# *Supplementary Material*

**1    Additional exploratory analyses**

**1.1    searchlight-based reconstructions and effect of coherence**

We also performed additional whole-brain searchlight decoding to identify regions potentially outside our pre-defined ROIs that might encode stimulus-related or choice-related information. For this, the single-subject BFCA maps resulting from the searchlight reconstructions were smoothed using a Gaussian kernel with a FWHM of 6 mm and spatially normalized with SPM12. In order to identify searchlights in which the reconstruction performance was significantly above chance (i.e. BFCA > 50%), we used two one-way factorial designs (one for stimulus and one for report) with coherence level as a within-subject factor. For each model, we specified three t-contrasts. In this way we could identify clusters with significant above-chance information for each coherence level. We expected the reconstruction performance of searchlights located in visual areas, to decrease as a function of decreasing coherence. We also expected to identify clusters of voxels carrying information about perceptual judgements in visual (Britten et al., 1996; Serences & Boynton, 2007; Hebart et al., 2012; Sousa et al., 2021) and parietal areas (Gold & Shadlen, 2007; Brincat et al., 2018; Hebart et al., 2012, 2016; Levine & Schwarzbach, 2017). We further predicted the coherence level to have an effect on report reconstruction performance as well. We evaluated the effect of coherence on the reconstruction performance by inclusively masking the voxels that showed an average effect of reconstruction across coherence levels ($p < 0.001$, uncorrected). This procedure was done separately for stimulus reconstruction and report reconstruction. Please note that this analysis is not circular because the test for an effect of coherence is orthogonal to the test for average reconstruction performance.

The whole-brain searchlight reconstruction at 100% coherence revealed stimulus information from voxel clusters located in the left ($FWE_c$, $p < 0.05$, K = 1677; cluster-defining threshold $p < 0.001$) and in the right occipital cortex ($FWE_c$, $p < 0.05$, K = 1024; cluster-defining threshold $p<0.001$). For the intermediate coherence condition and for the 0% coherence condition, we found no searchlights that were significantly predictive of the stimulus motion direction. Note that in the 0% coherence condition, the stimulus has no global motion direction, thus a chance-level reconstruction performance is to be expected. In the left and right occipital cortex, we also found clusters of voxels informative about participants' reports for the 100% coherence condition (left: $FWE_c$, $p < 0.05$, K = 1049; cluster-defining threshold $p < 0.001$; right: $FWE_c$, $p < 0.05$, K = 864; cluster-defining threshold $p < 0.001$) as well as for the intermediate coherence condition (left: $FWE_c$, $p < 0.05$, K = 201; cluster-defining threshold $p<0.001$; right: $FWE_c$, $p < 0.05$, K = 487; cluster-defining threshold $p<0.001$). For the 0% coherence condition, we were not able to identify clusters informative about participants' reports (Supplementary Figure 1).

Since neurons in early and extrastriate visual areas are tuned to motion directions (Albright et al. 1984; Movshon & Newsome, 1996; Nichols & Newsome, 2002), we reasoned that if population-level measurements of neural activity obtained from single voxels reflects this property  (Nevado et al., 2004; Haynes, 2015; Sprague et al., 2018), the stimulus reconstruction performance should be maximum in the 100% coherence condition, and progressively decrease at intermediate coherence. At 0% coherence instead, the stimulus has no net motion direction and the reconstruction

41performance should be at chance level. We therefore expected to identify a main effect of coherence
42on the stimulus reconstruction performance in visual areas. We identified a cluster of voxels located
43in the left occipital pole, where coherence level had an effect on the performance in stimulus
44reconstruction (FWE$_c$, p < 0.05, K = 262; cluster-defining threshold p<0.001). Similarly, we
45predicted that a possible effect of coherence on the report reconstruction might be present, and driven
46by the expected correlation between the stimulus identity and participants' report, when the stimulus
47is clearly visible. However, we did not find clusters that showed such effect for the report
48reconstruction performance (Supplementary Figure 2).

