

---

*Commentary | Cognition and Behavior*

## **Inverted Encoding Models Assay Population-Level Stimulus Representations, Not Single-Unit Neural Tuning**

**Thomas C. Sprague<sup>1</sup>, Kirsten C. S. Adam<sup>2</sup>, Joshua J. Foster<sup>2</sup>, Masih Rahmati<sup>1</sup>, David W. Sutterer<sup>2</sup> and Vy A. Vo<sup>3</sup>**

<sup>1</sup>*Department of Psychology, New York University, New York, NY 10003*

<sup>2</sup>*Department of Psychology and Institute for Mind and Biology, University of Chicago, Chicago, IL 60637*

<sup>3</sup>*Neurosciences Graduate Program, University of California, San Diego, La Jolla, CA 92093*

DOI: 10.1523/ENEURO.0098-18.2018

Received: 15 March 2018

Revised: 26 April 2018

Accepted: 3 May 2018

Published: 11 May 2018

---

**Author contributions:** TCS wrote the first draft, KCSA, JJF, MR, DWS and VAV edited the manuscript. All authors jointly developed ideas expressed in this commentary.

**Funding:** <http://doi.org/10.13039/100000053HHS> | NIH | National Eye Institute (NEI)  
F32-EY028438  
R01-EY016407

**Funding:** <http://doi.org/10.13039/100000001National Science Foundation> (NSF)  
Graduate Research Fellowship

**Funding:** <http://doi.org/10.13039/100000025HHS> | NIH | National Institute of Mental Health (NIMH)  
2R01-MH087214-06A1

**Conflict of Interest:** Authors declare no conflict of interest.

Supported by NEI F32-EY028438 (TCS); NSF Graduate Student Fellowship (VAV); NEI R01-EY016407 (MR); and NIMH 2R01-MH087214-06A1 (KCSA, JJF, and DWS).

K.C.S.A., J.J.F., M.R., D.W.S. and V.A.V. are contributed equally and listed alphabetically.

**Correspondence:** Thomas C. Sprague, 6 Washington Place, New York, NY, 10003; E-mail: [tsprague@nyu.edu](mailto:tsprague@nyu.edu)

**Cite as:** eNeuro 2018; 10.1523/ENEURO.0098-18.2018

**Alerts:** Sign up at [eneuro.org/alerts](http://eneuro.org/alerts) to receive customized email alerts when the fully formatted version of this article is published.

Accepted manuscripts are peer-reviewed but have not been through the copyediting, formatting, or proofreading process.

Copyright © 2018 Sprague et al.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

**Title:** Inverted encoding models assay population-level stimulus representations, not single-unit neural tuning

**Abbreviated title:** Assaying population-level stimulus representations

Thomas C. Sprague<sup>1\*</sup>, Kirsten C. S. Adam<sup>2\*</sup>, Joshua J. Foster<sup>2\*</sup>, Masih Rahmati<sup>1\*</sup>, David W. Sutterer<sup>2\*</sup>, and Vy A. Vo<sup>3\*</sup>

<sup>1</sup>Department of Psychology, New York University, New York, NY, 10003

<sup>2</sup>Department of Psychology and Institute for Mind and Biology, University of Chicago, Chicago, IL, 60637

<sup>3</sup>Neurosciences Graduate Program, University of California, San Diego, La Jolla, CA, 92093

+ authors contributed equally and are listed alphabetically

**Author contributions:** TCS wrote the first draft, KCSA, JJF, MR, DWS and VAV edited the manuscript. All authors jointly developed ideas expressed in this commentary.

**Correspondence:**

\* 6 Washington Place, New York, NY, 10003; [tsprague@nyu.edu](mailto:tsprague@nyu.edu) (corresponding author)

**Number of Figures:** 1

**Number of Tables:** 0

**Number of Multimedia:** 0

**Number of words, abstract:** n/a

**Number of words, significance statement:** 120/120

**Number of words, introduction:** n/a

**Number of words, discussion:** n/a

**Total number of words:** 1947

**Acknowledgments:** We wish to thank Clayton Curtis, Edward Ester, and John Serences for comments on early drafts of the manuscript and useful discussions.

**Conflict of interest:** the authors declare no conflict of interest

**Funding sources:** Supported by NEI F32-EY028438 (TCS); NSF Graduate Student Fellowship (VAV); NEI R01-EY016407 (MR); and NIMH 2R01-MH087214-06A1 (KCSA, JJF, and DWS).

