. The figure presents probability distributions of orientations contained within the pink noise patterns as a function of the dispersion coefficient. Right insets, showcase three levels of difficulty from non-oriented random noise (bottom) to the most oriented pattern used (top). **C.** Group performance ($n = 8$, two sessions each) obtained in the *attend-bar* (black) and *attend-fix* (gray) conditions. Error bars show \pm SEM and the dashed line shows the staircase convergence level ($\sim 79\%$ correct). **D.** In the *full screen* and *gaze center/left/right* conditions, a bar moved in different directions (black rectangles with arrows) interleaved with periods in which only the fixation bull's-eye was shown (empty black rectangles). **E.** To determine regions of interest (ROIs), we averaged *full screen* runs and fit a linear pRF model. **F.** Example V1 timeseries and its best explained pRF model and parameters. **G.** Single participant (*sub-005*) pRF angle and eccentricity maps projected on his inflated brain. Lines on the cortex delineate ROIs. Participants brains visualizations are available online (invibe.nohost.me/gazepRF/). **H.** Map of the effect of spatial attention obtained by comparing the explained variance of *attend-bar* and *attend-fix* conditions separately. Reddish colors of the inflated cortex illustrate the effect of spatial attention. **I.** Model explained variance (R^2) observed for the best-fitting voxels of each ROI in the *attend-bar* (black) and *attend-fix* (gray) conditions. Error bars show \pm SEM, asterisks show significant difference between conditions (two-sided $p < 0.05$).

1 Results

2 PRF modeling and attention effect

3 Participants were instructed to continuously fixate a bull's-eye on a black background while horizontal and
 4 vertical bars swiped across the full extent of the screen (Fig. 1A, *full screen*). They reported the orientation
 5 of noise patterns presented either within the moving bar (*attend-bar* condition) or within the fixation
 6 bull's-eye (*attend-fix* condition). We titrated the difficulty of the task with a staircase procedure which
 7 adjusted the orientation dispersion coefficient of pink noise patterns inside the presented stimuli (Fig. 1B).
 8 Participants maintained the expected staircase performance level (Fig. 1C, *full screen, attend-bar*: $74.06\% \pm 2.50$;
 9 *full screen, attend-fix*: $80.26\% \pm 0.40$), confirming their ability to attend to the noise patterns
 10 presented in the periphery or at fixation with equal performance. Apart from the attention task, the
 11 stimulus presentation sequence of *full screen* followed an identical sequence (Fig. 1D, top), allowing us to
 12 average signal timeseries across the attention tasks (Fig. 1E) into a single, high-SNR *fiducial* timeseries.

1 From each voxel's fiducial timeseries, we quantified the screen locations at which a stimulus evokes a
 2 visual response by estimating its population receptive field properties. Specifically, we fit a linear, isotropic
 3 Gaussian pRF model to find optimal location (pRF_x and pRF_y) and size (pRF_{size}) parameters that quantify
 4 the spatial tuning of the neural population within each voxel. Figure 1F shows an example of a V1 voxel
 5 timeseries and its best-fitting pRF model. From the position parameters, we derived pRF eccentricity and
 6 pRF angle maps that can be projected onto individual participants' cortical surfaces (Fig. 1G) to delineate
 7 visual field maps as regions of interests (see Methods). We next repeated this analysis using timeseries of
 8 *attend-bar* and *attend-fix* runs averaged separately (Fig. 1H, top), using the fiducial outcomes as starting
 9 parameters. Figure 1H show an inflated cortex of a participant displaying the change in explained variance
 10 (R^2) resulting from manipulating spatial attention. Using the 250 best-fitting voxels per ROI, we found a
 11 significant improvement of the explained variance when comparing the *attend-bar* to the *attend-fix*
 12 conditions for all ROIs (*attend-bar*: $0.51 > R^2 > 0.18$ vs. *attend-fix*: $0.43 > R^2 > 0.15$, $0.0117 > p > 0.0001$,
 13 two-sided p values) except area LO (*attend-fix*: $R^2: 0.38 \pm 0.03$ vs. *attend-bar* $R^2: 0.36 \pm 0.04$, two-sided $p =$
 14 0.24). This pattern of result was similar when including all voxels contained within the ROIs. Taken
 15 together, behavioral and functional imaging results illustrate classical effects of visual spatial attention: an
 16 improvement of orientation discrimination abilities^{19,20} correlated with an improvement of the
 17 corresponding fMRI BOLD signal-to-noise-ratio^{21,22}.

18

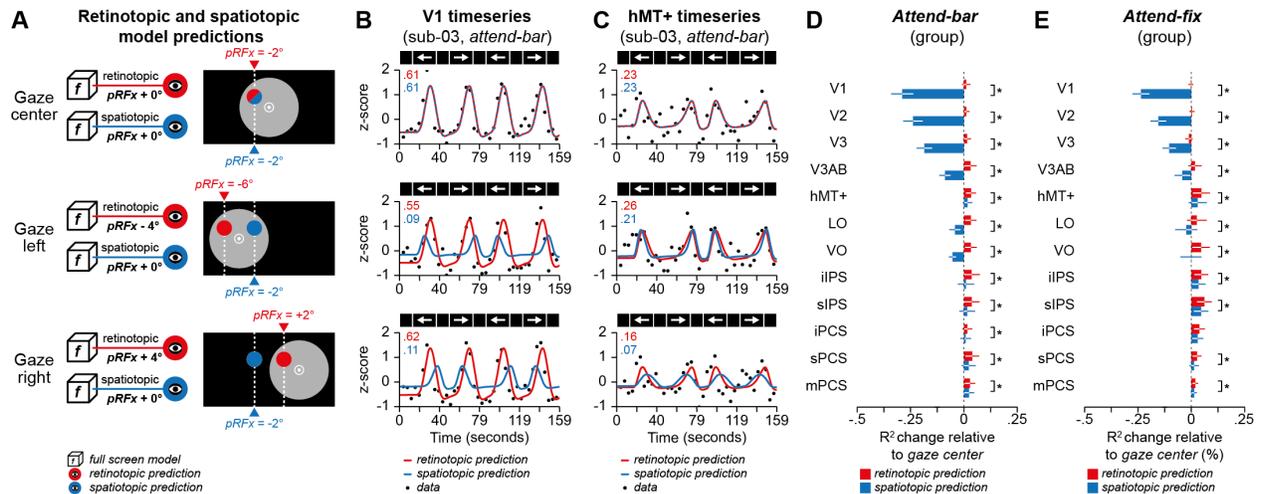


Figure 2. Out-of-sample retinotopic and spatiotopic predictions. **A.** Retinotopic predictions are obtained by adding the change in gaze direction to the *full screen* pRF_x parameters (e.g., $pRF_{x\text{gaze left}} = pRF_{x\text{full screen}} - 4^\circ$). The spatiotopic hypothesis dictates that the pRF_x should stay identical irrespective of the gaze direction (i.e., $pRF_{x\text{gaze left}} = pRF_{x\text{gaze right}} = pRF_{x\text{gaze center}} = pRF_{x\text{full screen}}$). **B-C.** Example retinotopic (red) and spatiotopic (blue) timeseries predictions for V1 (B) and hMT+ voxels (C) in the *gaze center* (top row), *gaze left* (middle row) and *gaze right attend-bar* runs of one participant (sub-003). Leftmost inset values show corresponding model R^2 . **D-E.** Retinotopic (red) and spatiotopic predictions (blue) change in explained variance between *gaze left/right* and *gaze center* conditions across ROIs. Note that only spatiotopic prediction of V1, V2 and V3 voxels display a significant change as compared to zero in both attention conditions (V1/V2/V3: $0.0156 > p > 0.0001$; two-sided p values), as well as V3AB (two-sided $p = 0.0156$) and VO in the attend-bar condition (two-sided $p = 0.0078$). Error bars show \pm SEM, asterisks show significant difference between retinotopic and spatiotopic predictions (two-sided $p < 0.05$).

19 **Out-of-sample predictions show retinotopic encoding of visual space**

20 We next aimed to use the pRF estimates from the *full-screen* condition to probe the reference frame of
 21 spatial vision. Across runs, participants were instructed to fixate their gaze either at the screen center (Fig.
 22 1A, *gaze center*), to the left (Fig. 1A, *gaze left*) or to the right of it (Fig. 1A, *gaze right*). The entire mapping
 23 stimulus sequence was centered on these fixation locations, and was restricted to vertical bars moving
 24 rightward or leftwards, vignettted by a circular aperture (Fig. 1A and Fig. 1D). Again, participants reported

1 the orientation of noise patterns presented within the bar (*attend-bar*) or within the fixation bull's-eye
 2 (*attend-fix*). They reached the staircase expected performance level in the *gaze center* (Fig. 1C, *attend-*
 3 *bar*: 76.27% ± 1.67; *attend-fix*: 76.95% ± 1.48), the *gaze left* (Fig. 1C, *attend-bar*: 76.59% ± 2.30; *attend-fix*:
 4 75.88% ± 2.49) and *gaze right* conditions (Fig. 1C, *attend-bar*: 78.09% ± 1.66; *attend-fix*: 80.61% ± 1.25).
 5 Using the *full screen* experiment's pRF parameters, we calculated predicted timeseries for the *gaze center*,
 6 *gaze left* and *gaze right* conditions in the retinotopic and spatiotopic reference frame. Specifically, to
 7 obtain retinotopic predictions we added the gaze direction change to the pRFx parameter (x-coordinate
 8 of the pRF) obtained in the corresponding *full screen* conditions (Fig. 2B, *retinotopic prediction*). The
 9 spatiotopic hypothesis supposed that the pRFx do not vary as a function of the gaze directions (Fig. 2A,
 10 *spatiotopic prediction*). We thus kept the pRFx parameter as observed in the corresponding *full screen*
 11 conditions. Figure 2B-C illustrate the predicted and measured timeseries of V1 and hMT+ voxels. We
 12 highlighted these regions to relate them to former studies reporting spatiotopic effects^{11,12}.
 13