### 49 1.2    Control analysis: stimulus and report reconstruction from eye-tracking data.

50The use of motion stimuli such as our RDK might trigger involuntary eye movement (Cohen et al.,
511977) that can be informative about the direction of perceived motion (Wilbertz et al., 2018). Eye
52movements also have an effect on brain activity measured with fMRI and can thus constitute a
53potential confound (Merriam et al., 2013), even when participants are specifically instructed to
54maintain fixation (Thielen et al., 2019). For this reason, we tested whether the recorded gaze position
55(x and y ordinates), was informative about the physical stimulus direction or participants' perceived
56direction. For this analysis we used the gaze position of 21 out of 23 subjects who participated in the
57main fMRI experiment (we couldn't record the traces of two participants for technical reasons). We
58reasoned that if the gaze position is systematically correlated with the presented motion direction or
59with participants' reports, the x and y ordinates should exhibit a specific position profiles, that can in
60turn be used to perform stimulus and report reconstruction. In order to check whether the pattern of
61participants' eye movements was related with the stimulus or the reported motion directions, we
62estimated such profiles with a cyclic version of the GPR (see *Materials and Methods* in the main
63manuscript) and performed stimulus and report reconstruction at each time point with a procedure
64similar to the one adopted for the main fMRI analysis (see *Materials and Methods*).

65The preprocessing pipeline employed for this analysis was different from the one performed for
66fixation control (see *Materials and Methods*). Blinks, detected by the provided software from Eye
67Link, were linearly interpolated using the approach described in Urai et al. (2017). The resulting
68traces were filtered for electronic noise using a Butterworth filter (low cut off 5Hz, high cut off
69100Hz) following Thielen et al. (2018). The complete trace of each session was linearly detrended to
70account for drifts that appear due to the long continuous recordings during each session (each
71approximately 1.5h). An additional linear detrending was performed separately on each run, to
72counterbalance slow drifts in head position in the scanner. Periods of interest (500 ms before stimulus
73onset together with 2000 ms stimulus period) were combined across the two recording sessions. Data
74from the period of interest were baseline corrected using the pre-stimulus interval (500 ms up to
75stimulus onset).

76After preprocessing, the x and y gaze positions of each subject were grouped by coherence levels
77(0%, medium, 100% coherence) resulting in a maximum of 160 trials per condition. To reduce
78computational time, we only considered time points during the stimulus period (2000 ms) and
79resampled the signal at 50Hz. For each time point we estimated two position-related profiles (one for
80each ordinate) by entering the trial-wise position value together with the corresponding stimulus
81motion direction $\theta_s$ or the reported direction $\theta_r$ into a cyclic version of the GPR. Please note that this
82procedure is very similar to the one previously described for the estimation of voxel-wise response
83profiles (see eq. 3 in the *Materials and Methods* section - where the parameter $\hat{\beta}_j$ represents now the
84trial-wise recorded position of each ordinate in a single time point). The estimation of the position

85profiles was performed with a leave-one-run-out cross-validation scheme, by only using trials in
86which participants were maintaining fixation (see *Materials and Methods*).

87The stimulus and report reconstructions were estimated using the same procedure described in the
88*Materials and Methods*. However, the estimated gaze position profiles instead of voxel response
89profiles, were used to predict the stimulus direction $\theta_s$ or the reported direction $\theta_r$ in a run-wise cross-
90validation procedure. Please note, that in this case it is not necessary to adopt regularization for
91estimation of the covariance matrix because the number of position profiles (one for x and one for y
92ordinates) does not exceed the number of trials across runs (see eq. 9 in the *Materials and Methods*
93section). The results were evaluated by testing if the averaged BFCA across subjects was above
94chance for each timepoint. Statistical analyses were corrected for multiple comparison by performing
95a cluster-based permutation test (Maris & Oostenveld, 2017).