32 **SIGNIFICANCE STATEMENT** (120/120 words): Inverted encoding models (IEMs) are a powerful tool  
33 for reconstructing population-level stimulus representations from aggregate measurements of neural  
34 activity (e.g., fMRI or EEG). In a recent report, Liu et al. (2018) tested whether IEMs can provide  
35 information about the underlying tuning of *single units*. Here, we argue that using stimulus  
36 reconstructions to infer properties of single neurons, such as neural tuning bandwidth, is an ill-posed  
37 problem with no unambiguous solution. Instead of interpreting results from these methods as  
38 evidence about single-unit tuning, we emphasize the utility of these methods for assaying *population-*  
39 *level* stimulus representations. These can be compared across task conditions to better constrain  
40 theories of large-scale neural information processing across experimental manipulations, such as  
41 changing sensory input or attention.

42

43 **MAIN TEXT**

44 Neuroscience methods range astronomically in scale. In some experiments, we record  
45 subthreshold membrane potentials in individual neurons, while in others we measure aggregate  
46 responses of thousands of neurons at the millimeter scale. A central goal in neuroscience is to bridge  
47 insights across all scales to understand the core computations underlying cognition (Churchland and  
48 Sejnowski, 1988). However, inferential problems arise when moving across scales: single-unit  
49 response properties cannot be inferred from fMRI activation in single voxels, subthreshold membrane  
50 potential cannot be inferred from extracellular spike rate, and the state of single ion channels cannot  
51 be inferred from intracellular recordings. These are all examples of an *inverse problem* in which an  
52 observation at a larger scale is consistent with an enormous number of possible observations at a  
53 smaller scale.

54 Recent analytical advances have circumvented challenges inherent in inverse problems by  
55 instead transforming aggregate signals from their native ‘measurement’ space (e.g., activation pattern  
56 across fMRI voxels) into a model-based ‘information space’ (e.g., activity level of modeled information  
57 channels). To make this inference possible, aggregate neural signals (fMRI voxel activation or EEG  
58 electrode activity) are modeled as a combination of feature-selective information channels, each with  
59 defined sensitivity profiles consistent with the single-unit literature (e.g., experimenter-defined tuning  
60 to a particular orientation, *Fig. 1A*; Brouwer and Heeger, 2009, 2011). When an aggregate neural  
61 signal is described with such an *encoding model*, it is possible to invert this model to infer the activity  
62 of each channel given a new pattern of neural activity (hence, these methods are often called  
63 ‘inverted encoding models’, IEM; Sprague et al., 2015). Importantly, rather than attempt to solve the  
64 inverse problem (how do single-units respond?), this method makes simplifying assumptions that  
65 enable transformation of one population-level measurement (aggregate neural signals in voxel- or

66 electrode space) into another (stimulus representations in ‘channel space’). These reconstructed  
67 ‘channel response functions’ enable visualization, quantification, and comparison of population-level  
68 stimulus representations across manipulations of task conditions (Brouwer and Heeger, 2011, 2013;  
69 Foster et al., 2017; Garcia et al., 2013; Scolari et al., 2012; Sprague and Serences, 2013).

70 Recently, Liu et al. (2018) examined whether an IEM applied to fMRI data can be used to  
71 unambiguously infer the underlying response properties of single units. To this end, they manipulated  
72 the contrast of orientated gratings, because contrast only affects the amplitude of single-unit  
73 orientation tuning functions, but not their tuning width (e.g. Sclar and Freeman, 1982). The authors  
74 reasoned that, if the width of single-unit tuning functions do not change with stimulus contrast, and if  
75 population-level feature reconstructions derived from aggregate neural signals can be used to make  
76 meaningful inferences about single-unit tuning, then manipulating contrast should not change the  
77 width of population-level channel-response functions.

78 To test this prediction, the authors used an IEM to reconstruct representations of grating  
79 orientations for two different contrast levels. The authors modeled voxel responses as a sum of neural  
80 channels tuned to different orientations based on known visual response properties (*Fig. 1A*). After  
81 extracting activation patterns from visual cortex, the authors split data from each contrast condition  
82 into a training set, used to estimate how each modeled neural channel contributes to each voxel (*Fig.*  
83 *1B*), and a testing set, which was used in conjunction with the best-fit model from the training set to  
84 compute channel response functions (*Fig. 1C*).