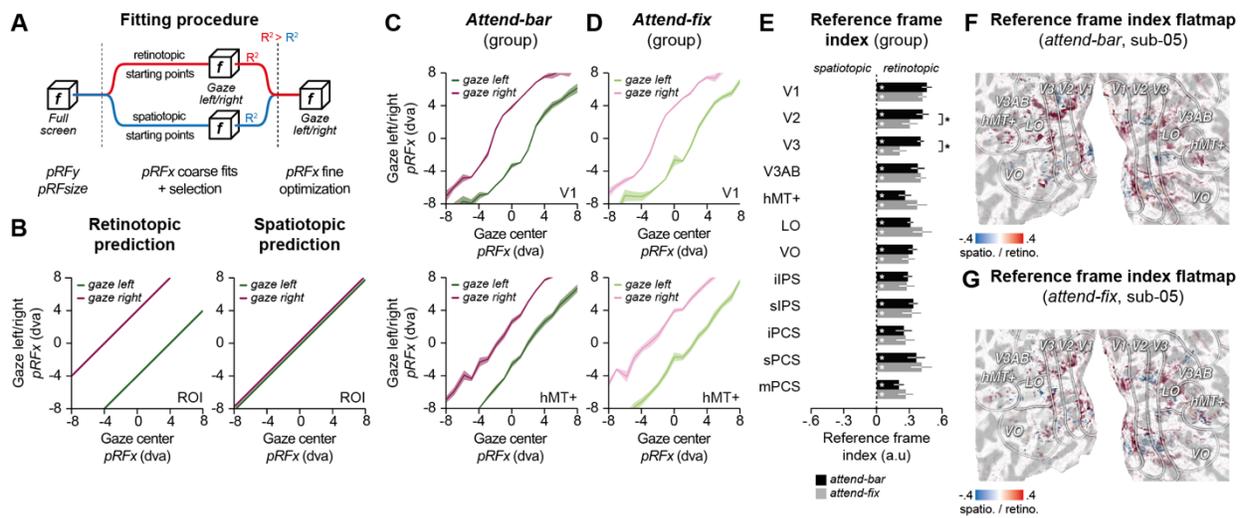


Figure 3. Fitting retinotopic and spatiotopic pRF models. **A.** Fitting procedure. We used a coarse-to-fine optimization fitting procedure in which we kept the pRFy and pRFsize parameters from the corresponding *full screen* runs. In order to avoid bias between the spatiotopic and retinotopic model, we first used two different sets of parameter starting points. These starting points were based on either the retinotopic (red) or the spatiotopic (blue) hypotheses. For a next optimization stage, we selected the parameters producing the highest fit quality. **B.** Retinotopic and spatiotopic predictions. The gaze center pRFx as a function of the gaze left/right pRFx will either shift (left panel, retinotopic prediction) or remain at the same position (right panel, spatiotopic prediction) when comparing the gaze center with the gaze left (green) or gaze right conditions (purple). **C-D.** Each panel shows the group average of 16 equal bins of the pRFx obtained in the gaze center runs as a function of the corresponding group averaged bins of the pRFx obtained in the gaze left (purple) and gaze right runs (green) for the *attend-bar* (C, dark colors) and *attend-fix* conditions (D, light colors). Error areas show ±SEM. **E.** Group average reference frame index (RFI) observed in the *attend-bar* (black) and *attend-fix* conditions (gray) for the best-fitting voxels of each ROI. Error bars show ±SEM, white asterisks show significant RFI as compared to a null effect (two-sided $p < 0.05$). Black asterisks indicate a significant difference between the *attend-bar* and *attend-fix* conditions (two-sided $p < 0.05$). **F-G.** Reference frame index maps obtained by projecting the RFI obtained in the *attend-bar* (F) and *attend-fix* conditions (G) on inflated cortex using a unimodal color scale (blue: spatiotopic vs. red: retinotopic) for a participant (sub-05, same as in Fig. 1).

14 Importantly, the gaze center condition can serve as a baseline: both the spatiotopic and retinotopic
 15 predictions necessarily produce identical pRFx parameters. We computed the change in the pRF model
 16 explained variance (R^2) between the gaze left and gaze right runs as compared to the gaze center runs. We
 17 found, again using the best 250 voxels across ROIs, significantly lower explained variance for the
 18 spatiotopic relative to the retinotopic prediction in both the *attend-bar* (Fig. 2D, spatiotopic prediction:

1 +2.36% > R^2 change > -28.91% vs. retinotopic prediction: +3.96 > R^2 change > +1.09%, 0.0195 > p s > 0.0001,
2 two-sided p values) and attend-fix conditions (Fig. 2E, spatiotopic prediction: +4.61% > R^2 change > -
3 23.59% vs. retinotopic prediction: +6.14 > R^2 change > -1.12%, 0.0352 > p s > 0.0001, two-sided p values).
4 Similar effects were found when including all voxels contained within the ROIs.

5 6 **Fitting pRF s positions shows retinotopic encoding of visual space**

7 The results presented above suggest that the visual system encodes location in a retinotopic reference
8 frame. However, our first analysis did not allow the pRF_x parameter to vary in between the retinotopic
9 and spatiotopic predictions, which could have provided evidence for spatiotopic representations of visual
10 space. Thus, we fit the *horizontal pRF position* parameter in the gaze conditions while keeping the other
11 pRF parameters constant (pRF_y and pRF_{size}). To avoid any bias due to the starting points of the fitting
12 procedure, we initiated parameters from distinct *horizontal pRF position* values referenced against either
13 the gaze position (Fig. 3A, retinotopic starting points) or the screen center (Fig 3A, spatiotopic starting
14 points). In a subsequent stage, we optimized the *horizontal pRF position* parameter starting from the
15 highest explained variance model (R^2) of the first stage. We formulated two distinct predictions based on
16 our two competing hypotheses. We used the *gaze center* condition as a reference for both hypotheses,
17 since they produce identical predictions. In the retinotopic reference frame hypothesis, *horizontal pRF*
18 *position* in the *gaze left* and *gaze right* conditions shifts based on gaze direction (Fig. 3B, retinotopic
19 prediction). Conversely, in the spatiotopic reference frame hypothesis, *horizontal pRF position* shouldn't
20 differ from that observed when participants look straight ahead (Fig. 3B, spatiotopic prediction). Figure 3C
21 and Figure 3D show these comparisons across participants for best-fitting voxels of V1 and hMT+ using
22 data modeled from the *attend-bar* and *attend-fix* runs respectively. In accordance with previous studies^{11–}
23 ¹³, results suggest that V1 uses a retinotopic frame of reference irrespective of where participant deploy
24 their spatial attention (Figure 3C-D, top). Our results suggest that hMT+ also uses a retinotopic framework,
25 regardless of spatial attention's focus (Figure 3C-D, bottom).

26 To evaluate quantitatively which reference frame best explains our results at the level of individual
27 voxels across all ROIs, we computed a reference frame index (RFI, see Methods), computed as in previous
28 studies^{11,13}. Briefly, it quantifies the difference in explained variance between the two model classes on a
29 scale between -1 and 1, where pure spatiotopy corresponds to -1, pure retinotopy to +1, and noise results
30 in an RFI value of 0. Figure 3E shows RFI obtained with the best-fitting voxels of each ROI for the *attend-*
31 *bar* and *attend-fix* conditions. Across participants, RFI were significantly positive for all ROIs in the *attend-*
32 *bar* (0.46 > RFI > 0.20, 0.0312 > p s > 0.0001, two-sided p values) and *attend-fix* conditions (0.42 > RFI >
33 0.21, 0.0312 > p s > 0.0001, two-sided p values). Moreover, RFI values did not differ between conditions in
34 all ROIs (*attend-bar*: 0.93 > p s > 0.38, two-sided p values), apart from V2 (*attend-bar* RFI: 0.42 ± 0.06 vs.
35 *attend-fix* RFI: 0.30 ± 0.08, two-sided p < 0.01) and V3 (*attend-bar* RFI: 0.40 ± 0.03 vs. *attend-fix* RFI: 0.21 ±
36 0.07, two-sided p < 0.01). Similar statistics were observed when including all voxels. This analysis suggests
37 that at the level of individual voxels, visual location is encoded in a retinotopic reference frame irrespective
38 of the deployment of attention. Figure 3F and 3G illustrate these effects by projecting the RFI onto the
39 cortical surface of a single participant when he attended stimuli presented within the bar (Fig. 3F) or within
40 the fixation bull's-eye (Fig. 3G).

41 42 **Bayesian decoding of visual field positions shows retinotopic encoding of visual space**

43 Even if local populations of neurons sampled by single voxels do not show evidence of spatiotopic location
44 coding, more dispersed effects across the entire visual field map might conceivably provide the means for
45 the brain to represent world-centered coordinates. Such a mechanism might only be evident when
46 inferring the locations of visual stimuli from the pattern of BOLD responses across an entire visual region,
47 while taking into account the covariance structure between voxels. Thus, we next aimed at decoding, for

1 each individual ROI, the position of the moving bar for different gaze conditions. We adapted a Bayesian
2 decoding method developed for estimating, from a pattern of BOLD responses, a full posterior distribution
3 along a single stimulus dimension, namely visual orientation¹⁶, to the multiple stimulus dimensions of
4 visual space. Our novel method decodes bar spatial locations on a TR-by-TR basis (see *Methods*), based on
5 pRF estimates and residual covariance structure from the *full screen* runs (Fig. 4A). Contrary to techniques
6 which decode a single stimulus value most consistent with the observed data, Bayesian decoding also
7 produces uncertainty estimates for the decoded features. This decoded uncertainty has been shown to
8 relate to sensory uncertainty and decision confidence^{16,23}. Our method determines, for every TR, the full
9 posterior probability along spatial bar location parameters based on the per-TR BOLD pattern in the *gaze*
10 *center*, *gaze left* or *gaze right* conditions (Fig. 4B). Because this posterior distribution lives in visual space
11 it can be compared to the ground truth of what was actually shown on the screen to participants.

