

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>

DOI:<http://dx.doi.org/10.1523/ENEURO.0098-18.2018>

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

**Key words:** cognitive vision; computational neuroimaging; fMRI; inverted encoding model

## Significance Statement

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

Neuroscience methods range astronomically in scale. In some experiments, we record subthreshold membrane

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

Recent analytical advances have circumvented challenges inherent in inverse problems by instead transforming aggregate signals from their native “measurement space” (e.g., activation pattern across fMRI voxels) into a model-based “information space” (e.g., activity level of modeled information channels). To make this inference possible, aggregate neural signals (fMRI voxel activation

Received March 15, 2018; accepted May 3, 2018; First published May 11, 2018.

The authors declare no competing financial interests.

Author contributions: T.C.S., K.C.S.A., J.J.F., M.R., D.W.S., and V.A.V. wrote the paper.

This work was supported by the National Eye Institute (NEI) Grant F32-EY028438 (to T.C.S.), a National Science Foundation Graduate Student Fellowship (to V.A.V.), the NEI Grant R01-EY016407 (M.R.), and the National Institute of Mental Health Grant 2R01-MH087214-06A1 (K.C.S.A., J.J.F., and D.W.S.).

\*K.C.S.A., J.J.F., M.R., D.W.S., and V.A.V. contributed equally to this work.

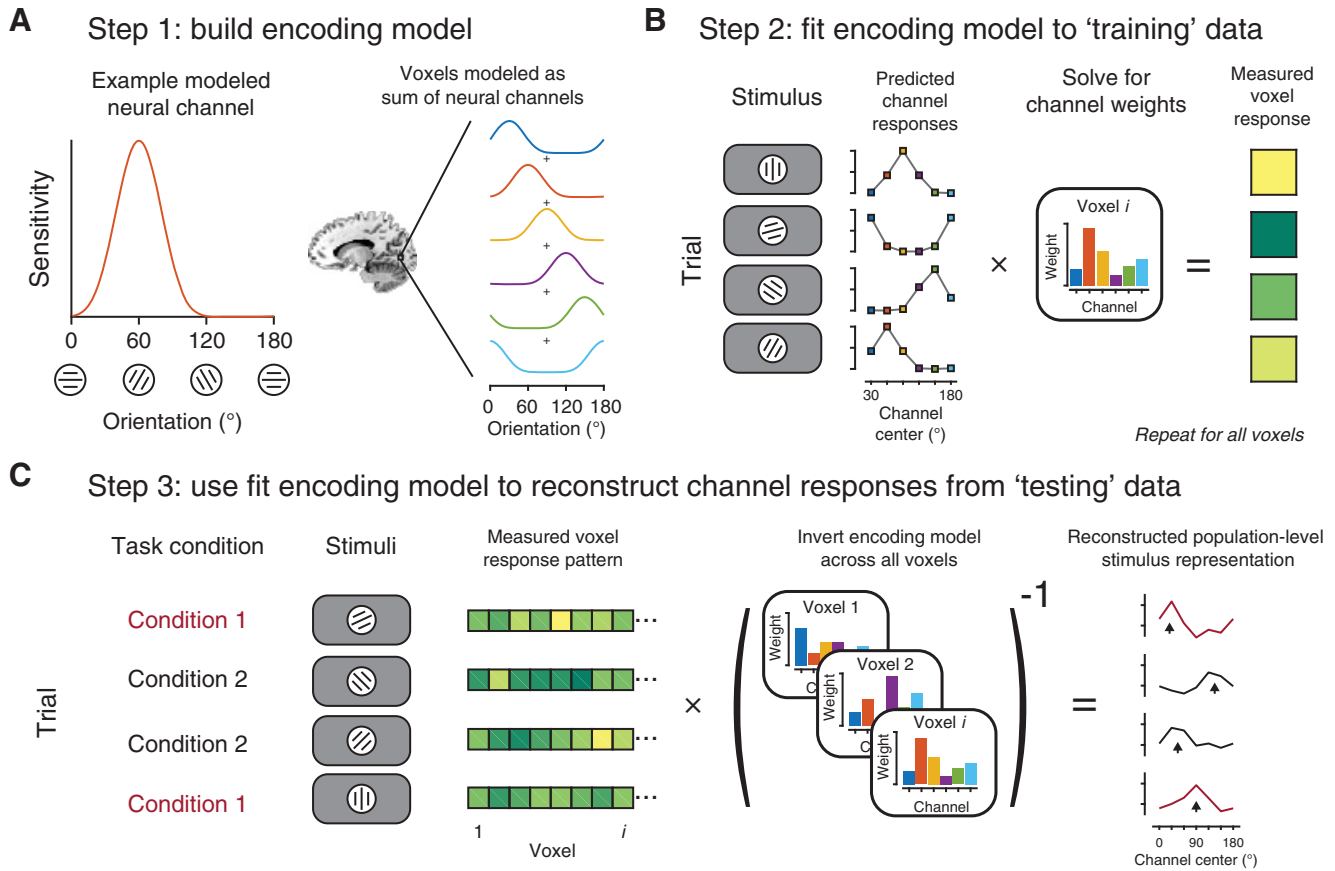
Acknowledgments: We thank Clayton Curtis, Edward Ester, and John Serences for comments on early drafts of this manuscript and useful discussions.

