

¹ Reverse Double-Dipping: When Data Dips You, Twice—Stimulus-Driven ²

³ Information Leakage in Naturalistic ⁴

⁵ Neuroimaging

⁵ **Seung-Goo Kim¹**

*For correspondence:

seung-goo.kim@ae.mpg.de (SGK)

⁶ ¹Research Group Neurocognition of Music and Language, Max Planck Institute for
⁷ Empirical Aesthetics, Grüneburgweg 14, 60322 Frankfurt, Germany

⁸

⁹ **Abstract** This article elucidates a methodological pitfall of cross-validation for evaluating
¹⁰ predictive models applied to naturalistic neuroimaging data—namely, ‘reverse double-dipping’
¹¹ (RDD). In a broader context, this problem is also known as ‘leakage in training examples’, which is
¹² difficult to detect in practice. RDD can occur when predictive modeling is applied to data from a
¹³ conventional neuroscientific design, characterized by a limited set of stimuli repeated across trials
¹⁴ and/or participants. It results in spurious predictive performances due to overfitting to repeated
¹⁵ signals, even in the presence of independent noise. Through comprehensive simulations and
¹⁶ real-world examples following theoretical formulation, the article underscores how such
¹⁷ information leakage can occur and how severely it could compromise the results and conclusions
¹⁸ when it is combined with widely spread informal reverse inference. The article concludes with
¹⁹ practical recommendations for researchers to avoid RDD in their experiment design and analysis.

²⁰

²¹ Introduction

²² Recent advancement of ‘naturalistic neuroimaging’ has opened up new exciting developments in
²³ various fields of human neuroscience (**Sonkusare et al., 2019; Nastase et al., 2020; Hamilton and**
²⁴ **Huth, 2020**). The key idea of naturalistic neuroimaging is that experiments with naturalistic (i.e.,
²⁵ real-world) stimuli are essential when investigating complex human behaviors. This latest rein-
²⁶ carnation of ecological psychology (**Brunswik, 1943; Gibson, 1978**) has gained wide popularity in
²⁷ psychology and human neuroscience. In particular, this idea has been widely adopted in domains
²⁸ where high-order cognitive and/or affective processes are involved, and a simple contrastive ex-
²⁹ perimental approach can explain only little. For example, comparing brain responses to *music*
³⁰ vs. *non-music* to find neural correlates of “music perception” may be under an overly reduction-

31 ist assumption (i.e., “music-as-fixed-effect” fallacy; *Kim, 2022*) that the human brain is governed by
32 simple, interpretable rules that can extrapolate to explain complex behaviors (for more discussion,
33 see *Nastase et al., 2020*).

34 Currently, one of the most popular frameworks in analyzing naturalistic neuroimaging data is
35 known as ‘linearized encoding analysis’ (for a comprehensive review of the encoding models in the
36 music domain, see *Kim, 2022*). This is a method that identifies how of a time-invariant system (e.g.,
37 the human brain or an artificial neural network [ANN]) transforms information from a stimulus to
38 a response, by separating separates the whole transformation into non-linear mapping and linear
39 mapping (*Wu et al., 2006; Naselaris et al., 2011*). Often, the non-linear mapping (‘linearization’) is
40 given by the researcher’s hypothesis (i.e., “*the human language system utilizes information X and Y but*
41 *not Z*”) and the linear mapping is found from the given linearization and human responses. This
42 approach has been more widely used in the era of deep neural networks. For example, studies
43 have shown the potentials of transfer learning of image classification models (*Allen et al., 2022*) or
44 large-language models (*Caucheteux et al., 2023*) to explain human behaviors and neural activity.

45 Besides encoding analysis, model-free approaches have been popularized. For example, the
46 synchrony of neural activity over time across multiple participants when watching an identical nat-
47 uralistic stimulus (i.e., the same movie) has been suggested to quantify stimulus-driven effects (*Has-*
48 *son et al., 2004*). Without explicit modeling of encoded information, the similarity in responses over
49 multiple repetitions can attribute repeated responses to the repeated stimuli—no matter which
50 information is being processed. The analysis of neural synchrony draws conclusions from strong
51 correlations between the neural responses of different participants or trials while presenting iden-
52 tical stimuli that the certain (but unspecified) information in the stimulus evoked the time-locked
53 neural responses.

54 More recently, data-driven segmentation algorithms such as hidden markov models have been
55 proposed to model high-level concepts such as *event boundaries* in narratives (*Baldassano et al.,*
56 *2018*). For model-free approaches, the repetition of identical stimuli is the key reference of the anal-
57 ysis. Thus, many naturalistic neuroscience experiments inspired by such approaches presented
58 identical stimuli to multiple participants.

59 More broadly, stimulus repetition remains a cornerstone of experimental design in neuroscience.
60 From the long tradition of event-related potential (ERP) experiments and the block-design in func-
61 tional magnetic resonance imaging (fMRI), repetition has been the gold standard for cancelling out
62 non-time-locked noise and isolating the time-locked signal.

63 However, in encoding analysis, time-locked responses to repeated stimuli, even with indepen-
64 dent noise, could introduce **information leakage**, disabling regularization, ultimately leading to a
65 false conclusion where irrelevant information is mistaken as relevant to human neural processes.
66 This paper explains how such a fallacy could occur in the Theory section, identifies contributing fac-
67 tors based on simulations in the Simulation section, and demonstrates real-world cases in the Real
68 Data section. Finally, implications for future analyses and experiments are discussed and practical
69 recommendations are given in the Discussion section.

70 **Theory**

71 **Types of information leakage**

72 In essence, information leakage is a circular fallacy. Leakage in data mining and machine learning
73 (or more generally, in predictive modeling) has been categorized into multiple cases: (i) *leaking
74 features*, (ii) *leakage by design decisions*, (iii) *leakage in training examples* (*Kaufman et al., 2012*).

75 In data-mining competitions, the training data (a set of feature-target pairs) and the test (hold-
76 out) data (an independent set of feature-target pairs) are separated by the organizer and only the
77 training data and test features are provided to contenders. Contenders create their predictive
78 models and submit their predictions on the test targets so that the organizer can compare their
79 hold-out validation performances objectively, under the assumption that the curated data is highly
80 representative of real-world problems.

81 In cross-validation, one individual researcher plays both the role of an *organizer* (i.e., partition-
82 ing data at hand into training and test sets) and that of a *contender* (i.e., modeling associations
83 between feature and response in the training set and predicting responses from features in the
84 test set). The *leaking features* could be information that is seemingly random but contains highly
85 predictive information by accidental association (or overlooked failure of dissociation) that exist
86 in both the training and test sets. Although in natural science no human organizer will create a
87 dataset, mistaking confounding features as predictive or even causal is not unheard of (*Nuzzo,
88 2015*).

89 An once prevalent—but still not uncommon—case of the circular fallacy in neuroscience is a
90 selective analysis based on the same data, which is widely known as *double-dipping* (*Kriegeskorte
91 et al., 2009*). This can be understood as *leakage by design decisions*. Double-dipping can also occur
92 in the context of cross-validation. For example, selecting features from the entire dataset prior to
93 splitting (i.e., using information from both of a training set and a test set) would leak information
94 (including noise) from the test set to the training set. In this case, it is you who dips the data into
95 analysis twice.

96 The last kind, *leakage in training exam-
97 ples*, occurs due to an incidental similarity
98 between training and test sets. This is also
99 known as *twinning*: i.e., having one individ-
100 ual participant in a training set and their
101 twin in a test set. If a researcher fails
102 to recognize the true relationship between
103 the twins, any random association stemming
104 from their shared characteristics may spu-
105 riously appear predictive. In this case, *it
106 is the data that dips you, twice*—thus, it can
107 be called **reverse double-dipping** (RDD; *Fig-
108 ure 1*). In the following, I will formally illus-
109 trate when RDD occurs, especially in the con-

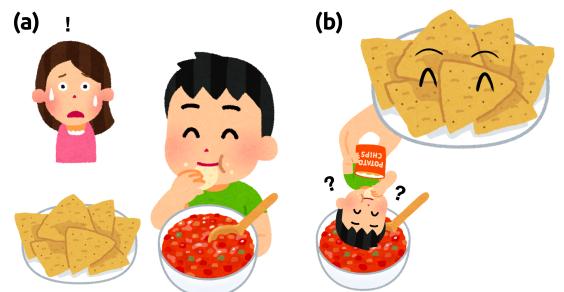


Figure 1. (a) In double-dipping, you dip an identical data point twice (the same stimulus and identical noise). (b) In reverse double-dipping, the dataset dips you twice, unbeknownst to you, with non-identical data points (the same stimulus but independent noise).

110 text of systems neuroscience.

111 Finite impulse response model

112 A finite impulse response (FIR) model is commonly used to identify a time-invariant linear system.
113 Please see **Table 1** for the definitions of variables and notations used in this article. Let us consider
114 an FIR model:

$$y = \mathbf{X}\mathbf{b} + \mathbf{e}, \quad (1)$$

115 where $y \in \mathbb{R}^{T \times 1}$ is a response vector over T time points of a single response unit (e.g., a channel
116 or a voxel), $\mathbf{X} \in \mathbb{R}^{T \times FD}$ is a nonzero Toeplitz design matrix of F features¹ with D delays such that
117 $\|\mathbf{x}_{(i)}^\top \mathbf{x}_{(i)}\| > 0$ for all $i \in \{1, \dots, FD\}$, $\mathbf{x}_{(i)}$ is the i -th column of \mathbf{X} , $\|\cdot\|$ denotes the l_2 -norm, which leads
118 to $\|\mathbf{X}^\top \mathbf{X}\|_F > 0$ where $\|\cdot\|_F$ denotes the Frobenius norm. $\mathbf{b} \in \mathbb{R}^{FD \times 1}$ is an unknown weight vector
119 (often called a temporal response function; or more generally a transfer function), $\mathbf{e} \in \mathbb{R}^{T \times 1}$ is a
120 zero-mean, unit-variance Gaussian noise vector $\mathbf{e} \sim \mathcal{N}_T(\mathbf{0}, \mathbf{I}_T)$ where $\mathbf{I}_T \in \mathbb{R}^{T \times T}$ is an identity matrix.
121 For convenience, we further assume that we standardize predictors and response variables prior
122 to analysis so that their sample means are zero and sample variances are one.

123 Here, let us assume that we have access to the true weights, which are nonzero (i.e., $\|\mathbf{b}\| > 0$).
124 With this model, we can generate a training dataset and a test dataset using the identical weights
125 \mathbf{b} but independent predictors and noise:

$$\mathbf{y}_i = \mathbf{X}_i \mathbf{b} + \mathbf{e}_i = \mathbf{s}_i + \mathbf{e}_i, \quad (2)$$

126 where $(\cdot)_i$ denotes the i -th independent partition in cross-validation with $i = 1$ for a training set and
127 $i = 2$ for a test set. The true signal is denoted as $\mathbf{s}_i \equiv \mathbf{X}_i \mathbf{b} \in \mathbb{R}^{T \times 1}$.

128 To avoid overfitting to noise, regularization such as l_2 -norm penalty (i.e., ridge penalty) is often
129 used.

$$\hat{\mathbf{b}} = (\mathbf{X}_1^\top \mathbf{X}_1 + \mathbf{L})^{-1} \mathbf{X}_1^\top \mathbf{y}_1, \quad (3)$$

130 where $\mathbf{L} \in \mathbb{R}^{FD \times FD}$ is a Tikhonov regularization matrix. In the usual case of ridge regression (i.e., a
131 single penalty applied to all predictors), $\mathbf{L} = \lambda \mathbf{I}_{FD}$. In the general case of multi-penalty ridge (**Hoerl**
132 and **Kennard**, 1970), predictor-delay-wise penalties can be defined: $\mathbf{L} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{FD})$. For
133 now, we assume a single penalty applies to all predictors, except for the intercept, which remains
134 unregularized.

135 In practice, the hyperparameter (i.e., λ) is typically optimized using the third, independent par-
136 tition of data² (**Hastie et al.**, 2009) as:

$$\hat{\mathbf{b}} = (\mathbf{X}_1^\top \mathbf{X}_1 + \mathcal{L}(\{\mathbf{X}_1, \mathbf{y}_1\}; \{\mathbf{X}_3, \mathbf{y}_3\}))^{-1} \mathbf{X}_1^\top \mathbf{y}_1, \quad (4)$$

137 where $\mathcal{L}(\cdot; \cdot)$ is an optimizer that finds the *optimal* regularization matrix \mathbf{L}^* minimizing prediction
138 error for given pairs of design and response variables $\{\mathbf{X}_i, \mathbf{y}_i\}$ and $\{\mathbf{X}_j, \mathbf{y}_j\}$. With this, a regularized

¹While *feature* and *predictor* are often interchangeably used, in this article a feature refers to a variable that describes the characteristics of interest of the input, while a predictor refers to a feature with a specific delay (i.e., each column of a design matrix). That is, F features and D delays make $P = FD$ predictors.

²An independent set for optimization is known as a *validation set* in the statistical learning literature (**Hastie et al.**, 2009). However, in some machine learning literature (**Varoquaux**, 2018), this set is referred to as a *test set*, and the test set is referred to as a *validation set*. To mitigate the confusion, this set for optimization is referred to as an *optimization set*.

¹³⁹ prediction for \mathbf{y}_2 based on \mathbf{X} is:

$$\hat{\mathbf{y}}_{2,\mathbf{X}} = \mathbf{X}_2 \hat{\mathbf{b}} = \mathbf{X}_2 (\mathbf{X}_1^\top \mathbf{X}_1 + \mathbf{L}^*)^{-1} \mathbf{X}_1^\top \mathbf{y}_1 = \mathbf{P}_{\mathbf{X}} \mathbf{y}_1, \quad (5)$$

¹⁴⁰ where $\mathbf{P}_{\mathbf{X}} \in \mathbb{R}^{T \times T}$ is a regularized projection matrix based on \mathbf{X} . Because of the standardization
¹⁴¹ ($\|\mathbf{y}\| = 1$) and the nonnegativity of the denominator, the expected value of a prediction accuracy
¹⁴² metric (e.g., Pearson correlation) over random noise can be seen as proportional to the expected
¹⁴³ inner product of the prediction and the response as:

$$\mathbb{E} [\text{corr}(\hat{\mathbf{y}}_{2,\mathbf{X}}, \mathbf{y}_2)] = \mathbb{E} \left[\frac{(\hat{\mathbf{y}}_{2,\mathbf{X}})^\top \mathbf{y}_2}{\|\hat{\mathbf{y}}_{2,\mathbf{X}}\| \|\mathbf{y}_2\|} \right] \propto \mathbb{E} [(\hat{\mathbf{y}}_{2,\mathbf{X}})^\top \mathbf{y}_2]. \quad (6)$$

¹⁴⁴ This can be further expanded as:

$$\begin{aligned} \mathbb{E} [(\hat{\mathbf{y}}_{2,\mathbf{X}})^\top \mathbf{y}_2] &= \mathbb{E} [\mathbf{y}_1^\top \mathbf{P}_{\mathbf{X}}^\top \mathbf{y}_2] \\ &= \mathbb{E} [(\mathbf{s}_1 + \mathbf{e}_1)^\top \mathbf{P}_{\mathbf{X}}^\top (\mathbf{s}_2 + \mathbf{e}_2)] \\ &= \mathbf{s}_1^\top \mathbf{P}_{\mathbf{X}}^\top \mathbf{s}_2 + \mathbb{E} [\mathbf{e}_1^\top] \mathbf{P}_{\mathbf{X}}^\top \mathbf{s}_2 + \mathbf{s}_1^\top \mathbf{P}_{\mathbf{X}}^\top \mathbb{E} [\mathbf{e}_2] + \mathbb{E} [\mathbf{e}_1^\top] \mathbf{P}_{\mathbf{X}}^\top \mathbf{e}_2, \end{aligned} \quad (7)$$

¹⁴⁵ where only the first term remains nonzero since $\mathbb{E} [\mathbf{e}] = \mathbf{0}$. That is,

$$\mathbb{E} [\text{corr}(\hat{\mathbf{y}}_{2,\mathbf{X}}, \mathbf{y}_2)] \propto \mathbf{s}_1^\top \mathbf{P}_{\mathbf{X}}^\top \mathbf{s}_2. \quad (8)$$

¹⁴⁶ With a sufficiently strong signal $\|\mathbf{b}\| \gg 0$, the optimal regularization approaches zero: $\mathbf{L}^* \approx \mathbf{0}$
¹⁴⁷ (*Hastie et al., 2009*), which allows for approximating the projection matrix as:

$$\mathbf{P}_{\mathbf{X}} = \mathbf{X}_2 (\mathbf{X}_1^\top \mathbf{X}_1 + \mathbf{L}^*)^{-1} \mathbf{X}_1^\top \approx \mathbf{X}_2 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top. \quad (9)$$

¹⁴⁸ Consequently, **Equation 8** can be approximated as:

$$\mathbb{E} [\text{corr}(\hat{\mathbf{y}}_{2,\mathbf{X}}, \mathbf{y}_2)] \propto \mathbf{s}_1^\top \mathbf{P}_{\mathbf{X}}^\top \mathbf{s}_2 \stackrel{(9)}{\approx} (\mathbf{X}_1 \mathbf{b})^\top \left\{ \mathbf{X}_2 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \right\}^\top (\mathbf{X}_2 \mathbf{b}). \quad (10)$$

¹⁴⁹ Due to the symmetry of the inverse covariance $\left\{ (\mathbf{X}_i^\top \mathbf{X}_i)^{-1} \right\}^\top = \left\{ (\mathbf{X}_i^\top \mathbf{X}_i)^\top \right\}^{-1} = (\mathbf{X}_i^\top \mathbf{X}_i)^{-1}$ and the
¹⁵⁰ nonzero assumption $\|\mathbf{X}_i^\top \mathbf{X}_i\|_F > 0$, this can be further simplified as:

$$\begin{aligned} \mathbb{E} [\text{corr}(\hat{\mathbf{y}}_{2,\mathbf{X}}, \mathbf{y}_2)] &\propto \mathbf{s}_1^\top \mathbf{P}_{\mathbf{X}}^\top \mathbf{s}_2 \approx \mathbf{b}^\top \mathbf{X}_1^\top \left[\mathbf{X}_1 \left\{ (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \right\}^\top \mathbf{X}_2^\top \right] \mathbf{X}_2 \mathbf{b} \\ &= \mathbf{b}^\top (\mathbf{X}_1^\top \mathbf{X}_1) (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_2^\top \mathbf{X}_2 \mathbf{b} \\ &= \mathbf{b}^\top \mathbf{X}_2^\top \mathbf{X}_2 \mathbf{b} = \|\mathbf{X}_2 \mathbf{b}\| = \|\mathbf{s}_2\| > 0. \end{aligned} \quad (11)$$

¹⁵¹ That is, with a sufficiently strong signal, the prediction accuracy is expected to be positive.

¹⁵² Red Team: null prediction

¹⁵³ In the above, we assumed that we have access to the true predictors \mathbf{X} and true weights \mathbf{b} unlike
¹⁵⁴ many real-world scenarios. Please note that the predictors in the current setting are not objectively
¹⁵⁵ observable conditions (e.g., the presence or absence of sounds) but a selected set of features of
¹⁵⁶ the stimulus that are hypothesized to be relevant to the neural responses by the researchers (e.g.,
¹⁵⁷ spectrotemporal modulation). In practice, while we clearly know which stimulus we presented,
¹⁵⁸ defining predictors requires knowledge of which information is encoded in the human brain, which

159 is ultimately unknown and often is the very aim of the study (e.g., “*Is information X encoded in the*
160 *brain region A or not?*”).

161 Now, let us consider an imaginary scenario where an independent group of researchers (“Red
162 Team” as in adversarial testing) tries to demonstrate that our analysis method is vulnerable to any
163 null predictors. That is, the red team wants to show that our analysis method can yield significant
164 results even with random predictors, not only the predictors of our choice. They receive our data
165 \mathbf{y} only, but not predictors \mathbf{X} . For simplicity, let us assume that we also inform the Red Team about
166 the true delays. In this case, a reasonable course of action for the Red Team could be to generate
167 a set of random vectors $\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}, \sigma\mathbf{I})$ and to delay them to create a Toeplitz matrix $\mathbf{U}_i \in \mathbb{R}^{T \times FD}$ to
168 make their prediction for \mathbf{y}_2 as:

$$\hat{\mathbf{y}}_{2,\mathbf{U}} = \mathbf{U}_2 \hat{\mathbf{b}}_0 = \mathbf{U}_2 (\mathbf{U}_1^\top \mathbf{U}_1 + \mathcal{L}(\{\mathbf{U}_1, \mathbf{y}_1\}; \{\mathbf{U}_3, \mathbf{y}_3\}))^{-1} \mathbf{U}_1^\top \mathbf{y}_1 = \mathbf{P}_{\mathbf{U}} \mathbf{y}_1. \quad (12)$$

169 Note that actual predictors that researchers would use based on theories, prior evidence, and
170 intuitions should be much more informative than the Red Team’s random numbers. That is, \mathbf{U}
171 is expected to perform worse than any reasonable predictors. Put differently, if the Red Team
172 somehow *magically* makes *significant* predictions using the null predictors \mathbf{U} , it indicates there is
173 something critically flawed in our analysis.

174 Since we assume that the Red Team generates their null predictors \mathbf{U} independently from the
175 true predictors \mathbf{X} , a valid optimization process such as cross-validation (**Hastie et al., 2009**) should
176 lead to a strong regularization of non-informative predictors (i.e., \mathbf{U}). This simplifies the projection
177 matrix to:

$$\mathbf{P}_{\mathbf{U}} = \mathbf{U}_2 (\mathbf{U}_1^\top \mathbf{U}_1 + \lambda^* \mathbf{I})^{-1} \mathbf{U}_1^\top \stackrel{\lambda^* \gg 0}{\approx} \mathbf{U}_2 (\lambda^* \mathbf{I})^{-1} \mathbf{U}_1^\top = \frac{1}{\lambda^*} \mathbf{U}_2 \mathbf{U}_1^\top. \quad (13)$$

178 Thus, given the \mathbf{U}_i , the expected value of the prediction accuracy over random noise is given as:

$$\begin{aligned} \mathbb{E} [\text{corr}(\hat{\mathbf{y}}_{2,\mathbf{U}}, \mathbf{y}_2)] &\propto \mathbf{s}_1^\top \mathbf{P}_{\mathbf{U}}^\top \mathbf{s}_2 \\ &\stackrel{(13)}{\approx} (\mathbf{X}_1 \mathbf{b})^\top \left(\frac{1}{\lambda^*} \mathbf{U}_2 \mathbf{U}_1^\top \right)^\top (\mathbf{X}_2 \mathbf{b}) \\ &= \frac{1}{\lambda^*} \mathbf{b}^\top \mathbf{X}_1^\top \mathbf{U}_1 \mathbf{U}_2^\top \mathbf{X}_2 \mathbf{b}. \end{aligned} \quad (14)$$

179 which converges to zero as λ^* approaches infinity:

$$\lim_{\lambda^* \rightarrow \infty} \mathbb{E} [\text{corr}(\hat{\mathbf{y}}_{2,\mathbf{U}}, \mathbf{y}_2)] = 0. \quad (15)$$

180 That is, the expected prediction accuracy of the Red Team is null.

181 Repetition of stimulus

182 So far, we assumed that the three partitions (i.e., training, optimization, and test sets) are indepen-
183 dent of each other, with independent stimuli and independent noise, but only sharing the identical
184 weights. However, presenting multiple repetitions of an identical, short stimulus—from tens to
185 thousands of times—has been one of the most classical techniques in neuroscience to cancel out
186 random noise in the data and reveal time-locked neural responses that are consistently evoked by
187 the stimulus (e.g., event-related potential, time-locked BOLD response). More recently, for investi-
188 gating the representation of naturalistic stimuli, a design to present an identical set of stimuli to

¹⁸⁹ multiple participants has been popularized in order to reveal stimulus-driven responses in terms
¹⁹⁰ of inter-subject correlation (*Hasson et al., 2004*) or to find a common functional coordinate via
¹⁹¹ hyperalignment (*Haxby et al., 2020*).

¹⁹² A problem occurs when the encoding analysis is naïvely applied to such data. To illustrate the
¹⁹³ point, let us consider an ideal design to reveal such a time-locked response, where the underlying
¹⁹⁴ signal is identical across partitions but the noise is independent: $\mathbf{s}_1 = \mathbf{s}_2 = \mathbf{s}_3$ but $\mathbf{e}_1 \neq \mathbf{e}_2 \neq \mathbf{e}_3$. Then,
¹⁹⁵ our projection matrix with $\lambda^* \approx 0$ will be:

$$\mathbf{P}_{\mathbf{X}} = \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1 + \mathbf{L}^*)^{-1} \mathbf{X}_1^T \approx \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T, \quad (16)$$

¹⁹⁶ when $\mathbf{X}_1^T \mathbf{X}_1$ is invertible. This simplifies the expected prediction accuracy to:

$$\begin{aligned} \mathbb{E} [\text{corr}(\hat{\mathbf{y}}_{2,\mathbf{X}}, \mathbf{y}_2)] &\propto \mathbf{s}_1^T \mathbf{P}_{\mathbf{X}}^T \mathbf{s}_1 \stackrel{(16)}{\approx} (\mathbf{X}_1 \mathbf{b})^T \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T (\mathbf{X}_1 \mathbf{b}) \\ &= \mathbf{b}^T \mathbf{X}_1^T \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_1 \mathbf{b} \\ &= \mathbf{b}^T \mathbf{X}_1^T \mathbf{X}_1 \mathbf{b} = \|\mathbf{s}_1\| > 0. \end{aligned} \quad (17)$$

¹⁹⁷ Equivalently, *Equation 17* being positive can be shown based on that the covariance matrix
¹⁹⁸ $\mathbf{X}_1^T \mathbf{X}_1$ is symmetric and that \mathbf{X}_1 is a rectangular matrix with independent columns, which means
¹⁹⁹ $\mathbf{X}_1^T \mathbf{X}_1$ is positive definite. By definition, $\mathbf{b}^T \mathbf{X}_1^T \mathbf{X}_1 \mathbf{b} > 0$ for $\mathbf{b} \neq \mathbf{0}$. This leads to a rather unsurprising
²⁰⁰ conclusion—with a sufficiently strong signal, the expected prediction accuracy will be positive, also
²⁰¹ with the repeated signals³.

²⁰² In the case of the Red Team, however, a repeated strong signal (i.e., $\|\mathbf{s}_1\| = \|\mathbf{s}_2\| = \|\mathbf{s}_3\| \gg 0$), even
²⁰³ unbeknownst to the Red Team, can alter the optimization process, disabling proper regularization
²⁰⁴ (see *Appendix 1*). When unregularized (i.e., $\mathbf{L}^* \approx \mathbf{0}$), the projection matrix can be approximated as:

$$\mathbf{P}_{\mathbf{U}} = \mathbf{U}_1 (\mathbf{U}_1^T \mathbf{U}_1 + \mathbf{L}^*)^{-1} \mathbf{U}_1^T \approx \mathbf{U}_1 (\mathbf{U}_1^T \mathbf{U}_1)^{-1} \mathbf{U}_1^T, \quad (18)$$

²⁰⁵ Thus, similarly to *Equation 17*, the expected null prediction accuracy of the Red Team can be ap-
²⁰⁶ proximated as:

$$\mathbb{E} [\text{corr}(\hat{\mathbf{y}}_{2,\mathbf{U}}, \mathbf{y}_2)] \propto \mathbf{s}_1^T \mathbf{P}_{\mathbf{U}}^T \mathbf{s}_1 \stackrel{(18)}{\approx} \mathbf{s}_1^T \mathbf{U}_1 (\mathbf{U}_1^T \mathbf{U}_1)^{-1} \mathbf{U}_1^T \mathbf{s}_1. \quad (19)$$

²⁰⁷ $\mathbf{U}_1^T \mathbf{U}_1$ is positive definite due to its symmetry and independence of the columns in \mathbf{U}_1 . Since
²⁰⁸ the inverse operation preserves the signs of eigenvalues of a square matrix, its inversion $(\mathbf{U}_1^T \mathbf{U}_1)^{-1}$
²⁰⁹ is also positive definite. A substitution ($\mathbf{d} \equiv \mathbf{U}_1^T \mathbf{s}_1 \neq \mathbf{0}$) can clarify that *Equation 19* is positive by
²¹⁰ definition:

$$\mathbf{s}_1^T \mathbf{P}_{\mathbf{U}}^T \mathbf{s}_1 \approx \mathbf{d}^T (\mathbf{U}_1^T \mathbf{U}_1)^{-1} \mathbf{d} > 0. \quad (20)$$

²¹¹ This may surprise some readers; however, *Equation 20* implies that the null prediction of the
²¹² Red Team is expected to be greater than zero. Depending on the signal-to-noise ratio, this may
²¹³ result in Type-I (false positive) errors. This is due to the circular fallacy introduced by the *leakage in*
²¹⁴ *training examples*, i.e., RDD.

²¹⁵ In short, the mechanism of RDD can be summarized in the following steps:

³Notably, here the expected prediction accuracy is proportional to the square of the signal in the training set, rather than the test set (*Equation 9*). This may affect generalization performance in hold-out validation (where the performance is evaluated only once on a separate test set), but in cross-validation the performance will be averaged out across sets. Either way, repetition of stimulus across sets leads to a spurious inflation of performance estimation.

- 216 1. The repeated signal disables the regularization of a seemingly valid optimization process even
 217 with independent noise.
 218 2. As the regularization hyperparameter approaches zero, the projection matrix becomes posi-
 219 tive definite, which was null when properly regularized.
 220 3. Because of the repeated signal, the expected prediction accuracy over random noise is pro-
 221 portional to a bilinear form involving a nonzero vector and a positive-definite square matrix,
 222 i.e., $\mathbf{d}^T \mathbf{M} \mathbf{d} > 0$.

223 **Toy example**

224 For a graphical illustration, a simple case of small-scale simulation (i.e., a toy example) is shown
 225 in *Figure 2* and *Figure 3*. See Simulation methods for details of how the simulations were created.
 226 In this example, two features with a high temporal autocorrelation were generated without repe-
 227 titions (*Figure 2a*) and with repetitions (*Figure 3a*). The Red Team created their own null features
 228 (*Figure 2m*). Thus, all coefficients were highly regularized for the Red Team (*Figure 2r*; pink). Conse-
 229 quently, the Red Team's predictions were mostly flat (*Figure 2o,p*; dotted lines) and expectedly, the
 230 Red Team's prediction accuracies were around $r = 0$ (*Figure 2l*; pink). However, with the stimulus
 231 repeated across sets (i.e., identical signals in the training and test sets), the regularization for the
 232 Red Team was much smaller, almost close to the true features (*Figure 3r*; pink vs. lime green) and
 233 the Red Team's prediction accuracies were around $r = 0.5$ (*Figure 3l*; pink). That is, due to the rep-
 234 etition of the stimulus, even with independent noise (*Figure 3e,k*), the prediction accuracies based
 235 on the null features were falsely inflated.

236 Having shown how RDD could occur, it is important to note that multiple assumptions were
 237 made to arrive here. In practice, various factors—such as signal strength, temporal and spatial
 238 correlation structures, feature collinearity, similarity across partitions, and the flexibility of the FIR
 239 model—might interact to inflate Type-I errors.

240 **Simulation**

241 In this section, I highlight major factors that worsen the Type-I error due to RDD. To keep the num-
 242 ber of combinations manageable, univariate models ($F = 1, V = 1$) were first considered. Then,
 243 while iteratively pruning out irrelevant factors, models with multivariate features (F) and multivari-
 244 ate responses (V) were considered. Methodological details of the simulation are described in the
 245 Simulation methods section. The parameters of simulation are summarized in *Table 2*.

246 **Univariate-feature, univariate-response**

247 The first batch of simulations was restricted to a univariate feature (the number of features $F = 1$)
 248 and a univariate response (the number of variates $V = 1$). The explored parameter levels were:
 249 $D \in \{1, 3, 5, 7, 9, 11\}$, $S \in \{-10, 0, 10\}$, $\phi_X \in \{0, 0.5, 1\}$, $\phi_U \in \{0, 0.5, 1\}$, $\phi_E \in \{0, 0.5, 1\}$, $\phi_B \in \{0, 0.5, 1\}$,
 250 $\text{IsRep} \in \{0, 1\}$. With these parameter levels, total 2,916 combinations were created, each sampled
 251 1,000 times.