12 Crucially, we can define distinct spatiotopic and retinotopic predictions for the reference frame in
13 which the decoded bar position is hypothesized to reside. As our decoder was trained on the *full screen*
14 dataset in which participants fixated straight ahead, the spatiotopic hypothesis predicts that our decoder
15 will sense the change of gaze as a shift of the bar position in the opposite direction. Conversely, the
16 retinotopic hypothesis predicts that the decoded position will be centered on the coordinates of the *full*
17 *screen* stimulus, irrespective of gaze direction. Figure 4C shows the decoded position of the bar (i.e., the
18 peak of the bar position distribution) as observed when considering V1 voxels from a single run of the *gaze*
19 *center*, the *gaze left* and the *gaze right* conditions. Figure 4D-E shows for a participant the same effect
20 averaged across runs and bar passes, for V1 (Fig. 4D) and hMT+ voxels respectively (Fig. 4E). To quantify
21 decoding performance, we computed the correlation between the decoded position and the retinotopic
22 ground truth of the moving bar (Fig. 4F). Correlations in the *attend-bar* condition were significantly positive
23 for all ROIs ($0.77 > r > 0.13$, $0.0391 > ps > 0.0001$, two-sided p values) with the strongest correlations for
24 low-level visual areas (e.g., V1, $r = 0.73 \pm 0.04$, $p < 0.01$; hMT+, $r = 0.70 \pm 0.05$, $p < 0.0001$, two-sided p
25 values). Correlations were significantly positive in the *attend-fix* condition across all ROIs ($0.67 > r > 0.17$,
26 $0.0078 > ps > 0.0001$, two-sided p values), except for sIPS ($r = 0.02 \pm 0.03$, two-sided $p = 0.49$) and iPCS (r
27 $= 0.07 \pm 0.05$, two-sided $p = 0.21$). Furthermore, when participants directed spatial attention to the moving
28 barn, correlations were significantly increased when compared to *attend-fix* condition in all ROIs (*attend-*
29 *bar* vs. *attend-fix*, $0.0391 > ps > 0.0001$, two-sided p values), apart from V3 (two-sided $p = 0.09$), hMT+
30 (two-sided $p = 0.06$) and iPCS (two-sided $p = 0.21$).

31 To answer our central question concerning visual reference frames, we again computed a RFI based on
32 these decoding results (see *Methods*) for each participant and ROI (Fig. 4E, see *Methods*). Across
33 participants, RFIs were significantly positive across all ROIs, indicating a predominant retinotopic reference
34 frame in both the *attend-bar* ($0.71 > RFI > 0.51$, $0.0078 > ps > 0.0001$, two-sided p values) and *attend-fix*
35 conditions ($0.64 > RFI > 0.47$, $0.0078 > ps > 0.0001$, two-sided p values). Deployment of attention to the
36 bar stimulus increased the RFI for V1, V3AB, iIPS, sPCS and mPCS ($0.0391 > ps > 0.0078$, two-sided p values),
37 indicating that increased decoding fidelity increases the strength of retinotopic representations in both
38 low and high levels of the visual hierarchy. Lastly, we inspected whether the decoder's uncertainty perhaps
39 provides residual indications for any role of spatiotopic spatial coding of visual locations. Specifically, the
40 mismatch between train and test set due to differences in gaze location should increase uncertainty for
41 gaze left and gaze right conditions relative to the gaze center condition. We quantified the uncertainty of
42 the decoder by calculating the dispersion of the posterior along the horizontal position dimension. In the
43 intermediate layers of the visual system, specifically in areas like V3AB, hMT+, LO, VO, and iIPS, we
44 observed distinct outcomes of spatial attention. Specifically, when attention is directed towards the bar
45 stimulus, there is a noticeable reduction in decoder uncertainty, aligning with the previously explained
46 impacts on decoding fidelity (Supp. Fig. 1A). The impact of decoder uncertainty on attention was
47 significantly more pronounced than any influence of gaze direction (Supp Fig. 1B).

1 Altogether these effects indicate that visual areas represent visual locations using a retinotopic
 2 reference frame irrespective of the deployment of attention, also when taking into account the region-
 3 level covariance structure of BOLD responses.
 4

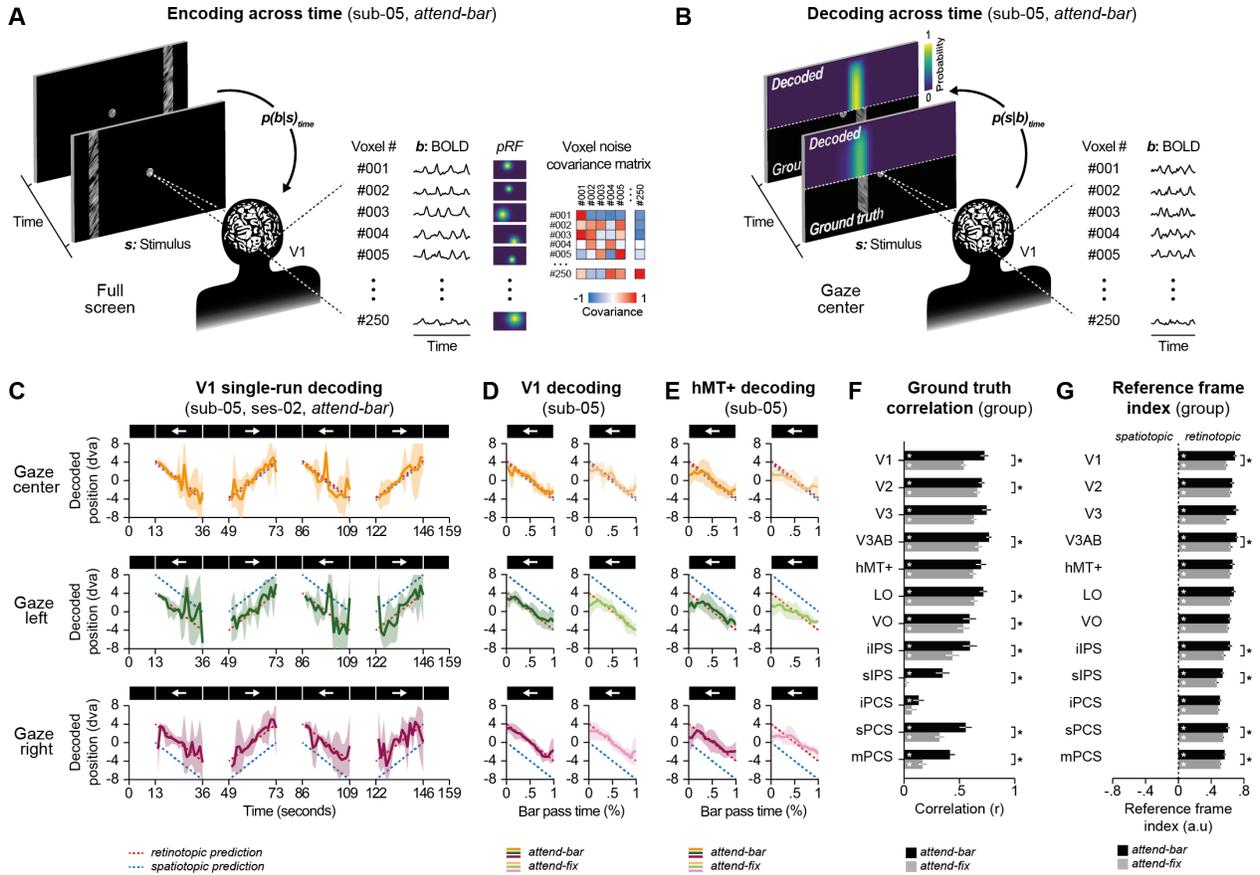


Figure 4. Bayesian decoding. **A.** Our encoding procedure consists of determining a likelihood $p(b|s)_{time}$ defined as the probability of observing fMRI BOLD signal (b) as a function of the stimulus (s) presented in the *full screen* condition. For this analysis we use the spatial tuning and noise correlations of the 250 best-fitting voxels per ROI. **B.** The decoding procedure illustrated here computes the posterior probability over time $p(s|b)_{time}$ using independent recorded BOLD signals for example from the *gaze center* condition. This procedure results in a direct comparison of decoded and actual bar position over time (see comparison of posterior distribution sample and ground truth). **C.** Decoded bar position as function of the time in the gaze center (top, orange), the gaze left (middle, green) and the gaze right (bottom, purple) conditions. Graphs present decoding from individual runs (from the second session) using V1 voxels activity (*attend-bar* condition) of a participant (sub-05). Predictions are illustrated as red and blue dashed lines for the retinotopic and spatiotopic predictions, respectively. Error areas show the posterior distribution STD. **D-E.** Average decoded bar position across bar passes and sessions for best-fitting voxels of V1 (D) and hMT+ (E) respectively, for *gaze center* (top), *gaze left* (middle) and *gaze right* conditions (bottom), in the *attend-bar* (left panels, dark colors) and *attend-fix* conditions (right panels, light colors). Error areas show \pm SEM. **F.** Group averaged correlation between the decoded bar position and the ground truth. Conventions are as in Figure 3E. **G.** Group average reference frame index observed in the *attend-bar* (black) and *attend-fix* conditions (gray) for the best-fitting voxels of each ROI. Conventions are as in Figure 3E.