85 The authors found that reconstructed channel response functions in visual cortex were ‘broader’  
86 for low-contrast gratings than for high-contrast gratings (Liu et al. *Fig. 2*), which they suggest could be  
87 interpreted as evidence that single-unit orientation tuning width depends on stimulus contrast.  
88 However, because this observation conflicts with demonstrations from single-unit physiology that  
89 orientation tuning is contrast-invariant, Liu et al. (2018) sought to resolve this discrepancy using  
90 simulations.

91 The authors simulated cortical fMRI data under different conditions to assess how changes in  
92 single-unit responses might be reflected in reconstructed channel response functions. Each simulated  
93 voxel’s response was modeled as a noisy weighted sum of orientation-tuned neurons, each with a  
94 different orientation preference (Liu et al., *Fig. 3*). Across runs of their simulations, the authors  
95 manipulated simulated response properties, like orientation tuning width of constituent model neurons  
96 and signal-to-noise ratio (SNR) of the voxel response. The authors found that by decreasing the  
97 response amplitude of each simulated neuron (thus, decreasing SNR) without changing the tuning  
98 width, they could almost exactly reproduce the broadening in the width of the channel response  
99 function when stimulus contrast was decreased (Liu et al., *Fig. 4*). Interestingly, they also found that

100 changes in modeled neural tuning width could alter the width of channel response functions. However,  
101 because such broadening is consistent with either a change in SNR or a change in neural tuning  
102 width, the authors conclude that it remains impossible to conclusively infer how changes in channel  
103 response functions relate to changes in neural tuning. Since it is plausible that low-contrast stimuli  
104 evoke weak, noisy responses relative to high-contrast stimuli, the authors argue this is a more  
105 parsimonious explanation for their observed data than overturning well-characterized results from the  
106 animal physiology literature and inferring that single-unit tuning properties change with contrast.  
107 Accordingly, the authors concluded that “changes in channel response functions do not necessarily  
108 reflect changes in underlying neural selectivity” (Liu et al., 2018, pg 404).

109 This report makes an important contribution in its dissection of how model-based analysis  
110 methods can be sensitive to features of the data that might vary across conditions (e.g., SNR), and  
111 clearly demonstrates that changes in population-level channel response functions cannot and should  
112 not be used to infer changes in unit-level neural tuning properties. However, we would like to  
113 emphasize that this is not the intended purpose of the IEM approach, which is designed to assess  
114 *population-level* stimulus-representations. Any inferences made about single-unit tuning from channel  
115 response functions are plagued by the same pitfalls encountered when attempting reverse inference  
116 about single-unit neural signals from aggregate measurements.

117 These issues are not unique to the IEM technique. For example, they also complicate  
118 interpretation of results from popular voxel receptive field (vRF) techniques. In these experiments,  
119 stimuli traverse the entire visual display while experimenters measure fMRI responses. Then, they fit  
120 a receptive field model that best describes how each voxel responds given the visual stimulus  
121 (Dumoulin and Wandell, 2008; Wandell and Winawer, 2015). Recent studies have demonstrated that  
122 changing task demands (e.g., locus of spatial attention) can change the shape and preferred position  
123 of vRFs (Kay et al., 2015; Klein et al., 2014; Sheremata and Silver, 2015; Sprague and Serences,  
124 2013; Vo et al., 2017). While it is tempting to infer that single-neuron RFs change accordingly, it could  
125 instead be the case that each neuron maintains a stable RF, but different neurons are subject to  
126 different amounts of response gain, altering the voxel-level spatial sensitivity profile measured with  
127 these techniques. Moreover, because aggregate measurements like fMRI pool over neurons of  
128 different types (excitatory vs. inhibitory), selectivity widths (narrow vs. broad), and cortical layers (e.g.,  
129 Layer IV vs. Layer II/III), the ability to make inferences about single-unit encoding properties is further  
130 limited.