252 To illustrate the most distinctive effects, prediction accuracies averaged across 1,000 random
 253 sampling are shown in *Figure 4*. Expectedly, the SNR increased the prediction accuracy based on

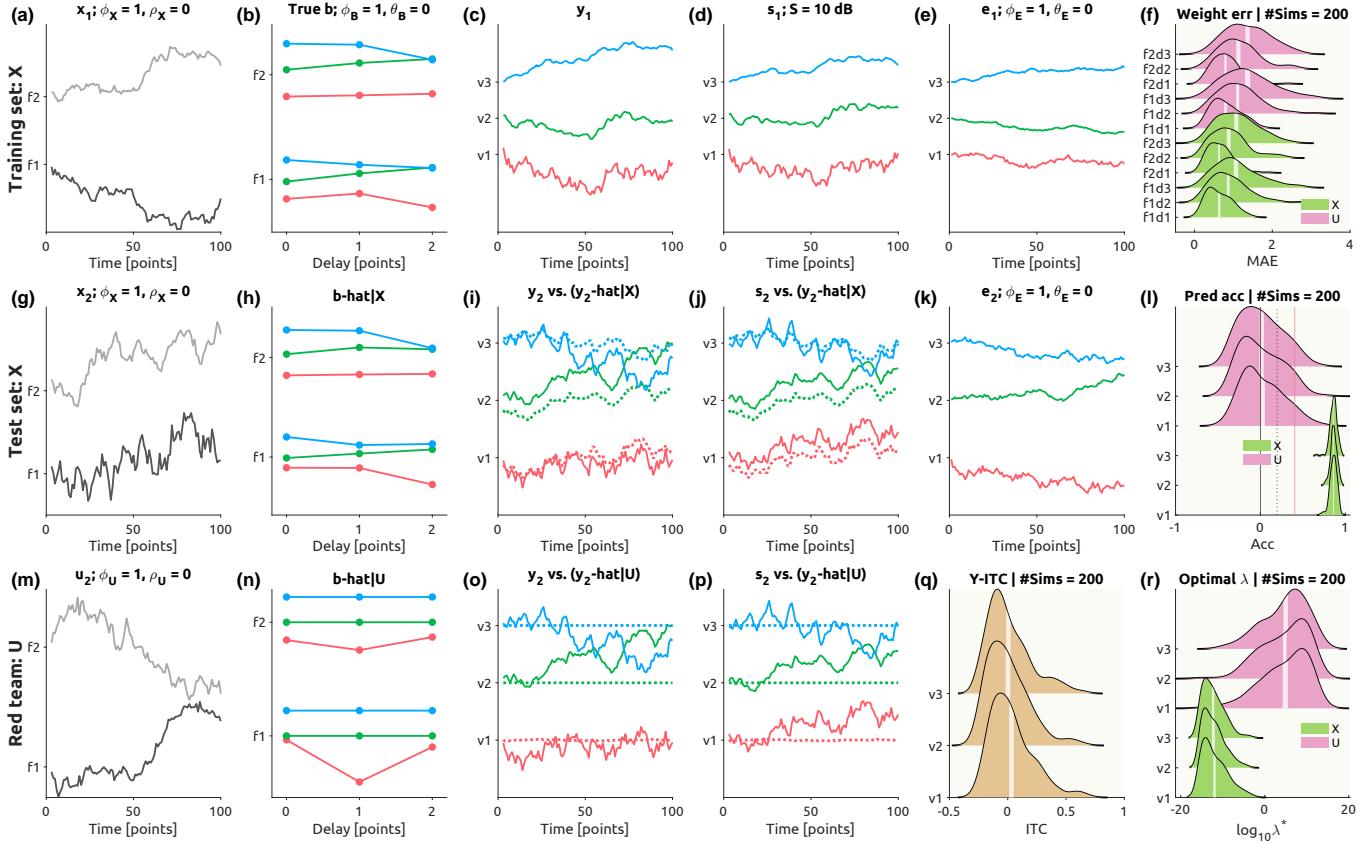


Figure 2. A toy example without stimulus repetition. Apart from the right-most panels in pale beige (**f, l, r**), the panels in the first row (**a–e**) correspond to the training set with the true features \mathbf{X}_1 , the panels in the second row (**g–k**) correspond to the test set with the true features \mathbf{X}_2 , and the panels in the third row (**m–p**) correspond to the Red Team's null features \mathbf{U}_2 . The panels in the first column (**a, g, m**) show the true features \mathbf{x} or null features \mathbf{u} in gray scale. The panels in the second column show the true weights \mathbf{b} (**b**) or estimation based on the true features (**h**) or null features (**n**) where the RGB color represents one of three response variates. The panels in the third column (**c, i, o**) show the response \mathbf{y} (solid lines) and prediction based on true or null features \mathbf{y}_hat (dashed lines), and the panels in fourth column (**d, j, p**) show the true signal \mathbf{s} (solid lines) and the prediction (dashed lines). The two panels in the fifth column (**e, k**) show true noise \mathbf{e} . In the pale beige panels (**f, l, q, r**), distributions from 200 simulations are shown in ridgeline plots with the 95% confidence interval of the mean shown in white strips: (**f**) mean absolute error (MAE) of weight estimation based on the true features (\mathbf{X} ; lime green) or null features (\mathbf{U} ; pink), (**l**) prediction accuracies in Pearson's correlation coefficients with a black vertical line for an absolute zero, a gray vertical line for an uncorrected $P < 0.05$ (assuming independent time points), and a red vertical line for a Bonferroni-corrected $P < 0.05$ adjusted for the number of variates. (**q**) inter-trial correlation (ITC) of the responses (brown), (**r**) exponents of the geometrically averaged optimal λ 's. Parameters to generate this simulation set are: $T = 100$, $V = 3$, $F = 2$, $D = 3$, $S = 1\text{dB}$, $\phi_X = 1$, $\rho_X = 0$, $\phi_B = 1$, $\theta = 0$, $\phi_E = 1$, $\theta_E = 0$, $\text{IsRep} = 0$. See **Table 2** for the explanation of parameters.

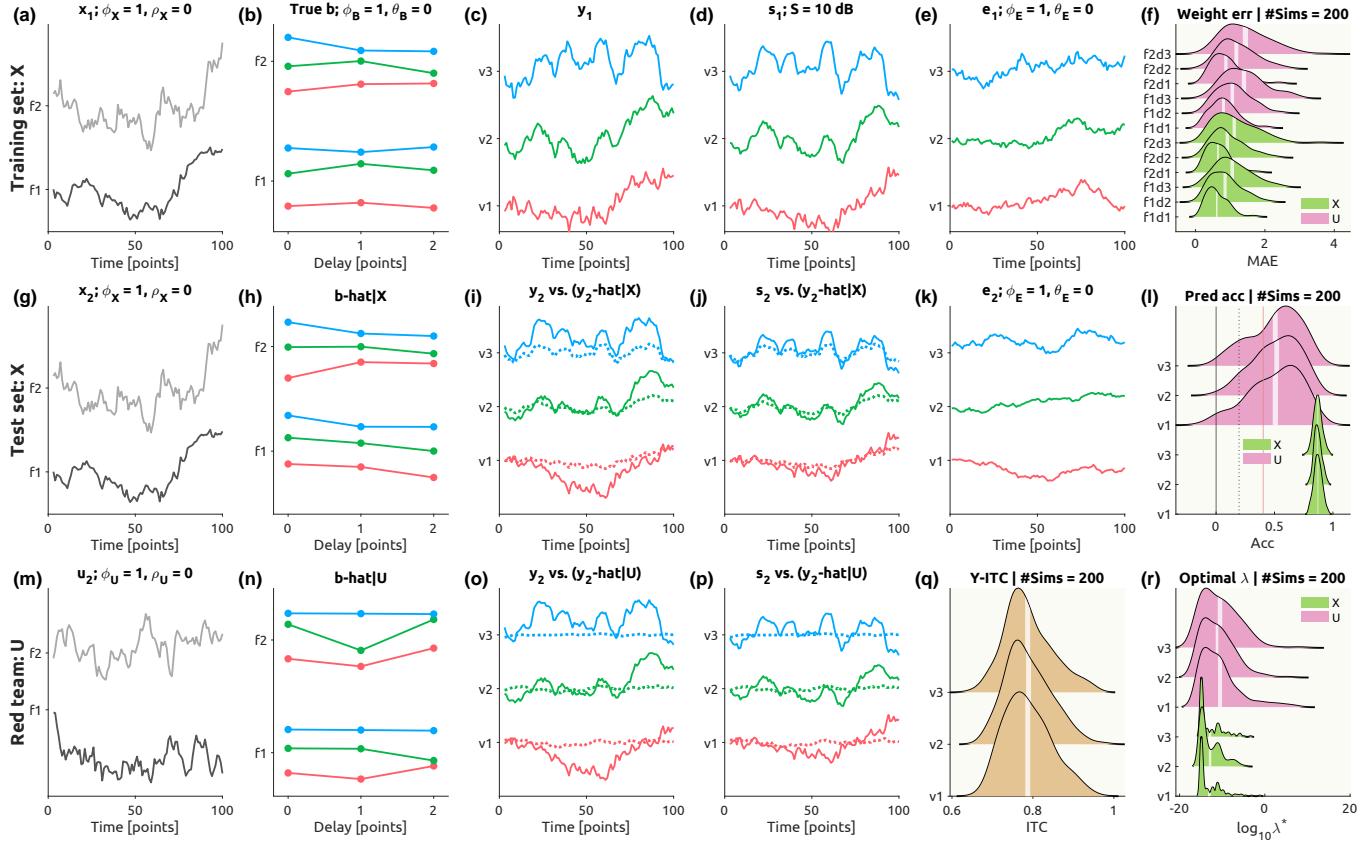


Figure 3. A toy example with stimulus repetition. The visualization scheme is identical to **Figure 2**. Parameters to generate this simulation set are identical to those in **Figure 2**, except IsRep = 1.

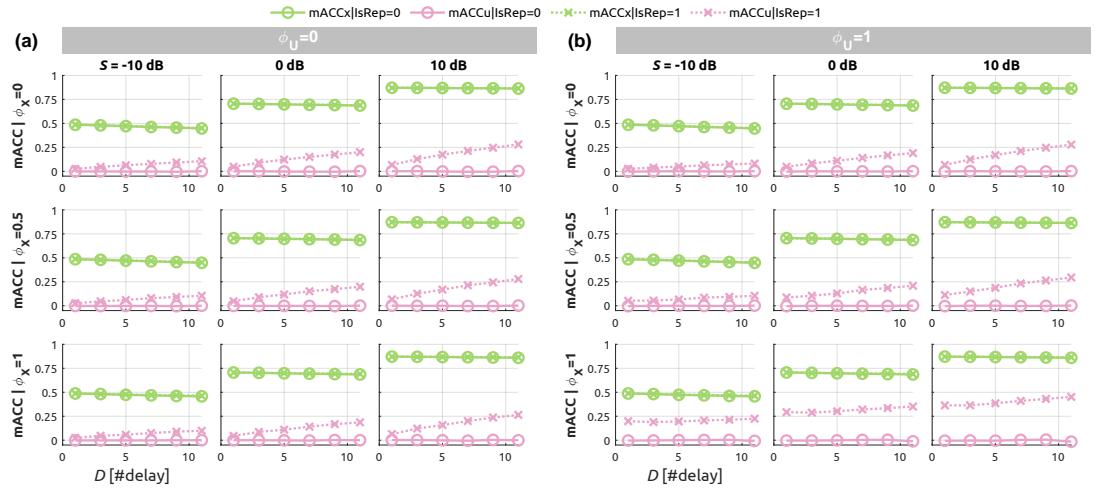


Figure 4. Simulations of univariate-feature, univariate-response models. Mean prediction accuracies are plotted over the number of delays D when the temporal autocorrelation **(a)** $\phi_U = 0$ and **(b)** $\phi_U = 1$. Each marker corresponds to the averaged Pearson correlation coefficients of 1000 simulations. Predictions based on true features (\mathbf{X}) are shown in lime green and null features (\mathbf{U}) in pink. Open circles with solid lines indicate accuracies when the true signals were not repeated. Crosses with dashed lines represent cases where the true signals were repeated. Each column corresponds to a specified SNR level. $\phi_X = 0$ in the top row and $\phi_X = 1$ in the bottom row. Other parameters were as follows: $K = 1000$, $\phi_B = 0$, $\phi_E = 0$. Other combinations of parameters showed generally similar patterns. See *Table 3* and *Table 4* for effect sizes.

the true features (green markers). The prediction based on the null features remained zero when the stimuli were independent across the cross-validation (CV) partitions (pink circles). However, when the stimuli were identical across the partitions, the null prediction accuracy increased over the number of delays and the SNR (pink crosses). In particular, the effect of temporal autocorrelation ϕ_X when $\phi_U = 1$ substantially increased the null prediction with the repeated stimuli.

To comprehensively assess the possible effects, a full factorial model on mean prediction accuracies was fitted with all seven variables as in Wilkinson notation:

$$r \sim 1 + D * S * \phi_X * \phi_U * \phi_E * \phi_B * \text{IsRep}, \quad (21)$$

where r is the mean prediction accuracy averaged for each set of 1000 simulations either based on true or null predictors, 1 represents an intercept, and $*$ denotes factor crossing as in $a * b = a+b+a : b$ with $:$ denoting an interaction. This yielded linear models with 128 terms from the intercept to the seven-way interaction, 2,916 observations, and adjusted $R^2 = 0.985$ for \mathbf{X} and 0.957 for \mathbf{U} , which are reasonable given that the accuracy metrics are averaged within each combination of parameters. Due to the large number of observations, P -values were not highly selective ($P_{\text{Bonferroni}} < 0.0001$ for many of contrasts; 9 for \mathbf{X} , 27 for \mathbf{U} out of 127). Thus, on top of $P_{\text{Bonferroni}} < 0.0001$, only effects with moderate effect sizes ($\eta_p^2 \geq 0.16$) were considered for further discussion (*Table 3*, *Table 4*).

As expected, the SNR consistently increased the prediction accuracy ($\eta_p^2[S] = 0.985$ for \mathbf{X} ; $\eta_p^2[S] = 0.660$ for \mathbf{U}). Most importantly, the interaction of the stimulus repetition (IsRep) with the signal strength (SNR, S), and the autocorrelation of the true and null features (ϕ_X , ϕ_U) were markedly found on the prediction accuracy based on the null features ($\max \eta_p^2 = 0.661$, *Table 4*). Put differently, the RDD effect (i.e., a false inflation of the null prediction accuracy by the repetition of stimulus

274 across CV partitions) was most pronounced when the true underlying signal was strong and the
275 true features were temporally autocorrelated. This finding is consistent with [Equation 17](#) above,
276 showing that the RDD effect depends on the nonzero strength of the underlying signal. The clearer
277 the signal, the more likely the RDD effect will occur when analyzed in a circular CV design.

278 In addition, the RDD effect strongly interacted with the complexity of the underlying (and fitted)
279 models (i.e., `IsRep : D`; [Table 4](#)). While in a correct CV design (`IsRep = 0`), the model complexity
280 did not increase the prediction accuracy as this was regularized by independent optimization and
281 evaluation. However, in a circular CV design (`IsRep = 1`), the model complexity increased the pre-
282 diction accuracy (hence the RDD effect) since the model was fitted to the identical signal across CV
283 partitions. Naturally, this effect further interacted with the strength of the signal (`IsRep : D : S`;
284 [Table 4](#)).

285 The effect of the autocorrelation of the weight time series ϕ_B was found to be negligible for
286 both **X** and **U** ($\eta_p^2 < 0.009$). Guided by these results, we fixed the autocorrelation of the weight time
287 series to zero ($\phi_B = 0$) in the following simulations.

288 Multivariate-feature, univariate-response

289 Since most often we are interested in multivariate features (e.g., motion energies, spectrograms,
290 deep ANN embeddings), it is of interest how the dimensionality and multicollinearity of features
291 influence RDD artifact.

292 The explored parameter levels were: $D \in \{1, 5, 9\}$, $F \in \{5, 10, 15, 20\}$, $S \in \{-10, 0, 10\}$, $\phi_X \in$
293 $\{0, 0.5, 1\}$, $\phi_U \in \{0, 0.5, 1\}$, $\phi_E \in \{0, 0.5, 1\}$, $\phi_B = 0$, $\rho_X \in \{0, 0.5, 1\}$, $\rho_U \in \{0, 0.5, 1\}$, `IsRep` $\in \{0, 1\}$.
294 With these levels, the full combinations amounted to 17,496, for each of which, once more, 1,000
295 random samplings were carried out.

296 [Figure 5](#) displays the simulated effects. The number of features F reduced the true models' pre-
297 diction accuracies without repetitions (green solid lines), and more so with more complex response
298 functions (e.g., $D = 9$). In contrast, increasing the number of features led to higher prediction accu-
299 racies in null models with stimulus repetition (pink dashed lines). Moreover, at a higher SNR (e.g.,
300 10 dB), the null prediction accuracies (pink dashed lines) were even higher than the true predi-
301 cation accuracies (green solid lines) with many independent predictors (e.g., 9 delays $\times \geq 15$ features
302 when $\rho_X = 0$; [Figure 5a](#)). This suggests that, in a bad combination, the null prediction accuracy can
303 go beyond the noise ceiling (i.e., green solid lines; assuming we know the true predictors), which
304 is the plausibly highest prediction accuracy bounded by the noise level of the signal. This crossing
305 of the true and null prediction accuracies was attenuated when features were highly correlated
306 ($\rho_X = 1$; [Figure 5b](#)), which reduced the effective degrees of freedom.

307 To quantify the observed effects, once again, a full factorial model on mean prediction accura-
308 cies was fitted with all nine variables:

$$r \sim 1 + D * F * S * \phi_X * \phi_U * \phi_E * \rho_X * \rho_U * \text{IsRep}. \quad (22)$$

309 A full factorial model with 17,496 observations and 511 terms resulted in adjusted $R^2 = 0.975$
310 for **X** and 0.923 for **U** ([Table 5](#), [Table 6](#)). A strong effect of the number of features F was found for
311 **X** ($\eta_p^2 = 0.276$) but only an intermediate effect for **U** ($\eta_p^2 = 0.102$). Additionally, its interaction with

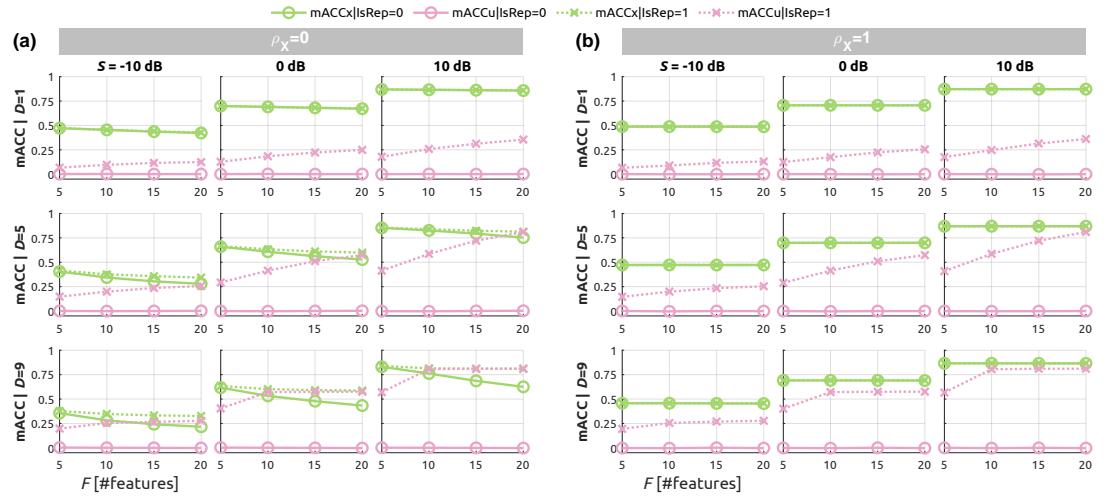


Figure 5. Simulations of multivariate-feature, univariate-response models. Mean prediction accuracies are plotted over the number of features F when the multicollinearity **(a)** $\rho_X = 0$ and **(b)** $\rho_X = 1$. Marker styles and colors match to those in **Figure 4**. Each column represents a specified SNR level. The number of delays is $D = 1$ in the top rows and $D = 9$ in the bottom rows. Other parameters were as follows: $K = 1000$, $\phi_U = 0$, $\phi_B = 0$, $\phi_E = 0$, $\rho_U = 0$. See **Table 5** and **Table 6** for effect sizes.

Figure 5—figure supplement 1. Additional cases with $\rho_U = 0$ and $\rho_U = 1$ when $\rho_X = 0$

stimulus repetition (IsRep) for the null features was only intermediate U ($\eta_p^2 = 0.101$), suggesting the dimension of the features alone was not a strong factor for the RDD effect.

However, the interaction between the stimulus repetition and the multicollinearity of the null features ($\rho_U : \text{IsRep}$) was strong ($\eta_p^2 = 0.442$; **Table 6**), where the RDD effect was more pronounced when the null features exhibited lower autocorrelation (**Figure 5—figure Supplement 1**). That is, when the null features have greater effective degrees of freedom (i.e., lower autocorrelation), the model could more flexibly fit the repeated signal, inflating the RDD effect.

In addition, an interesting finding was that the true prediction accuracy decreased as the number of features F and the number of delays D increased. This was due to the limited number of samples ($T = 100$) as compared to the high dimensionality of the feature space, which made the linear model ill-posed. While regularization makes fitting feasible, the inherent limitation exist. Put differently, these results suggest that a sufficient number of samples is required to faithfully estimate the transfer function of the high-dimensional feature space.

Multivariate-feature, multivariate-response

Finally, it was tested whether the dimensionality and spatial autocorrelation of responses (variates) affect the RDD artifact. First, it is worth noting that the encoding model is typically a univariate-response model (e.g., “voxel-wise” or “channel-wise”). The assumption of spatially (i.e., across response units) independent noise is similar to that of the classical general linear model (*Friston et al., 1994*), where a diagonal covariance structure across lattice sampling grids (e.g., voxels) is assumed⁴. Likewise, the encoding model is independently optimized for each response unit in the

⁴Therefore, the term ‘massive-univariate’ would be more fitting than ‘multivariate’.

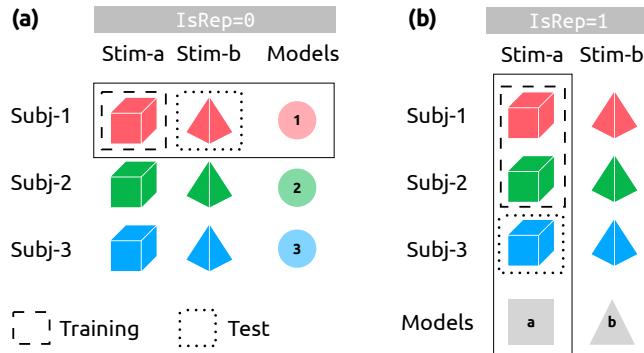


Figure 6. Schema of two cross-validation designs. For simplicity, let us say we have three subjects {1, 2, 3} (depicted in red, green, blue) and two stimuli $\{a, b\}$ (depicted as a cube and a tetrahedron). Note that we assume any repetitions of stimuli within a subject are averaged prior to the cross-validation. **(a)** $\text{IsRep}=0$. Subject-specific models (pale colored circles). Each model (e.g., pale red circle in solid rectangle) is trained with one stimulus (dashed rectangle) and tested on another stimulus (dotted rectangle). **(b)** $\text{IsRep}=1$. Stimulus-specific models (gray square and triangle). Each model (e.g., gray square in sold rectangle) is trained with two subjects (dashed rectangle) and tested on the other subject (dotted rectangle). The optimization set, which could be in the training set of the outer loop, is not marked for simplicity.

current paper⁵. Thus, while the spatial dependency may affect the family-wise error rates of multiple testing (since many correction methods exploit neighbouring supports), it is unexpected that the spatial dependency directly alters the unit-wise optimization and following weight estimation. However, for completeness, simulations were carried out with following parameters: $D \in \{1, 7, 11\}$, $F \in \{5, 10, 20\}$, $V \in \{5, 10, 20\}$, $S \in \{-10, 0, 10\}$, $\phi_X \in \{0, 0.5, 1\}$, $\phi_U \in \{0, 0.5, 1\}$, $\phi_E \in \{0, 0.5, 1\}$, $\phi_B = 0$, $\theta_B \in \{0, 0.5, 1\}$, $\theta_E \in \{0, 0.5, 1\}$, $\rho_X \in \{0, 0.5, 1\}$, $\rho_U \in \{0, 0.5, 1\}$, $\text{IsRep} \in \{0, 1\}$. With these parameter levels, there were 354,287 total combinations, as before, each sampled 1,000 times. A full factorial model, incorporating all twelve variables, was fitted using 354,287 observations and 4,096 terms. As expected, the spatial dependency θ_B and θ_E neither showed any marked main effects nor interactions ($\eta_p^2 < 0.160$; **Table 7**, **Table 8**).

Real Data

In this section, I present real-data examples where RDD spuriously inflated null prediction accuracies using open-access data where healthy participants listened to various musical excerpts while measuring neural activity (electroencephalography [EEG] or functional magnetic resonance imaging [fMRI]) or behavioral ratings (*Kaneshiro et al., 2020; Sachs et al., 2020*).

Based on the well-established encoding of the acoustic energy in the human auditory system, true features were the audio envelopes extracted from the musical stimuli using a cochlear model (*Chi et al., 2005*). Null features were either (a) the phase-randomized envelope (preserving spectral magnitudes and autocorrelation structures) as the most realistic one, (b) the normal noise, and (c) the uniform noise as the least realistic one. Details of the real data and analysis implementation are provided in the Materials and Methods section.

Importantly, the data were analyzed with two competing cross-validation (CV) designs (**Figure 6**):

⁵In some EEG studies, hyperparameters were averaged across response units (i.e., EEG channels). This practice introduces spatial dependency that leads to a suboptimal regularization for individual response units.

- 354 1. `IsRep` = 0: No stimulus was repeated across CV sets (training, optimization, and test). That is,
 355 subject-specific models were fit with stimuli assigned to CV sets.
 356 2. `IsRep` = 1: Identical stimuli, albeit with different noise realizations, were repeated in all three
 357 sets. Thus, stimulus-specific models were fit with participants partitioned into CV sets.

358 The magnitude of RDD artifact was estimated as the difference in null prediction accuracies
 359 between two CV schemes: $\text{RDD} = \bar{r}_{\text{stim}}(\mathbf{U}; \text{IsRep} = 1) - \bar{r}_{\text{subj}}(\mathbf{U}; \text{IsRep} = 0)$ where r_{stim} is a predic-
 360 tion accuracy of a stimulus-specific model and r_{subj} is that of a subject-specific model. $(\bar{\cdot})$) denotes
 361 averaging across null models ($M = 100$). That is, if there were no false inflation of prediction accu-
 362 racy due to RDD, the null predictions with and without stimulus repetitions should be equal (i.e.,
 363 $H_0 : \mathbb{E}(\text{RDD}) = 0$). Otherwise, the null prediction with stimulus repetitions is expected to be greater
 364 than the null prediction without repetitions ($H_A : \mathbb{E}(\text{RDD}) > 0$).

365 All data analyses were consistently done using an MATLAB package Linearized Encoding Analy-
 366 sis (LEA; <https://github.com/seunggookim/lea>).

367 **Electroencephalography**

368 Scalp electrical potential data were recorded in 48 healthy participants while listening to Western-
 369 style Indian pop music (i.e., Bollywood music; *Kaneshiro et al., 2020*). Using this EEG dataset, lin-
 370 earized encoding analysis was performed with the audio envelope as a true feature and its phase-
 371 randomized signals as null features.

372 Without stimulus repetition (`IsRep` = 0), a clear fronto-central topography is shown in the predic-
 373 tion accuracy ($\max r(X; 0) = 0.047$, **Figure 7a**) as well as in the ridge hyperparameter ($\min \log_{10} \lambda(X; 0) =$
 374 5.86, **Figure 7b**), reflecting the envelope encoding in the bilateral auditory cortices while listening
 375 to music. For this particular data, the estimated weights were stronger in the left than right fronto-
 376 central channels (**Figure 7c**). With the phase-randomized envelope, as expected, the null predic-
 377 tion accuracy was minimal ($\max r(U; 0) = 0.006$, **Figure 7d**; $\max \mathbb{E}[r(U; 0)] = 0.001$, **Figure 7g**) with
 378 all channels were highly regularized ($\min \log_{10} \lambda(U; 0) = 10.98$, **Figure 7e**; $\min \mathbb{E}[\log_{10} \lambda(U; 0)] = 12.11$,
 379 **Figure 7h**).

380 With stimulus repetition (`IsRep` = 1), true prediction accuracies were increased ($\max r(X; 1) =$
 381 0.085, **Figure 7j**). This is because, unlike the simulation, our feature (i.e., the audio envelope) was
 382 not the sole information that the human EEG data encode.

383 However, most strikingly, the null prediction accuracies with stimulus repetition showed an al-
 384 most identical topography to the actual encoding results (**Figure 7m,p**), with even higher values
 385 than the true prediction without repetition ($\max r(X; 0) = 0.047$, $\max r(U; 1) = 0.052$, $\max \mathbb{E}[r(U; 1)] =$
 386 0.056). Note that, by definition, the null feature (phase-randomized envelopes) should have not
 387 predicted anything in the EEG data. However, when the identical stimuli were repeated over CV
 388 partitions, the regularization was disabled—regardless of the given features—in channels where
 389 the stimulus-evoked response is strong (**Figure 7n,q**). Then, even random weights (**Figure 7o,r; Fig-**
 390 **ure 7—figure Supplement 9**; i.e., widely different from the true weights) could successfully predict
 391 the repeated signal. Because RDD artifact reflects the genuine biological signal that is repeatedly
 392 evoked by identical stimuli, the observed patterns of the prediction accuracy may appear indis-
 393 tinguishable from the true signal without close investigation of weights. This pattern of RDD was

394 observed, albeit weaker, even when the null features were unrealistic such as normal noise (*Figure 7—figure Supplement 1*) or uniform (*Figure 7—figure Supplement 2*). Moreover, the RDD effect
395 was consistent across different sets of delays (*Figure 7—figure Supplement 3*, *Figure 7—figure*
396 *Supplement 4*, *Figure 7—figure Supplement 5*, *Figure 7—figure Supplement 6*, *Figure 7—figure*
397 *Supplement 7*, *Figure 7—figure Supplement 8*).

398 Transfer function weights showed interesting patterns (*Figure 7—figure Supplement 9*, *Fig-*
399 *ure 7—figure Supplement 10*, *Figure 7—figure Supplement 11*). Spatially, the eigenvectors ex-
400 plaining the largest variance (i.e., PC1) of the phase-randomized envelope were similar to those of
401 the true envelope while the eigenvectors of the uniform or normal noise without autocorrelation
402 showed rather noisy (spatially high frequencies) patterns. Temporally, even though the uniform
403 or normal noise features had no autocorrelation, the eigenvariate time series showed smooth pat-
404 terns, reflecting the autocorrelation of the EEG data.

405 When comparing the null prediction accuracies between the CV schemes (e.g., *Figure 7g* vs.
406 *Figure 7p*), the RDD effects were found significant in all channels and displayed a fronto-central
407 topography that is highly plausible for the auditory cortical activity ($P_{FDR} < 0.01$; *Figure 8*). It is
408 noteworthy that the RDD effect is much stronger for the phase-randomized envelope, which pre-
409 serves the autocorrelation structure of the stimulus. Also, the number of delays seems to further
410 inflate the RDD effect as demonstrated in the simulation results (e.g., *Figure 4*).

412 **Functional magnetic resonance imaging**

413 Blood-oxygen-level-dependent (BOLD) data were acquired in 39 healthy participants while listening
414 to Western instrumental musical pieces that either evoke happiness or sadness as validated in
415 independent listeners (*Sachs et al., 2020*). As done for the EEG dataset, linearized encoding analysis
416 was performed with the fMRI data as responses, the audio envelope as a true feature, the phase-
417 randomized envelope as a null feature, and delays from 3 to 9 seconds (*Figure 9*). Similarly to
418 the EEG results, the phase-randomized envelope strikingly predicted the BOLD time series in the
419 bilateral auditory cortices including the Heschl's gyrus and planum temporale (*Figure 9m,p*) while
420 no consistent pattern in the transfer function weights was found over the phase randomizations
421 (*Figure 9r*), clearly demonstrating the RDD effect. Once again, the anatomical location and the
422 extent of the heightened null prediction accuracies precisely matched the true encoding results
423 (*Figure 9a,j*), which would seem 'highly convincing' to many human neuroscientists.

424 When analyzed with different noise models and delays (*Figure 9—figure Supplement 1*, *Fig-*
425 *ure 9—figure Supplement 2*, *Figure 9—figure Supplement 3*, *Figure 9—figure Supplement 4*, *Fig-*
426 *ure 9—figure Supplement 5*, *Figure 9—figure Supplement 6*, *Figure 9—figure Supplement 7*, *Fig-*
427 *ure 9—figure Supplement 8*), the similar RDD pattern was consistently observed (i.e., inflated pre-
428 dictions accuracies in the auditory cortices). The RDD effect was statistically significant in not only
429 the bilateral superior temporal gyri but the medial occipital cortices and the inferior frontal cortices,
430 where acoustic energy is not expected to be encoded (*Figure 10*). Consistently with the EEG results,
431 the RDD effect was stronger for the phase-randomized envelope than the normal or uniform noise
432 as well as for longer delay points than shorter ones.