5 Discussion

6 We determined the spatial reference frame of visual processing across the human visual hierarchy using a
 7 principled computational modeling approach. We estimated population receptive field parameters while
 8 participants directed their gaze to distinct points in outside-world space. This allowed us to directly
 9 juxtapose two competing hypotheses: that of retinotopic spatial coding moving with participants' gaze,
 10 and that of spatiotopic spatial coding fixed to the outside world or at least to the participant's head.

1 Confronting a conundrum of contradictory findings in the literature, we factorially manipulated both gaze
2 direction and the deployment of spatial attention. In brief, we find that all regions showing responses
3 tuned to visual spatial location, as gauged using the pRF model, are retinotopically organized. Moreover,
4 our results indicate that the deployment of spatial attention has no bearing on the reference frame of
5 visual location processing. We conclude that the human brain fundamentally encodes visual locations in a
6 retinotopic reference frame throughout the visual hierarchy. Our findings restrict solutions to the space
7 constancy problem, and increase the prominence of those centered on dynamic interactions between
8 action and perception, such as trans-saccadic remapping driven by re-afferent motor signals.

9 In the present study, we sought to leverage the precision afforded by the combination of ultra-high-
10 field fMRI's increased SNR and the quantitative pRF model of visual location tuning. This allowed us to
11 determine the reference frame of visual responses at the local neural population level, in single voxels.
12 We ensured that for all quantifications, differences in signal quality or pRF size would not bias our
13 outcomes toward the retinotopic or spatiotopic reference frame, following the logic outlined by Gardner
14 and colleagues¹³. We also took into account the fact that V1 neuron receptive fields and pRFs have a
15 nonlinear normalization field penumbra extending up to a factor of eight beyond their classical linear
16 extent²⁴. This means that visual information impinging on peripheral retinal areas, such as the visual
17 transients produced by eye movements, evoke measurable signal fluctuations in even the most central
18 portions of our retinotopic maps. Indeed, only experiments in full darkness will suffice to fully rule out the
19 misinterpretation of such visual signals originating from the periphery. Importantly, by employing static
20 gaze positions throughout our experimental runs and decreasing the strength of peripheral visual
21 stimulation, we explicitly prevent such spurious signals from producing false spatiotopic patterns in our
22 data⁷. We conceptualized experiments and analyses to be able to unequivocally adjudicate between
23 retinotopic and spatiotopic reference frames. A crucial step here was to explicitly formulate and directly
24 compare quantitative predictions for both hypotheses, at the level of single-voxel time series and across
25 populations of voxels within ROIs. Across subsequent analyses these direct comparisons between data and
26 hypotheses' predictions consistently favored the retinotopic reference frame hypothesis as the
27 fundamental spatial encoding scheme of visual cortex.

28 As for the role of attention, our results align with previous studies indicating that spatial attention
29 enhances the fidelity of spatial representations^{21,22}. Spatial attention is known to modulate pRF position in
30 accordance with the predictions of gain-field models^{25,26}. Our results indicate that enhanced signal fidelity
31 afforded by spatial attention renders visual cortex more, not less (Figs 3E, 4F-G), retinotopically
32 organized¹¹. We conclude that the differential deployment of spatial attention does not result in changes
33 in the retinotopic spatial reference frame of visual cortex.

34 Every day, we experience a stable visual world across our eye movements: an experience far removed
35 from the flutter of images projected onto our retiniae. Our results reframe and highlight the fundamental
36 question of visual stability across gaze changes, shifting the focus to mechanisms that achieve this stability
37 with a robustly retinotopic visual cortex. Different solutions to the space constancy problem have been
38 proposed^{27,28}, we here review these solutions through the lens of our results. A first category of solutions
39 to the space constancy problem relies on the use of proprioceptive information of the eye position in the
40 skull. An influential computational model²⁹ supported by numerous findings in both low-level^{30,31} and high-
41 level visual areas^{32,33} proposed that proprioceptive signal of gaze direction could be used to recover stimuli
42 position across eye movements. These findings suggest that the gain of retinotopically organized neurons
43 is modulated by a field which encodes the position of the animal's eye. It is known from
44 electrophysiological studies that although neighboring neurons encode similar locations in retinotopic
45 visual space, they experience wildly varying gain-field effects³¹, reflecting a salt-and-pepper organization.
46 We did not find an influence of gaze position at the level of individual voxels nor at the level of visual areas,
47 which is possibly in line with this organization at the single neuron level. Other than the potential absence
48 of structured gain field maps, other reasons, however, might have prevented us to find them. First, visual

1 display size used in fMRI studies only allows gaze direction changes smaller than what used in
2 electrophysiology, although the magnitude of our gaze shifts is in the range of most saccades made during
3 naturalistic viewing. Second, gain fields from animal studies involved modulation of neuronal activations
4 while our results rely on fMRI BOLD signals, which only is an indirect proxy of neurons firing rates^{34,35}.
5 Third, the timescales at which these eye position gain fields impact spatial processing may be faster than
6 those of our experimental runs. Thus, our results do not rule out gain field mediated solutions to space
7 constancy. Rather, our results support these accounts, while emphasizing the importance of the
8 retinotopic reference frame in their application³⁶.

9 A second category of solution involves the use of an efference copy³⁷ to correct the consequence of
10 gaze direction changes either before or during an eye movement^{38,39}. These corrections have been
11 proposed to be performed on retinotopically organized visual maps with the process being orchestrated
12 by visual attention^{40,41}. Remapping was observed at both the single neuron measures in distinct visual
13 maps⁴², fMRI studies⁴³⁻⁴⁵ and human behavior^{46,47}. Our results are compatible with this proposal but do
14 not test it. Indeed, it is important to note that in order to model voxel activity we asked participants to
15 maintain steady gaze throughout each acquisition. They change gaze directions only between runs, and
16 not across trials – an explicit design choice to avoid any peripheral visual stimulation arising from the eye
17 movements. Nevertheless, as participants laid still with the head and body fixed, our experiments left as
18 much space as possible for putative fundamental spatiotopic processes to occur (see also¹⁴). Our results
19 thus rule out space constancy solutions relying on an early encoding of space within craniotopic or
20 spatiotopic maps^{11,12}.

21 As we move up the visual hierarchy, responses become more spatially invariant. This invariance
22 abstracts neural responses away from their exact retinal origins, and reduces their adherence to the
23 retinotopic reference frame while allowing them to be influenced by factors such as attentional gain
24 fields^{25,26,48}. It is also clear from a host of canonical findings in the medial temporal lobe that the brain does
25 indeed encode the location of the agent in world-centered space. Clearly, a host of different information
26 sources such as self-motion⁴⁹⁻⁵¹, pictorial cues⁵², and navigational affordances⁵³⁻⁵⁵ play a role in generating
27 world-centered representations in our everyday interactions with the world. We argue that this process
28 must feature the continuous integration of low-level sensory and high-level world-centered information.
29 Interestingly, hippocampus was shown to encode locations in visual space in both humans and rodents⁵⁶⁻
30 ⁵⁸, indicating that world-centered and sensory-based spatial reference frames can coexist at higher levels
31 of processing. These findings demonstrate that V1-hippocampus interaction highlight the intimate links
32 between different sensory and world-centric reference frames. It is clear that oculomotor behavior plays
33 a foundational role in drawing these links⁵⁹.

34 35 **Conclusion**

36 Modeling ultra-high-field human cortical activity, we probed the reference frame of spatial vision. We
37 found, both the level of individual neural population and that of visual areas, that vision is organized
38 retinotopically, irrespective of where we look or attend. These results point to solutions to the space
39 constancy problem that use active correction mechanism by means of predictive or proprioceptive signals.

40 41 **Methods**

42 ***Ethics statement***

43 This experiment was approved by the Ethics Committee of Vrije Universiteit Amsterdam and conducted in
44 accordance with the Declaration of Helsinki. All participants gave written informed consent.

45 46 ***Participants***

47 Eight students and staff members of the Spinoza Centre for Neuroimaging participated in the experiment
48 (ages 21-40, 1 female, 4 authors). All except two authors were naïve to the purpose of the study and all

1 had normal or corrected-to-normal vision.

3 **MRI data acquisition**

4 Four T1-weighted (1.0 mm isotropic resolution) and one T2-weighted (1.0 mm isotropic resolution)
5 structural scans were acquired for each participant at The Spinoza Centre for Neuroimaging on a Philips
6 Achieva 3T scanner (Philips Medical Systems, Best, Netherlands) and a 32 channel receive-coil array with
7 a single channel transmit coil. Functional data were collected at the same center on a Philips Achieva 7T
8 scanner (Philips Medical Systems, Best, Netherlands) with a 32 channel receive-coil array with 8-channel
9 transmit coil (Nova Medical, Wilmington, MA). Functional data were collected using a 3D-EPI sequence at
10 a resolution of 1.8 mm isotropic with a 1.3 s volume acquisition time, 44 ms TR, 17 ms TE, consisting of 98
11 slices of 112 by 112 voxels covering the entire brain. SENSE acceleration was applied in both the Anterior-
12 Posterior (2.61-fold) and Right-Left (3.27-fold) directions. To estimate and correct susceptibility-induced
13 distortions, we acquired an identical EPI image with an opposite phase encoding direction, after each
14 functional scan. Transmit field homogeneity was improved by adjusting the 8-channel transmit output
15 based on a B1+ field population template (i.e., “universal pulse”⁶⁰).