131 Liu et al.’s (2018) report also highlights that it is important to consider how an encoding model is  
132 estimated when comparing channel response functions across conditions. In their work, Liu et al.  
133 estimated separate encoding models for each contrast condition (*Fig. 1B*). But because SNR likely  
134 differed between conditions, the observed differences between reconstructions may result from

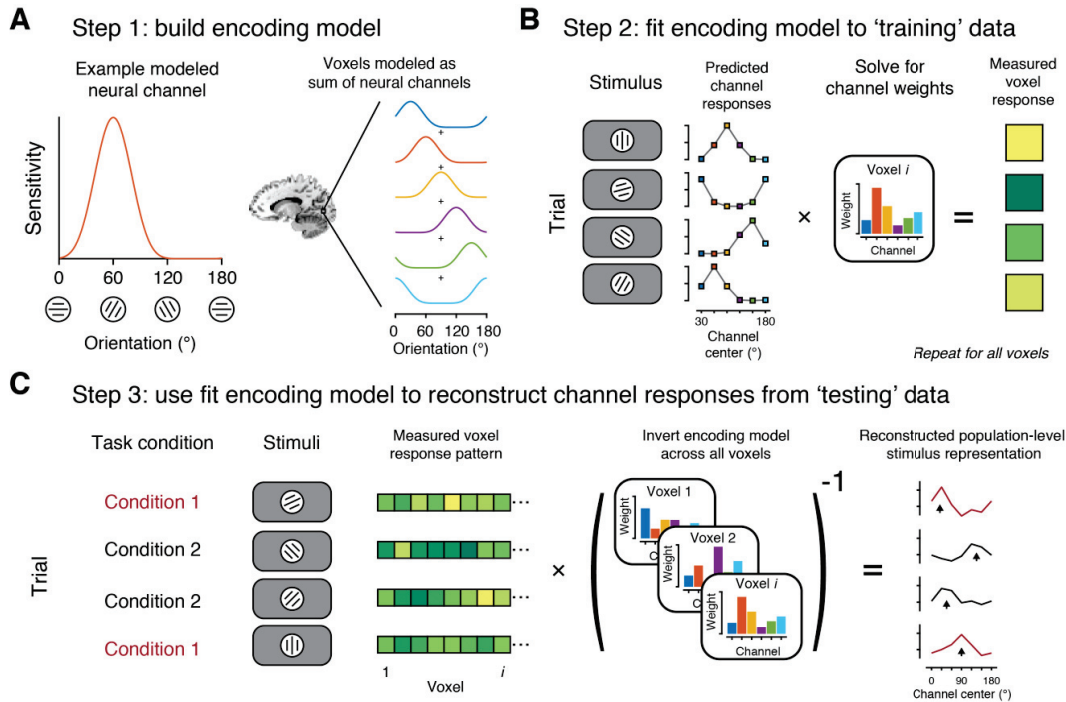
135 differences in the training sets (i.e. different model fits), or from differences in the testing sets (i.e.  
136 different reconstructed activation patterns), or from a combination of the two. More generally, this  
137 training scheme can pose a problem for researchers who wish to minimize the effect of known SNR  
138 differences between their conditions to study some other variable (e.g. the effect of attention) since it  
139 is not possible to unambiguously attribute changes in reconstructed channel response functions to  
140 changes in the quality of the model fit or the quality of the representation supported by the population  
141 activity pattern, which can both differ between conditions. This problem is roughly akin to reporting a  
142 change in a ratio, which can result in changes in the numerator, denominator, or both. One way that  
143 others have mitigated this issue is by estimating an encoding model (*Fig. 1B*) using an unbiased  
144 (equal numbers of trials from each relevant condition) or neutral (entirely separate task used solely for  
145 model estimation) set of data. They then apply that single 'fixed' encoding model to test data from  
146 multiple stimulus conditions to reconstruct stimulus representations from each condition. This  
147 implementation has the advantage that researchers can avoid problems with comparing channel  
148 outputs from different IEMs, so the only difference between conditions is the data used for stimulus  
149 reconstruction (*Fig. 1C*). We note that even with such a procedure the central result in Liu et al (2018)  
150 could remain true: reconstructions under a fixed encoding model could still broaden with lower  
151 contrast. But, as discussed above, this would reflect a change in the quality of the population-level  
152 representation rather than provide unambiguous evidence for a change in underlying tuning of  
153 individual units. When interpreting results from IEM analyses, it is always critical to consider how the  
154 model was estimated.