Correspondence should be addressed to Thomas C. Sprague, Department of Psychology, New York University, 6 Washington Place, New York, NY, 10003. E-mail: [tsprague@nyu.edu](mailto:tsprague@nyu.edu).

DOI:<http://dx.doi.org/10.1523/ENEURO.0098-18.2018>

Copyright © 2018 Sprague et al.

This is an open-access article distributed under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.



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

or EEG electrode activity) are modeled as a combination of feature-selective information channels, each with defined sensitivity profiles consistent with the single-unit literature (e.g., experimenter-defined tuning to a particular orientation; Fig. 1A; Brouwer and Heeger, 2009, 2011). When an aggregate neural signal is described with such an encoding model, it is possible to invert this model to infer the activity of each channel given a new pattern of neural activity [hence, these methods are often called inverted encoding models (IEMs); Sprague et al., 2015]. Importantly, rather than attempt to solve the inverse problem (how do single-units respond?), this method makes simplifying assumptions that enable transformation of one population-level measurement (aggregate neural signals in voxel or electrode space) into another (stimulus representations in “channel space”). These reconstructed “channel response functions” enable visualization, quan-

tification, and comparison of population-level stimulus representations across manipulations of task conditions (Brouwer and Heeger, 2011, 2013; Scolari et al., 2012; Garcia et al., 2013; Sprague and Serences, 2013; Foster et al., 2017).

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

should not change the width of population-level channel-response functions.

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

The authors found that reconstructed channel response functions in visual cortex were “broader” for low-contrast gratings than for high-contrast gratings (Fig. 2–4; Liu et al., 2018), which they suggest could be interpreted as evidence that single-unit orientation tuning width depends on stimulus contrast. However, because this observation conflicts with demonstrations from single-unit physiology that orientation tuning is contrast-invariant, Liu et al. (2018) sought to resolve this discrepancy using simulations.

The authors simulated cortical fMRI data under different conditions to assess how changes in single-unit responses might be reflected in reconstructed channel response functions. Each simulated voxel’s response was modeled as a noisy weighted sum of orientation-tuned neurons, each with a different orientation preference (Liu et al., 2018, their Fig. 3). Across runs of their simulations, the authors manipulated simulated response properties, like orientation tuning width of constituent model neurons and signal-to-noise ratio (SNR) of the voxel response. The authors found that by decreasing the response amplitude of each simulated neuron (thus, decreasing SNR) without changing the tuning width, they could almost exactly reproduce the broadening in the width of the channel response function when stimulus contrast was decreased (Liu et al., 2018, their Fig. 4). Interestingly, they also found that changes in modeled neural tuning width could alter the width of channel response functions. However, because such broadening is consistent with either a change in SNR or a change in neural tuning width, the authors conclude that it remains impossible to conclusively infer how changes in channel response functions relate to changes in neural tuning. Since it is plausible that low-contrast stimuli evoke weak, noisy responses relative to high-contrast stimuli, the authors argue this is a more parsimonious explanation for their observed data than overturning well-characterized results from the animal physiology literature and inferring that single-unit tuning properties change with contrast. Accordingly, the authors concluded that “changes in channel response functions do not necessarily reflect changes in underlying neural selectivity” (Liu et al., 2018, p 404).