433 The transfer function weights (*Figure 9—figure Supplement 9*, *Figure 9—figure Supplement 10*,

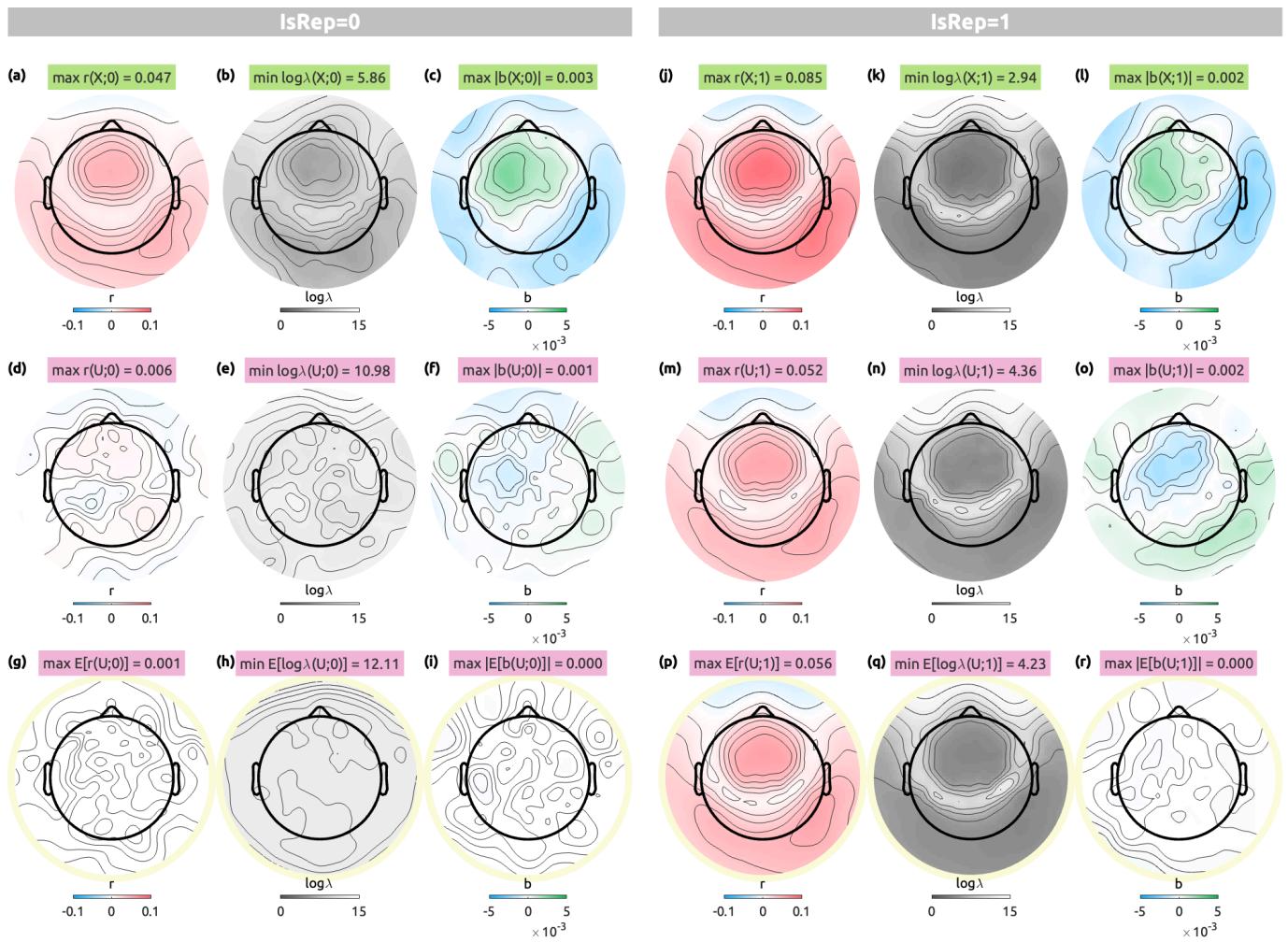


Figure 7. EEG linearized encoding analysis results with delays from 0 to 0.5 sec with an audio envelope (top row, (a-c, j-l)), a single case of a phase-randomized envelope (middle row, (d-f, m-o)), and an average of 100 phase-randomized envelopes (bottom row, circled in pale yellow, (g-i, p-r)). For each CV scheme ($\text{IsRep} = 0$, left panels, (a-i); $\text{IsRep} = 1$, right panels, (j-r)), prediction accuracy (r , blue to red, (a, d, g, j, m, p)), logarithmic ridge hyperparameter ($\log_{10} \lambda$, gray to white, (b, e, h, k, n, q)), transfer function weights that are summed over delays (b , blue to green, (c, f, i, l, o, r)) are shown along the columns. Note that stimulus repetition not only slightly inflated true prediction accuracies but even the predicted ‘brain activity pattern’ from null features (m,p) which was literally indistinguishable from true predictions (a, j).

Figure 7—figure supplement 1. EEG linearized encoding analysis with normal noise and delays from 0 to 0.5 sec

Figure 7—figure supplement 2. EEG linearized encoding analysis with uniform noise and delays from 0 to 0.5 sec

Figure 7—figure supplement 3. EEG linearized encoding analysis with the phase-randomized envelope and delays from 0 to 0.3 sec

Figure 7—figure supplement 4. EEG linearized encoding analysis with normal noise and delays from 0 to 0.3 sec

Figure 7—figure supplement 5. EEG linearized encoding analysis with uniform noise and delays from 0 to 0.3 sec

Figure 7—figure supplement 6. EEG linearized encoding analysis with the phase-randomized envelope and delays from 0 to 1 sec

Figure 7—figure supplement 7. EEG linearized encoding analysis with normal noise and delays from 0 to 1 sec

Figure 7—figure supplement 8. EEG linearized encoding analysis with uniform noise and delays from 0 to 1 sec

Figure 7—figure supplement 9. EEG transfer function weights with all three noise models and delays from 0 to 0.5 sec

Figure 7—figure supplement 10. EEG transfer function weights with all three noise models and delays from 0 to 0.3 sec

Figure 7—figure supplement 11. EEG transfer function weights with all three noise models and delays from 0 to 1 sec

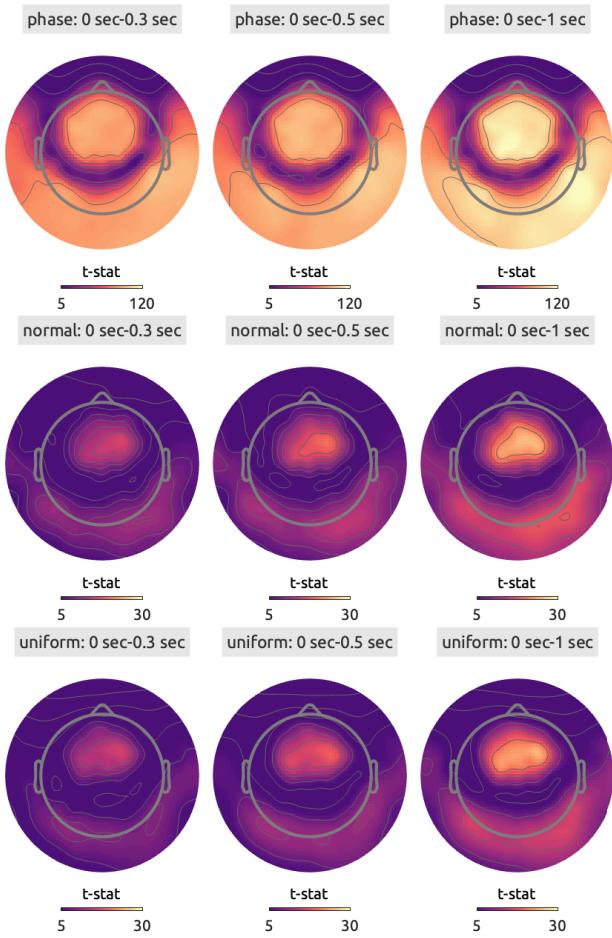


Figure 8. t -statistic maps comparing the null prediction accuracies between two CV schemes (e.g., [Figure 7g](#) vs. [Figure 7p](#)) to test the RDD effects in the EEG data for the phase-randomized envelope (top row), the normal noise (middle row), and the uniform noise (bottom row). All channels were found significant after FDR adjustment in all cases ($P_{FDR} < 0.01$).

434 **Figure 9—figure Supplement 11**) display strong “auditory components” even for normal and uniform noise in their eigenvectors while the corresponding eigenvariates were widely different from
435 the weights estimated by the true feature.

437 **Behavioral ratings**

438 Continuous ratings of music-evoked emotions were sampled from the same 39 healthy participants
439 who took part in the fMRI experiment above (*Sachs et al., 2020*). After the scanning session,
440 participants listened to the same musical pieces again and rated their Emotionality (how happy/sad
441 they felt) and Enjoyment (how much they enjoyed the piece) using a slider. A linearized encoding
442 analysis was performed with the ratings as responses, the audio envelope as a true feature, the
443 phase-randomized envelope as a null feature, and delays from 0 to 10 seconds (*Figure 11*). Once
444 more, while the true envelope predicted Emotionality to some degrees and Enjoyment to a greater
445 extent (*Figure 11a*), the phase-randomized envelope also predicted both scales well above zero
446 when the stimuli were repeated across CV partitions (*Figure 11m,p*) unlike when the stimuli were
447 not repeated (*Figure 11g*).

448 The RDD effect was significant also in the behavioral ratings consistently across all noise models
449 and delays (*Figure 12; Figure 11—figure Supplement 1, Figure 11—figure Supplement 2, Figure 11—figure Supplement 3, Figure 11—figure Supplement 4, Figure 11—figure Supplement 5, Figure 11—figure Supplement 6, Figure 11—figure Supplement 7, Figure 11—figure Supplement 8*).
450 Similarly to other modalities, the transfer function weights reflected the inherent autocorrelation
451 structure of the behavioral data (*Figure 11—figure Supplement 9, Figure 11—figure Supplement 10, Figure 11—figure Supplement 11*).

455 **Discussion**

456 The primary objective of cognitive neuroscience is to comprehend how the brain executes information-
457 processing operations (*Kay, 2018*). Linearized encoding analysis serves as a robust method to evaluate
458 a model (i.e., transfer function) that describes how the brain encodes sensory information
459 from the environment and processes this information further (*Naselaris et al., 2011*). Prediction
460 accuracy is crucial as it measures the model’s ability to generalize to unseen stimulus-response
461 pairs. This paper elucidates how information leakage in training examples (i.e., RDD) can artificially
462 inflate prediction accuracy and demonstrates this through extensive simulations and real
463 data analyses.

464 Firstly, it was mathematically shown that the expected prediction accuracy of the null feature
465 could exceed zero when the null features are identically repeated across CV partitions (*Equation 20*).
466 It is important to understand that RDD arises not from noise in the data but from the stimulus-
467 driven similarity. In particular, the similarity between training and optimization sets disables regu-
468 larization leading to an inflation of the null prediction accuracy.

469 Secondly, simulations showed that the RDD effect (i.e., the inflation of the prediction accuracy
470 due to RDD; an interaction with *IsRep*) is more pronounced with a higher signal-to-noise ratio,
471 greater flexibility (i.e., longer delay points or higher dimensional features), and more similar auto-
472 correlation structures between the true and null features (*Table 4; Table 6*).

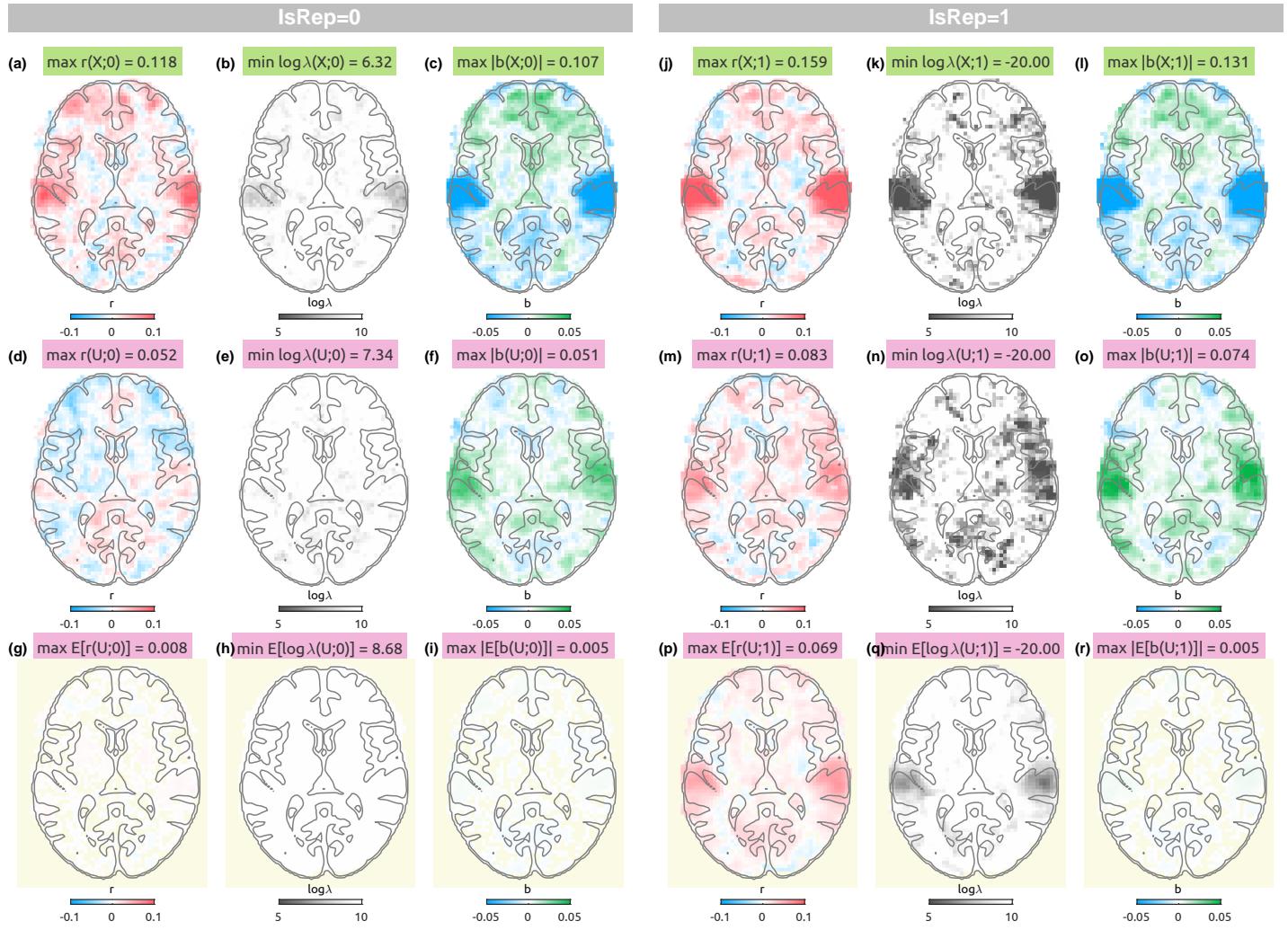


Figure 9. fMRI linearized encoding analysis results with delays from 3 to 9 sec with the audio envelope (top row, (a-c, j-l)), a single case of a phase-randomized envelope (middle row, (d-f, m-o)), and an average of 100 phase-randomized envelopes (bottom row, pale yellow background, (g-i, p-r)). For each CV scheme ($\text{IsRep} = 0$, left panels, (a-i); $\text{IsRep} = 1$, right panels, (j-r)), prediction accuracy (r , blue to red, (a, d, g, j, m, p)), logarithmic ridge hyperparameter ($\log_{10} \lambda$, gray to white, (b, e, h, k, n, q)), transfer function weights that are summed over delays (b , blue to green, (c, f, i, l, o, r)) are shown along the columns. The analysis was done in the 3-D space, but transverse slices (Montreal Neurological Institute [MNI]-coordinate $Z = 8$ mm) are chosen to display anatomical structures implicated in a meta-analysis on music-evoked emotions (Koelsch, 2020) such as the Heschl's gyrus, planum temporale, the inferior frontal cortex. The 3-D volumes can be viewed with the NeuroVault web viewer (<https://identifiers.org/neurovault.collection:19626>).

Figure 9—figure supplement 1. fMRI linearized encoding analysis with normal noise as null features and delays from 3 to 9 sec

Figure 9—figure supplement 2. fMRI linearized encoding analysis with uniform noise as null features and delays from 3 to 9 sec

Figure 9—figure supplement 3. fMRI linearized encoding analysis with the phase-randomized envelope and delays from 4 to 6 sec

Figure 9—figure supplement 4. fMRI linearized encoding analysis with normal noise as null features and delays from 4 to 6 sec

Figure 9—figure supplement 5. fMRI linearized encoding analysis with uniform noise as null features and delays from 4 to 6 sec

Figure 9—figure supplement 6. fMRI linearized encoding analysis with the phase-randomized envelope and delays from 0 to 12 sec

Figure 9—figure supplement 7. fMRI linearized encoding analysis with normal noise as null features and delays from 0 to 12 sec

Figure 9—figure supplement 8. fMRI linearized encoding analysis with uniform noise as null features and delays from 0 to 12 sec

Figure 9—figure supplement 9. fMRI transfer function weights with all three noise models and delays from 3 to 9 sec

Figure 9—figure supplement 10. fMRI transfer function weights with all three noise models and delays from 4 to 6 sec

Figure 9—figure supplement 11. fMRI transfer function weights with all three noise models and delays from 0 to 12 sec

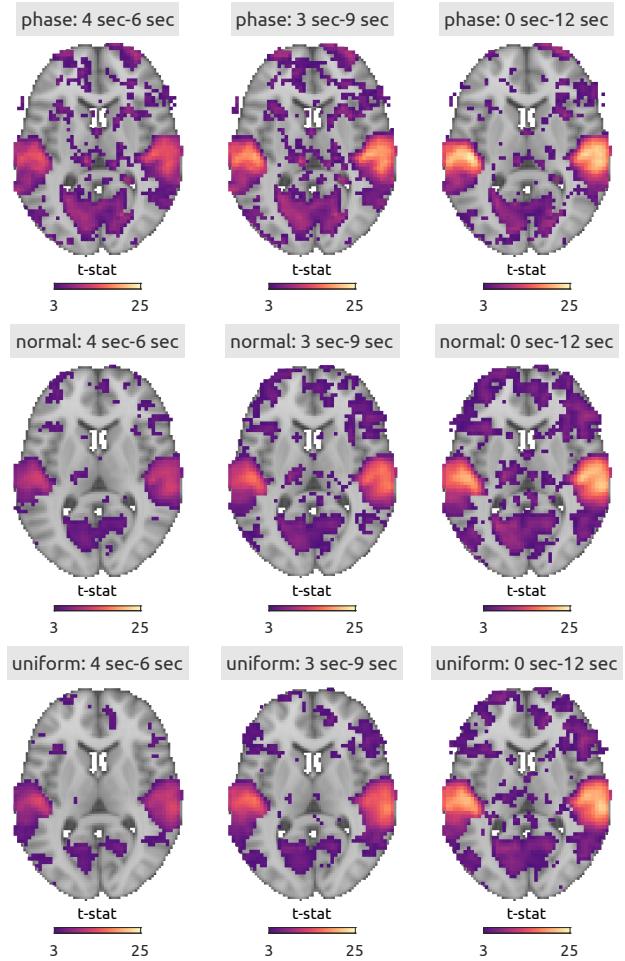


Figure 10. t -statistic maps on the transverse slices comparing the null prediction accuracies between two CV schemes (e.g., **Figure 9g** vs. **Figure 7p**) to test the RDD effects in the fMRI data with the phase-randomized envelope (top row), the normal noise (middle row), and the uniform noise (bottom row). Voxels were thresholded by statistical significance after FDR adjustment ($P_{FDR} < 0.01$). The background anatomical image is the MNI template included in FSL (MNI152_T1_2mm_brain.nii.gz). The 3-D volumes can be viewed with the NeuroVault web viewer (<https://identifiers.org/neurovault.collection:19626>).

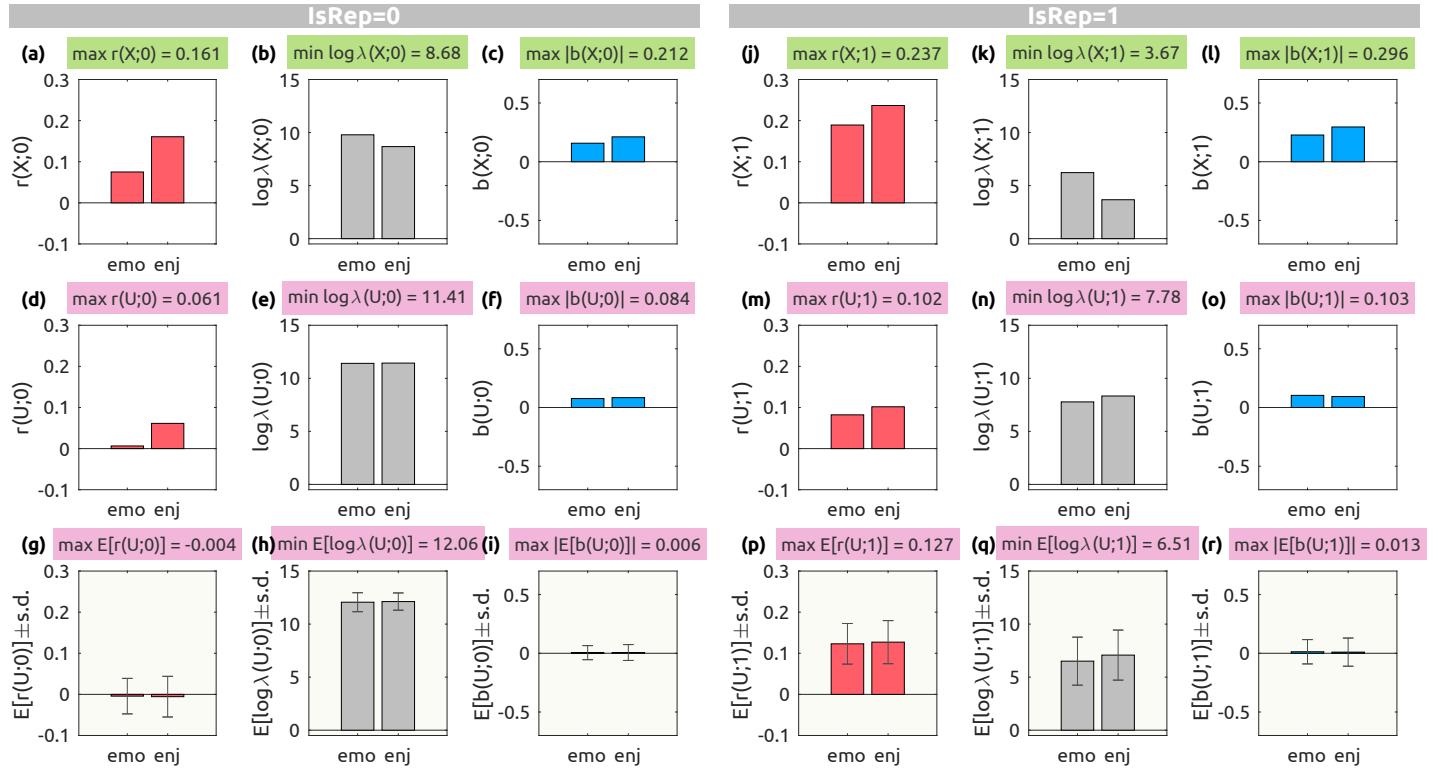


Figure 11. Behavioral linearized encoding analysis results with delays from 0 to 10 sec with an audio envelope (top row, (a-c, j-l)), a single case of the phase-randomized envelope (middle row, (d-f, m-o)), and an average of 100 phase-randomized envelopes (bottom row, pale yellow background, (g-i, p-r)). For each CV scheme (IsRep = 0, left panels, (a-i); IsRep = 1, right panels, (j-r)), prediction accuracy (r , red bars, (a, d, g, j, m, p)), logarithmic ridge hyperparameter ($\log_{10} \lambda$, gray bars, (b, e, h, k, n, q)), transfer function weights that are summed over delays (b , blue bars, (c, f, i, l, o, r)) are shown along the columns. For the averaged metrics (bottom row), the standard deviations are shown as error bars. Emo.: emotionality, Enj.: enjoyment.

Figure 11—figure supplement 1. Behavioral linearized encoding analysis with the normal noise and delays from 0 to 10 sec

Figure 11—figure supplement 2. Behavioral linearized encoding analysis with the uniform noise and delays from 0 to 10 sec

Figure 11—figure supplement 3. Behavioral linearized encoding analysis with the phase-randomized envelope and delays from 0 to 5 sec

Figure 11—figure supplement 4. Behavioral linearized encoding analysis with the normal noise and delays from 0 to 5 sec

Figure 11—figure supplement 5. Behavioral linearized encoding analysis with the uniform noise and delays from 0 to 5 sec

Figure 11—figure supplement 6. Behavioral linearized encoding analysis with the phase-randomized envelope and delays from 0 to 15 sec

Figure 11—figure supplement 7. Behavioral linearized encoding analysis with the normal noise and delays from 0 to 15 sec

Figure 11—figure supplement 8. Behavioral linearized encoding analysis with the uniform noise and delays from 0 to 15 sec

Figure 11—figure supplement 9. Behavioral transfer function weights with all three noise models and delays from 0 to 10 sec

Figure 11—figure supplement 10. Behavioral transfer function weights with all three noise models and delays from 0 to 5 sec

Figure 11—figure supplement 11. Behavioral transfer function weights with all three noise models and delays from 0 to 15 sec

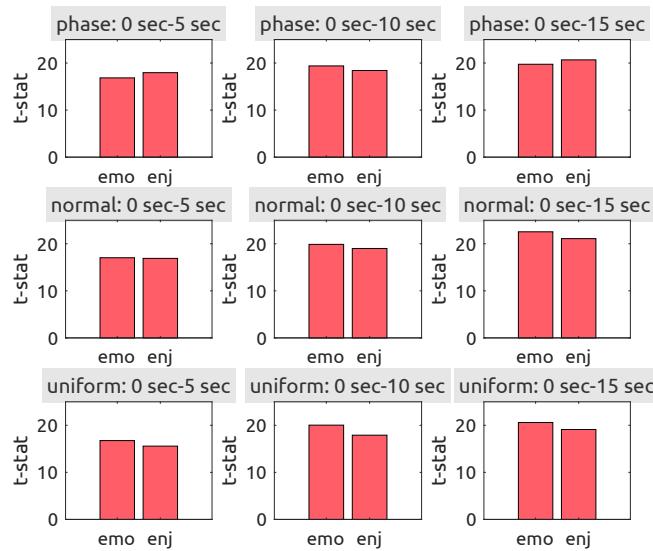


Figure 12. *t*-statistic bar plots comparing the null prediction accuracies between two CV schemes (e.g., [Figure 11g](#) vs. [Figure 7p](#)) to test the RDD effects in the behavioral data with the phase-randomized envelope (top row), the normal noise (middle row), and the uniform noise (bottom row). All effects were statistical significant after FDR adjustment ($P_{FDR} < 0.01$). Emo.: emotionality, Enj.: enjoyment.

473 Lastly, the RDD effect was consistently observed across popular data modalities in cognitive
 474 neuroscience ([Figure 8](#), [Figure 10](#), [Figure 12](#)). In particular, the inflated prediction accuracy exhibited
 475 highly plausible spatial patterns even when predicted by uniform noise as a null feature. It is
 476 essential to emphasize that these patterns are driven by time-locked neural responses to repeated
 477 stimuli (widely known as inter-trial-/subject synchrony), not by random noise in the data. Therefore,
 478 when combined with *informal reverse inference* (i.e., falsely inferring a mental process from a brain
 479 activity pattern without accounting for base rates), which is also a common logical fallacy in cog-
 480 nitive neuroscience ([Poldrack, 2006](#)), RDD can lead to completely incorrect conclusions (e.g., “*The*
 481 *auditory cortex was encoding this uniform random noise that was never presented to the participant.*”;
 482 [Figure 9—figure Supplement 2](#)).

483 **But my features are not just random noise!**

484 As clearly demonstrated in the current paper, RDD can lead to spurious findings in cognitive neuro-
 485 science research when combined with questionable research practices such as adaptation of novel
 486 analysis frameworks without considering the original design of their data and informal reverse in-
 487 ference.

488 I used random noise to demonstrate that the RDD effect can be observed with a feature that
 489 is not expected to be encoded in the real data. Certainly, no one would sincerely expect the audi-
 490 tory cortex to encode a random feature that cannot be extracted from the stimulus. In practice,
 491 researchers hypothesize that certain information extracted from the stimulus is encoded in the re-
 492 sponse based on some theoretical and/or empirical grounds. However, RDD can inflate the predic-
 493 tion accuracy of their hypothesized feature well above the chance level even when it is not encoded.
 494 This, in turn, can misguide future research and contribute to contamination of the literature.

495 Also, note that the RDD effect was greater for the phase-randomized envelope than the normal
496 or uniform noise (*Figure 8*, *Figure 10*, *Figure 12*). This is because the phase-randomized envelope
497 preserves the autocorrelation structure of the stimulus, making it more similar to true features
498 than random noise. Since the hypothesized feature is extracted from actual stimuli, it necessarily
499 bears some resemblance to true features, regardless of the nonlinear operations involved. There-
500 fore, the risk of RDD is greater when the hypothesized feature is derived from the stimulus rather
501 than from random noise.

502 Often nested regression models are compared to determine the unique predictive contribution
503 of a feature of interest. For example, in our previous study (*Leahy et al., 2021*), the prediction ac-
504 curacy of a model with an audio envelope (“reduced model”) was subtracted from the prediction
505 accuracy of a model with the envelope and musical beats (“full model”) to test the encoding of mu-
506 sical beats⁶. Even in such cases, RDD can still inflate the prediction accuracy of the hypothesized
507 feature (or any additional random features) as shown in the simulation where additional random
508 features moderately increased the RDD effect (*Table 8—source data 1*, the interaction between
509 the number of features and stimulus repetition on null prediction: $\eta_p^2 = 0.101$). Once again, in a
510 proper cross-validation without information leakage, irrelevant features would have been regu-
511 larized. However, in the presence of RDD, the regularization is disabled, and the irrelevant features
512 are not penalized, thus leading to an inflated prediction accuracy.

513 **Is RDD really a different type of leakage than double-dipping?**

514 All types of information leakage are essentially a circular fallacy (*Kaufman et al., 2012*). However,
515 RDD is different from the more widely known type of circular fallacy—double dipping (*Kriegeskorte
516 et al., 2009*)—in terms of the locus of the circularity. In double-dipping, it is the identical noise that
517 is repeated in a selective analysis that uses the same dataset twice. The identical random noise
518 appears once in setting up a selective analysis (e.g., channels or voxels of interest) and once again
519 in testing the selective analysis, thus yielding spurious findings.

520 In reverse-double-dipping, it is the underlying signal that is repeated in a predictive analysis
521 that uses apparently different datasets. Because the data points were acquired at different times
522 with independent noise (but with the same stimuli), an identical underlying signal may seem less
523 obvious. Because of the independence of noise, even researchers who are well-aware of the issue
524 of double-dipping might not notice the circularity in their analyses.

525 **Is RDD relevant to other analyses?**

526 Given the deceptive nature of RDD, readers may wonder whether the displayed problem is specific
527 to linearized encoding analysis or also relevant to other popular analyses in cognitive neuroscience.
528 Here, I discuss the relevance of RDD to other analyses.

⁶While the comparison of nested models is a standard practice in Ordinary Least Squares (OLS) regression, it can be complicated with regularized regression such as ridge. The major problem is over-regularization of relevant features due to irrelevant features in the full model when a single penalty hyperparameter is used for all features (feature spaces). A multi-penalty model can address this issue better (e.g., *La Tour et al., 2022; Kim et al., 2024*).

529 Beta image encoding

530 Unlike the FIR model that estimates transfer function weights for each time point, the beta image
 531 encoding model is on top of the classical GLM that estimates ‘beta’ weights for each stimulus. This
 532 approach is popular in visual fMRI experiments where the gazing and visual attention of partici-
 533 pants is difficult (or unnecessary) to resolve in time (*Kay et al., 2008*) or auditory fMRI experiments
 534 with short (1-2 seconds) stimuli (*Moerel et al., 2018*). Thus, instead of directly handling the autocor-
 535 related BOLD time series, an average activation amplitude (i.e., beta weight) during a short trial is
 536 first estimated using a GLM, either using a theoretical transfer function called the ‘canonical hemo-
 537 dynamic response function (cHRF)’ (*Henson et al., 1999*) or by fitting a regularized FIR model on a
 538 split data (*Prince et al., 2022*).

539 For more general noise covariance structures, the beta estimation can be described as a Gen-
 540 eralized Least Squares (GLS) problem (*Lage-Castellanos et al., 2019*):

$$\xi \equiv \hat{\beta} = (\Phi^T \Omega^{-1} \Phi)^{-1} \Phi^T \Omega^{-1} \mathbf{y}, \quad (23)$$

541 where $\xi \in \mathbb{R}^{M \times 1}$ is the estimated stimulus-response vector for M stimuli, $\Phi \in \mathbb{R}^{T \times M}$ is the cHRF-
 542 convolved design matrix for T time points, $\Omega \in \mathbb{R}^{T \times T}$ is the autocovariance matrix of the noise in
 543 the BOLD time series, and $\mathbf{y} \in \mathbb{R}^{T \times 1}$ is the BOLD time series. Then, the estimated beta weights are
 544 subjected to an encoding model:

$$\xi = \mathbf{F}\mathbf{g} + \mathbf{e}, \quad (24)$$

545 where $\mathbf{F} \in \mathbb{R}^{M \times F}$ is a matrix that describes F features for M presented stimuli, $\mathbf{g} \in \mathbb{R}^{F \times 1}$ is a feature-
 546 response vector of the voxel, and $\mathbf{e} \in \mathbb{R}^{M \times 1}$ is unknown noise.

547 Let us consider a case where two sets of beta images with M identical stimuli were partitioned
 548 into CV folds (e.g., even runs vs. odd runs where all M stimuli were presented in each run in
 549 randomized orders). Then, the expected null prediction accuracy by the Red Team would be also
 550 positive:

$$\mathbb{E} [\text{corr}(\hat{\xi}_{2,\mathbf{H}}, \xi_2)] \stackrel{(18)}{\approx} \mathbf{s}_i^T \mathbf{H}_1 (\mathbf{H}_1^T \mathbf{H}_1)^{-1} \mathbf{H}_1^T \mathbf{s}_i \geq 0, \quad (25)$$

551 where $\mathbf{s}_i = \mathbf{F}_i \mathbf{g} \in \mathbb{R}^{M \times 1}$ is the underlying signal pattern (a “response profile”) for M stimuli in the i -th
 552 CV partition ($i = 1$, training; $i = 2$, testing), $\mathbf{H}_i \in \mathbb{R}^{M \times F}$ is the null feature matrix.