96The group-level average reconstruction performance for the stimulus and the report labels are
97depicted in Supplementary Figure3. We were not able to identify stimulus-related information in any
98of the three coherence levels. Instead, the evaluation of the report model indicates that the pattern of
99eye movements was informative about participants' report in the 0% coherence condition. More
100precisely we were able to identify clusters of above-chance reconstruction performance, peaking after
1011000 ms. Eye movement were not predictive of participant's choices for the intermediate or 100%
102coherence levels.

## 1032 The problem of feature continuous accuracy with unbalanced labels

104In order to evaluate the performance of our GPR-based reconstructions, we implemented a balanced
105version of FCA (BFCA –  see eq. 15 in the *Materials and Methods* section). Our goal was to obtain a
106measure of performance that could be intuitively compared to a standard accuracy measure with
107values distributed between 0% and 100%. FCA is derived by rescaling the continuous values of the
108absolute angular deviation (see eq. 1 and 2 in the Materials and Methods section); see also Pilly &
109Seitz, 2009), to evaluate the reconstruction performance. The need for a *balanced* version of FCA
110was due to participants' responses in the 0% coherence condition being unbalanced (Supplementary
111Figure 4) as is often the case for reports, even despite the use of our sensory matching approach that
112minimizes such biases (Töpfer et al., 2022). In case of a standard classification analysis, training and
113testing a classifier with unbalanced labels make accuracy an unreliable measure of performance
114(Japkowicz & Stephen, 2002). More specifically, when the performance of a classifier is tested on an
115imbalanced dataset, it might lead to the misleading finding of significant above-chance performance
116of the classifier (Brodersen et al., 2010), simply because the classifier tends to reproduce the
117distribution of the training dataset.

## 1183 Comparison between FCA and BFCA

### 1193.1 Simulation analysis

120In order to illustrate how BFCA and FCA are related with each other, as well as with the underlying
121independent variable distribution, we here show a simulation performed on synthetic data. In order to
122match the features of our experimental design, we simulated a $t \times 1$ vector of trial-wise parameter
123estimates $\hat{\beta}_j$ for a total of 1000 voxels where $t = 160$ total trials were generated across 10 runs. We
124also generated a vector $\theta$ corresponding to the independent variable (the stimulus or the report
125direction).

126For the current simulation we distinguished four alternative scenarios:

4

127  1) $\theta$ **modulates** $\hat{\beta}_j$ **and the distribution of** $\theta$ **is balanced**. This situation corresponds to the
128      hypothesized behavior of voxels sensitive to motion directions (as the stimulus directions are
129      balanced across runs in our experimental design);

130  2) $\theta$ **modulates** $\hat{\beta}_j$ **and the distribution of** $\theta$ **is unbalanced**. This scenario corresponds to the
131      hypothesized behavior of voxels sensitive to participants' reports in our experimental design
132      (as participants' reports are unbalanced, especially at 0% coherence level);

133  3) $\theta$ **does not modulate** $\hat{\beta}_j$ **and the distribution of** $\theta$ **is balanced**. This should be the case for
134      voxels insensitive to the stimulus direction. Such voxels should not produce spurious above-
135      chance FCA when combined for searchlight-based reconstruction;

136  4) $\theta$ **does not modulate** $\hat{\beta}_j$ **and the distribution of** $\theta$ **is unbalanced**. We assume that this
137      scenario could possibly produce spurious above-chance FCA when the voxels are combined
138      for searchlight-based reconstruction.

139 We applied the same analyses described in the manuscript (see *Materials and Methods*) to estimate
140 voxel-wise response profiles using GPR and to perform the searchlight-based reconstruction using
141 MLE. The simulated searchlights consisted of 241 voxels. We finally evaluated the reconstruction
142 performance by using averaged FCA and BFCA.

