

Representational models: A common framework for understanding encoding, pattern-component, and representational-similarity analysis

Jörn Diedrichsen¹ & Nikolaus Kriegeskorte²

1. Brain and Mind Institute, Department for Computer Science, Department for Statistical and Actuarial Science, Western University, Canada

2. Cognitive and Brain Sciences Unit, Cambridge University, UK

Address correspondence:

Jörn Diedrichsen

Brain Mind Institute

Natural Science Center

Western University

London, Ontario, N6A 5B7

Canada

Email: jdiedric@uwo.ca

1 Abstract

2 Representational models specify how activity patterns in populations of neurons (or, more
3 generally, in multivariate brain-activity measurements) relate to sensory stimuli, motor responses,
4 or cognitive processes. In an experimental context, representational models can be defined as
5 hypotheses about the distribution of activity profiles across experimental conditions. Previous
6 studies have used three different methods to test such hypotheses: encoding analysis, pattern
7 component modeling (PCM), and representational similarity analysis (RSA). Here we develop a
8 common mathematical framework for understanding the relationship of these three methods, which
9 all share one core commonality: all three evaluate the second moment of the distribution of activity
10 profiles, which determines how well any feature can be linearly decoded from population activity.
11 Using simulated data for three different experimental designs, we compare the power of the
12 methods to adjudicate between competing representational models. PCM implements a likelihood-
13 ratio test and therefore provides the most powerful test if its assumptions hold. However, the other
14 two approaches – when conducted appropriately – can perform similarly. In encoding analysis, the
15 linear model needs to be appropriately regularized, which effectively imposes a prior on the activity
16 profiles. With such a prior, an encoding model specifies a well-defined distribution of activity
17 profiles. In RSA, the unequal variances and statistical dependencies of the dissimilarity estimates
18 need to be taken into account to enable near-optimal inference. The three methods render different
19 aspects of the information explicit (e.g. single-response tuning in encoding analysis and population-
20 response representational dissimilarity in RSA) and have specific advantages in terms of
21 computational demands, ease of use, and extensibility. The three methods are properly construed as
22 complementary components of a comprehensive data-analytical toolkit for understanding neural
23 representations on the basis of multivariate brain-activity data.

24 Author Summary

25 Modern neuroscience can measure activity of many neurons or the local blood oxygenation
26 of many brain locations simultaneously. As the number of simultaneous measurements grows, we
27 can better investigate how the brain represents and transforms information, to enable perception,
28 cognition, and behavior. Recent studies go beyond showing *that* a brain region is involved in some
29 function. They use representational models that specify *how* different perceptions, cognitions, and
30 actions are encoded in brain-activity patterns. In this paper, we provide a general mathematical
31 framework for such representational models, which clarifies the relationships between three
32 different methods that are currently used in the neuroscience community. All three methods
33 evaluate the same core feature of the data, but each has distinct advantages and disadvantages.
34 Pattern component modelling (PCM) implements the most powerful test between models, and is
35 analytically tractable and expandable. Representational similarity analysis (RSA) provides a highly
36 useful summary statistic (the dissimilarity) and enables model comparison with weaker
37 distributional assumptions. Finally, encoding models characterize individual responses and enable
38 the study of their layout across cortex. We argue that these methods should be considered
39 components of a larger toolkit for testing hypotheses about the way the brain represents
40 information.

41 Introduction

42 The measurement of brain activity is rapidly advancing in terms of spatial and temporal resolution,
43 and in terms of the number of responses that can be measured simultaneously [1]. Modern electrode
44 arrays and calcium imaging enable the recording of hundreds of neurons in parallel.
45 Electrophysiological signals that reflect summaries of the population activity can be recorded using
46 both invasive (e.g. the local field potential, LFP) and non-invasive techniques (e.g. scalp
47 electrophysiological measurements) at increasingly high spatial resolution. Modern functional

48 magnetic resonance imaging (fMRI) enables us to measure hemodynamic activity in hundreds of
49 thousands of voxels across the entire human brain at sub-millimeter resolution.

50 In order to translate advances in brain-activity measurement into advances in computational
51 theory [2], researchers increasingly seek to test representational models that capture both what
52 information is represented in a population of neurons, and how it is represented. Knowing the
53 content and format of representations provides strong constraints for computational models of brain
54 information processing. We refer to hypotheses about the content and format of brain
55 representations as *representational models*. We address here the important methodological question
56 of how to best test representational models.

57 Referring to an activity pattern as a “representation” of some variable (such as a perceptual
58 property, some cognitive content, or an action parameter) implies (1) that there is information about
59 the variable present in the pattern of activity and (2) that this information *serves a functional*
60 *purpose* for further processing [3]. Even though observational methods such as fMRI or single-cell
61 recording can generally not be used to test (2) directly, we can test whether the variable is encoded
62 in such a form that makes it easily accessible to read out. This criterion leads to our working
63 definition of an *explicit neural representation*, on which our specific formulation of a
64 “representational model” is based: A variable is *explicitly represented* in an area, if it can be
65 *linearly* decoded from the neural activity pattern in the area [4]. To understand why the limitation to
66 linear decoding makes sense, consider the following example: When looking at a face, the firing
67 pattern of retinal neurons contains *information* about the identity of the person, as it can be read out
68 from using a highly non-linear decoder. However, the face identity is *implicit* in the retinal
69 representation, requiring the decoder, in essence, to perform face recognition. For a population code
70 to constitute an *explicit representation*, another area must be able to read out the represented
71 variable directly, for example as a weighted sum of the neural activity pattern across neurons in the
72 region [2, 4, 5]. Note that this definition does not restrict representations to highly localized codes,

73 such as the “grandmother neuron” [6], but encompasses any distributed code, that can be read out
 74 by computing a linear combination of neuronal activities.

75 This definition of explicit representation has motivated researchers to use linear decoding
 76 methods to study neural representations [7-9]. Decoding analysis can establish that a given feature
 77 is explicitly represented in the activity pattern. Representational models, as considered here, go one
 78 step further: they fully characterize the representational geometry, which defines *all* explicitly
 79 represented features in a region, how strongly each of them is represented (relative signal to noise
 80 ratio), and how the activity patterns associated with different features relate to each other.
 81 Representational models therefore fully specify the explicit representational content of an area.

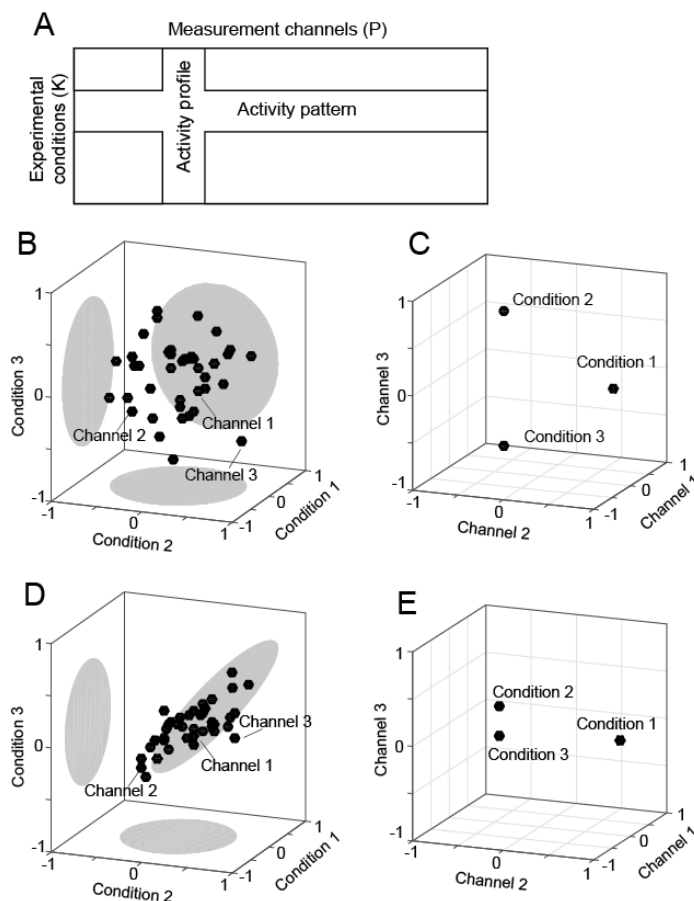


Figure 1. Two complementary perspectives on population activity.