553 Although it is standard practice to randomize the presentation order of stimuli across runs
 554 and participants, the beta image estimation process can reorganize the response profiles, aligning
 555 them in a consistent order across all runs (e.g., sequentially from the first to the M -th stimuli).
 556 Therefore, the risk of RDD in the beta image encoding analysis remains a concern.

557 Stimulus reconstruction

558 Reconstruction of unseen stimuli (e.g., images or sounds) based on neural data demonstrates the
 559 remarkable potential of neuroimaging techniques for “mind-reading” (*Kay et al., 2008; Santoro
 560 et al., 2017; Han et al., 2019*). In practice, a set of linear models decodes features from neural data
 561 (e.g., beta images), followed by a reconstruction step where a simple classifier or a deep neural
 562 network (such as variational autoencoder) synthesizes the stimulus from the decoded features.

563 For instance, a multivariate decoding model based on beta images can be described as:

$$\mathbf{F} = \mathbf{W}\Xi + \mathbf{e} \quad (26)$$

564 where $\mathbf{F} \in \mathbb{R}^{M \times F}$ is an F -dimensional feature matrix for M stimuli, $\mathbf{W} \in \mathbb{R}^{V \times F}$ is a spatial filter to
 565 decode the features, $\Xi \in \mathbb{R}^{M \times V}$ is the matrix of vectorized beta images for M stimuli and V voxels.
 566 Since the decoding model is also a linear model like the encoding model, in principle, RDD can also
 567 occur. However, the goal of the analysis makes the requirement of “unseen stimuli” more explicitly.
 568 For example, if one fit a linear model with $N - 1$ subjects’ responses to Stimulus A, and then tries to
 569 recover, once again, Stimulus A from a response of the i -th subject, the problem of the circularity in
 570 this design would be more visible. Having said that, it is still possible that latent similarity of stimuli
 571 could be overlooked. Especially if there are too many stimuli to manually inspect, it is possible
 572 that certain stimuli are indeed non-identical but still highly similar (e.g., utterances by the same
 573 speaker; repetitions of musical phrases). Checking inter-stimulus correlation for all features prior
 574 to partitioning them into training and test sets will prevent such cases of hidden ‘twinning’.

575 **Multivariate classification**

576 Multivariate classification analysis has been widely used in the cognitive neuroscience commu-
 577 nity, commonly known as multivoxel (or multivariate) pattern analysis (MVPCA; *Kriegeskorte et al.*,
 578 **2006**) or single-trial classification in electrophysiological data such as EEG, MEG, and ECoG (*Müller-*
 579 *Gerking et al., 1999; Pistoohl et al., 2012; Quandt et al., 2012*), either on the whole set of response
 580 units (e.g., whole-brain classification; *Ryali et al., 2010*) or a striding set of local neighbors (e.g., a
 581 “searchlight”; *Kriegeskorte et al., 2006*). While a highly accurate classifier is necessary to build a
 582 brain-computer interface system, classification analysis can be useful even with low but significant
 583 accuracy—often the case in cognitive neuroscience research—as a multivariate model that tests
 584 the existence of certain information in the brain.

585 In a simple case of two classes, a linear classifier can be constructed (*Hastie et al., 2009*) as:

$$C = \text{sign}(\mathbf{a}^T \mathbf{w} + w_0), / \quad (27)$$

586 where $\text{sign} : \mathbb{R} \rightarrow \{-1, +1\}$ is a sign function, $\mathbf{a} \in \mathbb{R}^{V \times 1}$ is an activation pattern (e.g., an M/EEG
 587 topography or beta values in a region of interest) for an unknown class instance, $\mathbf{w} \in \mathbb{R}^{V \times 1}$ is a
 588 classification weight vector, w_0 is a scalar bias term, and $C \in \{-1, +1\}$ is a class label.

589 Classification can be seen as model-free as compared to model-based encoding or decoding
 590 analysis because the classification does not require a definition of a ‘model’ that describes which
 591 feature of the stimulus contributes to which response to what extent. Note that the trained weight
 592 vector only represents a separation of given training examples in a functional space, regardless
 593 of the features of the stimuli. For example, one may try to classify ‘happy music’ vs. ‘sad music’
 594 based on the EEG responses. However, if the chosen exemplars of the classes are imbalanced
 595 (e.g., all exemplars of ‘happy music’ naturally happened to be faster and louder than ‘sad music’),
 596 classification could be strongly driven by different acoustics, not necessarily due to perceived or
 597 evoked emotions from the music. That is, without carefully matching the training examples for all
 598 relevant features, the interpretation of classification analysis may remain unclear.

599 Moreover, in principle, the repetition of stimuli (or at least classes) across data points is not just
 600 unavoidable but in fact necessary for the classification. In order to train a classifier, not only linear
 601 but also general, balanced training and testing examples of all classes are required (*Hastie et al.*,

602 2009). In other words, it is impossible to train a classifier for an ‘unseen class’ while it can be trained
603 for an ‘unseen instance’ of a known class. Therefore, it can be noted that RDD does not concern
604 classification analysis although the forward RDD (i.e., double-dipping) greatly concerns the feature
605 selection process of classification.

606 **Representation similarity analysis**

607 Finally, let us consider a method that evaluates the second-order isomorphism across representa-
608 tional systems (*Kriegeskorte et al., 2008*), as widely known as representational similarity analysis
609 (RSA). This flexible method can define a ‘model’ distance matrix between stimuli either based on
610 class labels (as in classification analysis) or feature descriptors (as in encoding analysis). The origi-
611 nal formulation of RSA does not involve cross-validation as the RSA is neither a predictive model nor
612 classification (*Kriegeskorte et al., 2008*). A later extension introduced a cross-validated, squared
613 Mahalanobis distance estimator—known as “crossnobis”—to enhance both reliability and inter-
614 pretability (*Diedrichsen and Kriegeskorte, 2017*). When estimating a crossnobis distance between
615 two conditions in a leave-one-out cross-validation (LOOCV) scheme, it is assumed that the num-
616 ber of responses for both conditions remains balanced across all CV partitions. Thus, if identical
617 stimuli are repeated across partitions as in case of RDD, then the crossnobis distance for the brain
618 RDM would be biased by the stimulus-specific (not condition-specific) activity. However, even so,
619 null descriptors would define a model RDM that is irrelevant to the brain RDM. Therefore, a false
620 conclusion that “a brain region represents null information” cannot be wrongly supported by the
621 association between the model and brain RDMs.

622 RSA primarily focuses on second-order association across RDMs, typically assessed using linear
623 (non-negative) regression. However, this association has been seldomly tested for its generalizabil-
624 ity. For instance, two model RDMs—one based on object identity (e.g., a human face vs. a house)
625 and another based on pixel intensity—and a brain RDM from trials where a subject views these
626 images. In a training set, the association weights can be estimated as

$$\mathbf{d}_{\text{brain}} = \beta_0 + \mathbf{d}_{\text{identity}}\beta_1 + \mathbf{d}_{\text{intensity}}\beta_2 \quad (28)$$

627 where \mathbf{d} is a vector of flattened upper triangular elements excluding the diagonal of an RDM. In
628 a test set, the relationship between RDMs for ‘unseen’ pictures can be predicted by the weights
629 obtained from the training set, similarly to any predictive regression models (i.e., the third-order
630 similarity). In such a case (i.e., a replication of RSA findings), the repetition of stimuli across train-
631 ing and test sets (i.e., ‘unseen’ pictures were actually ‘seen’ pictures) could inflate the prediction
632 accuracy. However, introducing null features would disrupt the model RDMs and would greatly de-
633 crease the weights (i.e., $\hat{\beta}_i$) already in the training set, making it unlikely to observe a high prediction
634 accuracy. In conclusion, the risk of RDD in RSA seems minimal.

635 **How do we detect and prevent RDD?**

636 The linearized encoding model has been utilized in various neuroimaging studies for many years
637 (*Kay et al., 2008*), so it may be surprising that this pitfall has gone underrecognized for so long. In
638 fact, many leading groups in this field have employed a hold-out validation design rather than cross-
639 validation. In this approach, the testing dataset is deliberately acquired separately from the training

640 set by the experiment design. That is, training stimuli are repeated and averaged, while different
641 stimuli are presented only once for training (e.g., *Nishimoto et al., 2011; Huth et al., 2016; Han et al.,*
642 *2019*). As a result, model predictions are evaluated only once per subject. This design prevents
643 accidental stimulus repetition across CV partitions and thus avoids the risk of RDD. However, when
644 the data is collected for a different purpose (e.g., intersubject synchrony) and researchers later
645 apply encoding analysis, the likelihood of encountering RDD increases. This risk is particularly high
646 for novice researchers who may misinterpret the valid recommendation to construct a ‘subject-
647 independent model’ (*Crosse et al., 2021*), which involves cross-validation across subjects. Without
648 careful implementation, they may inadvertently create a “stimulus-specific” model, as illustrated
649 in *Figure 6*, leading to RDD pitfall. Thus, in this section, I provide practical recommendations for
650 detecting and preventing RDD.

651 **Inter-trial correlation diagnosis**

652 As repeatedly shown throughout the paper, RDD occurs when the identical stimulus is presented
653 more than once across CV partitions. Thus, a simple way to test a CV design for the risk of RDD is
654 to compute the correlation across data epochs before partitioning the data into training, optimiza-
655 tion, and test sets. If the correlation is high, RDD is likely to occur. This was already illustrated in
656 the toy example. When no stimulus was repeated across CV partitions the inter-trial correlation
657 (ITC) was on average zero across 200 random samplings (*Figure 2q*). However, when the stimulus
658 was repeated across CV partitions, the ITC was on average about 0.8 (*Figure 3q*). Because this cor-
659 relation can be cheaply computed prior to costly optimization and modeling fitting, this can be a
660 useful diagnostic tool to assess risk for RDD.

661 Checking ITC is also useful to detect latent similarity across stimuli as well. For example, two
662 audio files are named differently but contain similar music, the researcher would not know about
663 the leakage in the training examples. The ITC can be a useful tool to detect such latent similarity.

664 For users’ convenience, an automatic validation test for a given CV design based on feature-
665 ITC and response-ITC is implemented as a default option in the MATLAB package for Linearized
666 Encoding Analysis (LEA; <https://github.com/seunggookim/lea>).

667 **Hold-out validation**

668 The most straightforward but also most expensive way to prevent RDD is to use a hold-out valida-
669 tion design. In this design, the testing dataset is deliberately acquired separately from the training
670 set by the experiment design. For example, the training stimuli are repeated and averaged, while
671 different stimuli are presented only once for training. This design prevents accidental stimulus
672 repetition across training and test partitions and thus avoids the risk of RDD.

673 However, the hold-out validation requires more data points and higher SNR to work as com-
674 pared to cross-validation. In cross-validation, the out-sample prediction performance is repeatedly
675 assessed and averaged across different CV partition schemes to achieve an accurate estimation in
676 the presence of sampling variance. In hold-out validation, the out-sample prediction performance
677 is assessed only once. Thus, the sampling variance needs to be already suppressed to achieve a
678 similarly accurate estimation of the true prediction performance.

679 Single-use stimulus

680 Another straightforward design to prevent RDD is not to re-use a stimulus during the whole study.
681 That is, a stimulus is presented to only one participant, only once, and never used again (i.e., a
682 single-use stimulus). This design prevents any accidental stimulus repetition across CV partitions
683 and thus avoids the risk of RDD. That is, any CV partitioning (across stimulus or participants or
684 both) cannot occur RDD.

685 However, this design may be not practical in most cases as it requires a large number of stimuli
686 to be presented to each participant in order to ensure the sampling variance due to stimuli is
687 sufficiently reduced. With a small number of stimuli, the attribution of the observed difference
688 between participants would be ambiguous—to what extent should it be attributed to inter-subject
689 variability and inter-stimulus variability?

690 Group-level modeling

691 In practice, acquiring neuroimaging data of high quality is still highly time-consuming and expen-
692 sive. In particular, acquiring extensive data from vulnerable populations (e.g., patients, young chil-
693 dren, elderly people) poses not only financial and logistical but also ethical concerns. Thus, many
694 researchers may be motivated to combine the limited data across multiple subjects to further re-
695 duce the sampling variance.

696 As discussed, if an identical set of stimuli were used (as a standard procedure) across subjects,
697 RDD could arise. In such a case, one option is to simply average all data across subjects within
698 each group to test the group difference by cross-encoding (see below). If only a single long stim-
699 ulus was used, multiple pseudo-trial segments should be created with appropriate gaps to avoid
700 carry-over effect (e.g., for fMRI data, a gap of at least six seconds is recommended to account for
701 the hemodynamic delay; for M/EEG data, a gap of one second or more), and design a CV across
702 pseudo-trial segments. To detect potential repetition across segments, checking ITC of features is
703 also recommended. While a one-size-fits-all hard threshold cannot be recommended because the
704 temporal autocorrelation alters the expected variance of the null correlation, a strong correlation
705 across segments (e.g., $r > 0.5$) may indicate a risk of RDD. In that case, design a different CV scheme
706 or a different segmentation scheme.

707 A group-level inference is typically performed by testing the null hypothesis that the mean
708 subject-wise prediction accuracy is indifferent between groups (e.g., clinical cases vs. healthy con-
709 trols), using either a parametric test (e.g., a t -test based on the t -distribution) or a non-parametric
710 test (e.g., a permutation test based on a null distribution generated by swapping group labels).
711 When analyzing a group difference after averaging data within each group, a group difference can
712 be tested by cross-encoding (i.e., predicting a mean response to the ‘unseen’ segment B in the
713 ‘unseen’ patient group by the weights fitted to a mean response to the ‘seen’ segment A in the
714 control group) with bootstrapping across subjects (i.e., resampling with replacement when creat-
715 ing the mean response). If the prediction is at or below chance level, the group difference can be
716 concluded.

717 **Limitations**

718 While the current paper provides a comprehensive analysis of RDD, some limitations should be
719 acknowledged.

720 First, only a single-penalty ridge regression was considered here. However, other types of reg-
721 ularization (e.g., multi-penalty ridge, LASSO, elastic net) and different modeling techniques (e.g.,
722 support vector regression, non-linear kernel ridge regression) were not explored. While the key
723 mechanism of RDD (i.e., disabling regularization due to the similarity of the underlying signal be-
724 tween the training and optimization sets) is likely to be present in these other methods, the extent
725 to which the RDD effect generalizes to these other methods remains an open question.

726 Also, no neural spike data or intracranial recordings—which is gradually becoming more ac-
727 cessible for human patients—were analyzed in this paper. The RDD effect was only shown in the
728 context of continuous-valued data acquired from non-invasive methods (e.g., EEG, fMRI, and behav-
729 ioral ratings). Sparse spike data or discrete firing rate data may behave differently from continuous-
730 valued data. However, given that the RDD effect is driven by the similarity of the underlying signal,
731 RDD is likely to occur when the data is transformed in such a way that stimulus-evoked, time-locked
732 response is often demonstrated (e.g., high gamma power envelope of a local population of neu-
733 rons; *Ray et al., 2008; Jacobs and Kahana, 2009*).

734 **Conclusion**

735 The current paper shows that RDD is a critical but underrecognized risk in encoding analysis and
736 stimulus reconstruction. This circular fallacy—which I named *reverse double-dipping*—may lead to
737 spurious findings that irrelevant information is encoded in neural signal. When carefully designed,
738 the model-based approach will remain a powerful tool for information-based neuroimaging.

739 **Materials and Methods**

740 **Simulation methods**

741 **Predictors**

742 True features were sampled from an F -dimensional multivariate Gaussian distribution with a cor-
743 relation between adjacent parameters as: $\mathbf{x} \sim \mathcal{N}_F(\mathbf{0}, \Sigma_X)$ where the covariance is defined as:

$$\Sigma_X = \begin{bmatrix} 1 & \rho_X & \rho_X^2 & \cdots & \rho_X^{F-1} \\ \rho_X & 1 & \rho_X & \cdots & \rho_X^{F-2} \\ \rho_X^2 & \rho_X & 1 & \cdots & \rho_X^{F-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_X^{F-1} & \rho_X^{F-2} & \rho_X^{F-3} & \cdots & 1 \end{bmatrix}. \quad (29)$$

744 The f -th feature $\mathbf{x}^{(f)} = [x_1^{(f)}, x_2^{(f)}, \dots, x_T^{(f)}]^T$ is given with an AR(1) temporal autocorrelation ϕ_X as:

$$x_t^{(f)} = x_t^{(f)} + \phi_X x_{t-1}^{(f)}. \quad (30)$$

745 Null features \mathbf{u} were created in the same way as \mathbf{x} , but independently. Only causal (i.e., nonneg-
746 ative) delays were considered from $\{0\}$ to $\{0, \dots, D - 1\}$ in creating a design matrix by horizontally

747 concatenating Toeplitz matrices as:

$$\mathbf{X} = \begin{bmatrix} x_1^{(1)} & \cdots & x_{1-(D-1)}^{(1)} & x_1^{(2)} & \cdots & x_{1-(D-1)}^{(2)} & \cdots & x_1^{(F)} & \cdots & x_{1-(D-1)}^{(F)} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \cdots & \vdots & \ddots & \vdots \\ x_T^{(1)} & \cdots & x_{T-(D-1)}^{(1)} & x_T^{(2)} & \cdots & x_{T-(D-1)}^{(2)} & \cdots & x_T^{(F)} & \cdots & x_{T-(D-1)}^{(F)} \end{bmatrix}. \quad (31)$$

748 For convenience, all predictors were standardized to zero-mean and unit variance.

749 Responses

750 Responses $\mathbf{Y}_i \in \mathbb{R}^{T \times V}$ were generated by plugging in a design matrix \mathbf{X}_i , true weights $\mathbf{B} \in \mathbb{R}^{FD \times V}$,
 751 and noise $\mathbf{E}_i \in \mathbb{R}^{T \times V}$ to the FIR model **Equation 1** for the i -th partition as $\mathbf{Y}_i = \mathbf{X}_i \mathbf{B} + \mathbf{E}_i$ where
 752 $\mathbf{Y}_i = [\mathbf{y}_i^{(1)}, \mathbf{y}_i^{(2)}, \dots, \mathbf{y}_i^{(V)}]$, $\mathbf{B} = [\mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \dots, \mathbf{b}^{(V)}]^T$, and $\mathbf{E}_i = [\mathbf{e}_i^{(1)}, \mathbf{e}_i^{(2)}, \dots, \mathbf{e}_i^{(V)}]$.

753 True weights of the v -th variate $\mathbf{b}^{(v)}$ (a column of a response variable) were created from the
 754 multivariate Gaussian distribution with an AR(1) temporal correlation ϕ_B along the delays and a
 755 covariance σ_b across predictors corresponding to the rows of the design matrix as:

$$\mathbf{b}^{(v)} = [b_{1,1}^{(v)}, \dots, b_{1,D}^{(v)}, b_{2,1}^{(v)}, \dots, b_{2,D}^{(v)}, \dots, b_{F,1}^{(v)}, \dots, b_{F,D}^{(v)}]^T \quad (32)$$

756 where $b_{f,d}^{(v)}$ is a scalar coefficient for the v -th variate, the f -th feature, and the d -th delay. Noise $\mathbf{e}^{(v)}$ at
 757 the v -th variate was also created as Gaussian noise with an AR(1) temporal correlation ϕ_E as $e_t^{(v)} =$
 758 $e_t^{(v)} + \phi_E e_{t-1}^{(v)}$. In case of multivariate models ($V \geq 2$), weights and noise have spatial autocorrelation:
 759 $\mathbf{b}^{(v)} = \mathbf{b}^{(v)} + \theta_B \mathbf{b}^{(v-1)} + \theta_B \mathbf{b}^{(v+1)}$ and $\mathbf{e}^{(v)} = \mathbf{e}^{(v)} + \theta_E \mathbf{e}^{(v-1)} + \theta_E \mathbf{e}^{(v+1)}$ where a variate v is a neighbor of $v - 1$
 760 and $v + 1$ except for boundaries ($v = 1$ and $v = V$).

761 Before summing the signal \mathbf{XB} and error \mathbf{E} , the mean variance of error across variates was
 762 scaled so that it achieved the intended signal-to-noise ratio as $S = 10 \log_{10} \sigma_S^2 / \sigma_E^2$. Then, all response
 763 variates were standardized to zero-mean and unit-variance.

764 Optimization and evaluation

765 Ridge parameters were optimized via a grid search (λ -grid = $10^{[-10, -9, \dots, 10]}$) for each variate indepen-
 766 dently. To implement a nested 4-by-3-fold cross-validation (CV), in total 4 “trials” were generated
 767 for each random sampling. For the 4-fold-outer-loop, 3 trials (outer-training) vs. 1 trial (outer-test)
 768 were partitioned. Then, within each 3-fold-inner-loop with the 3 trials, 2 trials (inner-training) vs.
 769 1 trial (inner-optimization) were partitioned. In total 12 different partitions (each was called a “CV-
 770 fold”) were used for training (50%), optimization (25%), and test (25%) models. Prediction accuracy
 771 of Pearson correlation was averaged across the CV-folds. Random sampling was repeated for 1000
 772 times for each combination of model parameters.

773 Computational considerations

774 To minimize the number of inversion operations, a general linear model (GLM) formulation was
 775 used. That is, for each possible λ , a regularized covariance matrix was inverted only once for all
 776 variates (i.e., voxels or channels). The sum of squared errors was then temporally stored, allowing
 777 the optimal λ to be determined for each variate. This approach is more efficient than a naïve
 778 method of inverting the covariance matrix separately for each variate, which would redundantly
 779 repeat the same calculation for all variates. Also, when computing Predictions, variates with the
 780 same optimal λ were grouped together into GLM models to reduce the number of inversions.

781 Computation was carried out using an in-house high-performance computing (HPC) server,
782 where a user is allowed to utilize up to 192 CPUs of Intel Xeon Gold 6130 [2.10 GHz] in parallel.
783 The actual utilization varied between 32–192 CPUs, depending on the demands of other users. An
784 individual job of 1,000 random samplings took 100–400 seconds of CPU time. A total of 45,899
785 jobs, which amounted to approximately 7 months of CPU time, was completed in about one week
786 on the HPC server.

787 **Real data methods**

788 For clarity, the original studies (*Kaneshiro et al., 2020; Sachs et al., 2020*) focused on inter-subject
789 synchrony in neural responses (EEG and fMRI) without specific hypotheses about the encoded
790 features. Thus, these studies did not suffer from RDD.

791 **EEG data**

792 **Data source** The original study investigated the electroencephalographic correlates of temporal
793 structure and beat in Western-style Indian pop music with Hindi lyrics (i.e., Bollywood music)
794 (*Kaneshiro et al., 2020*). The raw dataset of 48 healthy participants was downloaded from
795 the Stanford Digital Repository (<https://purl.stanford.edu/sd922db3535>).

796 **Data acquisition** Scalp electrical potential data were acquired using a 128-channel Geodesic EEG
797 System 300 (Electrical Geodesics, Inc., Oregon, USA) at a sampling rate of 1 kHz with vertex
798 reference and electrode impedances less than 60 kΩ.

799 **Preprocessing** The authors' preprocessed data (CleanEEG_*) from the repository were used with
800 the channel location file (GSN-HydroCel-125.sfp). As explained in the original publication
801 (*Kaneshiro et al., 2020*), the authors' preprocessing steps involved bandpass filtering (0.3–50
802 Hz), downsampling (125 Hz), ocular artifacts removal using independent component analy-
803 sis, bad channel interpolation, average re-referencing, and epoching. In addition, based on
804 previous EEG encoding analyses where the low-frequency bands of the EEG signal mostly car-
805 ried the audio envelope encoding (*Di Liberto et al., 2015, 2020*), we further bandpass filtered
806 the EEG data to δ and θ bands (1–8 Hz) using the FIR filter in the MATLAB Signal Processing
807 Toolbox (R2022b, RRID:SCR_001622).

808 **Data dimensions** The analyzed EEG data were comprised of 125 channels, 32,878–33,982 time
809 points (263.02–271.86 seconds at 125 Hz), 96 runs (48 subjects \times 2 run) per stimulus with
810 12 stimuli (3 versions of 4 songs). Two runs with the same stimuli were averaged for each
811 participant to increase the signal-to-noise ratio of the evoked response. In the original study
812 (*Kaneshiro et al., 2020*), a participant was randomly assigned to one of 4 versions of the 4 songs
813 (e.g., Intact-Song-A, Measure-shuffled-Song-B, Reserved-Song-C, Phase-scrambled-Song-D).
814 The current study only analyzed the 3 versions except for phase-scrambled version, which
815 did not evoke strong responses.

816 **fMRI data**

817 **Data source** The original study investigated the intersubject correlation in fMRI time series while
818 participants were listening to sad music (*Sachs et al., 2020*). The published data include 2 sad
819 musical pieces and 1 happy musical piece. The raw dataset of 39 healthy participants was

820 downloaded from OpenNeuro (RRID:SCR_005031, <https://openneuro.org/datasets/ds003085>/
821 versions/1.0.0).

822 **Data acquisition** Blood-oxygen-level-dependent (BOLD) signals were acquired using a 3-T Prisma
823 magnetic resonance imaging system with a 32-channel head coil (Siemens Healthineers, Er-
824 langen, Germany). Eight-fold accelerated multiband, T2*-weighted, gradient-echo, echo-planar
825 imaging sequence was used to acquire 40 transverse slices at the sampling rate of 1 Hz and
826 at the isotropic spatial resolution of 3 mm. Anatomical scans were also collected using a
827 T1-weighted contrast sequence at the isotropic spatial resolution of 1 mm.

828 **Preprocessing** After correcting for the susceptibility artifacts using the reversed phase-encoding
829 images with topup in FSL (v6.0.2; <https://fsl.fmrib.ox.ac.uk/fsl/>), SPM12 (v6225; <https://www.fil.ion.ucl.ac.uk/spm/software/spm12/>) was used for slice timing correction and realignment of
830 the functional 4-D time series data. Advanced Normalization Tools (v2.3.5; [https://github.com/](https://github.com/ANTsX/ANTs)
831 [ANTsX/ANTs](https://github.com/ANTsX/ANTs)) were used to perform symmetric diffeomorphic registration between individu-
832 als' anatomical 3-D images and the standard template 3-D image (`MNI152_T1_2mm_brain.nii.gz`)
833 from FSL as well as the coregistration (rigid-body affine transform) between the anatomical
834 3-D image and the temporally averaged functional 3-D image within each subject. The rigid-
835 body and diffeomorphic transformations were combined and applied to each volume of the
836 realigned 4-D functional image, which was resampled only once at the isotropic resolution of
837 3 mm. Thereafter, the functional images were spatially smoothed with an isotropic Gaussian
838 kernel (with a full width at half maximum of 6 mm). ICA-AROMA (v0.4.4.beta) was used to
839 automatically reject 'noise' components to attenuate head motion artifacts and non-BOLD
840 image intensity perturbations (<https://github.com/maartenmennes/ICA-AROMA>).

841 **Data dimensions** The analyzed fMRI data consisted of 62,062 "brain" voxels (based on a brain
842 mask created using `bet` in FSL), 178–525 time points (or seconds), 3 stimuli, and 39 subjects.
843 Due to limited field of views in some subjects, only 61,572 voxels had valid values in the lin-
844 earized encoding analysis results. The authors of the original study excluded 3 subjects from
845 their analysis for either motion artifacts and emotional ratings but included in the shared
846 data repository. In the current analysis, all 39 subjects were analyzed for no apparent mo-
847 tion artifacts after ICA-AROMA.

848

849 Behavioral data

850 **Data source** As part of the fMRI study (*Sachs et al., 2020*), behavioral data of 39 healthy partici-
851 pants were downloaded from OpenNeuro (<https://openneuro.org/datasets/ds003085/versions/1.0.0>).

852 **Data acquisition** After the fMRI session, participants listened to the same stimuli while rating their
853 evoked instantaneous emotions using a physical slider ("fader") in a silent room. Emotional
854 scales of *Emotionality* (the intensity of the evoked feelings of sadness or happiness, depending
855 on the intended emotion of each piece) and *Enjoyment* (the momentary feelings of enjoyment)
856 were rated, one at a time. In total, a participant listened to an identical stimulus for three
857 times including the fMRI session. The slider position values (an integer from 0 to 127) were
858 sampled at about 30.3 Hz.

859

860 **Preprocessing** The imported time series were linearly detrended and downsampled at 5 Hz after
861 an anti-aliasing low-pass filtering using the `resample` function in the MATLAB Signal Process-
862 ing Toolbox (R2022b).

863 **Data dimensions** The analyzed behavioral data comprised 2 emotional scales, 841–2,576 time
864 points depending on the stimulus (178–525 seconds at 5 Hz), 3 stimuli, and 39 subjects.

865 Linearized encoding analysis

866 **Features** True features \mathbf{X} were the audio envelope extracted from the music samples for the well-
867 established auditory response in various types of human brain data including fMRI (*Giraud et al., 2000; Harms et al., 2005; Overath et al., 2012*), EEG (*Aiken and Picton, 2008*), and in-
868 tracranial electrocorticography (*Kubanek et al., 2013*). The envelope was created by summing
869 the output of the “cochlear model” with a 128-channel filterbank, of which characteristic fre-
870 quencies ranged loglinearly from 180 to 7,040 Hz, as in the NSL Auditory-cortical MATLAB
871 Toolbox (v2001; <http://nsl.isr.umd.edu/downloads.html>), and then further downsampled to the
872 sampling rate of the human data (EEG: 125 Hz; fMRI: 1 Hz, behavioral ratings: 5 Hz). Null fea-
873 tures \mathbf{U} were either (a) uniform noise, (b) normal noise, or (c) the phase-randomized envelope.
874 Phase-randomization preserves the amplitude spectrum of the envelope, thus preserves the
875 temporal autocorrelation. This is necessary to correctly estimate the null correlation dis-
876 tribution of the autocorrelated noise. In previous studies (*Leahy et al., 2021; Kaneshiro et al., 2020*),
877 phase-randomization was used to generate an empirical null distribution for
878 non-parameteric statistical inference. The phase randomization was done using fast Fourier
879 transform (FFT), random rotation of phases while preserving complex conjugation for pos-
880 itive and negative frequencies, and followed by an inverse FFT. While this method ('unwin-
881 dowed Fourier transform' in *Theiler et al., 1992*) may introduce spurious high frequencies
882 for non-stationary signals, the envelope of stimuli that include zeros (i.e., silent periods) at
883 the beginning and end of the musical stimuli can be assumed as stationary (i.e., the values of
884 the last samples are similar to the first samples). To estimate the central tendency, 100 noise
885 realizations were created for all cases. All features (either true or null) were standardized
886 prior to fitting.

887 **FIR modeling** As in the simulations, an FIR model was regularized by ridge hyperparameters that
888 are specific for response units (i.e., EEG channels, fMRI voxels, rating scales). Accounting for
889 the inherent time scales of the measures, different delays of the features were used. To
890 demonstrate the effect of the model's flexibility, 3 cases of delays (common choice, shorter,
891 and longer) were used. For EEG, delays of [0 to 0.3 sec; 39 samples], [0 to 0.5 sec; 64 sam-
892 ples], and [0 to 1 sec; 126 samples] were used. For fMRI, [4 to 6 sec; 3 samples], and [3 to
893 9 sec; 7 samples], [0 to 12 sec; 13 samples]. For behaviors, [0 to 5 sec; 26 samples], [0 to
894 10 sec; 51 samples], and [0 to 15 sec; 76 samples]. To avoid transient onset/offset effects at
895 the boundaries of the musical stimuli, the first and last 15 seconds were excluded from the
896 analysis. Both stimuli and responses were standardized before the FIR modeling.

897 **Cross-validation** In all cases, the CV partition was 3-by-2 nested k-fold (i.e., 33%, 33%, 33% for
898 training, optimization, test sets) for the given structures of the real data (i.e., 3 stimuli per

900 participant). To compare CV schemes, prediction accuracies were averaged across CV folds
901 and models (either subject-specific or stimulus-specific).

902 **Statistical inference** The RDD effect is defined as the difference between null prediction accura-
903 cies: $\text{RDD} = \bar{r}_{\text{stim}}(\mathbf{U}; \text{IsRep} = 1) - \bar{r}_{\text{subj}}(\mathbf{U}; \text{IsRep} = 0)$. The null hypothesis is that the expected
904 RDD effect is zero $\mathcal{H}_0 : \mathbb{E}(\text{RDD}) = 0$ and the alterative hypothesis is that RDD effect is posi-
905 tive $\mathcal{H}_A : \mathbb{E}(\text{RDD}) > 0$. Non-parameteric P -values were computed by permutation test ($K =$
906 10,000). That is, for 200 (100 randomizations \times 2 CV-designs) vectors of prediction accura-
907 cies, the binary variable IsRep was randomly permuted ($\max = C(200, 100) > 10^{57}$), and the
908 two-sample t -statistics between two CV-designs were calculated for 10,000 times to form a
909 null distribution. Resulting one-sided P -values were further corrected for the multiple re-
910 sponse units using false discovery rate (FDR) adjustment (*Yekutieli and Benjamini, 1999*) to
911 control the family-wise error rate as $P_{\text{FDR}} < 0.01$.