17 **Stimuli and tasks**

18 The experiment consisted of 3 experimental sessions on different days. Participants started with two fMRI
19 sessions each composed of 10 consecutive runs together with different preparatory and field
20 inhomogeneity mapping scans (about 1h each). The last session was used to obtain multiple structural
21 images (about 45 min in total). Participants were trained on the behavioral task outside the scanner.
22 Stimuli were presented at a viewing distance of 225 cm, on a 32-inch LCD screen (BOLDscreen, Cambridge
23 Research Systems, Rochester, UK) situated at the end of the bore (17.8 dva horizontally by 10 dva
24 vertically) and viewed through a mirror. The screen had a spatial resolution of 1920 by 1080 pixels and a
25 refresh rate of 120 Hz. Button responses were collected using an MRI compatible button box (Current
26 Designs, Philadelphia, PA, USA). The experimental software controlling the display and the response
27 collection as well as eye tracking was implemented in Matlab (The MathWorks, Natick, MA, USA) using the
28 Psychophysics Toolbox^{61,62}.

29 Participants were instructed to continuously fixate a white bull’s-eye on a black background. The bull’s-
30 eye was composed of a central white dot of 0.04 dva radius surrounded by a white annulus of 0.4 dva
31 radius and a line width of 0.04 dva. They were instructed to attend visual noise contained either within
32 this central bull’s-eye (*attend-fix* condition, 1st, 3rd, 5th, 7th and 9th run) or within a moving bar (*attend-bar*
33 condition, 2nd, 4th, 6th, 8th and 10th run). Each functional session started with 4 runs in which the fixation
34 bull’s-eye was displayed at the screen center together with the moving bar traversing the entire screen
35 (*full-screen* condition, 1st, 2nd, 3rd and 4th run). In the next 6 runs, the moving bar stimulus was displayed
36 within a 4 dva radius aperture (0.04 dva width cosine edge between the visual noise and the background)
37 surrounding the fixation bull’s-eye. The bull’s-eye was displayed either 4 dva to the left of the screen center
38 (*gaze left* condition, 5th and 6th run), 4 dva to the right of the screen center (*gaze right* condition, 7th and
39 8th run) or at the screen center (*gaze center* condition, 9th and 10th run). *Full screen* runs lasted 195 seconds
40 (150 TRs). They were composed of 9 interleaved periods with and without the presentation of visual noise.
41 In periods without visual noise, the bull’s-eye was shown alone for 13 seconds (10 TRs). In periods with
42 the visual noise, the bar aperture movement direction was sequentially 180° (left), 270° (down), 0° (right)
43 and 90° (up). The bar aperture could either be horizontally or vertically oriented, in order to be
44 perpendicular to its movement direction. Its center moved in the bar direction on every TR by discrete
45 steps of 0.56 dva, such that 18 vertical and 32 horizontal steps were necessary to traverse the entire
46 screen. *Gaze left*, *gaze right* and *gaze center* runs lasted 158.6 seconds (122 TRs). They were composed of
47 the same nine interleaved periods with and without the presentation of visual noise. Contrary to the *full*
48 *screen* runs, the bar aperture movement direction was sequentially 180° (left), 0° (right), 180° (left) and 0°

1 (right) and the bar was always vertically oriented. Its center moved in the bar direction on every TR by 18
2 discrete steps of 0.22 dva to traverse the circular aperture centered on the bull's-eye. Visual noise
3 contained within the bar aperture (1 dva width, with 0.04 dva width cosine lateral edges) and the bull's-
4 eye was composed of distinct randomly generated pink noise (1/f) grayscale textures⁶³ (with individual
5 pixel of 0.04 dva width). To probe visual attention, we filtered the orientation contained within the
6 textures to generate clockwise and counterclockwise signal textures. We defined 15 different difficulty
7 levels by either keeping randomly generated noise or by filtering contained orientation by a von Misses
8 distribution with standard deviation values between 10^{-1} (large dispersion) to $10^{1.5}$ (narrow dispersion)
9 centered around $+45^\circ$ or -45° relative to the vertical axis. Each bar presentations started with 400 ms of
10 streams of unfiltered noise textures presented at 10 Hz in both the bull's-eye and the bar aperture. This
11 period was followed by 600 ms of clockwise or counterclockwise streams of oriented filtered noise, later
12 followed by 300 ms of unfiltered noise streams. The period with oriented filtered noise streams was
13 highlighted to participants with the bull's-eye central dot displayed in black. The orientation of the filtered
14 noise was selected randomly and independently every TR for both the bull's-eye and the bar aperture.
15 Participants reported the orientation of filtered noise streams presented within the bar aperture (*attend-*
16 *bar* condition) or the bull's-eye (*attend-fix* condition) with left (-45° or counter-clockwise) or right thumb
17 button presses (45° or clockwise). The difficulty of the task was titrated by a staircase procedure following
18 a 3 down 1 up rule adjusting for the attended stream, the orientation dispersion around $\pm 45^\circ$. Staircases
19 started at an intermediate difficulty value (level 10) on every run. Participants had until the end of the TR
20 to press a button if no response was registered the corresponding staircase wasn't modified.

21 **Anatomical data preprocessing**

22 T1-weighted (T1w) images were corrected for intensity non-uniformity⁶⁴. The T1w-reference was then
23 skull-stripped. Brain tissue segmentation of cerebrospinal fluid (CSF), white-matter (WM) and gray-matter
24 (GM) was performed on the brain-extracted T1w⁶⁵. A T1w-reference map was next computed after
25 registration of 4 T1w⁶⁶. Brain surfaces were reconstructed using freesurfer⁶⁷ and the brain mask estimated
26 previously was next refined⁶⁸, with the aid of the additional T2w acquisition. Finally, the surface
27 reconstruction was manually edited to correct for pial segmentation errors near the boundary of the
28 occipital lobe and the cerebellum.

29 **Functional data preprocessing**

30 For each of the 20 functional runs per participant (across all tasks and sessions), the following
31 preprocessing was performed. First, a reference volume and its skull-stripped version were generated
32 using a custom methodology of fMRIPrep⁶⁹. Head-motion parameters with respect to the BOLD reference
33 (transformation matrices, and six corresponding rotation and translation parameters) are estimated
34 before any spatiotemporal filtering⁷⁰. A fieldmap was estimated based on two echo-planar imaging (EPI)
35 references with opposing phase-encoding directions⁷¹. Based on the estimated susceptibility distortion, a
36 corrected EPI (echo-planar imaging) reference was calculated for a more accurate co-registration with the
37 anatomical reference. The BOLD reference was then co-registered to the T1w FreeSurfer⁷². Co-registration
38 was configured with twelve degrees of freedom to account for distortions remaining in the BOLD
39 reference. The BOLD timeseries were resampled onto their original, native space by applying a single,
40 composite transform to correct for head-motion and susceptibility distortions. Low frequency drift up to
41 0.01 Hz of the functional time series was removed using cosine drift regressors obtained from fMRIPrep.
42 Functional signal units were next converted to z-score.

43 **Regions of interest definition and pRF parameters**

44 We first averaged the timeseries of the 8 *full screen* runs, irrespective of whether participants attended to
45 the fixation bull's-eye or to the bar stimulus. We analyzed this *fiducial* averaged time series using an
46 isotropic Gaussian pRF model¹⁵. The model included a visual design matrix (converting the visual stimulus
47
48

1 in an on-off contrast map for each TR), the position (x, y), the size (standard deviation of the Gaussian), as
2 well as a signal amplitude and signal baseline as parameters. The time series were fitted in two parts. The
3 first part was a coarse spatial grid search of 40 linear steps in the 3 spatial model parameters: position (x ,
4 y), and size. A linear regression between the predicted and measured time series signal was used to
5 determine the baseline and amplitude parameters. The best-fitting parameters of the first step were next
6 used as the starting point of an optimization phase to produce finely tuned estimates. The visual stimulus
7 was down sampled by a factor 8 for the first of these fitting stages. The grid search position and size
8 parameters were distributed linearly within 20 dva from the center of the screen. In the fine search stage,
9 we used a gradient descent algorithm starting from the obtained grid search parameters, leaving the
10 parameter ranges essentially unconstrained. pRF polar angle maps were derived with the best explained
11 pRF position parameters. The polar angle was drawn using *Pycortex*⁷³ on an inflated and flattened cortical
12 visualization of individual participant anatomy. The explained variance of the model set as the coefficient
13 of determination (R^2) was used to determine color transparency. From this analysis we manually defined
14 for each participant 12 cortical visual regions of interests following pRF angle and eccentricity progression
15 on the cortex as well as anatomical references. These regions include V1, V2 and V3, V3AB which unites
16 V3A and V3B subregions, Lateral Occipital area (LO) which unites LO1 and LO2 subregions, Ventral Occipital
17 (VO) including selective voxels of the ventral path, hMT+ which unites both MT and MST subregions, the
18 inferior and superior part of the Intra Parietal Sulcus areas (iIPS and sIPS) and the inferior, superior and
19 medial parts of the Pre Central Sulcus areas (iPCS, sPCS and mPCS). For each ROI, we defined the best-fitting
20 voxels parameters as the 250 highest R^2 per ROI. We evaluated the effect of the task on the explained
21 variance of the pRFs by taking the task-related time series independently.