## 143 3.2  Simulation results

144 The results of the simulation are summarized in Supplementary Figure 5. We obtained an above-
145 chance reconstruction performance for cases 1 (mean FCA: 92.31%, SD: 1.57; mean BFCA: 91.73%,
146 SD: 1.9) and 2 (mean FCA: 92.72%, SD: 1.09; mean BFCA: 90.99%, SD: 1.91). Interestingly, for
147 case 3 the mean reconstruction performance is around chance for both measures (mean FCA: 50.36%
148, SD: 2.61; mean BFCA: 49.21%, SD: 2.56) whereas for case 4 the distribution of FCA is skewed
149 toward right (mean FCA: 56.23%, SD: 4.25) whereas BFCA values are not (mean BFCA:  46.98%,
150 SD: 6.06), as confirmed by two one-sample right tailed t-tests evaluating whether the mean values
151 were greater than 50% (FCA > 50% : t = 46.278; p < 0.001; BFCA > 50%: t = -15.747; p = 1).

## 152 3.3  Real data analysis

153 We computed the result the whole-brain searchlight analysis for all of the 23 subjects both with FCA
154 and BFCA as measures of reconstruction performance. The maps were obtained following the
155 procedure described in the *Materials and Methods* section. For the purpose of this analysis we only
156 considered three main conditions:

157  1) **Stimulus labels at 100% coherence**. In this condition, the stimulus motion directions are
158      balanced, therefore the reconstruction performance should be above chance only for the
159      searchlights with voxels sensitive to motion directions. Because the distribution of the
160      stimulus directions is balanced, we expect no difference between the reconstruction
161      performance computed with FCA and BFCA.

162  2) **Stimulus labels at 0% coherence**. In this condition, the stimulus had no real motion
163      direction, but each trial was assigned a motion direction, generated according to our
164      randomization scheme (see *Materials and Methods*) This results in a balanced label
165      distribution. Because of this, no searchlight should result in above-chance reconstruction

6

166 performance. Following the outcome of the simulation described above, we expect no
167 difference between the reconstruction performance computed with FCA and BFCA.

168 3) **Report labels at 0% coherence**. Here, the labels assigned to each trial correspond to the
169 motion directions reported by participants. Therefore, the distribution of the reports across
170 trials reflects the idiosyncratic biases of each subject. In this condition, because some
171 participants' choices lead to an unbalanced distribution of reported motion directions, we
172 suspect that FCA leads to spurious above-chance reconstruction performance. Based on the
173 outcome of the simulation, we hypothesize a difference in the reconstruction performance
174 computed with FCA and BFCA.

175 We then used SPM12 to compare the FCA and the BFCA maps of the 23 participants. We performed
176 three second-level paired t-tests to evaluate our hypotheses.
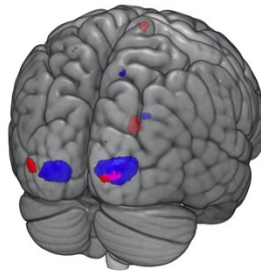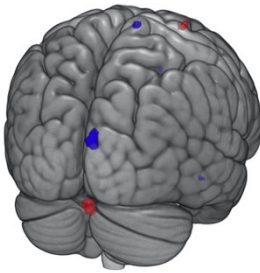
## 177 3.4 Real data results

178 The results are shown in Supplementary Figure 6. The two measures (FCA and BFCA) were not
179 significantly different when using the stimulus labels at 100% coherence and 0% coherence
180 (Appendix 2 - Figure 3). However, the FCA was significantly different from BFCA when using the
181 labels of the 0% reports in a large cluster covering various portions of the brain (FWE$_c$, $p < 0.05$, K =
182 667990; cluster-defining voxel threshold $p < 0.001$). The results remained consistent even when we
183 lowered the cluster-defining threshold to $p < 0.0001$ or $p < 0.00001$, with many significant clusters of
184 smaller size scattered throughout the brain reaching significance level (FWE$_c$, $p < 0.05$, smallest K =
185 70; cluster-defining voxel threshold $p < 0.0001$; FWE$_c$, $p < 0.05$, smallest K = 20; cluster defining
186 voxel threshold $p < 0.00001$).