912 **Software implementation** All encoding analyses of the multimodal datasets (EEG, fMRI, behav-
913 ior) were done using an MATLAB package named Linearized Encoding Analysis (LEA; <https://github.com/seunggookim/lea>), developed by the author.
914

915 Acknowledgments

916 The author declares that the research was conducted in the absence of any commercial or financial
917 relationships that could be construed as a potential conflict of interest. The author confirms being
918 the sole contributor of this work and approved it for publication. Max Planck Society supported
919 this work. The author thanks Dr. Daniela Sammler, Dr. Vincent Shi-chen Chien, Dr. Vincent K.
920 M. Cheung, and Mr. Seung-Cheol Baek for their constructive feedback and expert insights on an
921 earlier version of the manuscript. The author acknowledges the free image sources from www.irasutoya.com (Takashi Mifune), which reserves the copyrights of the source illustrations.
922

923 Data and code availability

924 All analysis code, simulation data, and real data analysis results are available at Zenodo: <https://doi.org/10.5281/zenodo.15100830>. An executable demo of simulation is available at CodeOcean: <https://codeocean.com/capsule/4591394/>. The EEG raw dataset is available at Standford Digital Repository: <https://purl.stanford.edu/sd922db3535>. The fMRI and behavioral raw datasets are available at Open-
925 Neuro: <https://openneuro.org/datasets/ds003085>. A MATLAB package for Linearized Encoding Analy-
926 sis (LEA) on multimodal data is available at Zenodo release: <https://doi.org/10.5281/zenodo.15107756>
927 and GitHub: <https://github.com/seunggookim/lea>. An executable demo of LEA on real data is avail-
928 able at MATLAB Online: <https://s.gwdg.de/7cOQmw>. fMRI analysis results in 3-D NIfTI format can
929 be viewed and downloaded at NeuroVault: <https://identifiers.org/neurovault.collection:19626>.
930

931 Declaration of AI-assisted technologies in the writing process

932 During the preparation of this work, the author, as a non-native English speaker, used ChatGPT-4o
933 (OpenAI) via GitHub Copilot to correct typographical and grammatical errors. The author reviewed
934 and edited the content as needed and takes full responsibility for the final manuscript.
935

937 References

- 938 Aiken SJ, Picton TW. Human cortical responses to the speech envelope. *Ear and hearing*. 2008; 29(2):139–157.
939 <https://doi.org/10.1097/AUD.0b013e31816453dc>.
- 940 Allen EJ, St-Yves G, Wu Y, Breedlove JL, Prince JS, Dowdle LT, Nau M, Caron B, Pestilli F, Charest I, et al. A
941 massive 7T fMRI dataset to bridge cognitive neuroscience and artificial intelligence. *Nature neuroscience*.
942 2022; 25(1):116–126. <https://doi.org/10.1038/s41593-021-00962-x>.
- 943 Baldassano C, Hasson U, Norman KA. Representation of real-world event schemas during narrative perception.
944 *Journal of Neuroscience*. 2018; 38(45):9689–9699. <https://doi.org/10.1523/jneurosci.0251-18.2018>.
- 945 Brunswik E. Organismic achievement and environmental probability. *Psychological review*. 1943; 50(3):255.
946 <https://doi.org/10.1037/h0060889>.
- 947 Caucheteux C, Gramfort A, King JR. Evidence of a predictive coding hierarchy in the human brain listening to
948 speech. *Nature human behaviour*. 2023; 7(3):430–441. <https://doi.org/10.1038/s41562-022-01516-2>.
- 949 Chi T, Ru P, Shamma SA. Multiresolution spectrotemporal analysis of complex sounds. *The Journal of the
950 Acoustical Society of America*. 2005; 118(2):887–906. <https://doi.org/10.1121/1.1945807>.
- 951 Crosse MJ, Zuk NJ, Di Liberto GM, Nidiffer AR, Molholm S, Lalor EC. Linear modeling of neurophysiological re-
952 sponses to speech and other continuous stimuli: methodological considerations for applied research. *Frontiers in neuroscience*.
953 2021; 15:705621.
- 954 Di Liberto GM, O'sullivan JA, Lalor EC. Low-frequency cortical entrainment to speech reflects phoneme-level
955 processing. *Current Biology*. 2015; 25(19):2457–2465. <https://doi.org/10.1016/j.cub.2015.08.030>.
- 956 Di Liberto GM, Pelofi C, Bianco R, Patel P, Mehta AD, Herrero JL, de Cheveigné A, Shamma S, Mesgarani N.
957 Cortical encoding of melodic expectations in human temporal cortex. *eLife*. 2020 mar; 9:e51784. <https://doi.org/10.7554/elife.51784>.
- 958 Diedrichsen J, Kriegeskorte N. Representational models: A common framework for understanding en-
959 coding, pattern-component, and representational-similarity analysis. *PLoS computational biology*. 2017;
960 13(4):e1005508.
- 961 Friston KJ, Holmes AP, Worsley KJ, Poline JP, Frith CD, Frackowiak RS. Statistical parametric maps in functional
962 imaging: a general linear approach. *Human brain mapping*. 1994; 2(4):189–210. <https://doi.org/10.1002/hbm.460020402>.
- 963 Gibson JJ. The Ecological Approach to the Visual Perception of Pictures. *Leonardo*. 1978; 11(3):227–235. <https://doi.org/10.2307/1574154>.
- 964 Giraud AL, Lorenzi C, Ashburner J, Wable J, Johnsrude I, Frackowiak R, Kleinschmidt A. Representation of the
965 temporal envelope of sounds in the human brain. *Journal of neurophysiology*. 2000; 84(3):1588–1598. <https://doi.org/10.1152/jn.2000.84.3.1588>.
- 966 Hamilton LS, Huth AG. The revolution will not be controlled: natural stimuli in speech neuroscience. *Language,
967 Cognition and Neuroscience*. 2020; 35(5):573–582. <https://doi.org/10.1080/23273798.2018.1499946>.
- 968 Han K, Wen H, Shi J, Lu KH, Zhang Y, Fu D, Liu Z. Variational autoencoder: An unsupervised model for encoding
969 and decoding fMRI activity in visual cortex. *NeuroImage*. 2019; 198:125–136.
- 970 Harms MP, Guinan Jr JJ, Sigalovsky IS, Melcher JR. Short-term sound temporal envelope characteristics deter-
971 mine multisecond time patterns of activity in human auditory cortex as shown by fMRI. *Journal of neuro-
972 physiology*. 2005; 93(1):210–222. <https://doi.org/10.1152/jn.00712.2004>.

- 977 **Hasson U**, Nir Y, Levy I, Fuhrmann G, Malach R. Intersubject Synchronization of Cortical Activity During Natural
978 Vision. *Science*. 2004; 303(5664):1634–1640. <https://doi.org/10.1126/science.1089506>.
- 979 **Hastie T**, Tibshirani R, Friedman JH, Friedman JH. The elements of statistical learning: data mining, inference,
980 and prediction, vol. 2. Springer; 2009. <https://doi.org/10.1007/978-0-387-84858-7>.
- 981 **Haxby JV**, Guntupalli JS, Nastase SA, Feilong M. Hyperalignment: Modeling shared information encoded in
982 idiosyncratic cortical topographies. *eLife*. 2020 jun; 9:e56601. <https://doi.org/10.7554/elife.56601>.
- 983 **Henson RNA**, Rugg MD, Shallice T, Josephs O, Dolan RJ. Recollection and Familiarity in Recognition Memory: An
984 Event-Related Functional Magnetic Resonance Imaging Study. *Journal of Neuroscience*. 1999; 19(10):3962–
985 3972.
- 986 **Hoerl AE**, Kennard RW. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*.
987 1970; 12(1):55–67. <https://doi.org/10.1080/00401706.1970.10488634>.
- 988 **Huth AG**, De Heer WA, Griffiths TL, Theunissen FE, Gallant JL. Natural speech reveals the semantic maps that
989 tile human cerebral cortex. *Nature*. 2016; 532(7600):453–458.
- 990 **Jacobs J**, Kahana MJ. Neural representations of individual stimuli in humans revealed by gamma-band electro-
991 corticographic activity. *Journal of neuroscience*. 2009; 29(33):10203–10214.
- 992 **Kaneshiro B**, Nguyen DT, Norcia AM, Dmochowski JP, Berger J. Natural music evokes correlated EEG responses
993 reflecting temporal structure and beat. *NeuroImage*. 2020; 214:116559. <https://doi.org/10.1016/j.neuroimage.2020.116559>.
- 995 **Kaufman S**, Rosset S, Perlich C, Stitelman O. Leakage in data mining: Formulation, detection, and avoidance.
996 *ACM Trans Knowl Discov Data*. 2012 dec; 6(4). <https://doi.org/10.1145/2382577.2382579>.
- 997 **Kay KN**. Principles for models of neural information processing. *NeuroImage*. 2018; 180:101–109.
- 998 **Kay KN**, Naselaris T, Prenger RJ, Gallant JL. Identifying natural images from human brain activity. *Nature*. 2008;
999 452(7185):352–355.
- 1000 **Kim SG**. On the encoding of natural music in computational models and human brains. *Frontiers in Neuro-
1001 science*. 2022; 16. <https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2022.928841>, doi:
1002 [10.3389/fnins.2022.928841](https://doi.org/10.3389/fnins.2022.928841).
- 1003 **Kim SG**, De Martino F, Overath T. Linguistic modulation of the neural encoding of phonemes. *Cerebral Cortex*.
1004 2024; 34(4):bhae155.
- 1005 **Koelsch S**. A coordinate-based meta-analysis of music-evoked emotions. *NeuroImage*. 2020; 223:117350.
- 1006 **Kriegeskorte N**, Goebel R, Bandettini P. Information-based functional brain mapping. *Proceedings of the
1007 National Academy of Sciences*. 2006; 103(10):3863–3868.
- 1008 **Kriegeskorte N**, Mur M, Bandettini PA. Representational similarity analysis-connecting the branches of sys-
1009 tems neuroscience. *Frontiers in systems neuroscience*. 2008; 2:249.
- 1010 **Kriegeskorte N**, Simmons WK, Bellgowan PSF, Baker CI. Circular analysis in systems neuroscience: the dangers
1011 of double dipping. *Nature Neuroscience*. 2009; 12(5):535–540. <https://doi.org/10.1038/nn.2303>.
- 1012 **Kubanek J**, Brunner P, Gunduz A, Poeppel D, Schalk G. The tracking of speech envelope in the human cortex.
1013 *PloS one*. 2013; 8(1):e53398. <https://doi.org/10.1371/journal.pone.0053398>.
- 1014 **La Tour TD**, Eickenberg M, Nunez-Elizalde AO, Gallant JL. Feature-space selection with banded ridge regression.
1015 *NeuroImage*. 2022; 264:119728.

- 1016 **Lage-Castellanos A**, Valente G, Formisano E, De Martino F. Methods for computing the maximum performance
1017 of computational models of fMRI responses. PLOS Computational Biology. 2019 03; 15(3):1–25.
- 1018 **Leahy J**, Kim SG, Wan J, Overath T. An analytical framework of tonal and rhythmic hierarchy in natural music
1019 using the multivariate temporal response function. Frontiers in Neuroscience. 2021; 15:665767.
- 1020 **Moerel M**, De Martino F, Kemper VG, Schmitter S, Vu AT, Uğurbil K, Formisano E, Yacoub E. Sensitivity and
1021 specificity considerations for fMRI encoding, decoding, and mapping of auditory cortex at ultra-high field.
1022 NeuroImage. 2018; 164:18–31.
- 1023 **Müller-Gerking J**, Pfurtscheller G, Flyvbjerg H. Designing optimal spatial filters for single-trial EEG classification
1024 in a movement task. Clinical Neurophysiology. 1999; 110(5):787–798.
- 1025 **Naselaris T**, Kay KN, Nishimoto S, Gallant JL. Encoding and decoding in fMRI. NeuroImage. 2011; 56(2):400–410.
- 1026 **Nastase SA**, Goldstein A, Hasson U. Keep it real: rethinking the primacy of experimental control in cog-
1027 nitive neuroscience. NeuroImage. 2020; 222:117254. <https://www.sciencedirect.com/science/article/pii/S1053811920307400>, doi: <https://doi.org/10.1016/j.neuroimage.2020.117254>.
- 1029 **Nishimoto S**, Vu AT, Naselaris T, Benjamini Y, Yu B, Gallant JL. Reconstructing visual experiences from brain
1030 activity evoked by natural movies. Current biology. 2011; 21(19):1641–1646.
- 1031 **Nuzzo R**. How scientists fool themselves –and how they can stop. Nature. 2015; 526(7572):182–185. <https://doi.org/10.1038/526182a>, doi: 10.1038/526182a.
- 1033 **Overath T**, Zhang Y, Sanes DH, Poeppel D. Sensitivity to temporal modulation rate and spectral bandwidth
1034 in the human auditory system: fMRI evidence. Journal of Neurophysiology. 2012; 107(8):2042–2056. <https://doi.org/10.1152/jn.00308.2011>, doi: 10.1152/jn.00308.2011, pMID: 22298830.
- 1036 **Pistohl T**, Schulze-Bonhage A, Aertsen A, Mehring C, Ball T. Decoding natural grasp types from human ECoG.
1037 NeuroImage. 2012; 59(1):248–260.
- 1038 **Poldrack RA**. Can cognitive processes be inferred from neuroimaging data? Trends in cognitive sciences. 2006;
1039 10(2):59–63.
- 1040 **Prince JS**, Charest I, Kurzawski JW, Pyles JA, Tarr MJ, Kay KN. Improving the accuracy of single-trial fMRI response
1041 estimates using GLMsingle. Elife. 2022; 11:e77599.
- 1042 **Quandt F**, Reichert C, Hinrichs H, Heinze HJ, Knight RT, Rieger JW. Single trial discrimination of individual finger
1043 movements on one hand: A combined MEG and EEG study. NeuroImage. 2012; 59(4):3316–3324.
- 1044 **Ray S**, Crone NE, Niebur E, Franaszczuk PJ, Hsiao SS. Neural correlates of high-gamma oscillations (60–200
1045 Hz) in macaque local field potentials and their potential implications in electrocorticography. Journal of
1046 Neuroscience. 2008; 28(45):11526–11536.
- 1047 **Ryali S**, Supekar K, Abrams DA, Menon V. Sparse logistic regression for whole-brain classification of fMRI data.
1048 NeuroImage. 2010; 51(2):752–764.
- 1049 **Sachs ME**, Habibi A, Damasio A, Kaplan JT. Dynamic intersubject neural synchronization reflects affective re-
1050 sponses to sad music. NeuroImage. 2020; 218:116512. <https://doi.org/10.1016/j.neuroimage.2019.116512>.
- 1051 **Santoro R**, Moerel M, Martino FD, Valente G, Ugurbil K, Yacoub E, Formisano E. Reconstructing the spectrotem-
1052 poral modulations of real-life sounds from fMRI response patterns. Proceedings of the National Academy
1053 of Sciences. 2017; 114(18):4799–4804.

- 1054 Sonkusare S, Breakspear M, Guo C. Naturalistic Stimuli in Neuroscience: Critically Acclaimed. Trends in Cognitive Sciences. 2019; 23(8):699–714. <https://www.sciencedirect.com/science/article/pii/S1364661319301275>, doi: <https://doi.org/10.1016/j.tics.2019.05.004>.
- 1057 Theiler J, Eubank S, Longtin A, Galdrikian B, Doyne Farmer J. Testing for nonlinearity in time series: the method
1058 of surrogate data. Physica D: Nonlinear Phenomena. 1992; 58(1):77–94. <https://www.sciencedirect.com/science/article/pii/016727899290102S>, doi: [https://doi.org/10.1016/0167-2789\(92\)90102-S](https://doi.org/10.1016/0167-2789(92)90102-S).
- 1060 Varoquaux G. Cross-validation failure: Small sample sizes lead to large error bars. NeuroImage. 2018; 180:68–
1061 77. <https://doi.org/10.1016/j.neuroimage.2017.06.061>.
- 1062 Wu MCK, David SV, Gallant JL. Complete functional characterization of sensory neurons by system identification.
1063 Annu Rev Neurosci. 2006; 29(1):477–505. <https://doi.org/10.1146/annurev.neuro.29.051605.113024>.
- 1064 Yekutieli D, Benjamini Y. Resampling-based false discovery rate controlling multiple test procedures for cor-
1065 related test statistics. Journal of Statistical Planning and Inference. 1999; 82(1):171–196. [https://doi.org/10.1016/S0378-3758\(99\)00041-5](https://doi.org/10.1016/S0378-3758(99)00041-5).

Table 1. Definition of variables, parameters, and notations.

Variables	Description
$\mathbf{X} \in \mathbb{R}^{T \times FD}$	True finite impulse response design matrix
$\mathbf{U} \in \mathbb{R}^{T \times FD}$	Null finite impulse response design matrix
$\mathbf{y} \in \mathbb{R}^{T \times 1}$	Response vector
$\mathbf{b} \in \mathbb{R}^{FD \times 1}$	Regression coefficient vector
$\mathbf{s} \in \mathbb{R}^{T \times 1}$	Signal time series
$\mathbf{e} \in \mathbb{R}^{T \times 1}$	Noise time strategies
$\mathbf{L} \in \mathbb{R}^{FD \times FD}$	Tikhonov regularization matrix
Parameters	Description
T	Number of time points
F	Number of features
D	Number of delays
P	Number of predictors ($= FD$)
ϕ	AR(1) temporal correlation parameter
θ	AR(1) spatial correlation parameter
ρ	Multicollinearity parameter
λ	regularization parameter
Notations	Description
$\ \cdot\ $	l_2 -norm of a vector
$\ \cdot\ _F$	Frobenius norm of a matrix
$(\cdot)^T$	Transposition
$(\cdot)_{(i)}$	i -th column vector of a matrix
$\text{diag}(\cdot)$	Diagonal matrix with the diagonal elements from the given vector
$\mathcal{L}(\cdot; \cdot)$	Optimization function
\cdot^*	Optimal solution
$\hat{\cdot}$	Estimate of a variable
$(\cdot)_i$	i -th cross-validation set
$\bar{\cdot}$	Mean
$\mathbb{E}[\cdot]$	Expectation of a random variable
\propto	Proportional to
\approx	Approximately equal to
\equiv	Equivalent to

Table 2. Parameters of the simulations

Category	Parameter	Notation
Complexity of model	Number of delays	D
	Number of features	F
Dimensionality of data	Number of variates	V
Strength of signal	Signal-to-noise ratio (SNR)	S
AR(1) temporal autocorrelation (if $T \geq 2$)	of true predictors	ϕ_X
	of null predictors	ϕ_U
	of true weight	ϕ_B
	of noise	ϕ_E
AR(1) spatial autocorrelation (if $V \geq 2$)	of true weight	θ_B
	of noise	θ_E
Multicollinearity (if $F \geq 2$)	of true predictors	ρ_X
	of null predictors	ρ_U
Presence of information leakage	Binary flag for the stimulus repetition	<code>IsRep</code>

Number of time points $T = 100$, λ -grid = $10^{[-10, -9, \dots, 15]}$, Number of samplings $K = 1000$

Table 3. Strong effects ($\eta_p^2 \geq 0.160$) on the prediction accuracies in univariate models with univariate responses using true predictors \mathbf{X}

Contrast	SS	F-stat	η_p^2
S	8.377×10^1	1.867×10^5	0.985
$\cdot \phi_X : \phi_E$	2.645×10^{-1}	5.895×10^2	0.175
$S : \phi_E$	2.037×10^{-1}	4.539×10^2	0.140
$S : \phi_X : \phi_E$	3.187×10^{-1}	7.101×10^2	0.203

$df_1 = 1$, $df_2 = 2788$, $P_{\text{Bonferroni}} \leq 0.0001$, $\eta_p^2 \geq 0.14$, SS: sum of squares

Table 3—source data 1. Full anova table: https://zenodo.org/records/15101528/files/uni-uni_anova-x.xlsx

Table 4. Strong effects ($\eta_p^2 \geq 0.160$) on the prediction accuracies in univariate models with univariate responses using null predictors \mathbf{U}

Contrast	SS	F-stat	η_p^2
IsRep	1.569×10^1	3.864×10^4	0.933
S	2.199	5.413×10^3	0.660
D	1.504	3.704×10^3	0.571
ϕ_X	4.312×10^{-1}	1.062×10^3	0.276
ϕ_U	3.607×10^{-1}	8.881×10^2	0.242
$S:\text{IsRep}$	2.209	5.438×10^3	0.661
$\text{IsRep}:D$	1.458	3.589×10^3	0.563
$\phi_X:\phi_U$	6.987×10^{-1}	1.720×10^3	0.382
$\phi_X:\text{IsRep}$	4.435×10^{-1}	1.092×10^3	0.281
$\phi_U:\text{IsRep}$	3.357×10^{-1}	8.266×10^2	0.229
$S:D$	2.164×10^{-1}	5.328×10^2	0.160
$\phi_X:\phi_U:\text{IsRep}$	7.277×10^{-1}	1.792×10^3	0.391
$S:\text{IsRep}:D$	2.408×10^{-1}	5.929×10^2	0.175

$df_1 = 1$, $df_2 = 2788$, $P_{\text{Bonferroni}} \leq 0.0001$, $\eta_p^2 \geq 0.14$, SS: sum of squares

Table 4—source data 1. Full anova table: https://zenodo.org/records/15101528/files/uni-uni_anova-u.xlsx

Table 5. Strong effects ($\eta_p^2 \geq 0.160$) on the prediction accuracies of multivariate models with univariate responses using true predictors \mathbf{X}

Contrast	SS	F-stat	η_p^2
S	5.729×10^2	5.984×10^5	0.972
ρ_X	2.374×10^1	2.479×10^4	0.593
D	1.258×10^1	1.314×10^4	0.436
IsRep	6.991	7.302×10^3	0.301
F	6.196	6.471×10^3	0.276
ϕ_X	4.010	4.188×10^3	0.198
ϕ_E	3.557	3.715×10^3	0.179
$\phi_X:\phi_E$	5.048	5.273×10^3	0.237
$\rho_X:D$	3.213	3.356×10^3	0.165
$\rho_X:\text{IsRep}$	2.841	2.967×10^3	0.149

$df_1 = 1$, $df_2 = 16984$, $P_{\text{Bonferroni}} \leq 0.0001$, $\eta_p^2 \geq 0.14$, SS: sum of squares

Table 5—source data 1. Full anova table: https://zenodo.org/records/15101528/files/mult-uni_anova-x.xlsx

Table 6. Strong effects ($\eta_p^2 \geq 0.160$) on the prediction accuracies of multivariate models with univariate responses using null predictors \mathbf{U}

Contrast	SS	F-stat	η_p^2
IsRep	4.161×10^2	1.057×10^5	0.862
S	6.946×10^1	1.765×10^4	0.510
ρ_U	5.282×10^1	1.342×10^4	0.441
D	3.429×10^1	8.715×10^3	0.339
$S:\text{IsRep}$	6.933×10^1	1.762×10^4	0.509
$\rho_U:\text{IsRep}$	5.296×10^1	1.346×10^4	0.442
IsRep: D	3.469×10^1	8.816×10^3	0.342

$df_1 = 1, df_2 = 16\,984, P_{\text{Bonferroni}} \leq 0.0001, \eta_p^2 \geq 0.14$, SS: sum of squares

Table 6—source data 1. Full anova table: https://zenodo.org/records/15101528/files/mult-uni_anova-u.xlsx

Table 7. Strong effects ($\eta_p^2 \geq 0.160$) on the prediction accuracies of multivariate models with multivariate responses using true predictors \mathbf{X}

Contrast	SS	F-stat	η_p^2
S	1.161×10^4	1.417×10^7	0.976
ρ_X	3.606×10^2	4.404×10^5	0.557
D	2.177×10^2	2.659×10^5	0.432
IsRep	8.921×10^1	1.090×10^5	0.237
F	8.851×10^1	1.081×10^5	0.236
ϕ_X	7.909×10^1	9.660×10^4	0.216
ϕ_E	7.119×10^1	8.696×10^4	0.199
$\phi_X:\phi_E$	1.057×10^2	1.292×10^5	0.269
$\rho_X:D$	5.227×10^1	6.384×10^4	0.154
$S:\phi_X:\phi_E$	4.999×10^1	6.106×10^4	0.148

$df_1 = 1, df_2 = 350\,191, P_{\text{Bonferroni}} \leq 0.0001, \eta_p^2 \geq 0.14$, SS: sum of squares

Table 7—source data 1. Full anova table: https://zenodo.org/records/15101528/files/mult-mult_anova-x.xlsx

Table 8. Strong effects ($\eta_p^2 \geq 0.160$) on the prediction accuracies of multivariate models on multivariate responses with null predictors \mathbf{U}

Contrast	SS	F-stat	η_p^2
IsRep	7.636×10^3	2.446×10^6	0.875
<i>S</i>	1.269×10^3	4.066×10^5	0.537
ρ_U	8.979×10^2	2.876×10^5	0.451
<i>D</i>	7.222×10^2	2.313×10^5	0.398
<i>S: IsRep</i>	1.270×10^3	4.068×10^5	0.538
$\rho_U : \text{IsRep}$	8.967×10^2	2.872×10^5	0.451
IsRep:D	7.226×10^2	2.315×10^5	0.398

$df_1 = 1$, $df_2 = 349\,964$, $P_{\text{Bonferroni}} \leq 0.0001$, $\eta_p^2 \geq 0.14$, SS: sum of squares

Table 8—source data 1. Full anova table: https://zenodo.org/records/15101528/files/mult-mult_anova-u.xlsx

1068
Optimization processes with identical signals

1069 In a scenario where the strong underlying signals are identical in the training and optimization sets ($\mathbf{s}_1 = \mathbf{s}_3$ and $\|\mathbf{s}_1\| \gg 0$), the optimization prediction accuracy for a given $\lambda \geq 0$ is:

$$1070 \begin{aligned} 1071 \mathbb{E} [\text{corr}(\hat{\mathbf{y}}_{3,\mathbf{X}}[\lambda], \mathbf{y}_3)] &= \mathbb{E} \left[\frac{(\hat{\mathbf{y}}_{3,\mathbf{X}}[\lambda])^\top \mathbf{y}_3}{\|\hat{\mathbf{y}}_{3,\mathbf{X}}\| \|\mathbf{y}_3\|} \right] \propto \mathbb{E} [(\hat{\mathbf{y}}_{3,\mathbf{X}}[\lambda])^\top \mathbf{y}_3] = \mathbb{E} [\mathbf{y}_1^\top \mathbf{P}_{\mathbf{X}}^\top[\lambda] \mathbf{y}_3] \\ 1072 &= \mathbb{E} [(\mathbf{s}_1 + \mathbf{e}_1)^\top \mathbf{P}_{\mathbf{X}}^\top[\lambda] (\mathbf{s}_1 + \mathbf{e}_3)] \\ 1073 &= \mathbf{s}_1^\top \left\{ \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1 + \lambda \mathbf{I})^{-1} \mathbf{X}_1^\top \right\}^\top \mathbf{s}_1 \end{aligned} \quad (\text{A1.1})$$

1074 If $\mathbf{X} \in \mathbb{R}^{T \times P}$ where $P = FD$ is orthonormal and full-rank, **Equation A1.1** can be simplified 1075 as a quadratic form of a scaled Ordinary Least Squares projection matrix and the signal as 1076 (Hastie et al., 2009, pp. 64):

$$1077 \mathbf{s}_1^\top \left\{ \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1 + \lambda \mathbf{I})^{-1} \mathbf{X}_1^\top \right\}^\top \mathbf{s}_1 = \frac{1}{1 + \lambda} \mathbf{s}_1^\top \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{s}_1 = \frac{\|\mathbf{s}_1\|}{1 + \lambda}, \quad (\text{A1.2})$$

1078 for which the optimal λ to maximize the prediction accuracy is zero ($\|\mathbf{s}_1\| > 0$):

$$1079 \lambda^* = \arg \max_{\lambda} \frac{\|\mathbf{s}_1\|}{1 + \lambda} = 0. \quad (\text{A1.3})$$

1080 More generally, the singular value decomposition of \mathbf{X} can be used; $\mathbf{X} = \mathbf{UDV}^\top$ where 1081 $\mathbf{U} \in \mathbb{R}^{T \times P}$ and $\mathbf{V} \in \mathbb{R}^{P \times P}$ are orthonormal, and the diagonal matrix $\mathbf{D} \in \mathbb{R}^{P \times P}$ contains the 1082 singular values: $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$. Then, **Equation A1.1** can be written as (Hastie et al., 1083 2009, Eq. 3.47):

$$1084 \mathbf{s}_1^\top \left\{ \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1 + \lambda \mathbf{I})^{-1} \mathbf{X}_1^\top \right\}^\top \mathbf{s}_1 = \mathbf{s}_1^\top \left(\sum_{j=1}^p \mathbf{u}_{(j)} \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_{(j)}^\top \right) \mathbf{s}_1. \quad (\text{A1.4})$$

1085 For any $d_j = 0$ (i.e., \mathbf{X} is rank-deficit), λ needs to be positive for the prediction accuracy 1086 value to be defined. Nonetheless, the prediction accuracy is maximized when $\lambda^* = \epsilon \approx 0$ for 1087 ϵ is the smallest positive value.

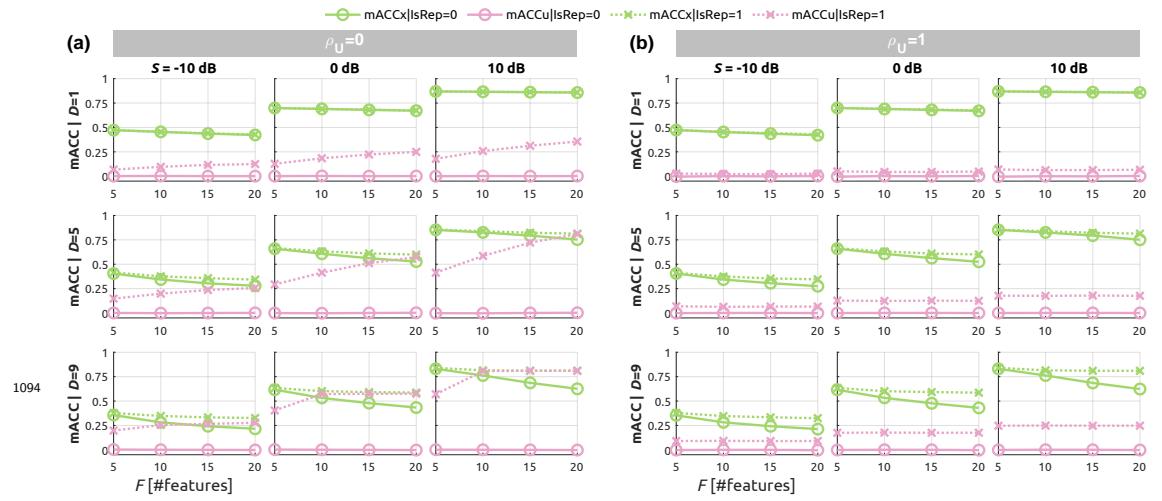


Figure 5—figure supplement 1. [TODO: FIXME] Simulations of multivariate-feature, univariate-response models when the multicollinearity of true feature is zero. Mean prediction accuracies are plotted over the number of features F when multicollinearity **(a)** $\rho_U = 0$ and **(b)** $\rho_U = 1$. Marker styles and colors are identical to those in **Figure 4**. Each column represents a specified SNR level. The number of delays is $D = 1$ in the top rows and $D = 9$ in the bottom rows. Other parameters were as follows: $K = 1000$, $\phi_U = 0$, $\phi_B = 0$, $\phi_E = 0$, $\rho_X = 0$.

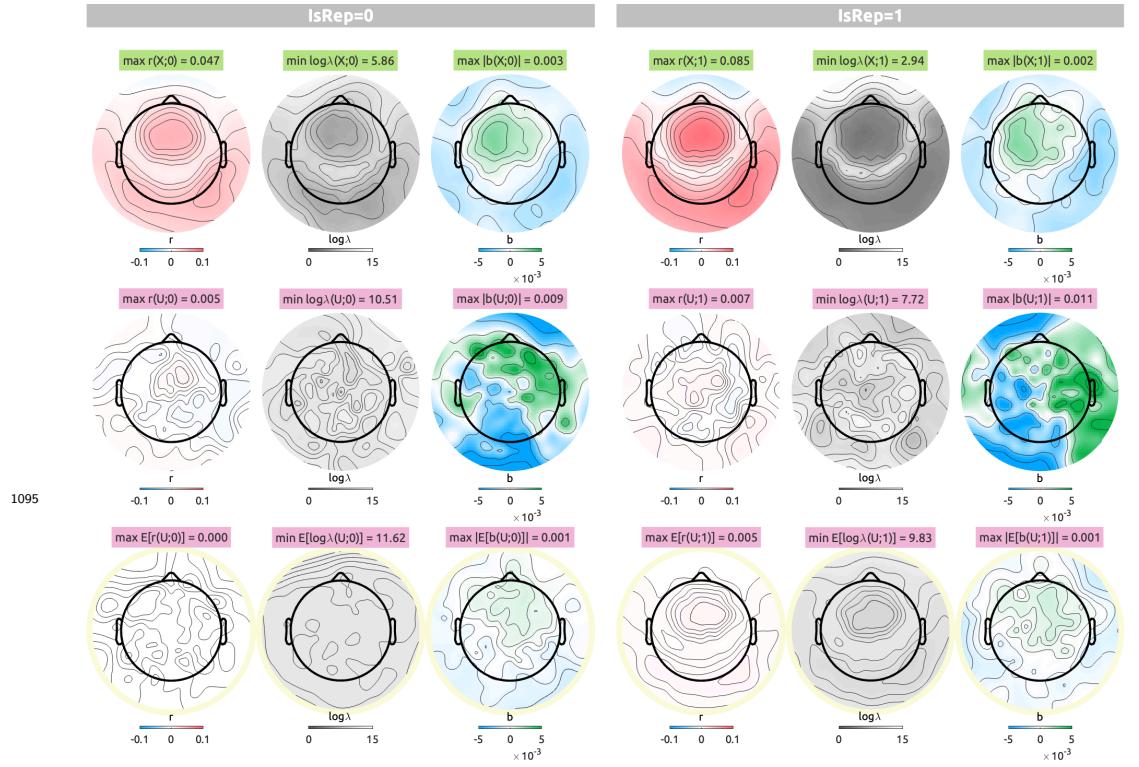


Figure 7—figure supplement 1. EEG linearized encoding analysis results with delays from 0 to 0.5 sec with an audio envelope (top row), a single case of the normal noise (middle row), and an average of 100 normal noises (bottom row; circled in pale yellow). For each CV scheme (IsRep = 0, left panels; IsRep = 1, right panels), prediction accuracy (r ; blue to red), logarithmic ridge hyperparameter ($\log_{10} \lambda$; gray to white), summed weights (b ; blue to green) are shown along the columns.