(A) The multivariate activity data can be viewed as a set of activity profiles (columns) or as a set of activity patterns (rows). An activity profile is a vector of responses of a single channel across experimental conditions. An activity pattern is a vector of responses across all channels for a single condition.

94 *Activity data can be visualized by plotting activity profiles as points in a space defined by the*
 95 *experimental conditions (B,D), or by plotting the activity patterns as points in a space defined*
 96 *by the measurement channels (C,E). (B) If the activities are uncorrelated between conditions,*
 97 *then (C) the corresponding activity patterns of all three conditions are equidistant to each*

98 *other, and can be equally well distinguished. (D) If the activities are positively correlated for*
99 *two conditions (conditions 2 and 3 here), then (E) the activity patterns for these conditions*
100 *are closer to each other and can be less well distinguished.*

101 To define representational models formally, we need to consider two complementary
102 perspectives on activity data, as illustrated in Fig. 1. The activity of many neurons, or more
103 generally *measurement channels* (neurons, electrodes, or fMRI voxels), can be measured across a
104 range of *experimental conditions* (stimuli, movements, or tasks). Thus, each channel will have an
105 *activity profile*, which can be plotted as a point in the space spanned by the experimental conditions
106 (Fig. 1b). A representational model specifies a *probability distribution of activity profiles* in the
107 space spanned by the experimental conditions – i.e. it treats the true activity profiles as a random
108 variable. In other words, a representational model predicts, for each possible activity profile, the
109 probability of observing a measurement channel exhibiting that profile. However, it does not predict
110 the activity profile for each individual channel actually measured. The motivation for this approach
111 derives from the idea that the computational function of a region does not depend on specific
112 neurons having specific response properties, but on the fact that certain features can be read out
113 from the population by downstream neurons. The probability distribution over activity profiles
114 determines which features can be linearly read out from the code and the signal-to-noise ratio of the
115 readout. By basing further analyses on the probability distribution of the activity profiles, we are
116 disregarding three aspects of the code: (1) which neuron fulfills which function, (2) where neurons
117 are located on the cortical sheet, and (3) the degree to which the information about a given
118 represented feature is concentrated in a few neurons (as in single-cell selectivity for a represented
119 feature) or spread out over the population. Ignoring these aspects may be viewed as an advantage or
120 a disadvantage, depending on the level of description that a researcher is interested in. We argue
121 that treating activity profiles as random vectors is a simplification that is useful for drawing
122 computational insights from population activity measurements.

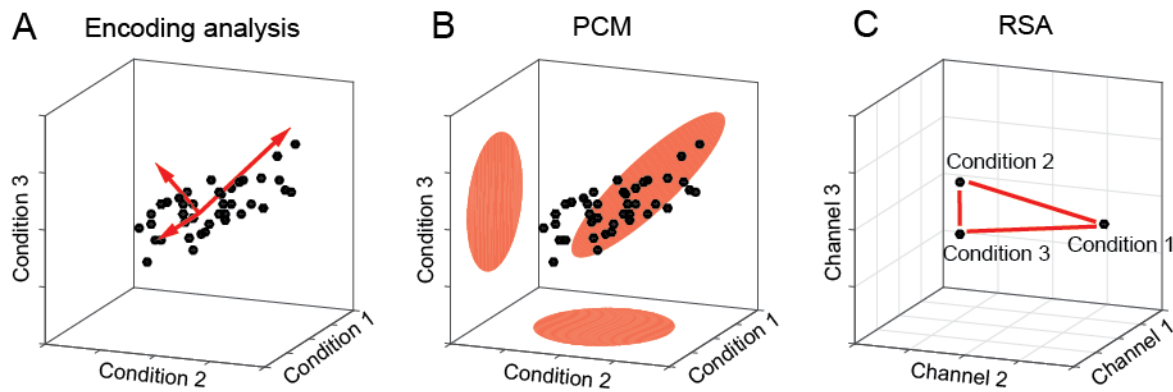
123 In this paper, we show that the signal-to-noise ratio with which any given feature can be
 124 linearly decoded is determined by a specific aspect of the distribution of the activity profiles,
 125 namely its multivariate *second moment*. We discuss three established methods for adjudicating
 126 between representational models: encoding analysis, pattern-component modeling (PCM) and
 127 representational similarity analysis (RSA, see Table 1). We show that these three techniques all
 128 exclusively rely on information contained in the second moment of the distribution of activity
 129 profiles. This core commonality enables us to consider them in the same formal framework.

130 **Table 1:** Comparison of encoding analysis with regularization, pattern component modelling
 131 (PCM), and representational dissimilarity analysis (RSA).

	Encoding analysis	PCM	RSA
Model definition	Model-feature matrix \mathbf{M} , regularization / prior	Predicted second-moments matrix (\mathbf{G})	Representational dissimilarity matrix (RDM)
First-level parameters (characterizing individual activity profiles)	One weight per feature and measurement channel	None; integrated out in the likelihood	None; integrated out when calculating dissimilarities
Second-level parameters (characterizing the distribution of activity profiles)	Regularization / Ridge coefficient (determined by noise / signal ratio)	Scale parameter s Noise variance	Scaling between predicted and observed distances (s)
Prediction target	Responses to test conditions	Distribution of measurement channels in activity-profile space	Dissimilarities among activity patterns
Training data required	always	not for fixed models, only if additional second-level parameters are to be fitted	not for fixed models, only if additional second-level parameters are to be fitted
Explicit likelihood for fitting additional model parameters	No – need to do nested within crossvalidation	Yes	Yes
Fitting algorithms for model parameters	-	EM Gradient descent Newton-Raphson	Linear and non-negative regression IRLS

132

133



134

135

136

137

138

139

140

141

142

143

144

145

146

Figure 2. Three approaches to testing representational models. (A) In encoding analysis, the distribution of activity profiles is described by the underlying features (red vectors). The direction of feature vector determines the associated activity profile, and the length the strength of the feature encoding in the representation. (B) PCM models the distribution of the activity profiles as a multivariate Gaussian. This model is parametrized by the second moment of the activity profiles, which determines at what signal-to-noise ratio any feature is linearly decodable from the population. (C) RSA uses the representational distances (or, more generally, dissimilarities) between activity patterns as a summary statistic to describe the underlying distribution. When the representational distances are estimates of the Mahalanobis distances (normalized by the error covariance), the distance matrix, like a second moment of the activity profiles, is a sufficient statistic that captures the linear decodability of any feature.

147

148

149

150

151

152

153

154

In *encoding analysis* [10, 11], representational models are defined in terms of the underlying features (Fig. 2A). Each activity profile can be characterized by a linear combination of such features. Examples include Gabor filters [12] (a model of low-level visual representation), abstract semantic dimensions [13] (a higher-level cognitive representation), and force, direction or hand position [14-16] (a movement representation). The prevalence of each feature in each channel is measured by a feature weight. Feature weights are considered first-level parameters in our framework, as they describe the individual activity profiles, as opposed to second-level parameters that describe the distribution of the activity profiles (Table 1). The large number of parameters

155 (number of features in the model times number of channels in the measurements) engenders a
156 danger of over-fitting. To account for effects of overfitting, encoding models are commonly
157 evaluated using cross-validation: The feature weights are estimated on a training set, and the model
158 is evaluated in terms of its performance at predicting left-out data [10]. The test data may consist in
159 a sample of experimental conditions not used in training, so as to test the model's generalization
160 performance [11, 12]. While many studies use simple linear regression to estimate the weights [11,
161 17], it is increasingly common to introduce a regularization penalty in the estimation of model
162 weights (for example the L2 norm of the vector of weights) [12, 13]. We will see that regularization
163 is not merely a technical trick used in fitting a given model. Instead, the regularization (and its
164 implicit distributional assumptions) are an essential part of the representational hypothesis that is
165 tested. Without it, encoding models do not specify a probability distribution with a defined second
166 moment and thus do not define the linear decodability of different features.