155 It would be a mistake to conclude from Liu et al. (2018) that the IEM technique is not useful in the  
156 context of its intended purpose: to assay properties of large-scale, population-level neural  
157 representations. The quality of these large-scale representations surely depends on myriad factors  
158 occurring at the single-unit level. It remains a fascinating question to evaluate how single  
159 measurement units, at either the neural or voxel level, change their response properties across visual  
160 and task manipulations, but the goal of the IEM approach is to assay the net effect of all these  
161 modulations on the superordinate population-level representation. Moreover, few behaviors are  
162 guided by single neurons in isolation, and so assaying the joint activity of many neurons, and the  
163 resulting population-level representations, is necessary to gain insight into the neural underpinnings of  
164 cognition (Graf et al., 2011; Jazayeri and Movshon, 2006; Ma et al., 2006). Indeed, IEMs have been  
165 used to assay the time course of covert attention (Foster et al., 2017), understand the consequences  
166 of attentional manipulations within WM (Rahmati et al., 2018; Sprague et al., 2016), evaluate how  
167 allocation of attention impacts the representation of irrelevant visual stimuli across the visual field  
168 (Sprague et al., 2018; Sprague and Serences, 2013; Vo et al., 2017), and probe the influence of top-  
169 down expectations on sensory stimulus representations (Kok et al., 2017; Myers et al., 2015).



170 We do not believe aggregate neural signals will ever be useful for unambiguously inferring single-  
171 unit response properties, including feature tuning. However, we see a bright future for collaborative  
172 efforts across labs studying similar questions in different model systems, such as human and  
173 macaque. When experiments are well-matched between species, both aggregate measurements in  
174 humans and single-unit responses in model systems can be used to inform our understanding of  
175 neural coding across different cognitive states. In bridging different levels of analysis, Liu et al. (2018)  
176 add to the growing literature using data-driven simulations to better understand the relationship  
177 between tuning properties and population-level feature representations (Kay et al., 2015; Sprague and  
178 Serences, 2013; Vo et al., 2017). Most importantly, their report underscores the importance of  
179 avoiding inferences about signal properties, such as single-unit neural feature tuning, that are  
180 fundamentally inaccessible via fMRI or EEG, even when using state-of-the-art acquisition and  
181 analysis techniques. We hope that future studies take these issues into account when interpreting  
182 findings from model-based analyses applied to aggregate measurement tools like fMRI and EEG.  
183 Finally, we remain optimistic that the IEM technique, when applied carefully and interpreted  
184 appropriately, will continue to reveal how experimental manipulations impact population-level  
185 representations of information.

186



187

188 **Figure 1. Inverted Encoding Model (IEM): use neural tuning as an assumption to estimate**  
 189 **population-level representations**

190 A. The IEM framework assumes that aggregate neural responses (e.g., voxels) can be modeled as a  
 191 combination of feature-selective information channels (i.e., orientation-selective neural populations).  
 192 Tuning properties of modeled information channels are experimenter-defined, and often based on  
 193 findings in the single-unit physiology literature. B. Once an encoding model (A) is defined, it can be  
 194 used to predict how each information channel should respond to each stimulus in the experiment.  
 195 These predicted channel responses are used to fit the encoding model to each voxel's activation  
 196 across all trials in a 'training' dataset, often balanced across experimental conditions, or derived from  
 197 a separate 'localizer' or 'mapping' task. C. By inverting the encoding models estimated across all  
 198 voxels (typically, within an independently-defined region), new activation patterns can be used to  
 199 compute the response of each modeled neural information channel. This step transforms activation  
 200 patterns from 'measurement space' (1 number per measurement dimension, e.g., voxel) to  
 201 'information space' (1 number per modeled information channel, (A)). These computed 'channel  
 202 response functions' can be aligned based on the known stimulus feature value on each trial (black  
 203 arrowheads), and quantified and compared across conditions (e.g., manipulations of stimulus  
 204 contrast, spatial attention, etc), especially when a fixed encoding model is used for reconstruction (as  
 205 schematized here). Cartoon data shown throughout figure.