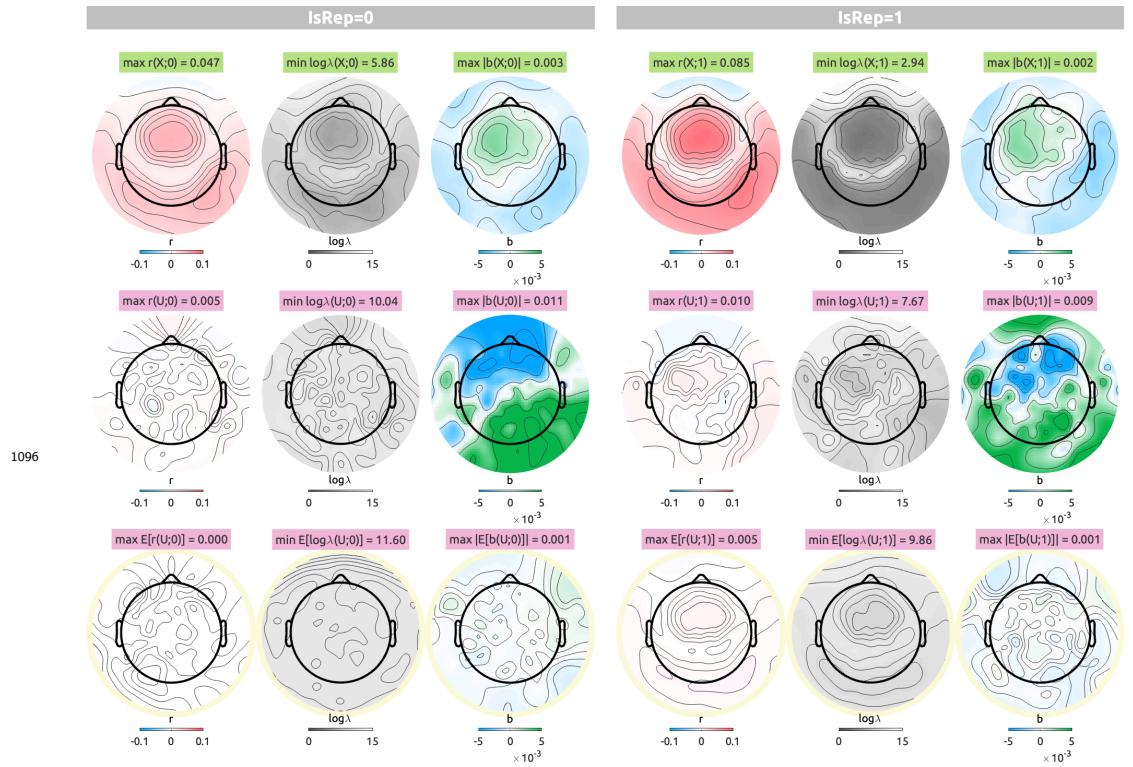


Figure 7—figure supplement 2. EEG topographies of linearized encoding analysis results with delays from 0 to 0.5 sec with the uniform noise as the null feature. The visualization scheme is identical to *Figure 7—figure Supplement 1*.

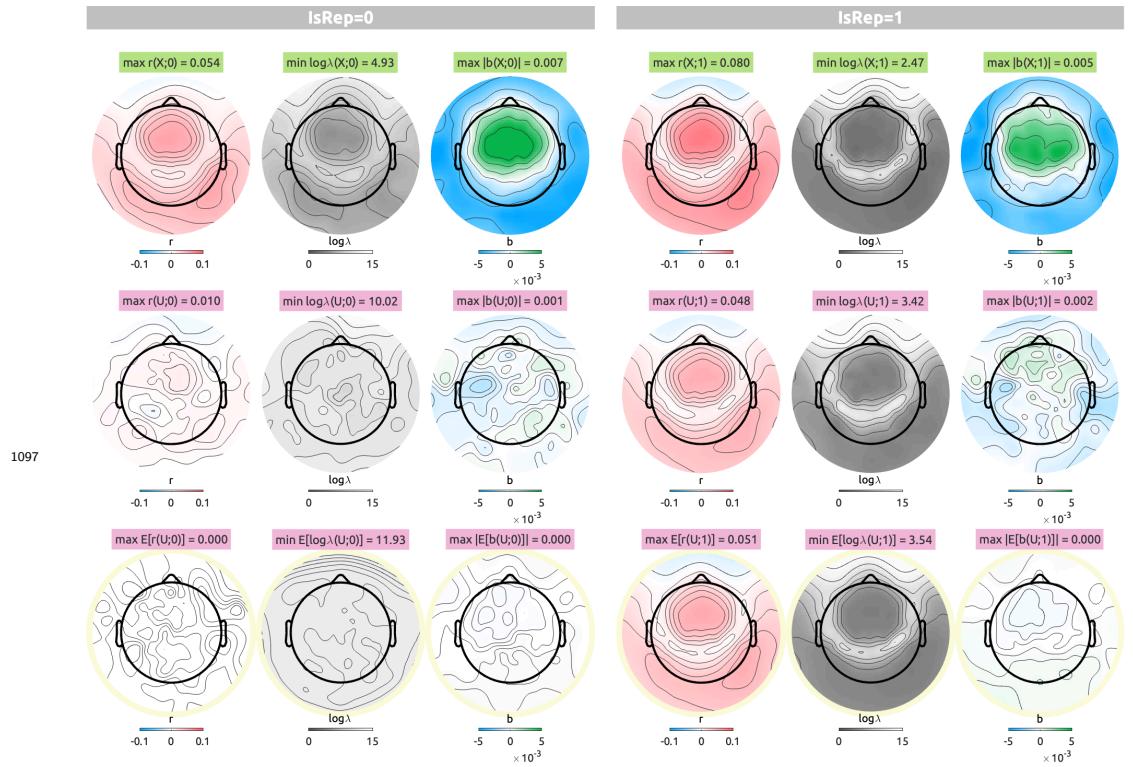


Figure 7—figure supplement 3. EEG topographies of linearized encoding analysis results with a delay from 0 to 0.3 sec with the phase-randomized envelope as the null feature. The visualization scheme is identical to *Figure 7—figure Supplement 1*.

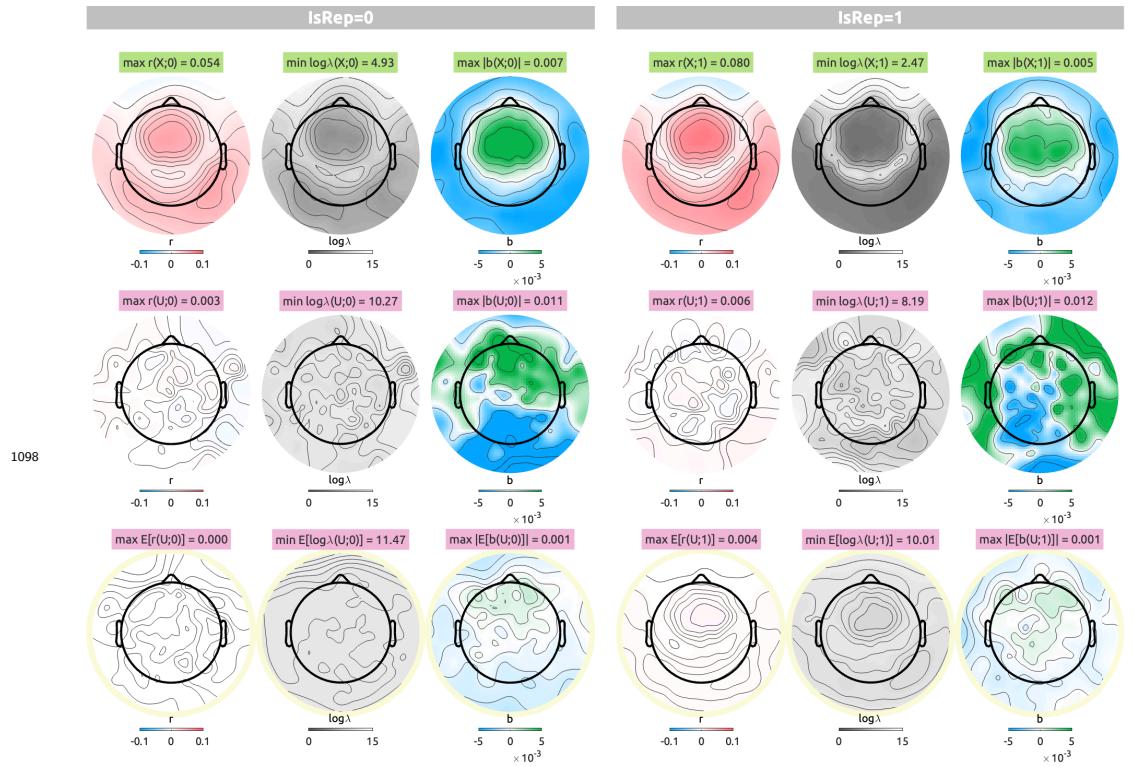


Figure 7—figure supplement 4. EEG topographies of linearized encoding analysis results with a delay from 0 sec to 0.3 sec with the normal noise as the null feature. The visualization scheme is identical to *Figure 7—figure Supplement 1*.

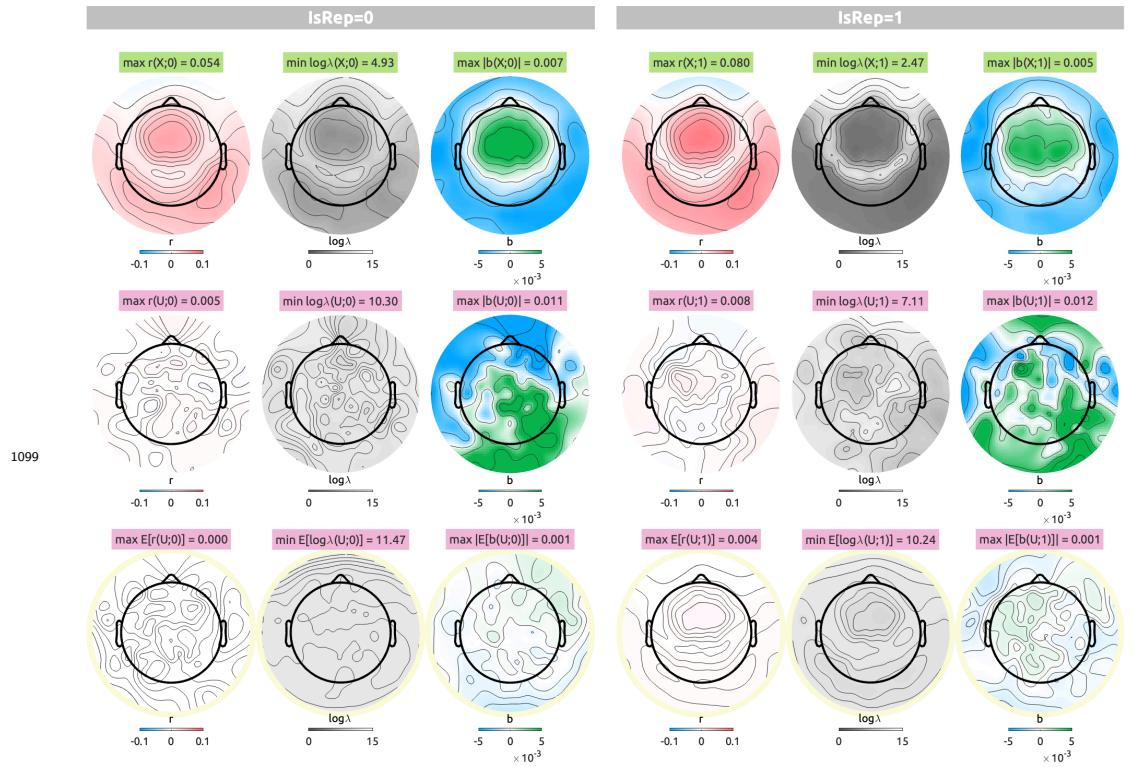


Figure 7—figure supplement 5. EEG topographies of linearized encoding analysis results with a delay from 0 to 0.3 sec with the uniform envelope as the null feature. The visualization scheme is identical to *Figure 7—figure Supplement 1*.

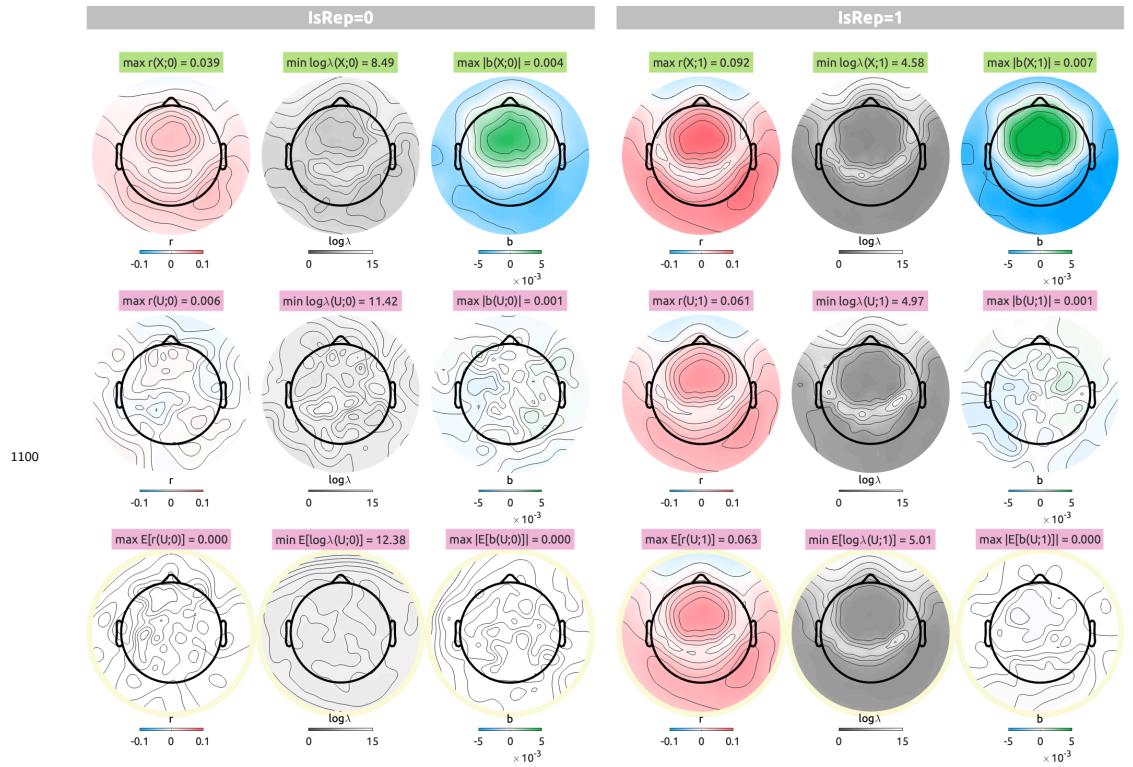


Figure 7—figure supplement 6. EEG topographies of linearized encoding analysis results with delays from 0 to 1 sec with the phase-randomized envelope as the null feature. The visualization scheme is identical to *Figure 7—figure Supplement 1*.

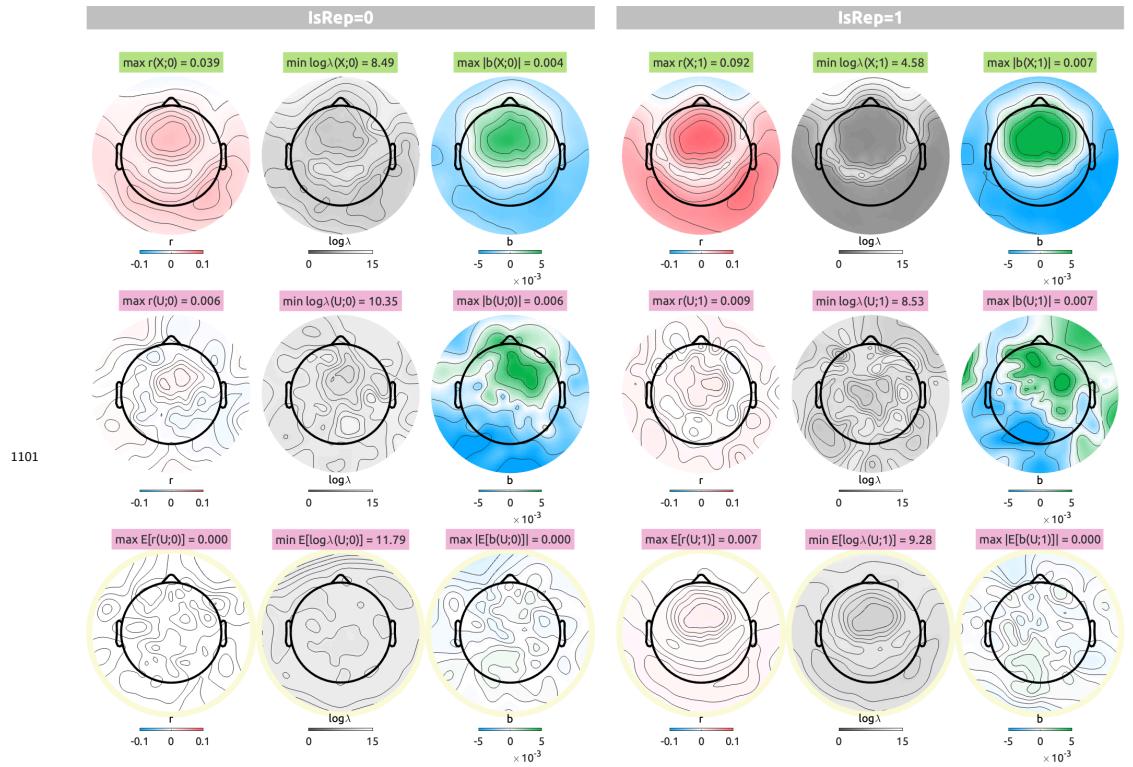


Figure 7—figure supplement 7. EEG topographies of linearized encoding analysis results with delays from 0 to 1 sec with the normal noise as the null feature. The visualization scheme is identical to *Figure 7—figure Supplement 1*.

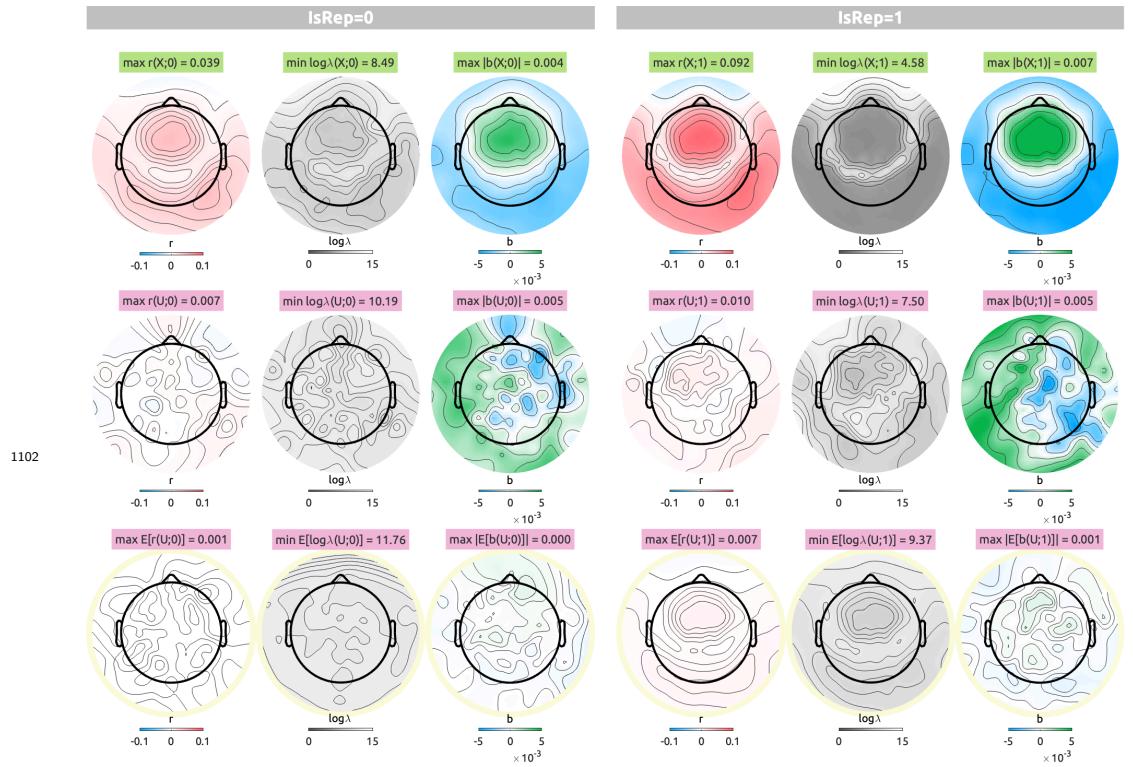


Figure 7—figure supplement 8. EEG topographies of linearized encoding analysis results with delays from 0 to 1 sec with the uniform noise as the null feature. The visualization scheme is identical to *Figure 7—figure Supplement 1*.

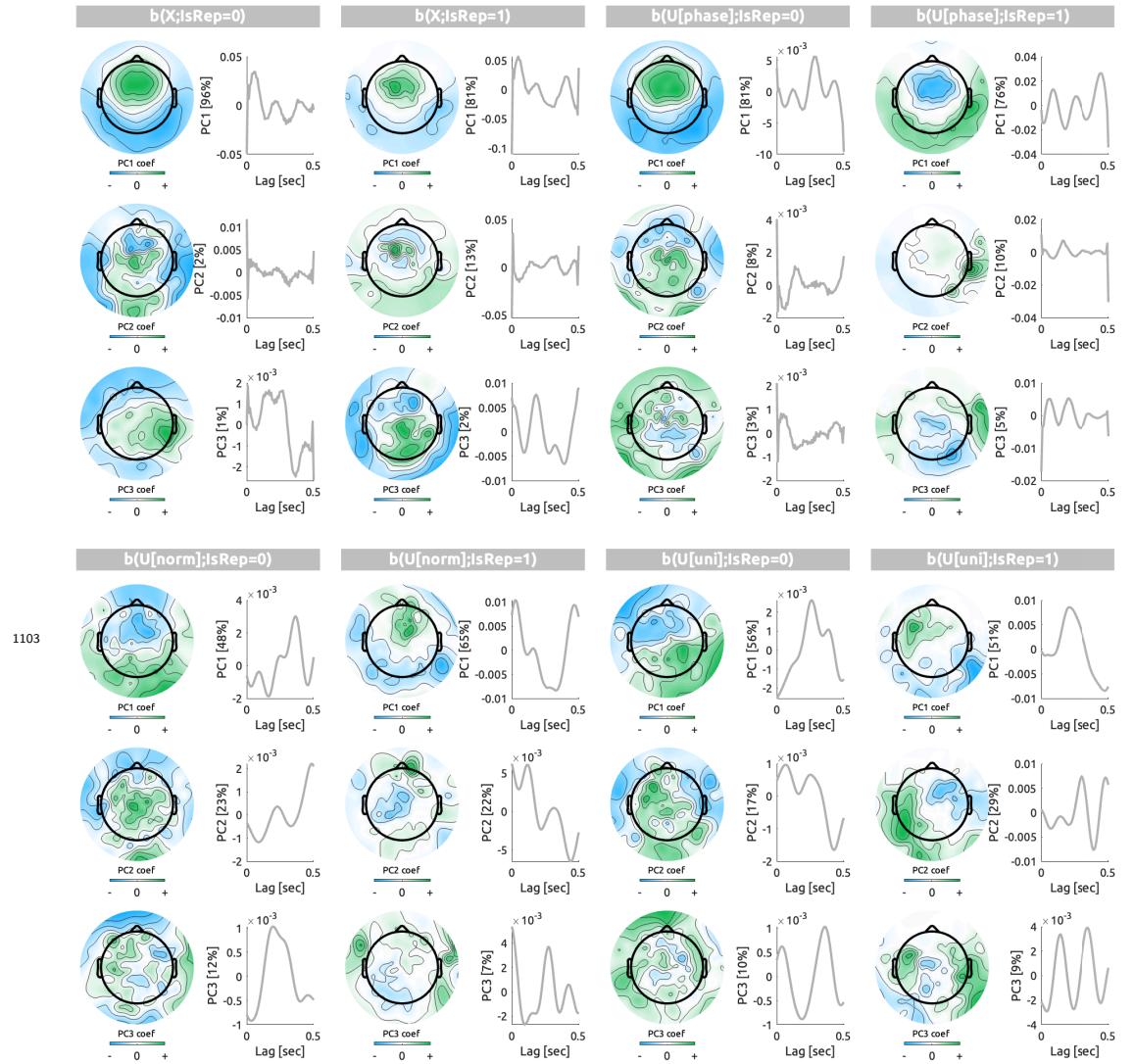


Figure 7—figure supplement 9. First three principal components (PCs) of the transfer function weights (b) of the models with delays from 0 to 0.5 sec are displayed by the topography of eigenvectors and the time series of eigenvariates with the explained variance noted. The weights are grouped by features (X , true audio envelope; $U[\text{phase}]$, phase-randomized envelope; $U[\text{norm}]$, normal noise; $U[\text{uni}]$, uniform noise) and the CV schemes ($\text{IsRep} = 0$, $\text{IsRep} = 1$).

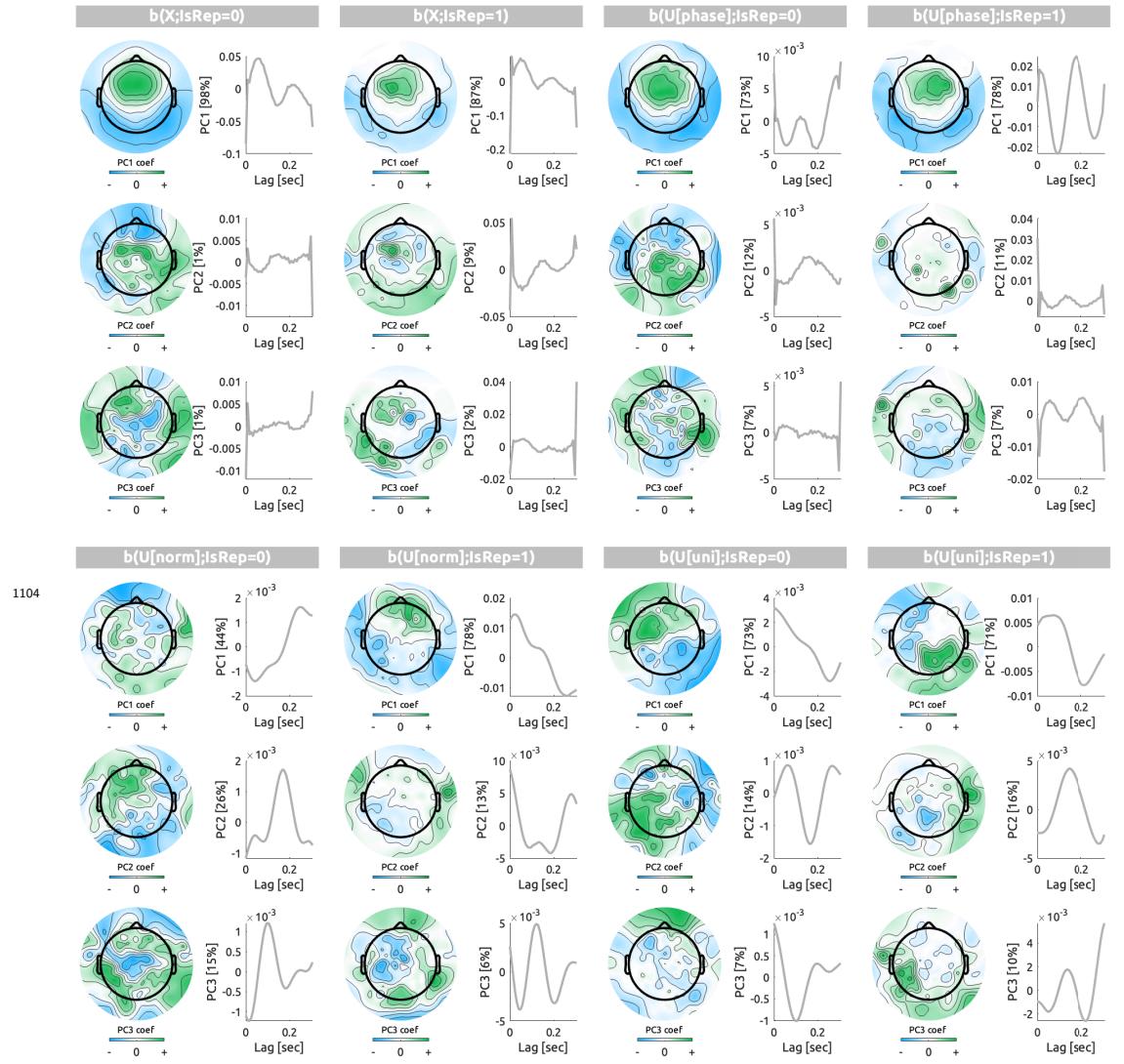


Figure 7—figure supplement 10. First three principal components (PCs) of the transfer function weights (b) of the models with delays from 0 to 0.3 sec are displayed. The visualization scheme is identical to *Figure 7—figure Supplement 9*.

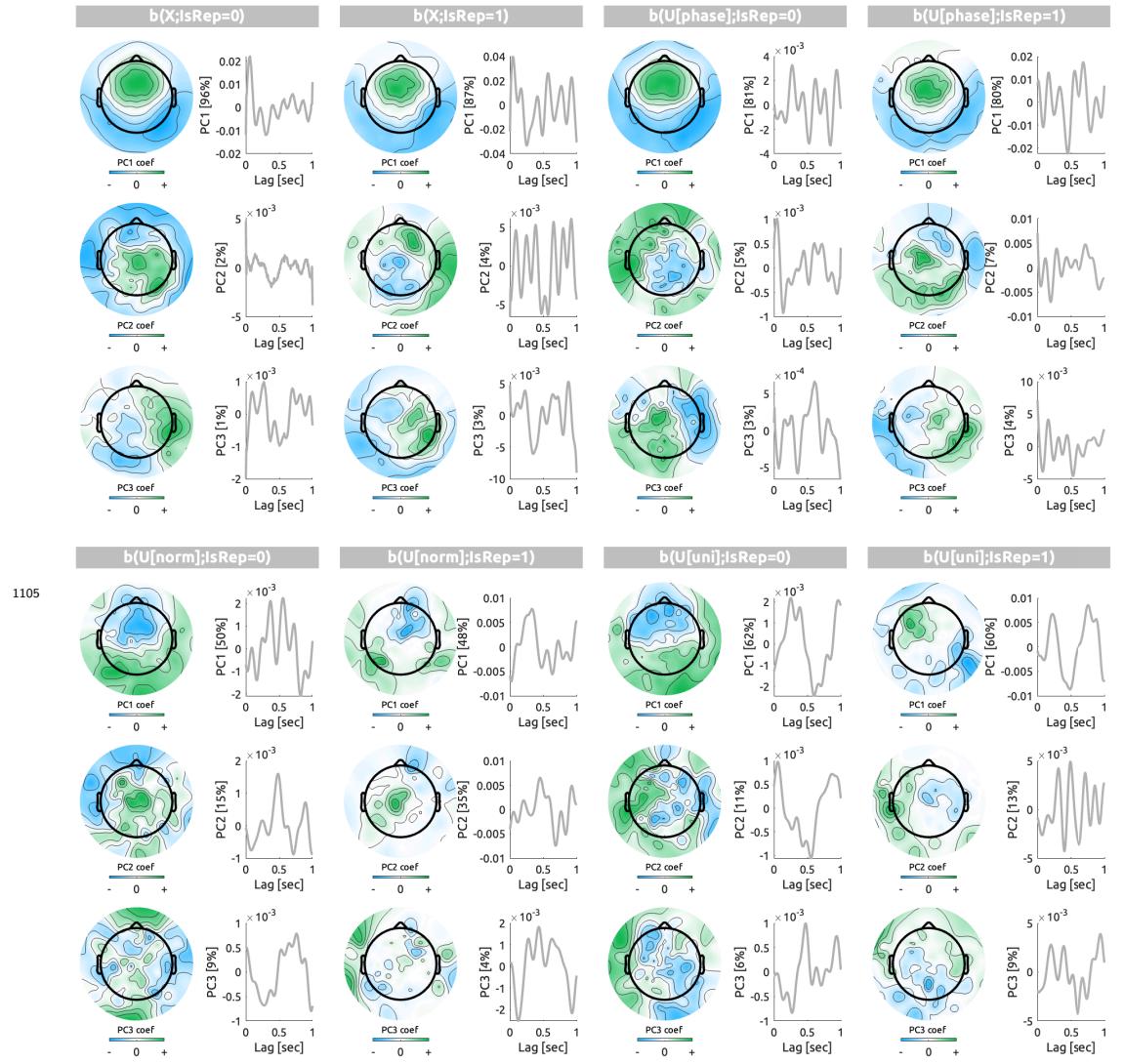


Figure 7—figure supplement 11. First three principal components (PCs) of the transfer function weights (b) of the models with delays from 0 to 1 sec are displayed. The visualization scheme is identical to *Figure 7—figure Supplement 9*.