167 *Pattern component modeling* [18] is based on an explicit generative model of the process
168 that produced the data and can be considered a Bayesian approach. The true activity profiles are
169 assumed to have a multivariate Gaussian distribution in the space spanned by the experimental
170 conditions (Fig. 2B). The activity measurements are assumed to be corrupted by zero-mean additive
171 Gaussian noise. This formulation enables us to evaluate the marginal likelihood of the observed
172 activity profiles under the probability distribution specified by the model. Thus, we do not fit any
173 first-level parameters (feature weights) and hence reduce the risk of overfitting. It thus enables us to
174 compare models with different numbers of features without having to correct for model complexity.
175 If the assumptions of the generative model hold, PCM implements the likelihood-ratio test between
176 models [19], which by the Neyman-Pearson lemma [20], is the most powerful test of its size. In
177 theory, therefore, PCM should yield more accurate inferences than any of its competitors, that is it
178 should be able to more sensitively adjudicate among competing models.

179 Finally, *representational similarity analysis* (RSA [21-23]) approaches the problem from a
180 complementary perspective. Rather than considering the activity profiles of the measurement

181 channels as points in the space spanned by the conditions (Fig. 1B,D), it considers the activity
182 patterns associated with the experimental conditions as points in the space spanned by the
183 measurement channels (Fig. 1C,E). RSA then uses the representational distances (Fig. 2C) between
184 the conditions as a summary statistic. We will see that these distances again exclusively depend on
185 the second moment of the distribution of activity profiles. Having obtained a matrix of
186 dissimilarities between activity patterns (the representational dissimilarity matrix, RDM), RSA then
187 tests models by comparing the observed distances to the distances predicted by each
188 representational model. In this paper, we consider different methods to perform this comparison,
189 including rank-based correlations [24], Pearson correlations [25], and a novel likelihood-based
190 approach that uses a multivariate normal approximation to the joint distribution of the cross-
191 validated Mahalanobis distances [26, 27]. As shown in our simulations here, this latter technique
192 provides a particularly powerful means to adjudicate between representational models.

193 In the remainder of the paper, we first introduce the second moment of the activity profiles
194 and explain why it is the sufficient statistic of linear decodability. We then define the three methods
195 in detail, using a common mathematical notation. Finally, using simulated data and models taken
196 from our fMRI work, we assess the statistical efficiency, i.e. how well these methods adjudicate
197 between two or more competing representational models given limited data. We also compare the
198 methods in terms of their computational efficiency.

199 Materials and Methods

200 *Basic definitions*

201 All symbols used in the following derivations are summarized in Table 2. First, we define \mathbf{U} to be
202 the matrix of noiseless activity profiles with K (number of experimental conditions) rows and P
203 (number of measurement channels) columns. Each row of this matrix is an activity pattern, the
204 response of the whole population to a single condition. Each column of this matrix is an activity
205 profile (Fig. 1a).

206 Table 2: Notation used. For non-scalars, the second column indicates the vector / matrix size.

K :		Number of conditions
M :		Number of independent partitions of the data (imaging runs)
P :		Number of measurement channels (voxels, electrodes, neurons)
N :		Overall number of measurements ($N_m \times M$)
Q :		Number of features in model
\mathbf{U} :	$K \times P$	Matrix of true activation patterns
$\mathbf{u}_{i,:}$:	$1 \times P$	Activation pattern for condition i ; i^{th} row of \mathbf{U}
$\mathbf{u}_{:,j}$:	$K \times 1$	Activation profile for measurement channel j ; j^{th} column of \mathbf{U}
$\hat{\mathbf{U}}^{(m)}$:	$K \times P$	Matrix of estimated activity patterns, based on data from partition m
$\tilde{\mathbf{U}}^{(\sim m)}$:	$K \times P$	Model prediction for activity patterns, based on data independent of m
\mathbf{M} :	$K \times Q$	Matrix of model features for all condition
\mathbf{W} :	$Q \times P$	Matrix of voxel weights for each feature
\mathbf{Y} :	$N \times P$	Matrix of brain measurements, concatenated activity estimates or time series data
\mathbf{Z} :	$N \times K$	Design matrix, indicating how measurements relate to activity patterns
\mathbf{X} :	$N \times R$	Design matrix containing n regressors of no-interest
\mathbf{G} :	$K \times K$	Second moment of \mathbf{U}
$d_{i,k}$:		Distance between condition i and k
J :		Number of distances, normally $K(K-1)/2$
\mathbf{D} :	$K \times K$	Representational dissimilarity matrix of all pairwise distances
\mathbf{d} :	$J \times 1$	Vector of all pairwise distances
$\tilde{\mathbf{d}}$:	$J \times 1$	Vector of predicted distances
\mathbf{C} :	$J \times K$	Contrast matrix, defining the J pairwise differences between conditions
Σ_P :	$P \times P$	Variance-covariance matrix between the P voxels
Σ_K :	$K \times K$	Variance-covariance matrix of the columns of $\hat{\mathbf{U}}^{(m)}$
\mathbf{V} :	$N \times N$	Variance-covariance matrix of \mathbf{Y}
\mathbf{S} :	$J \times J$	Variance-covariance matrix of all pair-wise distances

207 Because we are interested in the distribution of activity profiles, but not in the activity
 208 profiles per se, we consider the columns of \mathbf{U} to be a random variable. This is an essential step
 209 underlying our common framework, which is justified by the fact that, for the purpose of reading
 210 out information, the different measurement channels are exchangeable (see introduction). We
 211 assume that the activity profiles are repeatedly measured, with the data consisting of M independent
 212 partitions, each containing at least one activity measurement for each condition and measurement

213 channel. In the context of fMRI, a partition will consist of a separate phase of data acquisition, e.g.
214 a scanner run. The activity estimates $\hat{\mathbf{U}}^{(m)}$ of partition m are the true patterns \mathbf{U} plus noise $\mathbf{E}^{(m)}$.
215 The noise captures both neural trial-by-trial variability of the activity pattern in a single condition,
216 as well as measurement noise.

$$217 \quad \hat{\mathbf{U}}^{(m)} = \mathbf{U} + \mathbf{E}^{(m)}. \quad (\text{Eq. 1})$$

218 For the purposes of this paper, we assume that the noise is Gaussian, and independent and
219 identically distributed (i.i.d.) across conditions and partitions (homoscedasticity). Possible
220 dependence within each partition, however, can be easily accounted for [26, 28].

221 ***Dependence between measurement channels***

222 The discussion below further assumes that the noise is also i.i.d. across different measurement
223 channels (isotropicity). However, noise in fMRI, MEG, and even invasive electrophysiology
224 exhibits strong correlations between neighboring locations in the brain. To account for these
225 dependencies, we employ multivariate noise normalization (i.e. spatial prewhitening), which has
226 been shown to increase the reliability of inference [29]. Across all measurement channels, we
227 estimate the $P \times P$ variance-covariance matrix across trials, Σ_P and then regularize the estimate by
228 shrinking it towards a diagonal matrix [30]. In the context of fMRI, we can use the residual time
229 series from the fitting of the time-series model to estimate noise covariance [29] [31]. We then post-
230 multiply our activity estimates by $\hat{\Sigma}_P^{-1/2}$, rendering the model errors in the channels approximately
231 uncorrelated. If multivariate noise normalization is not performed or is incomplete, inference will
232 be suboptimal in all three methods (for details see [26]).

233 ***Second moment matrix and linear decodability***

234 In this section, we show that the *second moment* of the activity profiles fully characterizes
235 the linear decodability of any feature in the space spanned by the experimental conditions. This
236 provides a motivation, from the perspective of brain computation, for using the second moment
237 matrix as a summary statistic. While higher statistical moments may reveal some computationally

238 important features of the population activity (a point we will return to in the Discussion), we believe
239 that focusing on linear encoding (and hence the second moment) constitutes the natural starting
240 point for adjudicating among representational models.