206  
207**REFERENCES:**

- 208 Brouwer G, Heeger D (2011) Cross-orientation suppression in human visual cortex. *J Neurophysiol*  
209 106:2108–2119.
- 210 Brouwer G, Heeger D (2009) Decoding and Reconstructing Color from Responses in Human Visual  
211 Cortex. *J Neurosci* 29:13992–14003.
- 212 Brouwer GJ, Heeger DJ (2013) Categorical clustering of the neural representation of color. *J Neurosci*  
213 33:15454–65.
- 214 Churchland PS, Sejnowski TJ (1988) Perspectives on cognitive neuroscience. *Science* 242:741–5.
- 215 Dumoulin S, Wandell B (2008) Population receptive field estimates in human visual cortex.  
216 *Neuroimage* 39:647–660.
- 217 Foster JJ, Sutterer DW, Serences JT, Vogel EK, Awh E (2017) Alpha-Band Oscillations Enable  
218 Spatially and Temporally Resolved Tracking of Covert Spatial Attention. *Psychol Sci* 28:929–941.
- 219 Garcia J, Srinivasan R, Serences J (2013) Near-Real-Time Feature-Selective Modulations in Human  
220 Cortex. *Curr Biol* 23:515–522.
- 221 Graf ABA, Kohn A, Jazayeri M, Movshon JA (2011) Decoding the activity of neuronal populations in  
222 macaque primary visual cortex. *Nat Neurosci* 14:239–245.
- 223 Jazayeri M, Movshon JA (2006) Optimal representation of sensory information by neural populations.  
224 *Nat Neurosci* 9:690–696.
- 225 Kay KN, Weiner KS, Grill-Spector K (2015) Attention reduces spatial uncertainty in human ventral  
226 temporal cortex. *Curr Biol*.
- 227 Klein BP, Harvey BM, Dumoulin SO (2014) Attraction of position preference by spatial attention  
228 throughout human visual cortex. *Neuron* 84:227–37.
- 229 Kok P, Mostert P, de Lange FP (2017) Prior expectations induce prestimulus sensory templates. *Proc*  
230 *Natl Acad Sci U S A* 114:10473–10478.
- 231 Liu T, Cable D, Gardner JL (2018) Inverted Encoding Models of Human Population Response  
232 Conflate Noise and Neural Tuning Width. *J Neurosci* 38:398–408.
- 233 Ma WJ, Beck JM, Latham PE, Pouget A (2006) Bayesian inference with probabilistic population  
234 codes. *Nat Neurosci* 9:1432–8.
- 235 Myers NE, Rohenkohl G, Wyart V, Woolrich MW, Nobre AC, Stokes MG (2015) Testing sensory  
236 evidence against mnemonic templates. *Elife* 4:e09000.

- 237 Rahmati M, Saber GT, Curtis CE (2018) Population Dynamics of Early Visual Cortex during Working  
238 Memory. *J Cogn Neurosci* 30:219–233.
- 239 Sclar G, Freeman RD (1982) Orientation selectivity in the cat's striate cortex is invariant with stimulus  
240 contrast. *Exp brain Res* 46:457–61.
- 241 Scolari M, Byers A, Serences JT (2012) Optimal Deployment of Attentional Gain during Fine  
242 Discriminations. *J Neurosci* 32:1–11.
- 243 Sheremata SL, Silver MA (2015) Hemisphere-Dependent Attentional Modulation of Human Parietal  
244 Visual Field Representations. *J Neurosci* 35:508–517.
- 245 Sprague TC, Ester EF, Serences JT (2016) Restoring Latent Visual Working Memory  
246 Representations in Human Cortex. *Neuron* 91:694–707.
- 247 Sprague TC, Itthipuripat S, Vo VA, Serences JT (2018) Dissociable signatures of visual salience and  
248 behavioral relevance across attentional priority maps in human cortex. *J Neurophysiol*  
249 jn.00059.2018.
- 250 Sprague TC, Saproo S, Serences JT (2015) Visual attention mitigates information loss in small- and  
251 large-scale neural codes. *Trends Cogn Sci* 19:215–26.
- 252 Sprague TC, Serences JT (2013) Attention modulates spatial priority maps in the human occipital,  
253 parietal and frontal cortices. *Nat Neurosci* 16:1879–87.
- 254 Vo VA, Sprague TC, Serences JT (2017) Spatial Tuning Shifts Increase the Discriminability and  
255 Fidelity of Population Codes in Visual Cortex. *J Neurosci* 37:3386–3401.
- 256 Wandell BA, Winawer J (2015) Computational neuroimaging and population receptive fields. *Trends*  
257 *Cogn Sci* 19:349–357.
- 258