22
23 **Out-of-sample predictions**
24 We first evaluated the explained variance of retinotopic or spatiotopic pRF models by applying the
25 parameters obtained in the *full screen* conditions to the *gaze* conditions time series. To do so, we predicted
26 out-of-set timeseries of the gaze conditions visual designs knowing the pRF parameters of the attention
27 task-specific *full screen* condition to which we either subtracted ($pRFx_{gaze\ left} = pRFx_{full\ screen} - 4$) or added
28 ($pRFx_{gaze\ right} = pRFx_{full\ screen} + 4$) the gaze direction change to the $pRFx$ parameter. Signal amplitude (beta)
29 and signal baseline parameters were also adjusted to fit with the new timeseries. Change in explained
30 variance between condition was computed by subtracting the averaged R^2 observed in the *gaze left* and
31 *gaze right* condition from the R^2 observed in the *gaze center* condition.

32
33 **Fitting predictions**
34 We refine our analysis by determining anew the model parameters in the gaze conditions. To do so, we
35 kept the $pRFy$ and $pRFsize$ parameters obtained from the *full screen* runs, while fitting anew the $pRFx$
36 parameter as well as amplitude and the baseline. To avoid any bias toward the retinotopic or the
37 spatiotopic model, we started the coarse fitting procedure with x coordinate spatial grids centered either
38 on the gaze or on the screen coordinates. The best model (highest R^2) was next used in the optimization
39 phase to produce finely tuned parameters.

40
41 **Bayesian decoding**
42 We developed a novel method to probabilistically decode the position of our moving bar stimulus from
43 the multivariate fMRI data, within a Bayesian framework. The method is based on the standard pRF model
44 fitted on an independent training dataset from the same participant. Provided with a multivariate fMRI
45 timeseries $Y = y_{1..n,1..t}$ with n voxels and t time points, the method yields a posterior distribution
46 $p(x_{1..t}, h_{1..t} | Y)$ over possible stimulus positions x and heights h for every frame in the time series where
47 a stimulus was presented. The method allows to probe what bar position is represented on the level of a
48 given ROI, on a TR-to-TR basis.

1 By jointly sampling across all time points of the data concurrently, we took into account the
 2 hemodynamic delay and hemodynamic temporal smoothing in the estimation procedure. This is made
 3 computationally feasible by a state-of-the-art computational graph library developed for deep neural
 4 network training (Tensorflow), massively parallel graphics processing units for computation (GPUs), as well
 5 as Hamiltonian MCMC method⁷⁴, which exploits the exact gradient functions that the "autodiff"
 6 functionality that Tensorflow provide.

7 To obtain a posterior distribution over bar positions, given fMRI data, we first needed a likelihood
 8 function that describes the probability of the data given a stimulus time series S . Let's assume S a time
 9 series of 2D stimulus images. The likelihood function is the probability density function of a multivariate t
 10 distribution over the residuals of the pRF model with a covariance matrix Σ of a multivariate t-distribution
 11 with df degrees-of-freedom. Hence, $Y = pRF(S; \theta_{1..n}) + \epsilon$ where $\epsilon \sim t(0, \Sigma, df)$ and thus
 12 $p(Y|x_{1..t}, h_{1..t}; \theta_{1..n}, \Sigma, df) = t(Y - prf(S; \theta_{1..n}); \Sigma, df)$ where pRF is the output of a standard pRF
 13 model for the image time series S , as described in the previous section, with voxelwise position, dispersion,
 14 amplitude, and baseline parameters $\theta_{1..n}$ estimated using standard pRF fitting methods on an
 15 independent training dataset.

16 The pRF function includes a convolution step where "neural" time series are convolved with a canonical
 17 hemodynamic response function. This means that the pixel intensities on timepoint t influence the
 18 residuals and thereby likelihoods of multiple frames of neural activity data (roughly between timepoint t
 19 and 20 seconds later). This also means that different dimensions of the likelihood function that are close
 20 together in time are correlated, necessitating more advanced samplers than traditional Gibbs
 21 sampling^{75,76}.

22 To estimate the noise covariance matrix Σ , we couldn't use the sample covariance of the residuals of
 23 the training set, because these estimates would be highly unstable in our data regime of sparse data (a
 24 few hundred time points) and a large number of dimensions⁷⁷. Therefore, we estimated a regularized
 25 covariance matrix $\hat{\Sigma}$, using a method developed for decoding Gabor orientations from fMRI signals in early
 26 visual areas¹⁶. Specifically, the covariance matrix estimate $\hat{\Sigma}$ is a weighted sum of a diagonal covariance
 27 matrix $I \circ \tau\tau^T$, a perfectly correlated covariance matrix $\tau\tau^T$, and a matrix WW^T , where W is a $n \times m$
 28 matrix with n the number of voxels and m the total number of pixels in a single stimulus image. Each
 29 element of the matrix W describes the linear weight of its respective pixel to the activation of its respective
 30 voxel. More specifically, this value is the probability density of the x and y coordinates for the multivariate
 31 Gaussian described by the pRF parameters θ for that voxel, multiplied by its estimated amplitude.

32 Thus, we assume that the covariance of the residuals of two voxels that have particularly similar pRF
 33 should generally be higher than that of those voxels of which the pRF do not overlap¹⁶⁻¹⁸. The complete
 34 formula for the estimated residual covariance matrix $\hat{\Sigma}$ is: $\hat{\Sigma} = \rho I \circ \tau\tau^T + (1 - \rho)\tau\tau^T + \sigma^2 WW^T$ where
 35 the vector $\tau > 0$, scalar $0 \leq \rho < 1$, and scalar $\sigma^2 \geq 0$ are estimated using maximum likelihood estimation
 36 on the residuals $Y^{training} - pRF(S^{training}; \theta_{1..n})$ of the training data, using gradient descent
 37 optimization. However, the dimensionality of the stimulus time series S was so high (above 100,000) that
 38 it was infeasible to sample its posterior distribution. Moreover, its dimensions are not semantically
 39 meaningful. Therefore, we opted to drastically reduce the dimensionality of S using a two-parameter
 40 stimulus model $s(x, h)$ that returns a two-dimensional pixel image with a vertically centered bar of a fixed
 41 width, centered on the x coordinate and a height of h . We chose this stimulus model because both the x
 42 position and height of the bar varied over time in the gaze conditions. Note that, in the stimulus function,
 43 the pixel intensity rapidly but smoothly falls off as a sigmoidal function of distance to the border of the bar
 44 stimulus. Such a smooth functional form was necessary to keep the likelihood function differentiable and
 45 thereby Hamiltonian MCMC feasible. This way, the number of parameters of the likelihood function
 46 was heavily reduced to 2 times the number of timepoints: $p(Y|x_{1..t}, h_{1..t}; \theta) = t(Y -$
 47 $pRF(s(x_{1..t}, h_{1..t}); \theta_{1..n}); \Sigma, df)$. We reduced this number a bit further by not estimating stimulus

1 properties for frames where there was no stimulus bar on the screen such that the model assumed these
2 frames to be an empty screen with a fixation bull's-eye.

3 We used uniform priors on the estimated x coordinates, restricted between the leftmost and rightmost
4 coordinates on the screen as well as uniform priors on the height of the stimulus bar between 0 and the
5 height of the stimulus screen. This was achieved using bijective sigmoidal transformation functions. The
6 entire procedure was implemented in TensorFlow 2 and packaged as a freely available Python package
7 called Braincoder (<https://github.com/Gilles86/braincoder>).

8 As a first step, a 2D grid search (position by height) was performed for every time point separately,
9 optimizing the likelihood function with respect to single frame x and h parameters, neglecting the
10 covariance structure of the likelihood function. After this grid search, a maximum a posteriori (MAP)
11 estimate for the entire time series $\operatorname{argmax} p(x_{1..t}, h_{1..t} | Y)$ was estimated using gradient descent, starting
12 from the values found in the grid search. Finally, this MAP estimate was used as a starting point for the
13 NUTS MCMC sampler, a self-tuning Hamiltonian MCMC sampler⁷⁶, to obtain samples from the posterior
14 $p(x_{1..t}, h_{1..t} | Y)$. We used 4 independent chains, with 250 steps and 500 samples per chain. After some
15 experimenting, we settled on an acceptance probability of 0.3 to more effectively sample our highly
16 complex 144-dimensional posterior distribution.

17 **Reference frame index**

18 We devised a reference frame index (RFI), based on earlier work on the topic^{11,13}. This index contrasts the
19 ability of a spatiotopic versus a retinotopic model to explain the observed data. Specifically, it expresses
20 the difference in explained variance between the retinotopic and the spatiotopic model as a fraction of
21 the sum of both explained variances. Hence, the index ranges between -1 (completely and exclusively
22 consistent with spatiotopic reference frame) and 1 (completely and exclusively consistent with the
23 retinotopic reference frame):

$$24 \text{RFI} = \frac{\text{EV}_{\text{retinotopic}} - \text{EV}_{\text{spatiotopic}}}{\text{EV}_{\text{retinotopic}} + \text{EV}_{\text{spatiotopic}}}$$

25 Where EV is the reduction in variance from the raw data versus the residuals of the model fit $\text{EV} =$
26 $\text{Var}(Y - \hat{Y}) - \text{Var}(Y)$. Note that for the out-of-sample predictions (Figure 3) of individual PRF time series,
27 the variance is defined on real and predicted BOLD timeseries of individual voxels. For the decoding
28 analysis (Figure 4) the variance is defined on the actual and the decoded spatial positions of the bar
29 stimulus.

30 **Statistics**

31 For statistical comparisons we used SciPy⁷⁸ paired permutation test with 10,000 permutation
32 bootstrapped iterations used to approximate the null distribution⁷⁹. This procedure allows us to resample
33 our data to create a distribution with randomly rearranged labels of the conditions to compare for each
34 participant. We determined statistical significance by deriving two-tailed (see Results) p values for the
35 comparison of the distributions of the compared conditions.