241 The n^{th} moment of a scalar random variable u is $E(u^n)$, where $E()$ denotes the expected
242 value. Here we use a multivariate extension of the concept, with the second moment of the random
243 vector \mathbf{u} defined as the matrix $E(\mathbf{u}\mathbf{u}^T)$, the expected outer product of the activity profiles, where the
244 expectation is across the measurement channels. The second-moment matrix of the activity profiles
245 is given by

$$246 \quad \mathbf{G} \equiv \sum_{j=1}^P \mathbf{u}_{\cdot,j} \mathbf{u}_{\cdot,j}^T / P = \mathbf{U}\mathbf{U}^T / P. \quad (\text{Eq. 2})$$

247 Thus, each cell of this matrix contains the normalized inner product of two activity patterns.

248 Before calculating \mathbf{G} , some investigators subtract the mean activity across measurement
249 channels for each condition from the data. In this case, Eq. 2 becomes the variance-covariance
250 matrix of the activity profiles –the second moment around the mean activity profile. Here we do not
251 remove the mean, but use the second moment around zero. From the perspective of a neuron that
252 reads out the activity pattern of an area, any difference between activity patterns across conditions
253 can be used to encode information. Some features (for example, stimulus intensity) may be encoded
254 in the mean activity over all measurement channels. Other properties (for example, stimulus
255 identity) may be encoded in relative activity differences, with some measurement channels
256 responding more to one condition, and others to a different condition. The second moment (around
257 zero) captures both of these potentially meaningful types of difference.

258 To obtain a read-out for a certain feature y for condition i , for example a stimulus property
259 represented in the code, we would weight each channel's observed activity during condition i using
260 the $P \times I$ read-out vector \mathbf{v} ,

$$261 \quad \hat{y}_i = \hat{\mathbf{u}}_{i,\cdot} \mathbf{v}. \quad (\text{Eq. 3})$$

262 We would now like the estimate \hat{y} to have very different values for two conditions that differ in the
263 feature being read out, while showing small differences for trials that have the same feature value.

264 The feature to be read out can be defined by a specific contrast on the experimental conditions,
265 defined by a $K \times I$ vector \mathbf{f} . We are looking for the readout weight vector \mathbf{v} that maximizes the ratio
266 between feature variance and error variance, and thus the signal-to-noise ratio (S), of the readout:

$$267 \quad S = \frac{\mathbf{v}^T \mathbf{U}^T \mathbf{f} \mathbf{f}^T \mathbf{U} \mathbf{v}}{\mathbf{v}^T \mathbf{E}^T \mathbf{f} \mathbf{f}^T \mathbf{E} \mathbf{v}} . \quad (\text{Eq. 4})$$

268 The solution to this equation is commonly known as Fisher's linear discriminant [32], which, under
269 the assumption of homoscedastic Gaussian noise, is the best achievable linear decoder. If the noise
270 is isotropic (or the data is adequately pre-whitened), then $\mathbf{E}^T \mathbf{f} \mathbf{f}^T \mathbf{E} = \mathbf{I}b$, where b is a constant. The
271 denominator then depends only on the norm of the read-out vector \mathbf{v} , not on its direction, and can be
272 ignored when \mathbf{v} is constrained to have a norm of 1. The best readout vector \mathbf{v} is then given by the
273 first eigenvector of the matrix $\mathbf{U}^T \mathbf{f} \mathbf{f}^T \mathbf{U}$, and the quality of the best readout is determined by the
274 largest eigenvalue.

275 The non-zero eigenvalues (*eig*) of a square matrix are invariant to cyclic permutations of the
276 product order:

$$277 \quad \text{eig}(\mathbf{U}^T \mathbf{f} \mathbf{f}^T \mathbf{U}) = \text{eig}(\mathbf{f}^T \mathbf{U} \mathbf{U}^T \mathbf{f}) = P \text{eig}(\mathbf{f}^T \mathbf{G} \mathbf{f}) . \quad (\text{Eq. 5})$$

278 Therefore, the quality of the best linear decoder for *any* feature (as defined by \mathbf{f}) is fully
279 characterized by the second moment matrix \mathbf{G} of the pre-whitened activity patterns.

280 ***Representational analysis in the context of fMRI***

281 The methods in this paper were first developed in the context of fMRI data analysis, and our
282 examples will come from this domain. A simple way to apply the analyses to fMRI data is to use as
283 activity estimates ($\hat{\mathbf{U}}^{(m)}$) the regression coefficients, or “beta”-weights, from a first-level time series
284 analysis [33, 34]. The time-series model accounts for the hemodynamic lag and the temporal
285 autocorrelation of the noise. The activity estimates usually express the difference in activity during
286 a condition relative to rest. Activity estimates commonly co-vary together across fMRI imaging
287 runs, because all activity estimates within a partition are measured relative to the same resting
288 baseline. This positive correlation can be reduced by subtracting, within each partition, the mean

289 activity pattern (across conditions) from each activity pattern. This makes the mean of each
290 measurement channel (across condition) zero and thus centers the ensemble of points in activity-
291 pattern space that is centered on the origin.

292 Rather than using the concatenated activity estimates from different partitions, encoding
293 analysis and PCM can also be applied directly to time series data. As a universal notation that
294 encompasses both situations, we can use a standard linear mixed model [35]:

$$295 \quad \mathbf{Y} = \mathbf{Z}\mathbf{U} + \mathbf{X}\mathbf{B} + \mathbf{E} , \quad (\text{Eq. 6})$$

296 where \mathbf{Y} is an $N \times P$ matrix of all activity measurements, \mathbf{Z} the $N \times K$ design matrix, which relates the
297 activity measurements to the K experimental conditions, and \mathbf{X} is a second design matrix for
298 nuisance variables. \mathbf{U} is the $K \times P$ matrix of activity patterns (the random effects), \mathbf{B} are the
299 regression coefficients for these nuisance variables (the fixed effects), and \mathbf{E} is the matrix of
300 measurement errors. If the data \mathbf{Y} are the concatenated activity estimates, the nuisance variables
301 typically only model the mean pattern for each run. If \mathbf{Y} consists of time-series data, the nuisance
302 variables typically capture additional effects such as time-series drifts and residual head-motion-
303 related artifacts.

304 ***Representational analysis in the context of neurophysiological recordings***

305 All three methods can also be applied to recordings of single cell activity or neurophysiological
306 potentials [21, 22]. The activity estimates can then be firing rates estimated over a temporal window
307 for each trial, or other multivariate summaries of the spatiotemporal activity patterns measured, e.g.
308 the time-frequency response pattern. Because the trial-by-trial variability of firing rates will usually
309 increase with the mean firing rate, it is advisable to use the square root of firing rates to make the
310 data conform better to the assumption that the variance of the noise is independent of the signal
311 [36].

312 Here we focus on models that treat the activity patterns \mathbf{U} as static snapshots. To exploit the
313 temporal detail provided by electrophysiological recordings, the analyses can be either performed

314 using a sliding window over the time course of the trial [37-39], or by “stacking” the time series and
315 conditions, resulting in T (timepoints) \times K (condition) by P (neurons) activity matrix [e.g., 40].

316 ***Encoding analysis***

317 An encoding model characterizes the structure of the representation in terms of a set of features [10-
318 13]. The value of each feature for each experimental condition is coded in the *model matrix* \mathbf{M} (K
319 conditions by Q features). The *feature weight matrix* \mathbf{W} (Q features by P channels) then determines
320 how the different model features contribute to the activity profiles of different measurement
321 channels to produce the predicted activity patterns \mathbf{U} :

$$322 \quad \mathbf{U} = \mathbf{M}\mathbf{W}. \quad (\text{Eq. 7})$$

323 Geometrically, we can think of the features as the basis vectors of the subspace, in which the
324 activity profiles reside (Fig. 2A).

325 ***Encoding analysis without regularization***

326 To adjudicate among encoding models of different numbers of features – and hence different
327 numbers of free first-level parameters - most researchers use independent test sets [11-13]. A
328 training data set is used to estimate the feature weights for each channel, and the resulting
329 prediction is then evaluated on a held-out test data set. This can be implemented in a statistically
330 efficient manner by using cross-validation, which is usually performed by holding out a single
331 partition (e.g. fMRI imaging run) as a test set, and using the remaining $M-1$ partitions as the training
332 set. Each partition is held out as the test set once and prediction performance is averaged across the
333 M folds of cross-validation. Encoding models can also make predictions about conditions that are
334 not in the training set (Discussion). However, we focus our simulations on cases, in which training
335 and test sets include the same experimental conditions.