## 187 4 Supplementary Figures

7

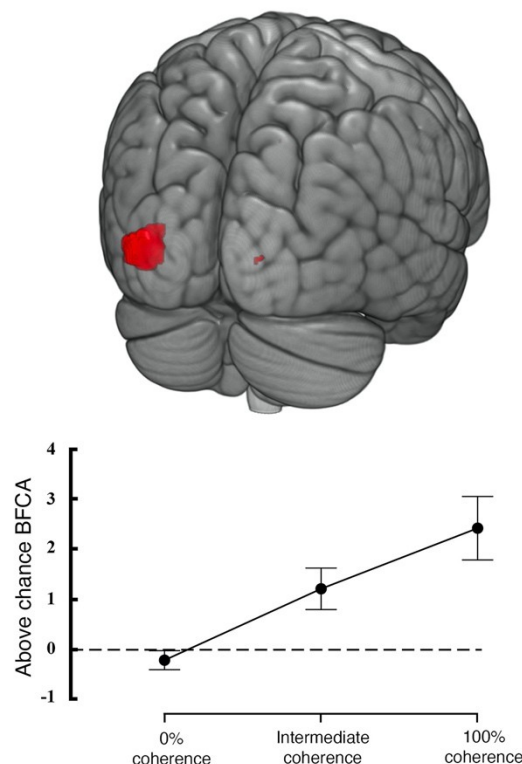0% Coherence          Intermediate Coherence          100% Coherence



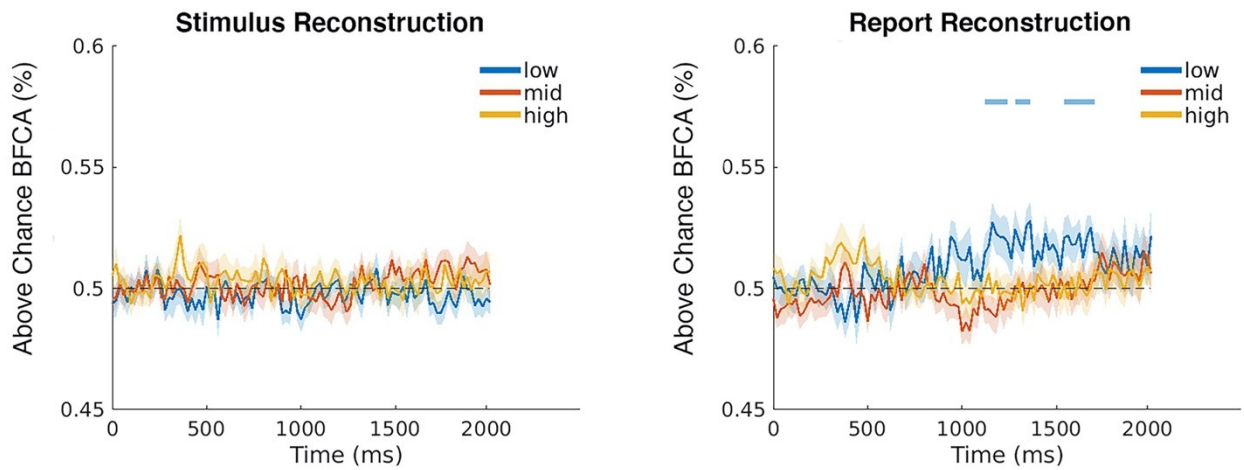● Stimulus          ● Report          ● Overlap

188

**Supplementary Figure 1.** Searchlight-based accuracy maps (plotted here for BFCA, see Methods). The images show results of the searchlight-based stimulus and report reconstructions for three coherence levels (left: 0%; middle: intermediate; right: 100%). The searchlights are mapped with different colors: red indicates significantly above-chance reconstruction performance for stimulus, blue indicate above-chance reconstruction performance for report, and purple indicate the overlap between the two. Please note that the maps are shown for display purposes (for visualization thresholded at $p < 0.001$ and not corrected for multiple comparisons).