36 **Data and code availability**

37 We make available online our imaging dataset, including all derivatives and individual participant and
38 group figures (openneuro.org/datasets/ds004091). Brain models with data analysis visualization are made
39 available online (invibe.nohost.me/gazeprf) together with the experimental and data analyses codes
40 (github.com/mszinte/gaze_prf).

41 **Acknowledgments**

42 We are grateful to the members of the Knapen and Dumoulin laboratories in Amsterdam for helpful
43 comments and discussions and to Alice and Clémence Szinte for their support. Centre de Calcul Intensif
44

1 d'Aix-Marseille is acknowledged for granting access to its high performance computing resources. This
2 research was supported by a Marie Skłodowska-Curie Action Individual Fellowship to M.S. (704537) and
3 an NWO-CAS (012.200.012) and ABMP (2015-7) grants to T.K.
4

5 **Author contributions**

6 Conceptualization, M.S., and T.K.; Methodology, M.S., G.H., M.A., and T.K.; Software, M.S., G.H., M.A., and
7 T.K.; Validation, M.S., G.H., M.A., and T.K.; Formal Analysis, M.S., G.H., M.A., and T.K.; Investigation, M.S.,
8 I.V., and T.K.; Resources, T.K.; Data Curation, M.S.; Writing – Original Draft, M.S., G.H., and T.K.; Writing –
9 Review & Editing, M.S., G.H., M.A., I.V., S.D., and T.K.; Visualization, M.S., G.H. and T.K.; Supervision, S.D.,
10 T.K.; Project Administration, T.K.; Funding Acquisition, M.S. and T.K.
11

12 **References**

- 13 1. Baylor, D. A., Lamb, T. D. & Yau, K. W. Responses of retinal rods to single photons. *J Physiology* **288**, 613–
14 634 (1979).
- 15 2. Sereno, M. I. *et al.* Borders of multiple visual areas in humans revealed by functional magnetic resonance
16 imaging. *Science* **268**, 889–893 (1995).
- 17 3. Wandell, B. A., Dumoulin, S. O. & Brewer, A. A. Visual field maps in human cortex. *Neuron* **56**, 366–383
18 (2007).
- 19 4. Zeki, S. M. Functional organization of a visual area in the posterior bank of the superior temporal sulcus
20 of the rhesus monkey. *J Physiology* **236**, 549–573 (1974).
- 21 5. Essen, D. C. V., Maunsell, J. H. R. & Bixby, J. L. The middle temporal visual area in the macaque:
22 Myeloarchitecture, connections, functional properties and topographic organization. *The Journal of*
23 *comparative neurology* **199**, 293–326 (1981).
- 24 6. Snyder, L. H., Batista, A. P. & Andersen, R. A. Coding of intention in the posterior parietal cortex. *Nature*
25 **386**, 167–170 (1997).
- 26 7. Fabius, J. H., Moravkova, K. & Fracasso, A. Topographic organization of eye-position dependent gain
27 fields in human visual cortex. *Nat Commun* **13**, 7925 (2022).
- 28 8. Andersen, R. A. & Mountcastle, V. B. The influence of the angle of gaze upon the excitability of the light-
29 sensitive neurons of the posterior parietal cortex. *The Journal of Neuroscience* **3**, 532–548 (1983).
- 30 9. Andersen, R. A., Snyder, L. H., Bradley, D. C. & Xing, J. Multimodal representation of space in the
31 posterior parietal cortex and its use in planning movements. *Annual Review of Neuroscience* **20**, 303–330
32 (1997).
- 33 10. Andersen, R. A., Essick, G. K. & Siegel, R. M. Encoding of spatial location by posterior parietal neurons.
34 *Science* **230**, 456–458 (1985).
- 35 11. Crespi, S. *et al.* Spatiotopic Coding of BOLD Signal in Human Visual Cortex Depends on Spatial Attention.
36 *PLoS ONE* **6**, (2011).
- 37 12. d'Avossa, G. *et al.* Spatiotopic selectivity of BOLD responses to visual motion in human area MT. **10**,
38 249–255 (2007).
- 39 13. Gardner, J. L., Merriam, E. P., Movshon, J. A. & Heeger, D. J. Maps of Visual Space in Human Occipital
40 Cortex Are Retinotopic, Not Spatiotopic. *The Journal of Neuroscience* **28**, 3988–3999 (2008).
- 41 14. Golomb, J. D. & Kanwisher, N. Higher Level Visual Cortex Represents Retinotopic, Not Spatiotopic,
42 Object Location. *Cereb Cortex* **22**, 2794–2810 (2012).
- 43 15. Dumoulin, S. O. & Wandell, B. A. Population receptive field estimates in human visual cortex.
44 *NeuroImage* **39**, 647–660 (2008).
- 45 16. Bergen, R. S. van, Ma, W. J., Pratte, M. S. & Jehee, J. F. M. Sensory uncertainty decoded from visual
46 cortex predicts behavior. *Nature neuroscience* **18**, 1728–1730 (2015).
- 47 17. Bergen, R. S. van & Jehee, J. F. M. Probabilistic Representation in Human Visual Cortex Reflects
48 Uncertainty in Serial Decisions. *J Neurosci* **39**, 8164–8176 (2019).

- 1 18. Bergen, R. S. van & Jehee, J. F. M. Modeling correlated noise is necessary to decode uncertainty.
2 *Neuroimage* **180**, 78–87 (2018).
- 3 19. Anton-Erxleben, K. & Carrasco, M. Attentional enhancement of spatial resolution: linking behavioural
4 and neurophysiological evidence. *Nature reviews. Neuroscience* **14**, 188–200 (2013).
- 5 20. Carrasco, M. Visual attention: The past 25 years. *Vision Research* **51**, 1484–1525 (2011).
- 6 21. Bressler, D. W. & Silver, M. A. Spatial attention improves reliability of fMRI retinotopic mapping signals
7 in occipital and parietal cortex. *Neuroimage* **53**, 526–533 (2010).
- 8 22. Silver, M. A., Ress, D. & Heeger, D. J. Topographic Maps of Visual Spatial Attention in Human Parietal
9 Cortex. *J Neurophysiol* **94**, 1358–1371 (2005).
- 10 23. Geurts, L. S., Cooke, J. R. H., Bergen, R. S. van & Jehee, J. F. M. Subjective confidence reflects
11 representation of Bayesian probability in cortex. *Nat Hum Behav* **6**, 294–305 (2022).
- 12 24. Aqil, M., Knapen, T. & Dumoulin, S. O. Divisive normalization unifies disparate response signatures
13 throughout the human visual hierarchy. *Proc National Acad Sci* **118**, e2108713118 (2021).
- 14 25. Klein, B. P., Harvey, B. M. & Dumoulin, S. O. Attraction of position preference by spatial attention
15 throughout human visual cortex. *Neuron* **84**, 227–237 (2014).
- 16 26. Es, D. M. van, Theeuwes, J. & Knapen, T. Spatial sampling in human visual cortex is modulated by both
17 spatial and feature-based attention. *eLife* **7**, 3771 (2018).
- 18 27. Sun, L. D. & Goldberg, M. E. Corollary Discharge and Oculomotor Proprioception: Cortical Mechanisms
19 for Spatially Accurate Vision. *Annual Review of Vision Science* **2**, 61–84 (2016).
- 20 28. Neupane, S., Guitton, D. & Pack, C. C. Perisaccadic remapping: What? How? Why? *Rev Neuroscience*
21 **31**, 505–520 (2020).
- 22 29. Andersen, R. A., Essick, G. K. & Siegel, R. M. Encoding of Spatial Location by Posterior Parietal Neurons.
23 *Science* **230**, 456–458 (1985).
- 24 30. Durand, J.-B., Trotter, Y. & Celebrini, S. Privileged processing of the straight-ahead direction in primate
25 area V1. *Neuron* **66**, 126–137 (2010).
- 26 31. Trotter, Y. & Celebrini, S. Gaze direction controls response gain in primary visual-cortex neurons.
27 *Nature* **398**, 239–242 (1999).
- 28 32. Andersen, R. & Mountcastle, V. The influence of the angle of gaze upon the excitability of the light-
29 sensitive neurons of the posterior parietal cortex. *J Neurosci* **3**, 532–548 (1983).
- 30 33. Morris, A. P. & Krekelberg, B. A Stable Visual World in Primate Primary Visual Cortex. *Curr Biol* **29**, 1471-
31 1480.e6 (2019).
- 32 34. Arthurs, O. J. & Boniface, S. How well do we understand the neural origins of the fMRI BOLD signal?
33 *Trends Neurosci* **25**, 27–31 (2002).
- 34 35. Sirotin, Y. B. & Das, A. Anticipatory haemodynamic signals in sensory cortex not predicted by local
35 neuronal activity. *Nature* **457**, 475–479 (2009).
- 36 36. Zhao, Z., Ahissar, E., Victor, J. D. & Rucci, M. Inferring visual space from ultra-fine extra-retinal
37 knowledge of gaze position. *Nat Commun* **14**, 269 (2023).
- 38 37. Sperry, R. W. Neural basis of the spontaneous optokinetic response produced by visual inversion.
39 *Journal of comparative and physiological psychology* **43**, 482–489 (1950).
- 40 38. Sommer, M. A. & Wurtz, R. H. Influence of the thalamus on spatial visual processing in frontal cortex.
41 *Nature* **444**, 374–377 (2006).
- 42 39. Duhamel, J. R., Colby, C. L. & Goldberg, M. E. The Updating of the Representation of Visual Space in
43 Parietal Cortex by Intended Eye-Movements. *Science* **255**, 90–92 (1992).
- 44 40. Cavanagh, P., Hunt, A. R., Afraz, A. & Rolfs, M. Visual stability based on remapping of attention pointers.
45 *Trends in Cognitive Sciences* **14**, 147–153 (2010).
- 46 41. Rolfs, M. & Szinte, M. Remapping Attention Pointers: Linking Physiology and Behavior. *Trends in*
47 *Cognitive Sciences* **20**, 399–401 (2016).
- 48 42. Wurtz, R. H. Neuronal mechanisms of visual stability. *Vision Research* **48**, 2070–2089 (2008).