336 The weights can be chosen to minimize the sum of squared errors on the training data, i.e.
337 using linear regression:

$$338 \quad \hat{\mathbf{W}} = (\mathbf{M}^T \mathbf{M})^{-1} \mathbf{M}^T \hat{\mathbf{U}}^{(\sim m)}, \quad (\text{Eq. 8})$$

339 where we define $\hat{\mathbf{U}}^{(\sim m)}$ to be the average activity estimates from all partitions except m . The
 340 prediction for the left-out test data of run m is

$$341 \quad \tilde{\mathbf{U}}^{(\sim m)} = \mathbf{M}\hat{\mathbf{W}}. \quad (\text{Eq. 9})$$

342 The accuracy of the prediction can be assessed by relating the residual sums-of-squares (SSR) of
 343 the prediction to the total sums-of-squares (SST) of the observed activities, summed over all
 344 partitions, conditions, and voxels

$$345 \quad R_{cv}^2 = 1 - \frac{SSR}{SST} = 1 - \frac{\sum_{m,i,j} [\hat{\mathbf{U}}_{i,j}^{(m)} - \tilde{\mathbf{U}}_{i,j}^{(\sim m)}]^2}{\sum_{m,i,j} \hat{\mathbf{U}}_{i,j}^{(m)2}}. \quad (\text{Eq. 10})$$

346 Alternatively, we can evaluate the prediction by correlating the predicted and observed activity
 347 patterns across all conditions and channels. Assuming that the mean of each channel across all
 348 conditions is zero (given mean pattern subtraction), the correlation is given by

$$349 \quad r = \frac{\sum_{m,i,j} \hat{\mathbf{U}}_{i,j}^{(m)} \cdot \tilde{\mathbf{U}}_{i,j}^{(\sim m)}}{\sqrt{\sum_{m,i,j} \hat{\mathbf{U}}_{i,j}^{(m)2} \cdot \sum_{m,i,j} \tilde{\mathbf{U}}_{i,j}^{(\sim m)2}}} \quad (\text{Eq. 11})$$

350 The correlation introduces an arbitrary scaling factor between prediction and observations and, in
 351 contrast to Eq. 10, allows the model to over- or under-predict the data by a scalar factor without
 352 penalty. Encoding analysis can also be applied directly to the time-series data without an
 353 intervening model (Eq. 6). In this case, the design matrix for the estimation of the weights (Eq. 8)
 354 and for the prediction of the left-out data (Eq. 9) becomes the product of the original first-level
 355 design matrix \mathbf{Z} and the model feature matrix \mathbf{M} .

356 To understand how encoding analysis adjudicates between models, consider the graphical
 357 representation of the estimation process for a single measurement channel in Figure 3. The training
 358 data set (black cross) is the activity profile of the measurement channel, which can be visualized as
 359 a point in activity-profile space. Regression analysis can be understood as the orthogonal projection
 360 of the measured activity profile onto the linear subspace spanned by the features of the model. The
 361 two models depicted in Fig. 3A and Fig. 3B have different features (blue arrows) that define
 362 different subspaces (planes with blue outlines). Therefore, the training data is projected onto two

363 different planes and the prediction for the test data differs between the two models. The model with
364 a subspace that better describes the cloud of activation profiles will make better predictions overall
365 across the measurement channels, show lower cross-validation error, and will hence be more likely
366 selected as the winning model.

367 Importantly, encoding analysis without regularization compares the subspaces of the
368 competing models, but not their probability distributions. For example, the model depicted in Fig.
369 3C predicts a different distribution than the one in Fig. 3A. The features of these two models,
370 however, span the same subspace. Therefore, without regularization, the predictions of these two
371 models are identical (black dots) and the models indistinguishable.

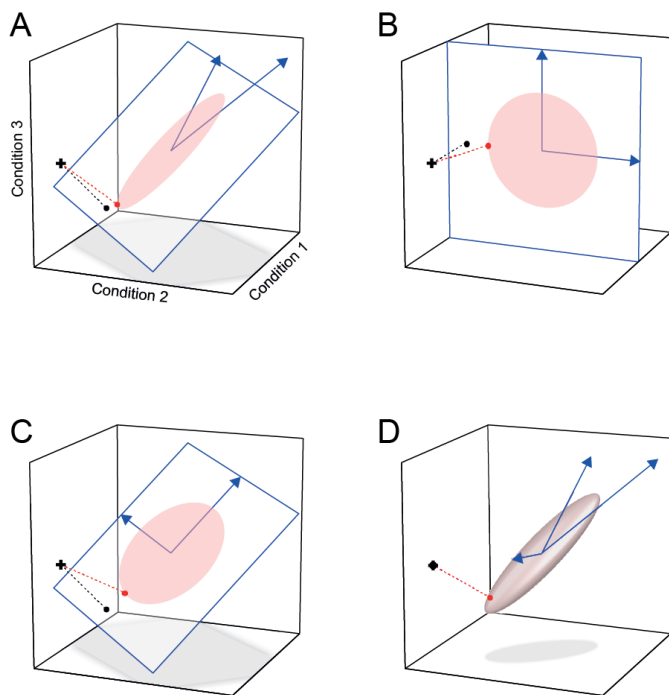


Figure 3. Adjudicating between encoding models. The axes of the three-dimensional space are formed by the response to three experimental conditions. The activity profile of each unit defines a point in this space. Models are defined by their features (blue arrows) and a prior distribution of the weights for these features. The features and the prior, together, define

382 a distribution of activity profiles (ellipsoids indicate an iso-probability-density contours of the
383 Gaussian distributions). To predict the activity profile of a single measurement channel, the
384 model is fitted to the training data set (cross). Simple regression finds the shortest projection
385 (black dot) onto the subspace defined by the features, whereas regression with regularization
386 (red dot) biases the prediction towards the predicted distribution. Two models (**A, B**) with
387 features that span different model subspaces are distinguishable using regression without
388 regularization. (**C**) This model spans the same subspace as model A. Unregularized

389 *regression results in the same projection as for model A, whereas regression with*
390 *regularization leads to a different projection. (D) A saturated model with as many features as*
391 *conditions. Unregularized regression can perfectly fit any data point (cross and black dot*
392 *coincide). With regularization, the prediction is biased towards the predicted distribution*
393 *(iso-probability-density ellipsoid).*

394 ***Encoding analysis with regularization***

395 When using regularized regression, encoding analysis evaluates models according to their
396 predicted distribution of activity profiles. From a Bayesian perspective, regularization can be
397 motivated by assuming a prior probability distribution on the weight vectors $\mathbf{w}_{\cdot,i}$ the columns of \mathbf{W} .
398 Specifically, L2-norm (Tikhonov) regularization can be derived by assuming a multivariate
399 Gaussian prior with zero mean and variance-covariance matrix $\mathbf{\Omega}$. Under this assumption, the
400 predicted second moment of the activity profiles is given by

$$\begin{aligned} 401 \quad \mathbf{G} &= \mathbf{M}\mathbf{W}\mathbf{W}^T\mathbf{M}^T/P \\ 402 \quad &= \mathbf{M}\mathbf{\Omega}\mathbf{M}^T. \end{aligned} \quad (\text{Eq. 12})$$

403 Thus, the model features together with the prior distributional assumption on the feature weights
404 define a probability distribution over activity profiles. For example, a representational model of
405 motor cortical activity could be defined by assuming that the features are individual units with
406 cosine-tuning for different movement directions [14], and that (as a prior) the preferred directions of
407 the units are uniformly distributed.

408 In practice, we allow a scalar factor, s , between the predicted and measured second moment.
409 This accounts for the fact that different subjects or regions will have different signal levels and that
410 hence the distribution of activity profiles have different widths. Under the assumption that the
411 feature weights come from a multivariate Gaussian distribution with variance $\mathbf{\Omega} \cdot s$, the best linear
412 unbiased predictor [BLUP, 41], i.e. the predictor that minimizes the squared error on the held-out
413 data is:

414
$$\widehat{\mathbf{W}} = (\mathbf{M}^T \mathbf{M} + \boldsymbol{\Omega}^{-1} s^{-1} \sigma_\epsilon^2)^{-1} \mathbf{M}^T \widehat{\mathbf{U}}^{(\sim m)}, \quad (\text{Eq. 13})$$

415 where σ_ϵ^2 is the noise variance on the observations [42]. The strength of regularization is
416 determined by the ratio of this noise variance and the variance of the signal $\boldsymbol{\Omega} \cdot s$, consistent with
417 Bayesian inference of the weights on the basis of the prior and the data.