196

**Supplementary Figure 2.** Effect of coherence on searchlight-based stimulus reconstruction performance. The picture on top shows clusters of voxels where coherence has a significant effect on stimulus reconstruction performance. The map is thresholded at $p < 0.001$, uncorrected for multiple comparisons. The plot on the bottom displays the averaged above-chance accuracy (BFCA minus baseline of 50%) extracted from the searchlights in which coherence had a significant effect on stimulus reconstruction performance, error bars are standard errors (N=23).
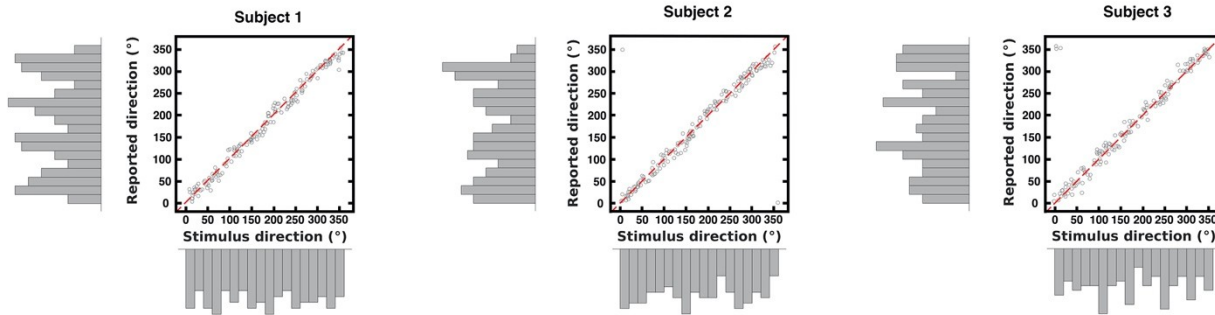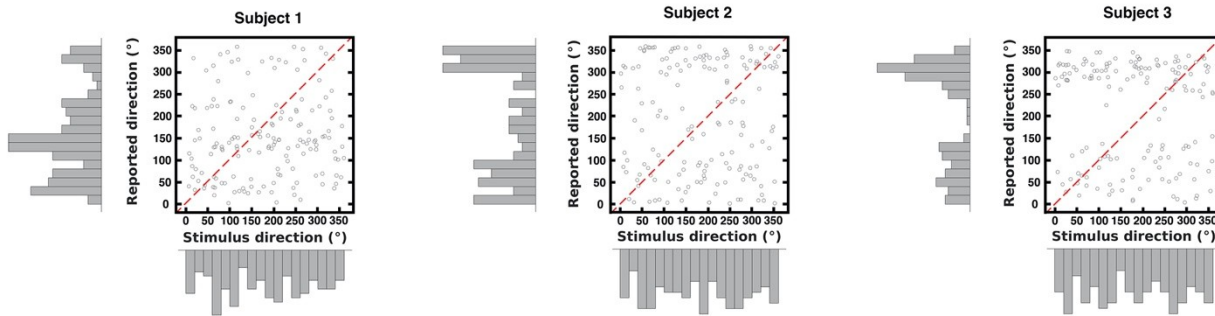
9