This report makes an important contribution in its dissection of how model-based analysis methods can be sensitive to features of the data that might vary across conditions (e.g., SNR), and clearly demonstrates that

changes in population-level channel response functions cannot and should not be used to infer changes in unit-level neural tuning properties. However, we would like to emphasize that this is not the intended purpose of the IEM approach, which is designed to assess population-level stimulus representations. Any inferences made about single-unit tuning from channel response functions are plagued by the same pitfalls encountered when attempting reverse inference about single-unit neural signals from aggregate measurements.

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

Liu et al. (2018)’s report also highlights that it is important to consider how an encoding model is estimated when comparing channel response functions across conditions. In their work, Liu et al. (2018) estimated separate encoding models for each contrast condition (Fig. 1B). But because SNR likely differed between conditions, the observed differences between reconstructions may result from differences in the training sets (i.e., different model fits), or from differences in the testing sets (i.e., different reconstructed activation patterns), or from a combination of the two. More generally, this training scheme can pose a problem for researchers who wish to minimize the effect of known SNR differences between their conditions to study some other variable (e.g., the effect of attention), since it is not possible to unambiguously attribute changes in reconstructed channel response functions to changes in the quality of the model fit or the quality of the representation supported by the population activity pattern, which can both differ between conditions. This problem is roughly akin to reporting a change in a ratio, which can result from changes in the numerator, denominator, or both. One way that others have mitigated this issue is by estimating an encoding model (Fig. 1B) using an unbiased (equal numbers of trials from each relevant condition) or neutral (entirely separate task used solely for model estimation) set of data. They then apply that single “fixed” encoding model to test data from multiple stimulus

conditions to reconstruct stimulus representations from each condition. This implementation has the advantage that researchers can avoid problems with comparing channel outputs from different IEMs, so the only difference between conditions is the data used for stimulus reconstruction (Fig. 1C). We note that even with such a procedure the central result in Liu et al. (2018) could remain true: reconstructions under a fixed encoding model could still broaden with lower contrast. But, as discussed above, this would reflect a change in the quality of the population-level representation rather than provide unambiguous evidence for a change in underlying tuning of individual units. When interpreting results from IEM analyses, it is always critical to consider how the model was estimated.

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

We do not believe aggregate neural signals will ever be useful for unambiguously inferring single-unit response properties, including feature tuning. However, we see a bright future for collaborative efforts across labs studying similar questions in different model systems, such as human and macaque. When experiments are well-matched between species, both aggregate measurements in humans and single-unit responses in model systems can be used to inform our understanding of neural coding across different cognitive states. In bridging different levels of analysis, Liu et al. (2018) add to the growing literature using data-driven simulations to better understand the relationship between tuning properties and population-level feature representations (Sprague and Serences, 2013; Kay et al., 2015; Vo et al., 2017). Most importantly, their report underscores the importance of avoiding inferences about signal properties, such as single-unit neural feature tuning, that are fundamentally inaccessible via fMRI or EEG, even when using state-of-the-art acquisition and analysis techniques. We hope that future studies take these issues into account when interpret-

ing findings from model-based analyses applied to aggregate measurement tools like fMRI and EEG. Finally, we remain optimistic that the IEM technique, when applied carefully and interpreted appropriately, will continue to reveal how experimental manipulations impact population-level representations of information.