418 After assuming a prior on the model weights, the two models depicted in Fig. 3A and 3C
419 predict different distributions of the activation profiles. When estimating the weights (Eq. 13), the
420 activity profiles are projected onto the space spanned by \mathbf{M} , but this time biased (red dot) towards
421 the denser part of the model-predicted distribution of activity profiles. As a result, the two models
422 make different predictions. An accurate prior will help the model generalize to the held-out data; an
423 inaccurate prior will hurt generalization performance. The model with the distribution that is closest
424 to the true distribution of activity profiles will yield the best cross-validation performance (as
425 measured by R^2 or r). When using regularized regression (Eq. 13), models can also have as many
426 features as conditions (Fig. 3D), or even more features than conditions. When using unregularized
427 regression, such *saturated* models are indistinguishable from each other. They become distinct only
428 after adding weight-distributional priors.

429 Because regularization is equivalent to imposing a prior on the feature weights, it is not just a
430 technical trick for estimation. Instead the prior is an integral part of the hypothesis being evaluated
431 as it co-determines (together with the features) the probability distribution over activity profiles that
432 the model predicts. Therefore, we will refer to encoding models evaluated using regularized
433 regression analysis in the following as “encoding models with a prior”.

434 One important consequence of Eq. 12 is that the same representational model can be defined
435 using different feature sets. Because a representational model is defined by its second moment, two
436 sets of features \mathbf{M}_1 and \mathbf{M}_2 , combined with corresponding second moment matrices of the weights,
437 $\boldsymbol{\Omega}_1$ and $\boldsymbol{\Omega}_2$, define the same representational model, if

438
$$\mathbf{G} = \mathbf{M}_1 \boldsymbol{\Omega}_1 \mathbf{M}_1^T = \mathbf{M}_2 \boldsymbol{\Omega}_2 \mathbf{M}_2^T. \quad (\text{Eq. 14})$$

439 Thus, an important caveat when using encoding models is that one does not compare different
440 feature sets per se – but rather different distributions (when using regularization) or different
441 subspaces of activity profiles (when not using regularization). The winning model in either case can
442 be equivalently re-expressed using a different feature set. Interpretation, therefore, must consider
443 the model-predicted distributions or subspaces of activity profiles, not the particular feature basis
444 set chosen (as the latter is not unique for any given representational model).

445 Technically, this also means that regression with a Gaussian prior can be implemented using
446 ridge regression [42]. The equivalence is established by scaling and rotating the model matrix \mathbf{M} in
447 such a way that $\mathbf{\Omega}$ becomes the identity matrix. Any representational model can be brought into this
448 diagonal form by setting the columns of \mathbf{M} to the eigenvectors of \mathbf{G} , each one multiplied by the
449 square root of the corresponding eigenvalue:

$$450 \quad \mathbf{M} = [\mathbf{v}_1\sqrt{\lambda_1} \quad \dots \quad \mathbf{v}_2\sqrt{\lambda_2}] \quad (\text{Eq. 15})$$

$$451 \quad \mathbf{G} = \mathbf{M}\mathbf{M}^T.$$

452 The strength of the regularization is determined by a scalar ridge coefficient defined by $s^{-1}\sigma_\epsilon^2$. For
453 an encoding model with regularization, the ridge coefficient still needs to be determined for each
454 cross-validation fold. This can be done again by nested cross-validation [12], generalized cross-
455 validation [43], or restricted maximum-likelihood estimation (Eq. 18). To save time, it is also
456 possible to use a constant regularization coefficient. For our simulations, we estimated the optimal
457 $s^{-1}\sigma_\epsilon^2$ by maximizing Eq. 18 for the training set (across all voxels). Generalized cross-validation
458 [43] yielded very similar results.

459 ***Pattern component modeling***

460 An alternative to cross-validation is to evaluate the likelihood of the measured activity profiles
461 under the representational model. This approach is taken in pattern-component modeling [18]. We
462 start with a generative model of the activity profiles (Eq. 6). We consider the activity profiles
463 (columns of \mathbf{U}) to come from a multivariate Gaussian distribution with zero mean and second-
464 moment matrix \mathbf{G} . To account for other nuisance effects (mean activity for each partition, low-

465 frequency drift, etc), the model also has some fixed-effects regressors (\mathbf{B}). We are not interested in
 466 fitting \mathbf{U} per se, but simply want to evaluate the likelihood of the data under different models,
 467 marginalized over all possible values of \mathbf{U} . The marginal distribution for each channel (columns of
 468 matrix \mathbf{Y}) takes the form of a multivariate normal:

$$\begin{aligned}
 469 \quad & \mathbf{y}_{\cdot,j} \sim \mathcal{N}(\mathbf{X}\mathbf{b}_{\cdot,j}, \mathbf{V}(\boldsymbol{\theta})) \\
 470 \quad & \mathbf{V}(\boldsymbol{\theta}) = \mathbf{Z}\mathbf{G}s\mathbf{Z}^T + \mathbf{I}\sigma_{\epsilon}^2 \quad (\text{Eq. 16}) \\
 471 \quad & \boldsymbol{\theta} = \{s, \sigma_{\epsilon}^2\}.
 \end{aligned}$$

472 The predicted covariance matrix of the activity measurements for each person is the function of the
 473 model (as encoded in the second-moment matrix \mathbf{G}) and two second-level parameters ($\boldsymbol{\theta}$): one that
 474 determines the strength of the signal (s) and one that determines the variance of the noise (σ_{ϵ}^2). In
 475 determining the likelihood, we remove the fixed effects using the residual forming matrix

$$476 \quad \mathbf{R} = \mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1}. \quad (\text{Eq. 17})$$

477 We need to then account for the removal of these fixed effects by evaluating the restricted
 478 likelihood $l(\mathbf{Y}|\mathbf{G},\boldsymbol{\theta})$ [44]:

$$\begin{aligned}
 479 \quad & l(\mathbf{Y}|\mathbf{G}, \boldsymbol{\theta}) = -\frac{NP}{2} \log(2\pi) - \frac{P}{2} \log|\mathbf{V}| \\
 480 \quad & -\frac{1}{2} \text{trace}(\mathbf{Y}^T\mathbf{R}^T\mathbf{V}^{-1}\mathbf{R}\mathbf{Y}) - \frac{P}{2} \log|\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X}|. \quad (\text{Eq. 18})
 \end{aligned}$$

481 To evaluate the fit of a model, the scaling and noise parameters need to be determined. For fMRI
 482 data, these two parameters can vary widely between different brain regions and individuals, and are
 483 not meaningful in themselves. We therefore replace $\boldsymbol{\theta}$ with point estimates that maximize Eq. 18 –
 484 i.e., the approach uses Empirical Bayes, or Type-II maximum likelihood for model comparison
 485 [42]. Because every model has the same two free second-level parameters, even models that are
 486 based on different numbers of features can be compared directly. An efficient implementation of
 487 this algorithm can be found in the open-source Matlab package for PCM [45].

488 ***Representational similarity analysis***

489 ***Relationship between representational dissimilarities and second-moment matrices***

490 In RSA, representational models are conceptualized in terms of the dissimilarities between the
 491 activity patterns elicited across channels by the experimental conditions (Fig. 2C). One important
 492 dissimilarity measure is the Euclidean distance, which is closely related to the second-moment
 493 matrix \mathbf{G} . The squared Euclidean distance between the true activity patterns for condition i and k
 494 (normalized by the number of measurement channels) is

495
$$d_{i,k} = (\mathbf{u}_{i,\cdot} - \mathbf{u}_{k,\cdot})(\mathbf{u}_{i,\cdot} - \mathbf{u}_{k,\cdot})^T / P = \mathbf{G}_{i,i} - 2\mathbf{G}_{i,k} + \mathbf{G}_{k,k} . \quad (\text{Eq. 19})$$

496 The Euclidean distance matrix is therefore a function the second moment of the activity profiles.
 497 The generalization of the Euclidean distances to non-isotropic noise is the Mahalanobis distance
 498 (see below). Correlation distances, another class of popular dissimilarity measures, can also be
 499 computed from the second-moment matrix. The cosine angle distance is defined as

500
$$d_{i,k} = 1 - \frac{\mathbf{u}_k \mathbf{u}_i^T}{\sqrt{\mathbf{u}_i \mathbf{u}_i^T \mathbf{u}_k \mathbf{u}_k^T}} = 1 - \frac{\mathbf{G}_{k,i}}{\sqrt{\mathbf{G}_{i,i} \mathbf{G}_{k,k}}} . \quad (\text{Eq. 20})$$

501 Here we focus on Euclidean and Mahalanobis distances, as they are independent of the resting
 502 baseline and generally easier to interpret [29].