203



204

**Supplementary Figure 3.** Group average (N=21) accuracy (expressed as BFCA, see *Materials and Methods*) at each time point relative to stimulus onset. The upper picture displays the stimulus reconstruction performance, the lower picture shows the report reconstruction performance for three coherence levels (blue: 0%; red: intermediate; orange: 100%). Shading around the individual curves indicates ±1 SEM. The light blue lines on top of the curves depict clusters of time points for which the reconstruction performance was greater than chance after correction for multiple comparisons. Please note that GPR estimated from eye-movements are predictive of the reports for the 0% coherence condition but not of those given at intermediate and 100% coherence. Such result together with those of our model consistency and model generalization analyses (see *Results*), suggest that eye movements were unrelated with the brain signals used to reconstruct participants' choices in the 0% coherence condition (see *Results*).
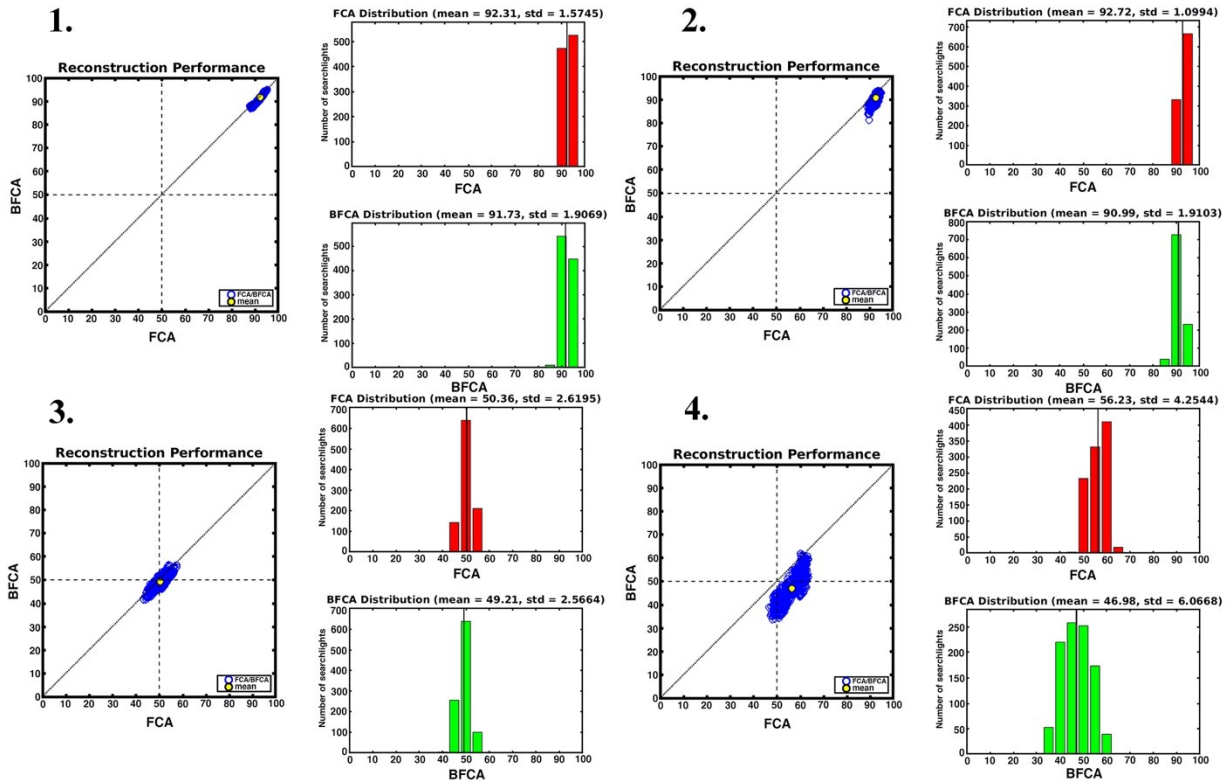
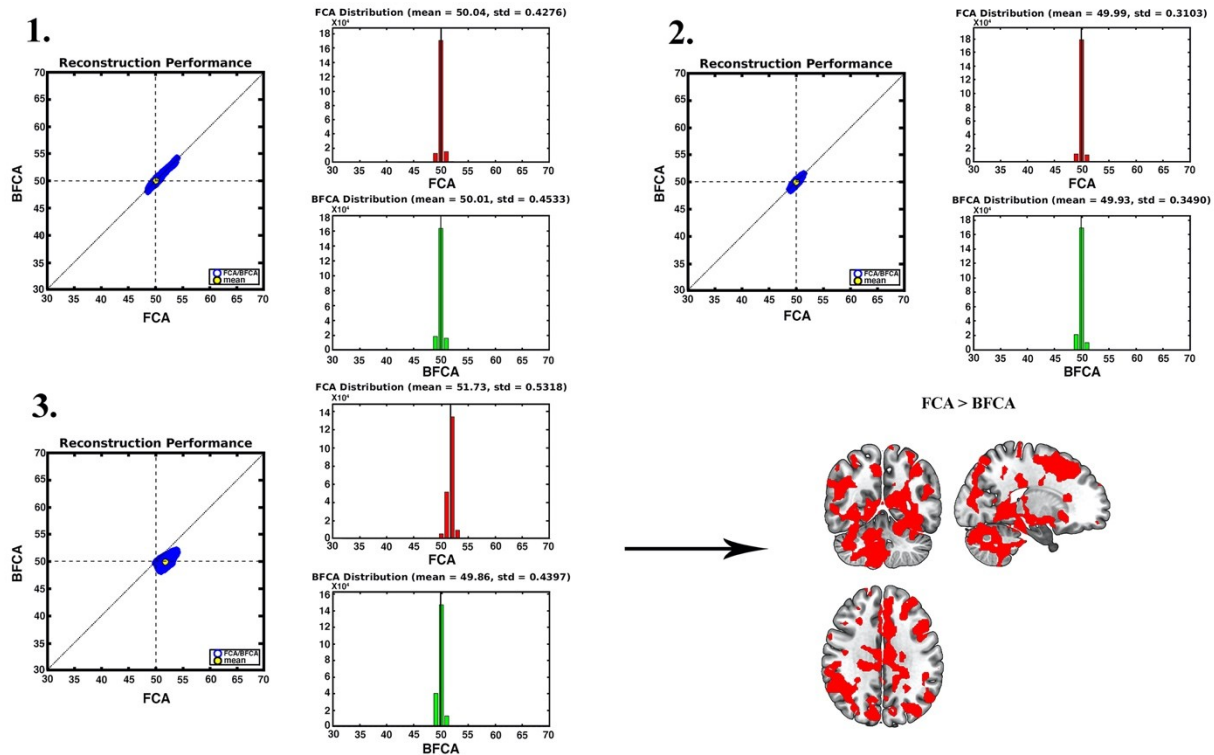216