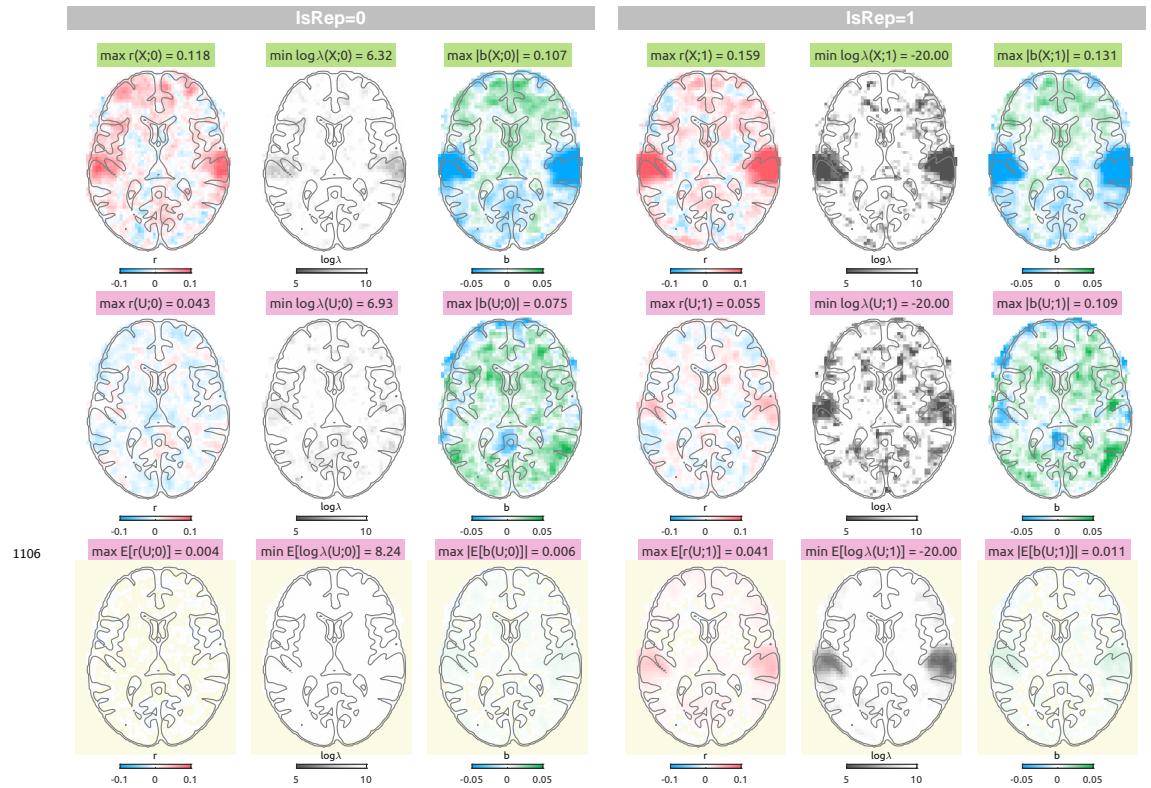


Figure 9—figure supplement 1. fMRI linearized encoding analysis results with delays from 3 to 9 sec with an audio envelope (top row), a single case of the normal noise (middle row), and an average of 100 phase-randomized envelopes (bottom row, pale yellow background). For each CV scheme ($\text{IsRep} = 0$, left panels; $\text{IsRep} = 1$, right panels), prediction accuracy (r , blue to red), logarithmic ridge hyperparameter ($\log_{10} \lambda$, gray to white), transfer function weights that are summed over delays (b , blue to green) are shown along the columns. The 3-D volumes can be viewed with the NeuroVault web viewer (<https://identifiers.org/neurovault.collection:19626>).

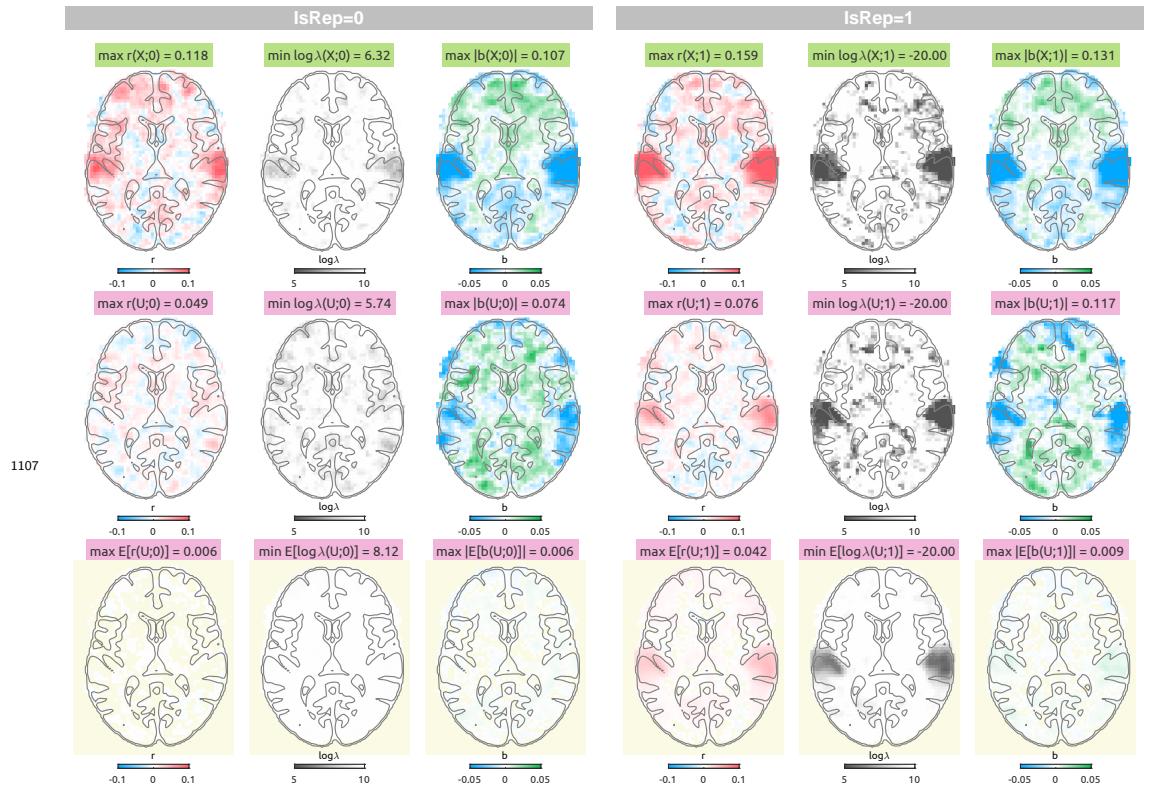


Figure 9—figure supplement 2. fMRI linearized encoding analysis results with delays from 3 to 9 sec with the uniform noise as the null feature. The visualization scheme is identical to **Figure 9—figure Supplement 1**.

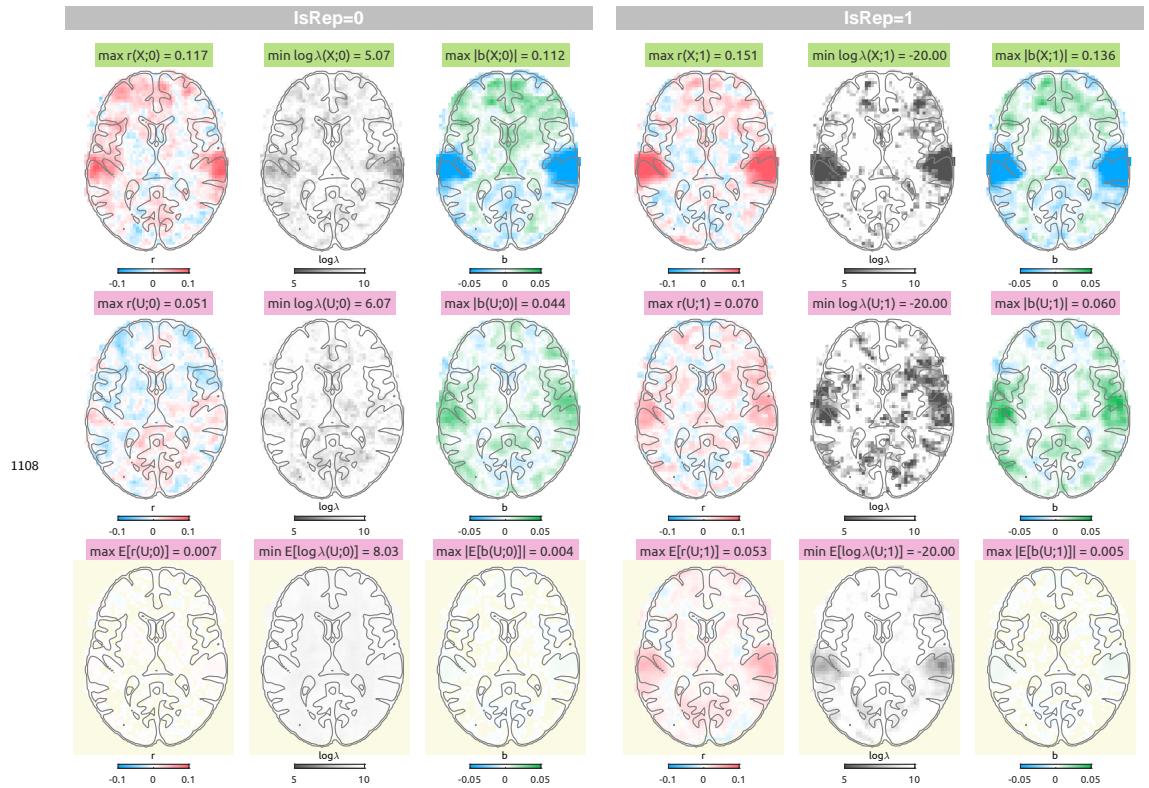


Figure 9—figure supplement 3. fMRI linearized encoding analysis results with delays from 4 to 6 sec with the phase-randomized envelope as the null feature. The visualization scheme is identical to *Figure 9—figure Supplement 1*.

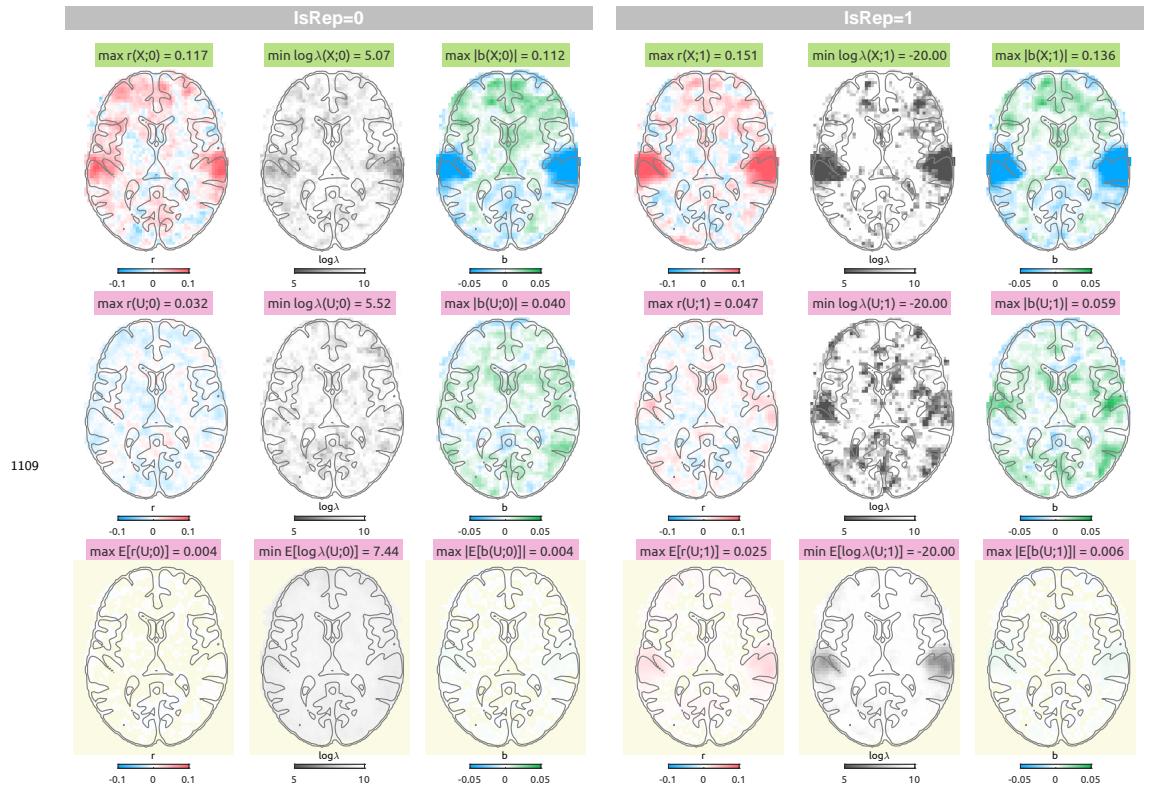


Figure 9—figure supplement 4. fMRI linearized encoding analysis results with delays from 4 to 6 sec with the normal noise as the null feature. The visualization scheme is identical to [Figure 9—figure Supplement 1](#).

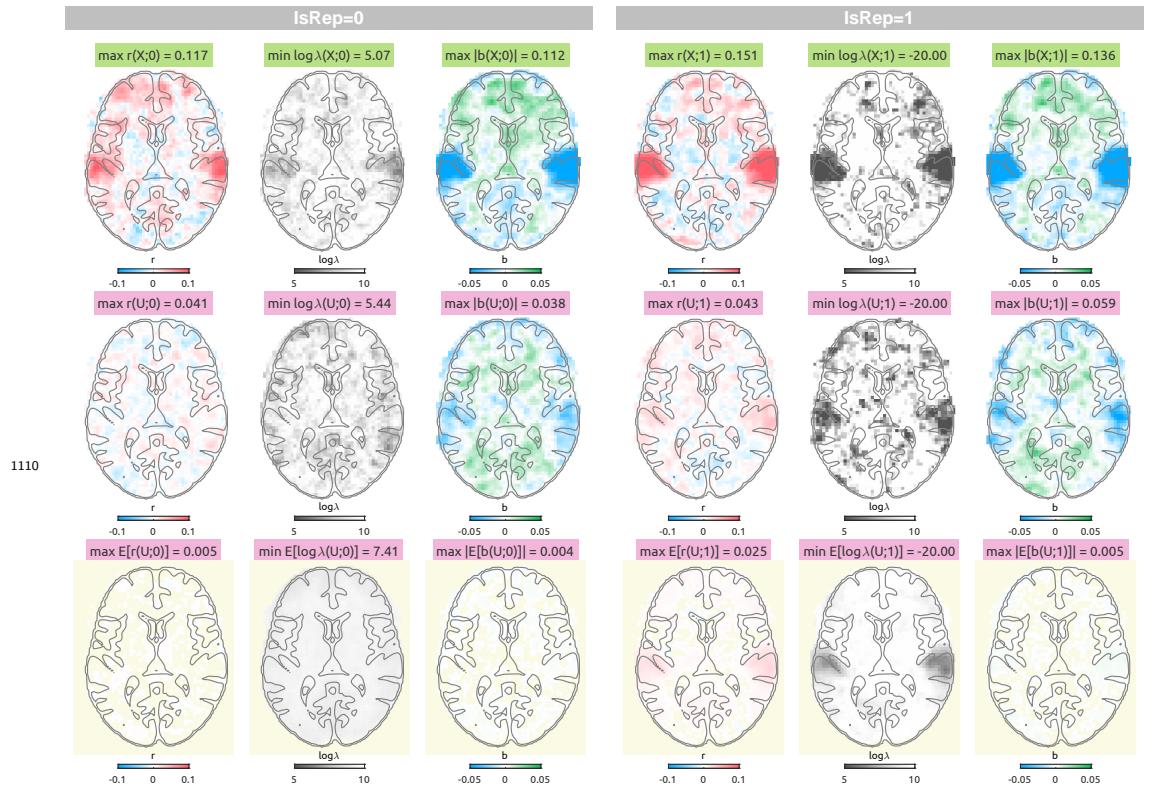


Figure 9—figure supplement 5. fMRI linearized encoding analysis results with delays from 4 to 6 sec with the uniform noise as the null feature. The visualization scheme is identical to **Figure 9—figure Supplement 1**.

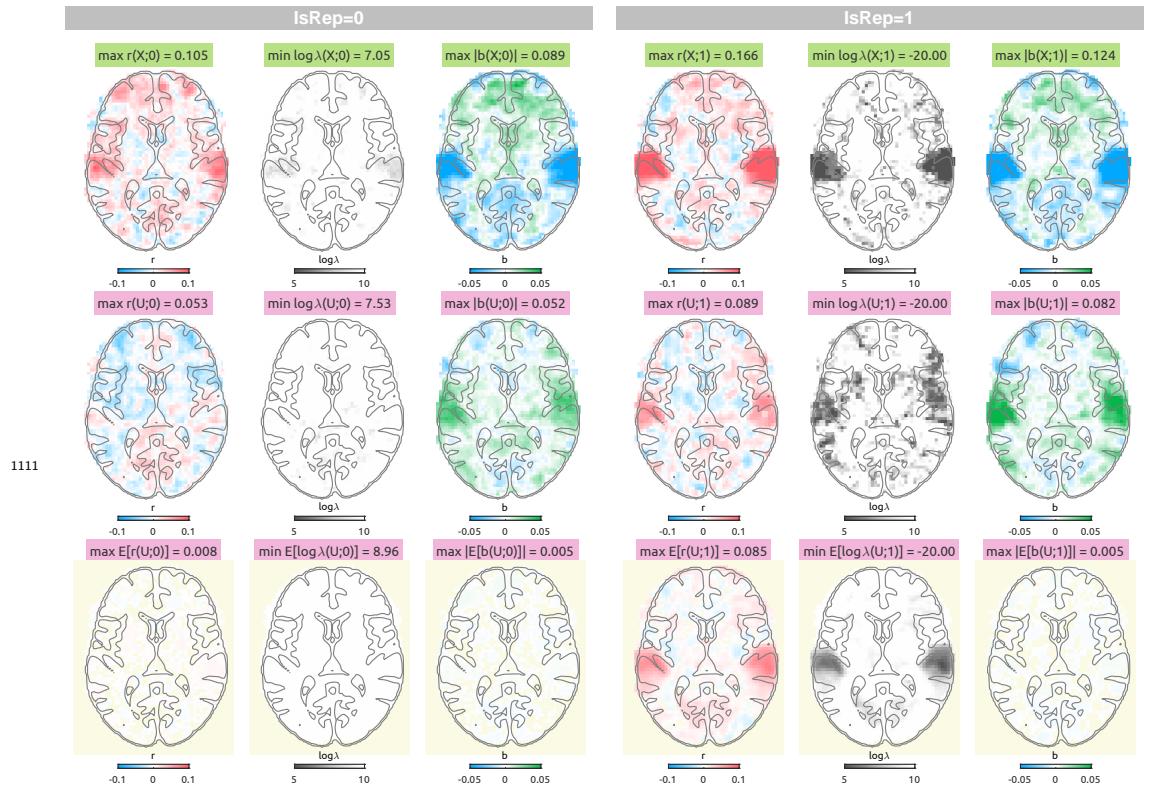


Figure 9—figure supplement 6. fMRI linearized encoding analysis results with delays from 0 to 12 sec with the phase-randomized envelope as the null feature. The visualization scheme is identical to *Figure 9—figure Supplement 1*.

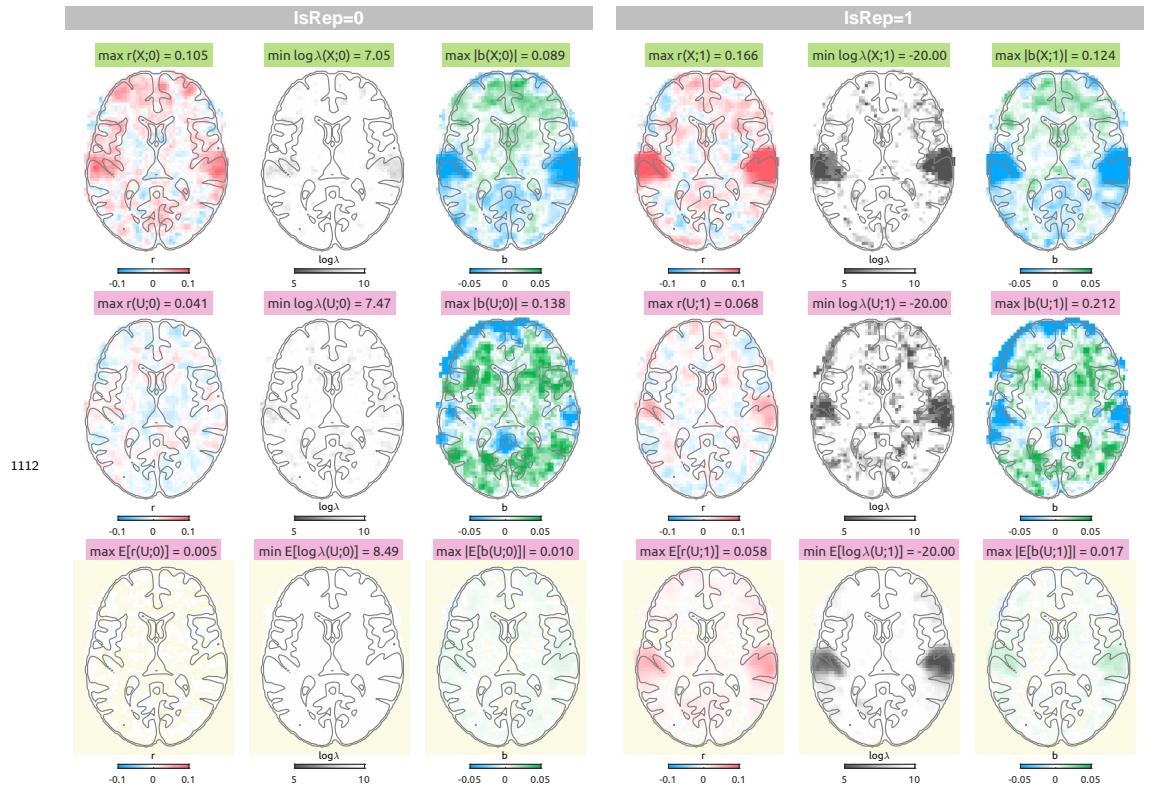


Figure 9—figure supplement 7. fMRI linearized encoding analysis results with delays from 4 to 6 sec with the normal noise as the null feature. The visualization scheme is identical to [Figure 9—figure Supplement 1](#).

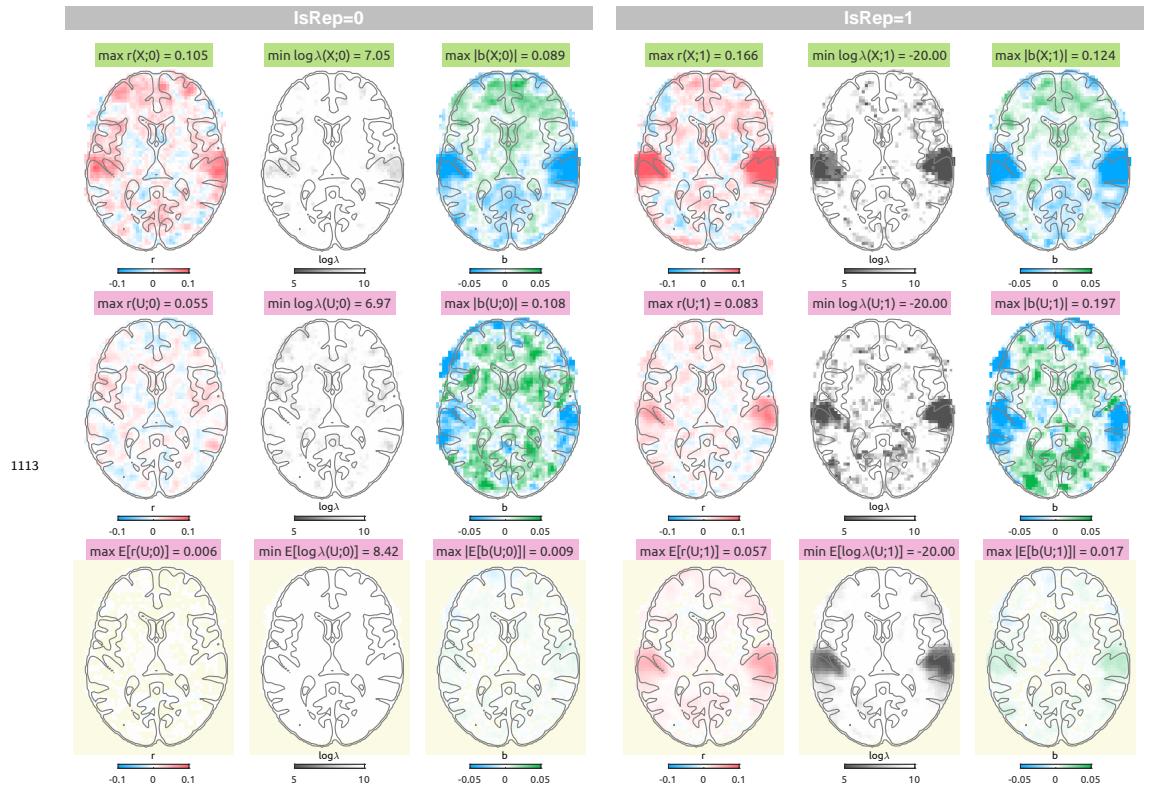


Figure 9—figure supplement 8. fMRI linearized encoding analysis results with delays from 0 to 12 sec with the uniform noise as the null feature. The visualization scheme is identical to [Figure 9—figure Supplement 1](#).

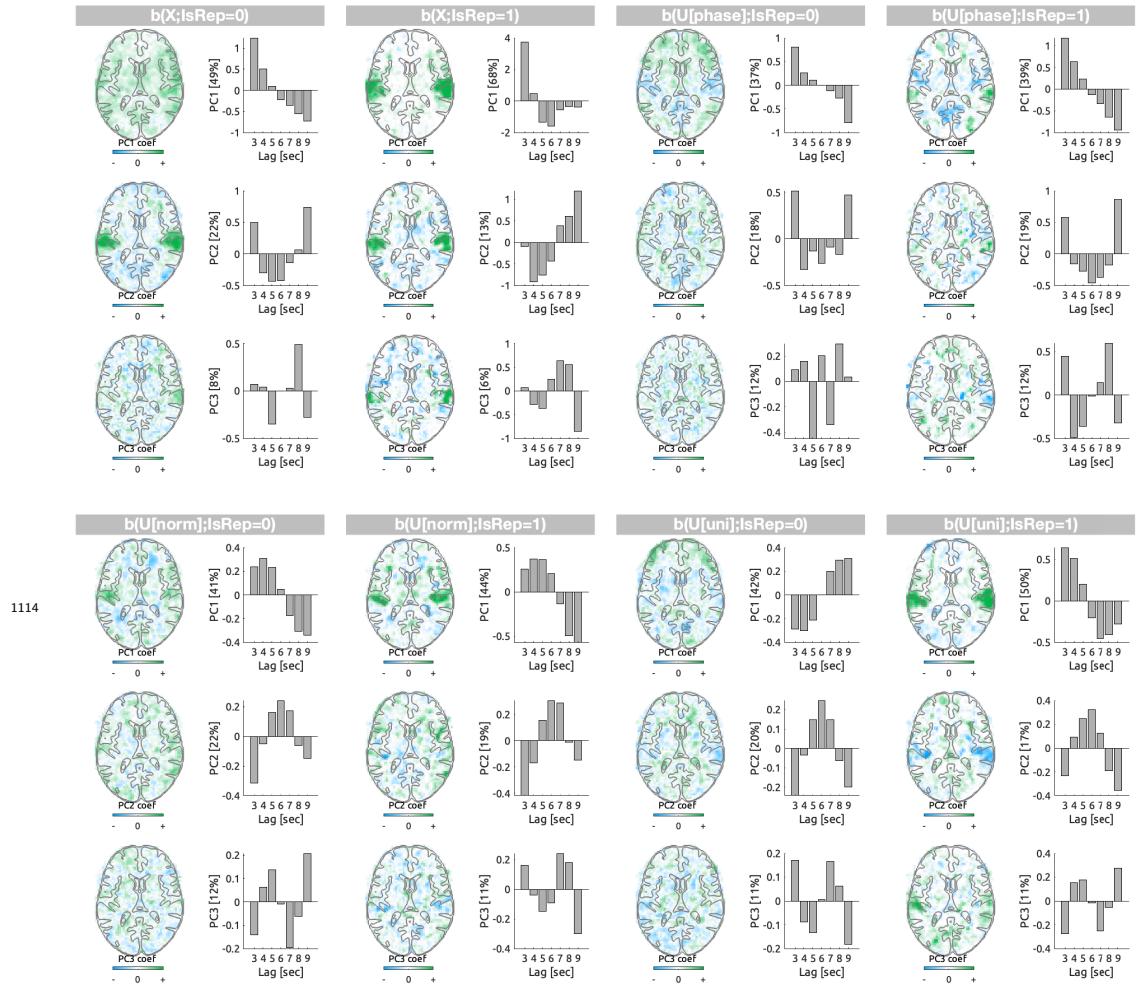


Figure 9—figure supplement 9. First three principal components (PCs) of the transfer function weights (b) of the models with delays from 3 to 9 sec are displayed in the transverse slice of eigenvectors and the time series of eigenvariates with the explained variance noted. The weights are grouped by features (X , true audio envelope; $U[\text{phase}]$, phase-randomized envelope; $U[\text{norm}]$, normal noise; $U[\text{uni}]$, uniform noise) and the CV schemes ($\text{IsRep} = 0$, $\text{IsRep} = 1$).

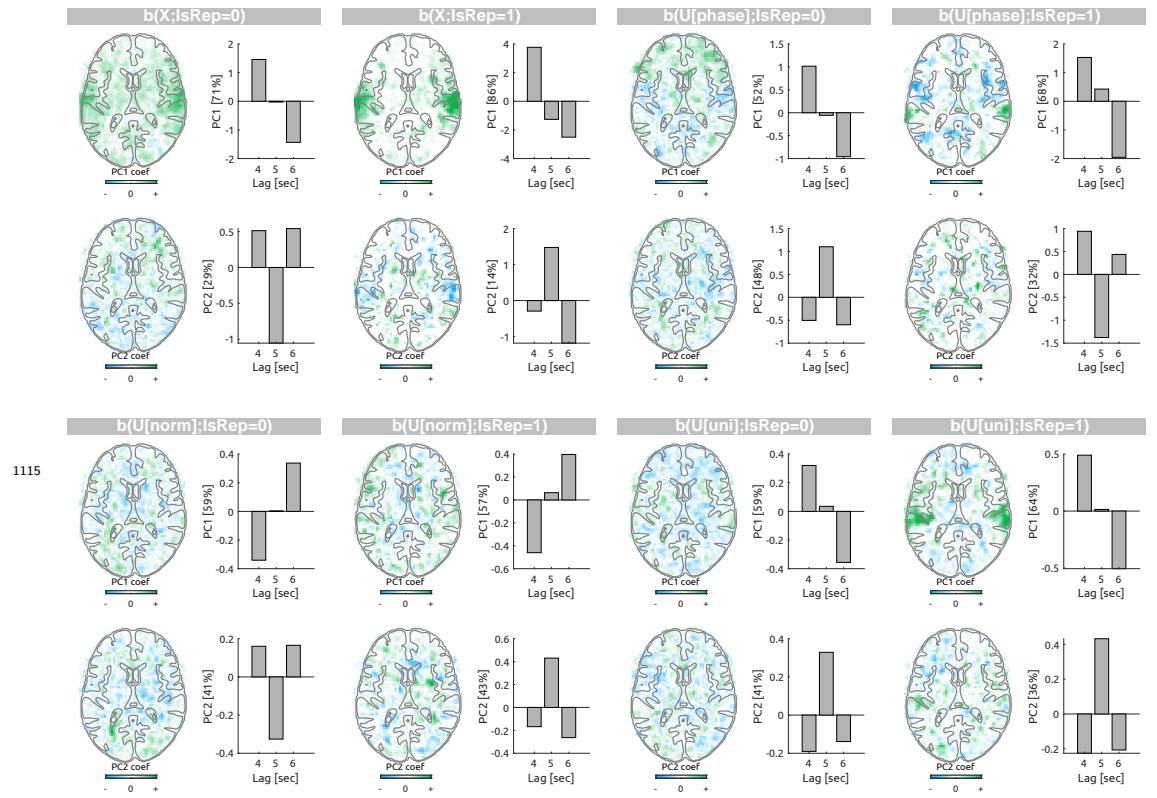


Figure 9—figure supplement 10. First two principal components (PCs) of the transfer function weights (b) of the models with delays from 4 to 6 sec are displayed. The visualization scheme is identical to **Figure 9—figure Supplement 9**.