503 In the following, we either represent these distances as a $K \times K$ representational dissimilarity
 504 matrix (RDM) \mathbf{D} , or a $K(K-1)/2$ vector \mathbf{d} that contains all unique pairwise dissimilarities (the lower
 505 triangular entries of \mathbf{D}). The vector of all pairwise dissimilarities can be obtained from \mathbf{G} by
 506 defining a contrast matrix \mathbf{C} , with each row encoding one of the pairwise contrasts, with a 1 and a
 507 -1 for the contrasted conditions and zeros elsewhere:

508
$$\mathbf{d} = \text{diag}(\mathbf{C}\mathbf{G}\mathbf{C}^T) . \quad (\text{Eq. 21})$$

509 The distances contain the same information as the second moment matrix – however, we are losing
 510 the distance of each pattern to the baseline, which was encoded on the diagonal of \mathbf{G} . Thus, in order
 511 to go from a distance matrix to a second-moment matrix, we need to re-set the origin of the
 512 coordinate system. An obvious choice is to define the mean activity pattern across all conditions to

513 be the baseline. This is equivalent making the sum of all rows and columns of \mathbf{G} zero, which can be
514 achieved by defining the centering matrix $\mathbf{H} = \mathbf{I}_K - \mathbf{1}_K/K$, with $\mathbf{1}_K$ being a square matrix of ones.
515 Under these conditions, \mathbf{G} can be computed from \mathbf{D} as

$$516 \quad \mathbf{G} = -\frac{1}{2}\mathbf{H}\mathbf{D}\mathbf{H}. \quad (\text{Eq. 22})$$

517 This yields the \mathbf{G} that would be obtained if the patterns in \mathbf{U} were centered about the origin, as can
518 be achieved by subtracting the mean pattern from each pattern.

519 *Multivariate noise normalization and cross-validation: the crossnobis estimator*

520 A particularly useful dissimilarity measure is the cross-validated, squared Mahalanobis distance
521 estimator (or crossnobis estimator for short). This estimator has superior characteristics in terms of
522 reliability and interpretability as compared to other dissimilarity measures [29].

523 The crossnobis estimator uses multivariate noise normalization (see section Spatial
524 dependence) to make the errors of different measurement channels approximately independent of
525 each other. Euclidean distances (Eq. 19) computed on these pre-whitened activity estimates are
526 equivalent to the Mahalanobis distance defined by the error-covariance matrix between channels
527 (for details see [26, 29]).

528 The crossnobis estimator is cross-validated to yield an unbiased estimate of the Mahalanobis
529 distance (assuming that the error covariance is correctly estimated). Conventional distances, which
530 are non-negative by definition, are positively biased when estimated on noisy data: When one
531 replaces the true activity patterns in Eq. 19 with their noisy estimates, the expected value of the
532 Euclidean distance will be always higher than the true distances, because the noise terms are
533 squared and summed. We can obtain an unbiased estimate of the true distance by computing the
534 difference vectors between the two activity patterns from two independent data partitions and taking
535 the inner product of the difference vectors. Thus, if we have M independent partitions, the
536 crossnobis estimator can be computed using a leave-one-out cross-validation scheme:

$$537 \quad d_{i,k} = 1/M \sum_{m=1}^M \left(\hat{\mathbf{u}}_{i,\cdot}^{(m)} - \hat{\mathbf{u}}_{k,\cdot}^{(m)} \right) \left(\hat{\mathbf{u}}_{i,\cdot}^{(\sim m)} - \hat{\mathbf{u}}_{k,\cdot}^{(\sim m)} \right)^T / P, \quad (\text{Eq. 23})$$

538 where $\hat{\mathbf{u}}_{i,\cdot}^{(m)}$ is the prewhitened pattern for condition i measured on partition m , and $\hat{\mathbf{u}}_{i,\cdot}^{(\sim m)}$ is same
539 activity pattern determined from the data of all other partitions. The expected value of this estimator
540 matches the true Mahalanobis distance [26, 29] (except for a multiplicative bias caused by
541 inaccuracies of the error covariance). In particular, if the patterns of two conditions only differ by
542 noise, then the expected value of this measure will be zero. We will see below that the interpretable
543 zero point can be advantageous for adjudicating among representational models.

544 ***Model comparison***

545 In RSA, different representational models are evaluated by comparing the predicted to the
546 observed dissimilarities. The overall magnitude of the Mahalanobis distances can vary considerably
547 between subjects. The inter-subject variation is caused by differences in physiological
548 responsiveness, physiological noise, and head movements – in short, by all the factors contributing
549 to signal strength or the noise distribution, by which the Mahalanobis distance is scaled. Therefore,
550 it is advisable to introduce a subject-specific scaling factor between observed and predicted
551 distances, relying on the ratios between distances to distinguish models.

552 The unknown scaling of the observed dissimilarities is usually accounted for by calculating
553 the correlation between the predicted and observed representational dissimilarity vectors (not to be
554 confused with the use of correlation distance as an activity-pattern dissimilarity measure, Eq. 20).

555 The most cautious approach is to assume that we can only predict the rank ordering of
556 distances [22]. It is then only appropriate to use Spearman correlation, or (in the case any of the
557 models predict equal ranks for different pairs of conditions) Kendall's τ_a [24]. Evaluating models
558 based on their ordinal dissimilarity predictions is conservative in terms of assumptions. However,
559 the lesser reliance on assumptions comes at the cost of reduced sensitivity to certain differences
560 between models. For more quantitative models, it may be appropriate to assume that distance
561 predictions can be made on an interval scale. The assumption of a linear relationship between the
562 predicted and measured distances motivates the use of Pearson correlation [25]. It may be justifiable
563 in certain cases and can increase our sensitivity to differences between representational models.

564 Both rank-based and linear correlation coefficients not only allow an arbitrary scaling factor
565 between observed and predicted distances, but also an arbitrary additive constant due the intercept
566 of regression. However, the crossnobis estimator has an interpretable zero point: If a model predicts
567 a zero distance for two conditions, then a brain region explained by the model should not be
568 sensitive to the difference between the two conditions. This is a very meaningful prediction, which
569 we can exploit to discriminate among models. Pearson and rank-based correlation coefficients
570 discard this information. This suggest the use of a normalized inner product, a quantity analogous to
571 a correlation coefficient, but in which the predictions and the data are not centered about their
572 mean:

$$573 \quad r_n = \mathbf{d}^T \tilde{\mathbf{d}} / \sqrt{\tilde{\mathbf{d}}^T \tilde{\mathbf{d}} \mathbf{d}^T \mathbf{d}} \quad (\text{Eq. 24})$$

574 This amounts to a linear regression model between the predicted and observed distances, where the
575 regression line is constrained to pass through the origin [46]:

$$576 \quad \mathbf{d} = \tilde{\mathbf{d}} \cdot s. \quad (\text{Eq. 25})$$

577 Here s is a scaling factor that is estimated from the data by minimizing the sum-of-squared errors
578 between predicted and observed values.