**A)    100 % coherence**



**B)    0 % coherence**



217

**Supplementary Figure 4.** The scatterplots display the trial-wise reported direction of 3 example participants against the trial-wise motion direction, with the corresponding marginal distributions. A) Each plot on the top row shows the data distribution obtained from 160 trials in the 100% coherence condition. B) The bottom row shows the data distribution for the 0% coherence condition. Note that in this case the motion direction labels were generated following the randomization scheme described in the manuscript (see the *Materials and Methods* section), as no real motion direction was present in the stimulus.

225

**Supplementary Figure 5.** Comparison of reconstruction performances obtained with simulated data. The picture illustrates the distribution of FCA (red) and BFCA (green) in the four scenarios examined in the simulation. 1. The relationship between FCA and BFCA for the condition in which $\theta$ modulates $\hat{\beta}_j$ and the distribution of $\theta$ is balanced. 2. Condition in which $\theta$ modulates $\hat{\beta}_j$ and the distribution of $\theta$ is unbalanced. 3. $\theta$ does not modulate $\hat{\beta}_j$ and the distribution of $\theta$ is balanced. 4. $\theta$ does not modulate $\hat{\beta}_j$ and the distribution of $\theta$ is unbalanced.

232

13

233

**Supplementary Figure 6.** The plots illustrate the difference between reconstruction performances obtained with two accuracy measures, FCA and BFCA. 1-3 show the relationship between FCA and BFCA for the 100% coherence stimulus reconstruction, for the 0% coherence stimulus reconstruction, and for the 0% coherence report reconstruction respectively. Results are plotted for each searchlight, averaged across subjects (N=23). The brain map on the bottom right is obtained from a 2[nd] level t-test evaluating searchlights where FCA was higher than BFCA in the 0% coherence report reconstruction (N=23). The map is thresholded at $p<0.001$, uncorrected for multiple comparisons.