- 1 43. Merriam, E. P., Genovese, C. R. & Colby, C. L. Spatial Updating in Human Parietal Cortex. *Neuron* **39**,
2 361–373 (2003).
- 3 44. Merriam, E. P., Genovese, C. R. & Colby, C. L. Remapping in Human Visual Cortex. *Journal of*
4 *Neurophysiology* **97**, 1738–1755 (2007).
- 5 45. Knapen, T., Swisher, J. D., Tong, F. & Cavanagh, P. Oculomotor Remapping of Visual Information to
6 Foveal Retinotopic Cortex. *Frontiers in Systems Neuroscience* **10**, 54 (2016).
- 7 46. Rolfs, M., Jonikaitis, D., Deubel, H. & Cavanagh, P. Predictive remapping of attention across eye
8 movements. *Nature neuroscience* **14**, 252–256 (2011).
- 9 47. Szinte, M., Jonikaitis, D., Rangelov, D. & Deubel, H. Pre-saccadic remapping relies on dynamics of spatial
10 attention. *eLife* **7**, e37598 (2018).
- 11 48. Womelsdorf, T., Anton-Erxleben, K. & Treue, S. Receptive field shift and shrinkage in macaque middle
12 temporal area through attentional gain modulation. *The Journal of Neuroscience* **28**, 8934–8944 (2008).
- 13 49. Campbell, M. G. *et al.* Principles governing the integration of landmark and self-motion cues in
14 entorhinal cortical codes for navigation. *Nat Neurosci* **21**, 1096–1106 (2018).
- 15 50. Mallory, C. S. *et al.* Mouse entorhinal cortex encodes a diverse repertoire of self-motion signals. *Nat*
16 *Commun* **12**, 671 (2021).
- 17 51. Arnoldussen, D. M., Goossens, J. & Berg, A. V. van den. Adjacent visual representations of self-motion
18 in different reference frames. *Proc National Acad Sci* **108**, 11668–11673 (2011).
- 19 52. Lescroart, M. D. & Gallant, J. L. Human Scene-Selective Areas Represent 3D Configurations of Surfaces.
20 *Neuron* **101**, 178-192.e7 (2019).
- 21 53. Bonner, M. F. & Epstein, R. A. Coding of navigational affordances in the human visual system. *Proc*
22 *National Acad Sci* **114**, 4793–4798 (2017).
- 23 54. Epstein, R. A. & Baker, C. I. Scene Perception in the Human Brain. *Annu Rev Vis Sc* **5**, 1–25 (2019).
- 24 55. Purandare, C. S. *et al.* Moving bar of light evokes vectorial spatial selectivity in the immobile rat
25 hippocampus. *Nature* **602**, 461–467 (2022).
- 26 56. Saleem, A. B., Diamanti, E. M., Fournier, J., Harris, K. D. & Carandini, M. Coherent encoding of subjective
27 spatial position in visual cortex and hippocampus. *Nature* **562**, 124–127 (2018).
- 28 57. Diamanti, E. M. *et al.* Spatial modulation of visual responses arises in cortex with active navigation.
29 *Elife* **10**, e63705 (2021).
- 30 58. Knapen, T. Topographic connectivity reveals task-dependent retinotopic processing throughout the
31 human brain. *Proc National Acad Sci* **118**, e2017032118 (2020).
- 32 59. Nau, M., Schröder, T. N., Bellmund, J. L. S. & Doeller, C. F. Hexadirectional coding of visual space in
33 human entorhinal cortex. *Nat Neurosci* **21**, 188–190 (2018).
- 34 60. Oliveira, Í. A. F., Roos, T., Dumoulin, S. O., Siero, J. C. W. & Zwaag, W. van der. Can 7T MPRAGE match
35 MP2RAGE for gray-white matter contrast? *Neuroimage* **240**, 118384 (2021).
- 36 61. Pelli, D. G. The VideoToolbox software for visual psychophysics: transforming numbers into movies.
37 *Spatial Vision* **10**, 437–442 (1997).
- 38 62. Brainard, D. H. The Psychophysics Toolbox. *Spatial Vision* **10**, 433-436. (1997).
- 39 63. Hanning, N. M., Deubel, H. & Szinte, M. Sensitivity measures of visuospatial attention. *Journal of Vision*
40 **19**, 17–17 (2019).
- 41 64. Tustison, N. J. *et al.* N4ITK: Improved N3 Bias Correction. *Ieee T Med Imaging* **29**, 1310–1320 (2010).
- 42 65. Zhang, Y., Brady, M. & Smith, S. Segmentation of Brain MR Images Through a Hidden Markov Random
43 Field Model and the Expectation-Maximization Algorithm. *Ieee T Med Imaging* **20**, 45 (2001).
- 44 66. Reuter, M., Rosas, H. D. & Fischl, B. Highly accurate inverse consistent registration: A robust approach.
45 *Neuroimage* **53**, 1181–1196 (2010).
- 46 67. Dale, A. M., Fischl, B. & Sereno, M. I. Cortical Surface-Based Analysis: I. Segmentation and Surface
47 Reconstruction. *NeuroImage* **9**, 179–194 (1999).
- 48 68. Klein, A. *et al.* Mindboggling morphometry of human brains. *Plos Comput Biol* **13**, e1005350 (2017).

1 69. Esteban, O. *et al.* fMRIPrep: a robust preprocessing pipeline for functional MRI. *Nature Methods* **16**,
2 111–116 (2019).
3 70. Jenkinson, M., Bannister, P., Brady, M. & Smith, S. Improved Optimization for the Robust and Accurate
4 Linear Registration and Motion Correction of Brain Images. *Neuroimage* **17**, 825–841 (2002).
5 71. Cox, R. W. & Hyde, J. S. Software tools for analysis and visualization of fMRI data. *Nmr Biomed* **10**, 171–
6 178 (1997).
7 72. Greve, D. N. & Fischl, B. Accurate and robust brain image alignment using boundary-based registration.
8 *Neuroimage* **48**, 63–72 (2009).
9 73. Gao, J. S., Huth, A. G., Lescroart, M. D. & Gallant, J. L. Pycortex: an interactive surface visualizer for
10 fMRI. *Frontiers in neuroinformatics* **9**, 162 (2015).
11 74. Betancourt, M. A Conceptual Introduction to Hamiltonian Monte Carlo. *Arxiv* (2017).
12 75. Turner, B. M., Sederberg, P. B., Brown, S. D. & Steyvers, M. A Method for Efficiently Sampling From
13 Distributions With Correlated Dimensions. *Psychol Methods* **18**, 368–384 (2013).
14 76. Hoffman, M. D. & Gelman, A. The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian
15 Monte Carlo. *Journal of Machine Learning Research* **15**, 1593–1623 (2014).
16 77. Ledoit, O. & Wolf, M. Improved estimation of the covariance matrix of stock returns with an application
17 to portfolio selection. *J Empir Financ* **10**, 603–621 (2003).
18 78. Virtanen, P. *et al.* SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature*
19 *Methods* **17**, 261–272 (2020).
20 79. Ernst, M. D. Permutation Methods: A Basis for Exact Inference. *Stat Sci* **19**, (2004).
21
22

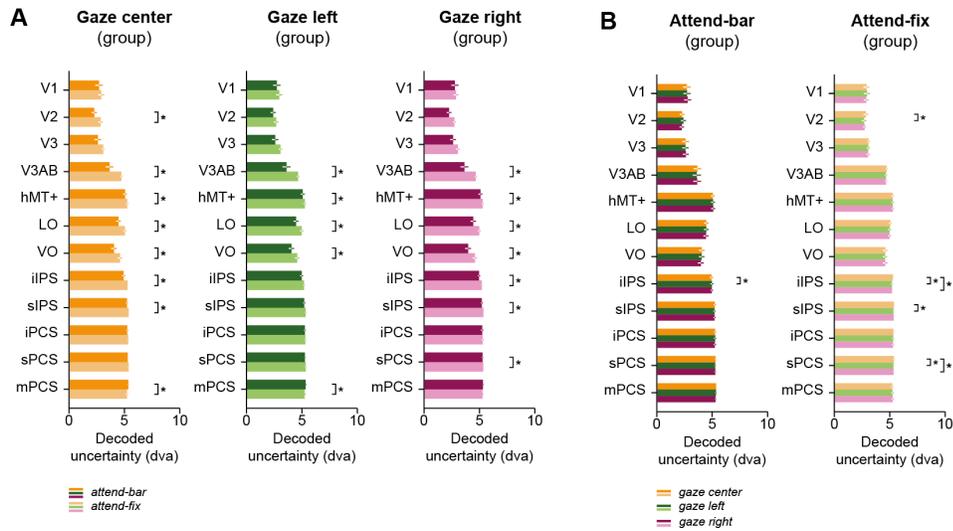


Figure S1. Decoding uncertainty. **A.** Average decoded uncertainty across ROIs and participants for the gaze center (left column), gaze left (middle column) and gaze right (right column) conditions. Black asterisks indicate a significant difference between the *attend-bar* and *attend-fix* conditions (two-sided $p < 0.05$). **B.** Same results for the *attend-bar* (left column) and *attend-fix* (right column) conditions. Black asterisks indicate a significant difference between the *gaze center*, *gaze left* and *gaze right* conditions (two-sided $p < 0.05$).