## References

- Brouwer G, Heeger D (2011) Cross-orientation suppression in human visual cortex. *J Neurophysiol* 106:2108–2119. [CrossRef Medline](#)
- Brouwer G, Heeger D (2009) Decoding and reconstructing color from responses in human visual cortex. *J Neurosci* 29:13992–14003. [CrossRef Medline](#)
- Brouwer GJ, Heeger DJ (2013) Categorical clustering of the neural representation of color. *J Neurosci* 33:15454–15465. [CrossRef Medline](#)
- Churchland PS, Sejnowski TJ (1988) Perspectives on cognitive neuroscience. *Science* 242:741–745. [Medline](#)
- Dumoulin S, Wandell B (2008) Population receptive field estimates in human visual cortex. *Neuroimage* 39:647–660. [CrossRef Medline](#)
- Foster JJ, Sutterer DW, Serences JT, Vogel EK, Awh E (2017) Alpha-band oscillations enable spatially and temporally resolved tracking of covert spatial attention. *Psychol Sci* 28:929–941. [CrossRef Medline](#)
- Garcia J, Srinivasan R, Serences J (2013) Near-real-time feature-selective modulations in human cortex. *Curr Biol* 23:515–522. [CrossRef Medline](#)
- Graf ABA, Kohn A, Jazayeri M, Movshon JA (2011) Decoding the activity of neuronal populations in macaque primary visual cortex. *Nat Neurosci* 14:239–245. [CrossRef Medline](#)
- Jazayeri M, Movshon JA (2006) Optimal representation of sensory information by neural populations. *Nat Neurosci* 9:690–696. [CrossRef Medline](#)
- Kay KN, Weiner KS, Grill-Spector K (2015) Attention reduces spatial uncertainty in human ventral temporal cortex. *Curr Biol* 25:595–600.
- Klein BP, Harvey BM, Dumoulin SO (2014) Attraction of position preference by spatial attention throughout human visual cortex. *Neuron* 84:227–237. [CrossRef Medline](#)
- Kok P, Mostert P, de Lange FP (2017) Prior expectations induce prestimulus sensory templates. *Proc Natl Acad Sci USA* 114:10473–10478. [CrossRef Medline](#)
- Liu T, Cable D, Gardner JL (2018) Inverted encoding models of human population response conflate noise and neural tuning width. *J Neurosci* 38:398–408. [CrossRef Medline](#)
- Ma WJ, Beck JM, Latham PE, Pouget A (2006) Bayesian inference with probabilistic population codes. *Nat Neurosci* 9:1432–1438. [CrossRef Medline](#)
- Myers NE, Rohenkohl G, Wyart V, Woolrich MW, Nobre AC, Stokes MG (2015) Testing sensory evidence against mnemonic templates. *Elife* 4:e09000. [CrossRef Medline](#)
- Rahmati M, Saber GT, Curtis CE (2018) Population dynamics of early visual cortex during working memory. *J Cogn Neurosci* 30:219–233. [CrossRef Medline](#)
- Sclar G, Freeman RD (1982) Orientation selectivity in the cat's striate cortex is invariant with stimulus contrast. *Exp Brain Res* 46:457–461. [CrossRef](#)
- Scolari M, Byers A, Serences JT (2012) Optimal deployment of attentional gain during fine discriminations. *J Neurosci* 32:1–11. [CrossRef](#)
- Sheremata SL, Silver MA (2015) Hemisphere-dependent attentional modulation of human parietal visual field representations. *J Neurosci* 35:508–517. [CrossRef Medline](#)
- Sprague TC, Serences JT (2013) Attention modulates spatial priority maps in the human occipital, parietal and frontal cortices. *Nat Neurosci* 16:1879–1887. [CrossRef Medline](#)

- Sprague TC, Saproo S, Serences JT (2015) Visual attention mitigates information loss in small- and large-scale neural codes. *Trends Cogn Sci* 19:215–226. [CrossRef](#)
- Sprague TC, Ester EF, Serences JT (2016) Restoring latent visual working memory representations in human cortex. *Neuron* 91:694–707. [CrossRef](#) [Medline](#)
- Sprague TC, Itthipuripat S, Vo VA, Serences JT (2018) Dissociable signatures of visual salience and behavioral relevance across attentional priority maps in human cortex. *J Neurophysiol*. Advance online publication. Retrieved May 16, 2018. [CrossRef](#)
- Vo VA, Sprague TC, Serences JT (2017) Spatial tuning shifts increase the discriminability and fidelity of population codes in visual cortex. *J Neurosci* 37:3386–3401. [CrossRef](#) [Medline](#)
- Wandell BA, Winawer J (2015) Computational neuroimaging and population receptive fields. *Trends Cogn Sci* 19:349–357. [Cross-Ref](#)