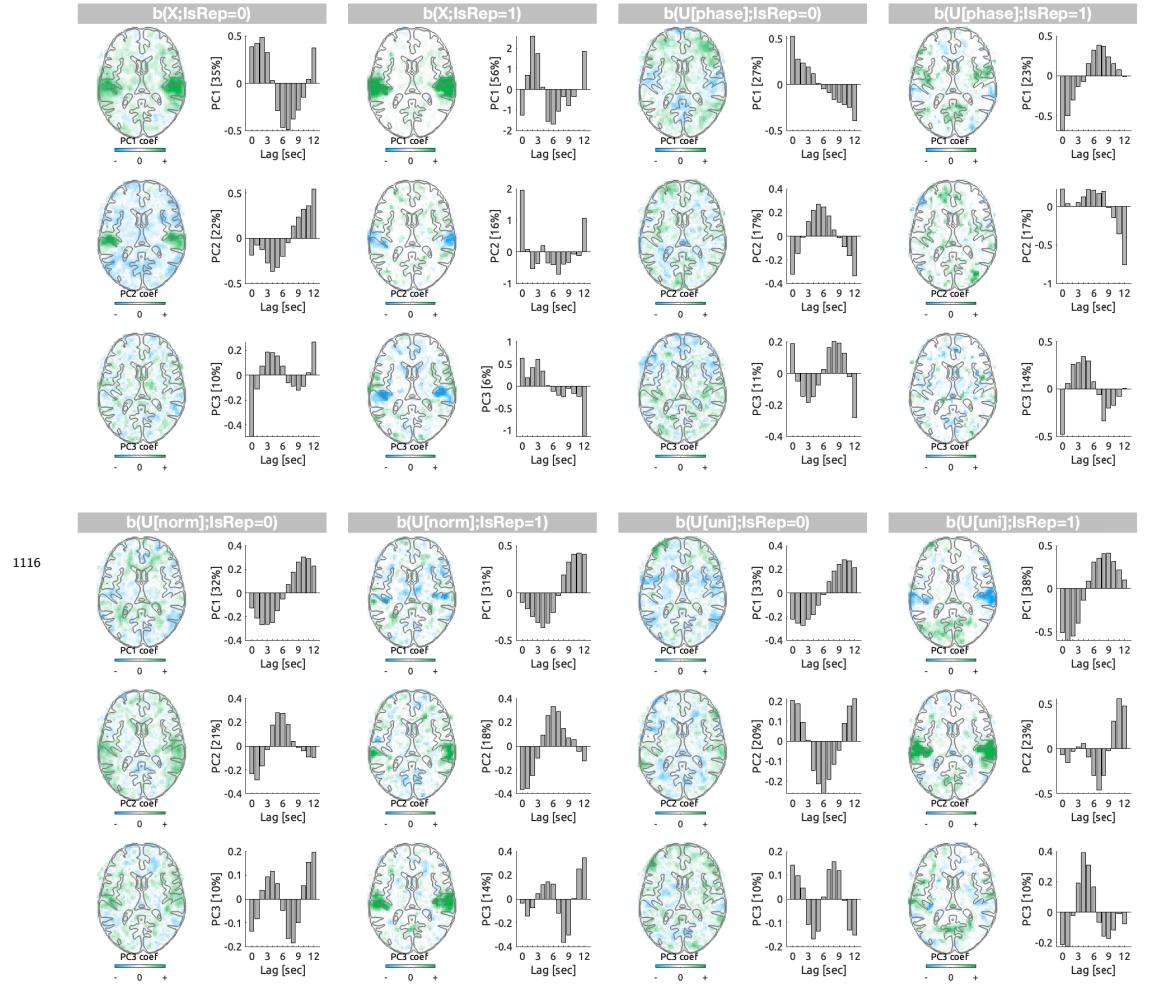


Figure 9—figure supplement 11. First three principal components (PCs) of the transfer function weights (b) of the models with delays from 0 to 12 sec are displayed. The visualization scheme is identical to **Figure 9—figure Supplement 9**.

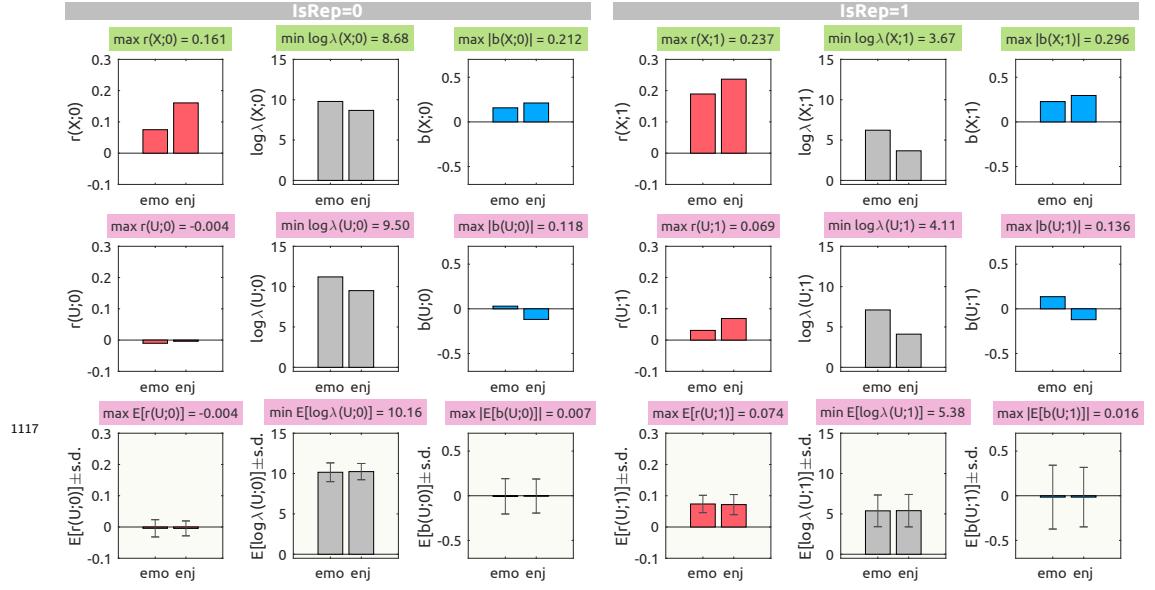


Figure 11—figure supplement 1. Behavioral linearized encoding analysis results with delays from 0 to 10 sec with an audio envelope (top row), a single case of a phase-randomized envelope (middle row), and an average of 100 phase-randomized envelopes (bottom row, pale yellow background). For each CV scheme ($\text{IsRep} = 0$, left panels; $\text{IsRep} = 1$, right panels), prediction accuracy (r , red bars), logarithmic ridge hyperparameter ($\log_{10} \lambda$, gray bars), transfer function weights that are summed over delays (b , blue bars) are shown along the columns.

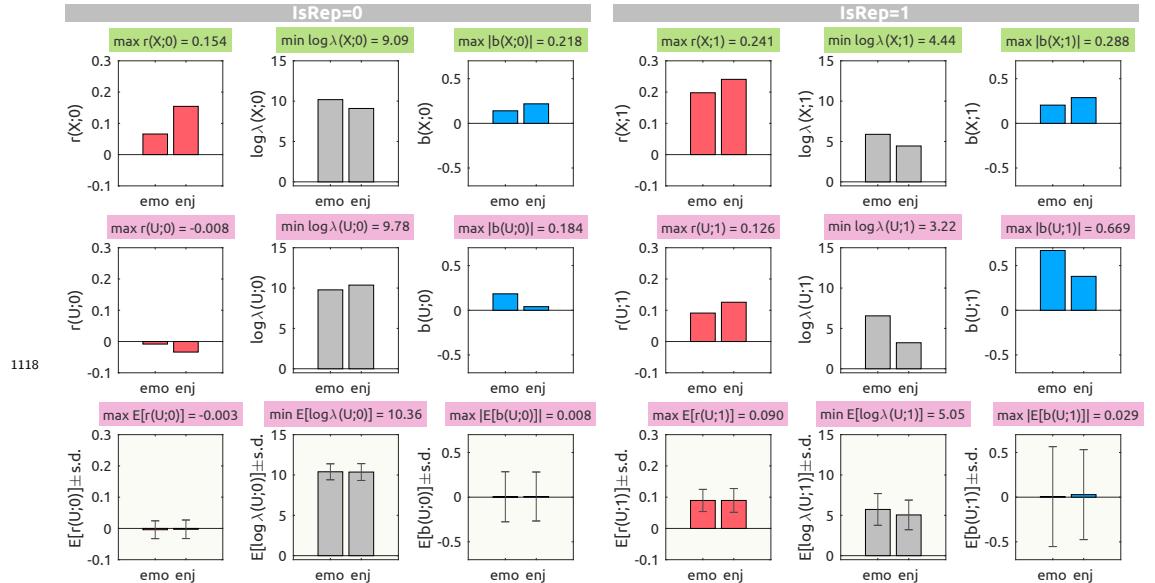


Figure 11—figure supplement 2. Behavioral linearized encoding analysis results with delays from 0 to 10 sec with the uniform noise as the null feature. The visualization scheme is identical to **Figure 11—figure Supplement 1**.

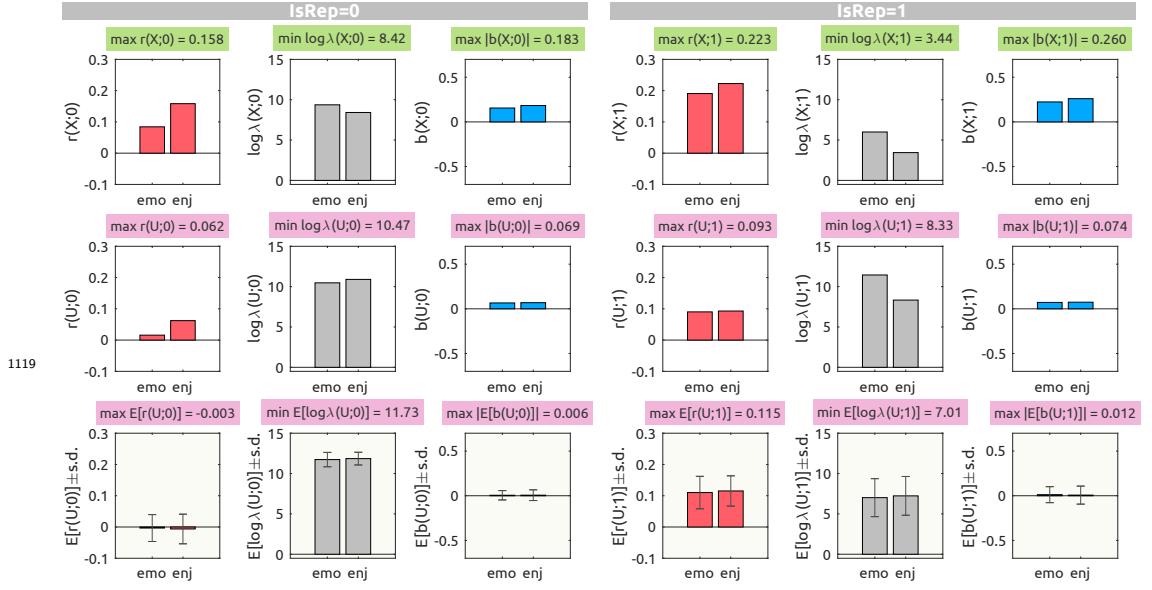


Figure 11—figure supplement 3. Behavioral linearized encoding analysis results with delays from 0 to 5 sec with the phase-randomized envelope as the null feature. The visualization scheme is identical to **Figure 11—figure Supplement 1**.

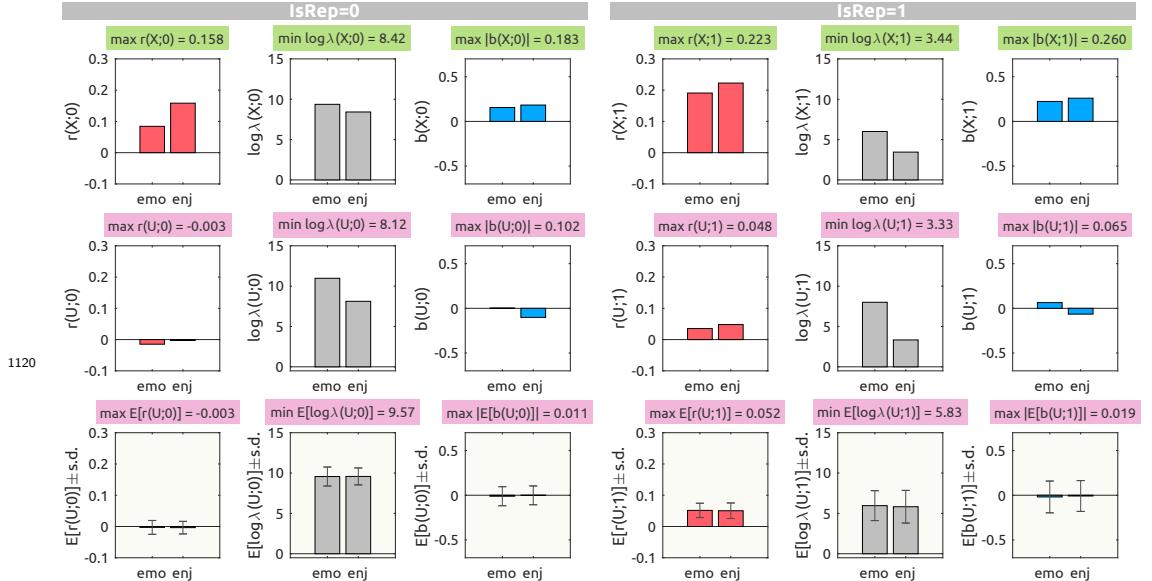


Figure 11—figure supplement 4. Behavioral linearized encoding analysis results with delays from 0 to 5 sec with the phase-randomized envelope as the null feature. The visualization scheme is identical to **Figure 11—figure Supplement 1**.

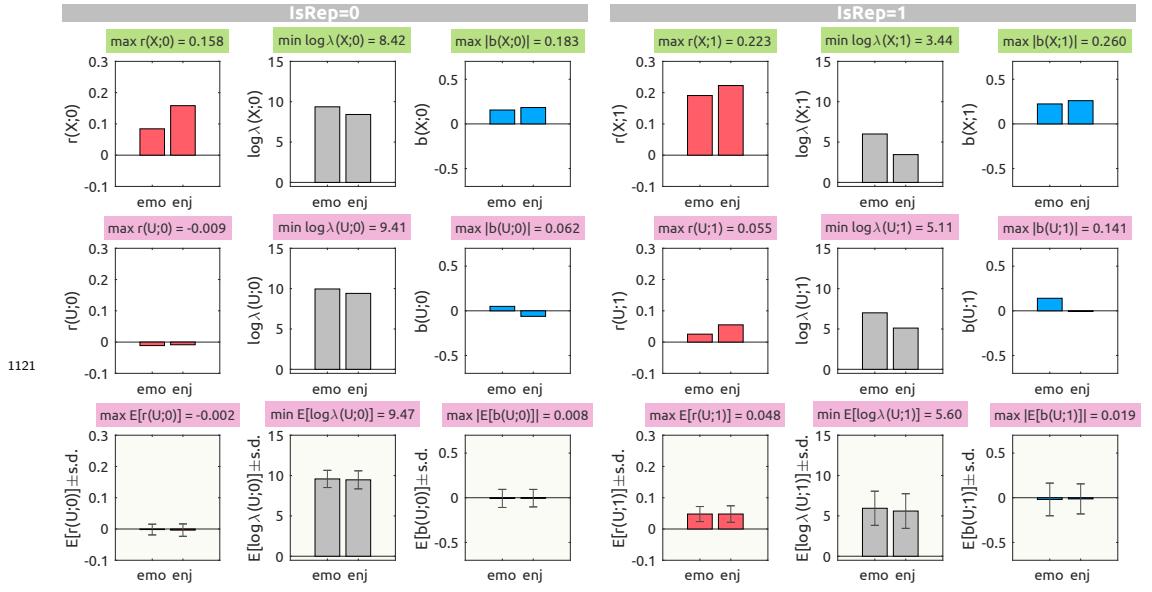


Figure 11—figure supplement 5. Behavioral linearized encoding analysis results with delays from 0 to 5 sec with the uniform noise as the null feature. The visualization scheme is identical to *Figure 11—figure Supplement 1*.

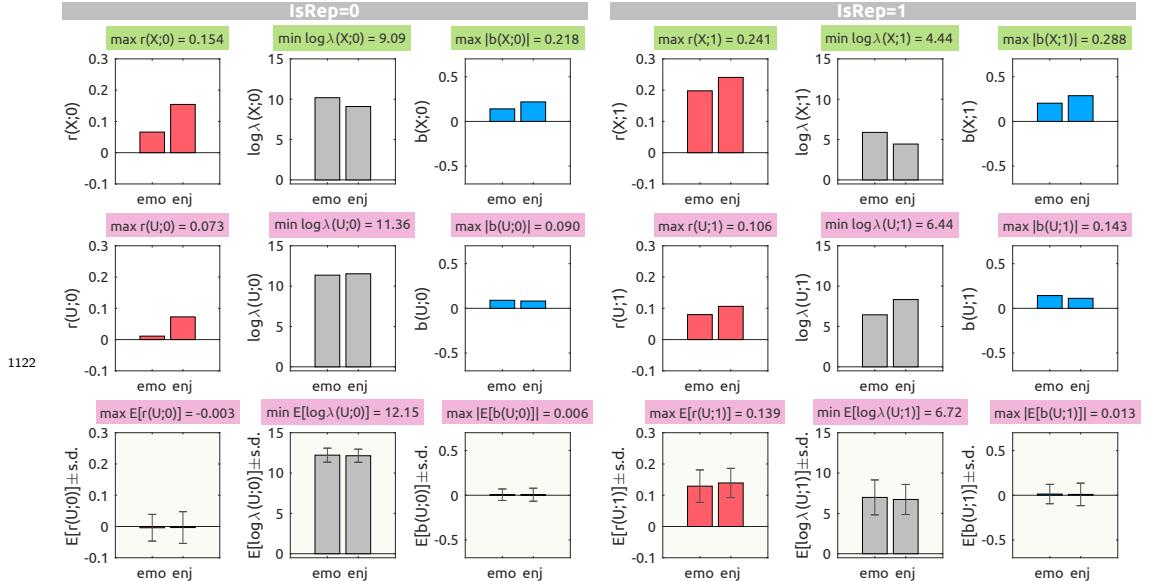


Figure 11—figure supplement 6. Behavioral linearized encoding analysis results with delays from 0 to 15 sec with the phase-randomized envelope as the null feature. The visualization scheme is identical to *Figure 11—figure Supplement 1*.

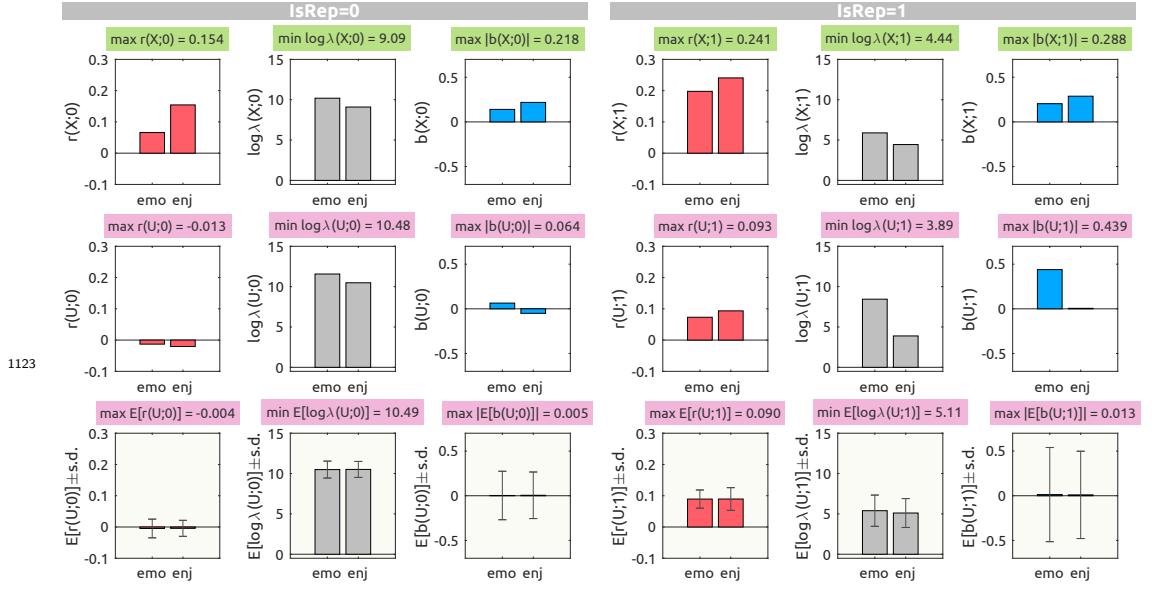


Figure 11—figure supplement 7. Behavioral linearized encoding analysis results with delays from 0 to 15 sec with the phase-randomized envelope as the null feature. The visualization scheme is identical to **Figure 11—figure Supplement 1**.

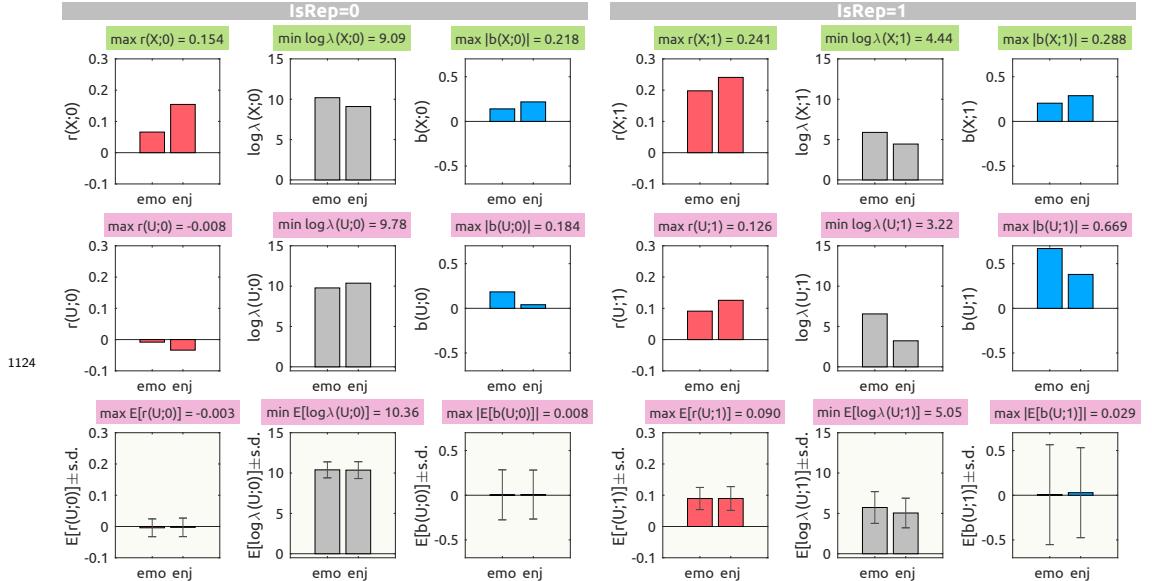


Figure 11—figure supplement 8. Behavioral linearized encoding analysis results with delays from 0 to 15 sec with the uniform noise as the null feature. The visualization scheme is identical to **Figure 11—figure Supplement 1**.

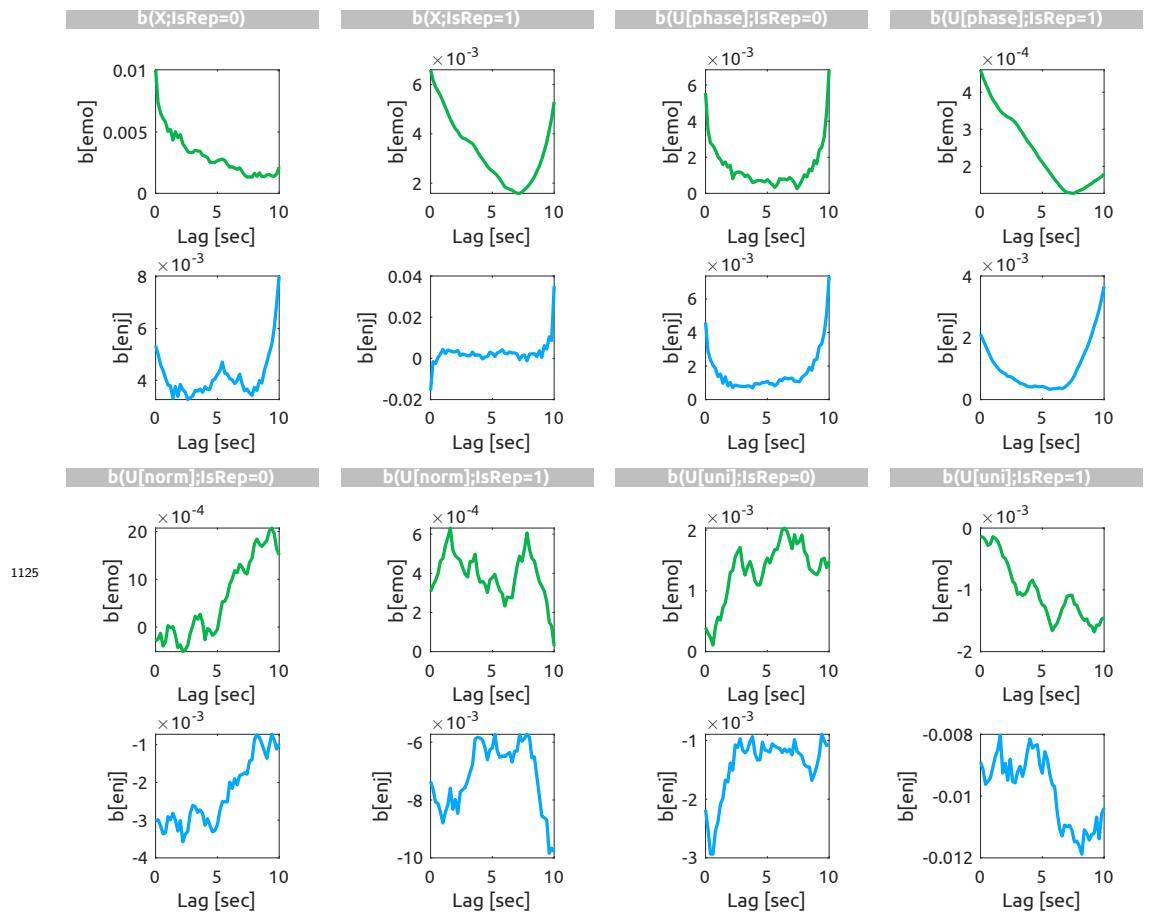


Figure 11—figure supplement 9. Transfer function weights (b) of Emotionality (green) and Enjoyment (blue) for the models with delays from 0 to 15 sec. The weights are grouped by features (X , true audio envelop; $U[\text{phase}]$, phase-randomized envelop; $U[\text{norm}]$, normal noise; $U[\text{uni}]$, uniform noise) and the CV schemes ($\text{IsRep} = 0$, $\text{IsRep} = 1$). Note that each weight time series is individually scaled.

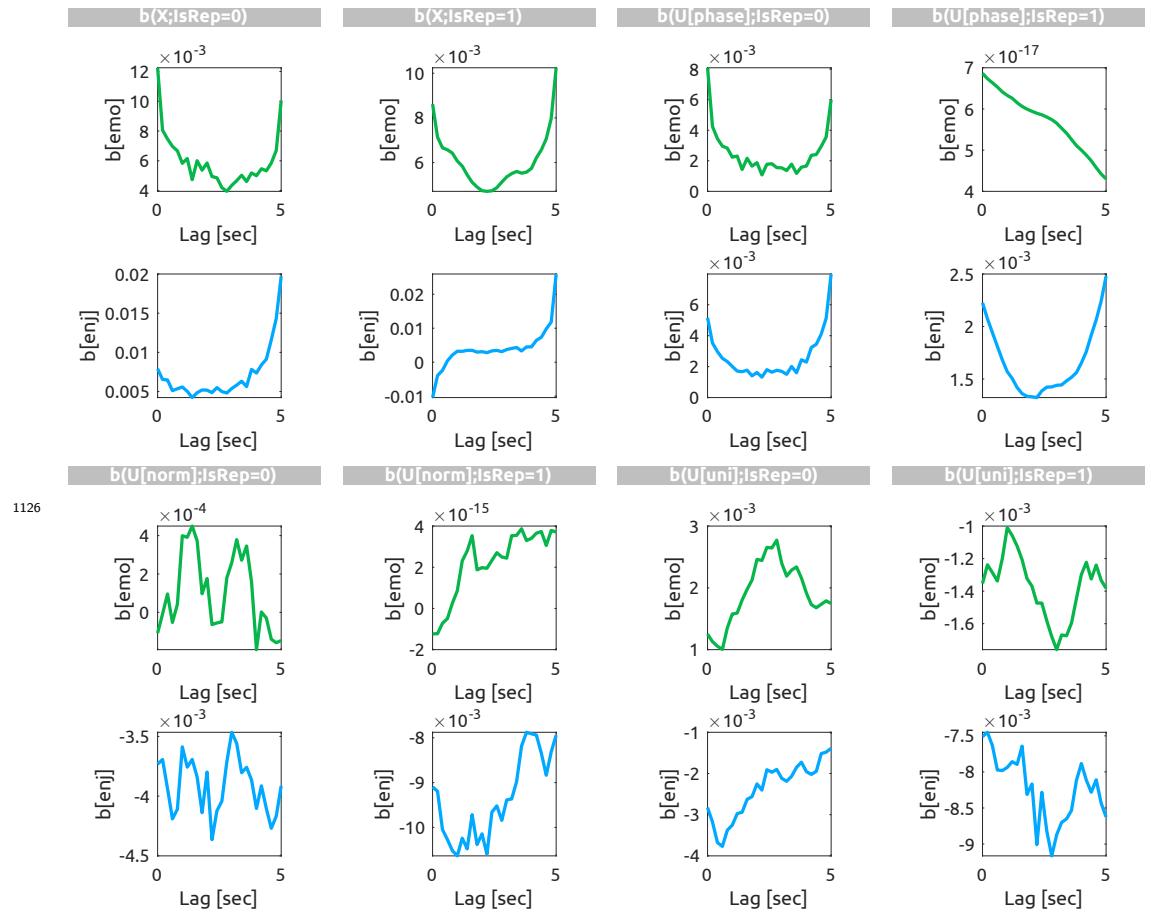


Figure 11—figure supplement 10. Transfer function weights (b) of the models with delays from 0 to 5 sec are displayed. The visualization scheme is identical to *Figure 11—figure Supplement 9*.

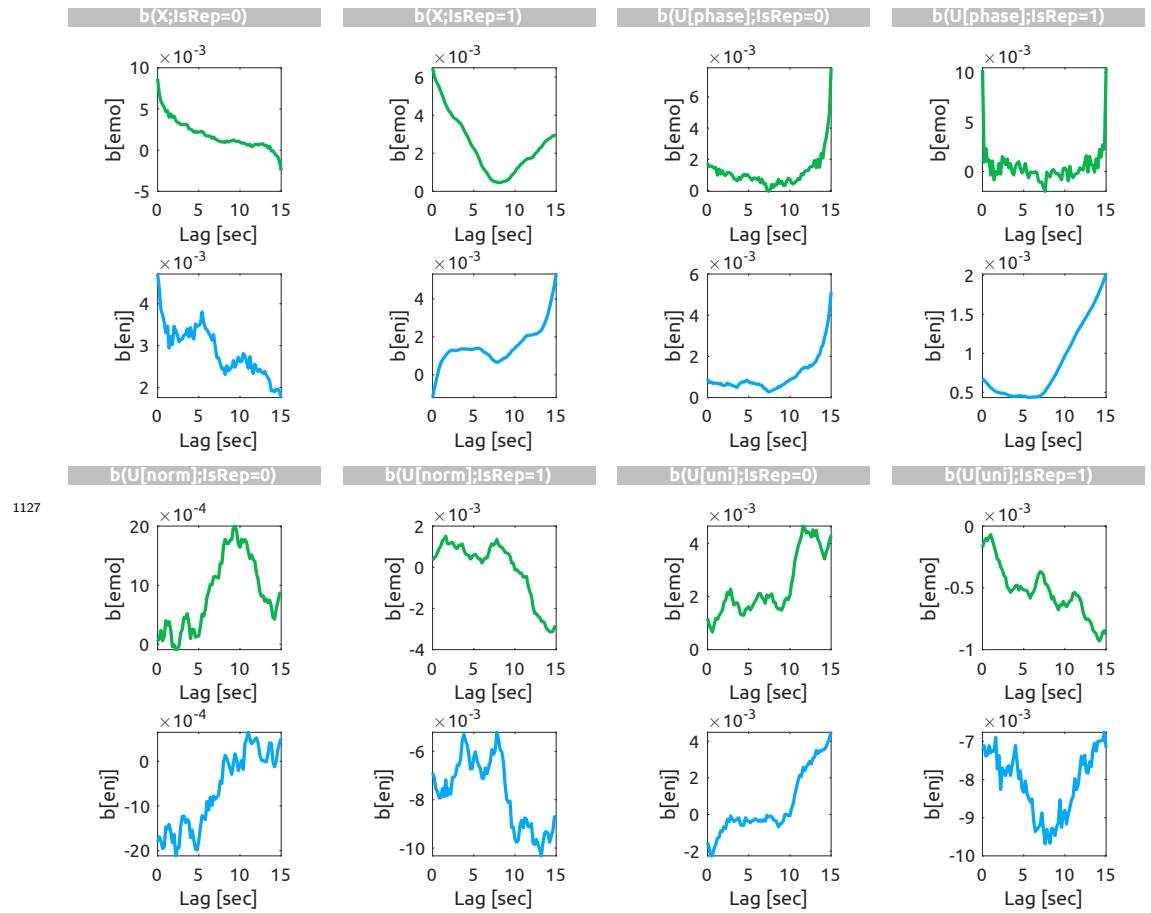


Figure 11—figure supplement 11. Transfer function weights (b) of the models with delays from 0 to 15 sec are displayed. The visualization scheme is identical to [Figure 11—figure Supplement 9](#).

1127