579 Eq. 24 would provide optimal inference, if all distances estimates were independent and of equal
580 variance. However, for the crossnobis estimator (and for most other dissimilarity measures), the
581 assumptions of independence and equal variance are both violated. Estimated squared distances
582 with larger true values are estimated with higher variability. Furthermore, the estimated distance
583 between conditions A and B is not independent from the estimated distances between A and C [26].
584 To account for these factors, we need to know the predicted probability distribution of
585 representational dissimilarity matrix estimates given a model. While the exact distribution of the
586 vector of $K(K-1)/2$ crossnobis estimates is difficult to obtain, we have shown that their distribution
587 is well approximated by a multivariate normal distribution [26]:

$$588 \quad \mathbf{d} \sim N(\tilde{\mathbf{d}}s, \mathcal{S}(\tilde{\mathbf{d}}s)). \quad (\text{Eq. 26})$$

589 The mean of the distribution is the true distance matrix, scaled by a parameter relating to the
 590 signal strength in this subject (s). In [26], we showed that that the variance-covariance matrix of \mathbf{d} is
 591 given by

$$592 \quad \mathbf{S}(\mathbf{G}, s, \mathbf{\Sigma}_K, \mathbf{\Sigma}_P) = \left[4 \frac{[\mathbf{sCGC}^T] \circ [\mathbf{C}\mathbf{\Sigma}_K\mathbf{C}^T]}{M} + 2 \frac{[\mathbf{C}\mathbf{\Sigma}_K\mathbf{C}^T] \circ [\mathbf{C}\mathbf{\Sigma}_K\mathbf{C}^T]}{M(M-1)} \right] \cdot \frac{\text{trace}(\mathbf{\Sigma}_P\mathbf{\Sigma}_P)}{P^2}. \quad (\text{Eq. 27})$$

593 Where \mathbf{G} is the predicted second-moment matrix of the patterns, \mathbf{C} the contrast matrix that
 594 transforms the second-moment matrix into distances, and \circ refers to the element-by-element
 595 multiplication of two matrices. $\mathbf{\Sigma}_K$ is the condition-by-condition covariance matrix of the estimates
 596 of the activation profiles across partitions, which can be estimated from the variability of the
 597 activity patterns around their mean ($\bar{\mathbf{U}}$):

$$598 \quad \hat{\mathbf{\Sigma}}_K = \frac{1}{M-1} \sum_m (\hat{\mathbf{U}}^{(m)} - \bar{\mathbf{U}})(\hat{\mathbf{U}}^{(m)} - \bar{\mathbf{U}})^T / P \quad (\text{Eq. 28})$$

599 $\mathbf{\Sigma}_P$ is the voxel-by-voxel correlation matrix of the activation estimates. If multivariate noise-
 600 normalization [29] was completely successful, then this would be the identity matrix. However,
 601 given the shrinkage of the noise-covariance matrix used for noise-normalization, some residual
 602 correlations will remain; for accurate predictions of the variance, these must be estimated and
 603 accounted for [26].

604 Based on this approximation we can now express the log-likelihood of the measured distances \mathbf{d}
 605 under the model predictions $\tilde{\mathbf{d}}$.

$$606 \quad l(\mathbf{d} | \tilde{\mathbf{d}} \cdot s) = -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{S}(\tilde{\mathbf{d}}s)| - \frac{1}{2} (\mathbf{d} - \tilde{\mathbf{d}}s)^T \mathbf{S}(\tilde{\mathbf{d}}s)^{-1} (\mathbf{d} - \tilde{\mathbf{d}}s) \quad (\text{Eq. 29})$$

607 To evaluate the likelihood, we first need to estimate the scaling coefficient between predicted and
 608 observed distances by choosing s to maximize the likelihood. This can be done efficiently using
 609 iteratively-reweighted least squares (IRLS): Given a starting estimate of \mathbf{S} , we can obtain the
 610 generalized least squares estimate of s ,

$$611 \quad s = (\tilde{\mathbf{d}}^T \mathbf{S}^{-1} \tilde{\mathbf{d}})^{-1} \tilde{\mathbf{d}}^T \mathbf{S}^{-1} \mathbf{d}, \quad (\text{Eq. 30})$$

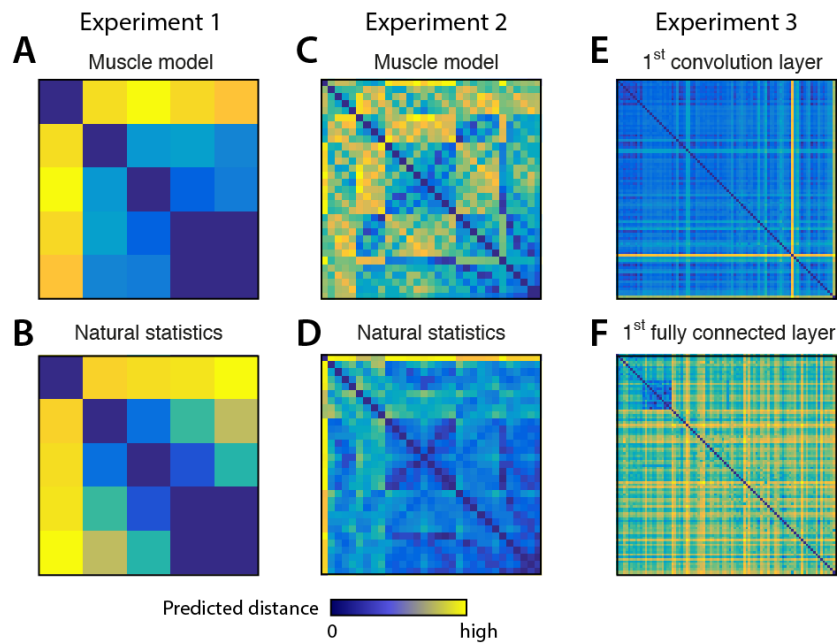
612 re-estimate \mathbf{S} according to Eq. 27, and iterate until convergence.

613 *Simulated example data sets*

614 We use simulated data sets here to evaluate and compare the three analysis techniques in a situation
615 where the ground-truth is known. The three simulated example data sets are inspired by real fMRI
616 studies. The first two examples are motivated by a paper investigating the representational structure
617 of finger movements in primary motor and sensory cortex [25]. The structure of the empirically
618 measured distances between movements of the five fingers was highly reliable across different
619 individuals. The main question was whether this invariant structure is best explained by the
620 correlations of finger movements in every-day life – i.e. the natural statistics of movement [47], or
621 by the patterns of muscle activity required for these movements. Rather than hypothesizing that
622 certain features form the basis set generating the activity profiles distribution, we could directly
623 predict the second-moment matrices, and hence the RDMs, from the correlations between naturally
624 occurring movements, or the correlations of muscle activity patterns. The predicted RDM for
625 individuated movements of the five fingers (Exp. 1) is shown in Fig. 3A, B. The second example
626 comes from experiment 3 in the same paper, this time looking at 31 different finger movements,
627 which span the whole space of possible “piano-chord” combinations (Fig. 3C, D).

628 The third example uses an experiment investigating the response of the human inferior
629 temporal cortex to 96 images, including animate and inanimate objects [21]. The model predictions
630 are derived from a convolutional deep neural network model – with each of the 7 layers providing a
631 separate representational model [48].

632 All data sets were simulated with 8 runs, 160 voxels, and independent noise on the
633 observations. The noise variance was set to $\sigma^2 = 1$. We first normalized the model predictions,
634 such that the norm of the predicted squared Euclidean distances was 1. We then derived the second
635 moment matrix (\mathbf{G}) for each model using Eq. 22 and created true activity patterns that were
636 normally distributed with second moment $\mathbf{U}\mathbf{U}^T/P = \mathbf{G} \cdot s$. The signal-strength parameter s was
637 varied systematically starting from 0 (pure noise data).



638

639 **Figure 4. Representational dissimilarity matrices (RDMs) for the models used in**
640 **simulation.** Each entry of an RDM shows the dissimilarity between the patterns associated
641 with two experimental conditions. RDMs are symmetric about a diagonal of zeros. Note that
642 while zero is meaningfully defined (no difference between conditions), the scaling of the
643 distances is arbitrary. For Experiment 1, the distance between the activity patterns for the
644 five fingers are predicted from the structure of (A) muscle activity and (B) the natural statistics
645 of movement. In Experiment 2 (C, D) the same models predict the representational
646 dissimilarities between finger movements for 31 piano-like chords. For Experiment 3 (E, F),
647 model predictions come from the activity of the seven layers of a deep convolutional neural
648 network in response to 96 visual stimuli. The 1st convolutional layer and the 1st fully
649 connected layer are shown as examples.

650 We generated 3,000 data sets for each experiment, parameter setting, and model. Each data
651 set was generated by one model (ground truth) and was analyzed so as to infer the data-generating
652 model, using each of the inference methods. To evaluate how well the methods adjudicated between
653 the models, we compared the fit of the true model (i.e. the model that generated that particular data
654 set) with each alternative model by counting the number of instances, in which the method decided

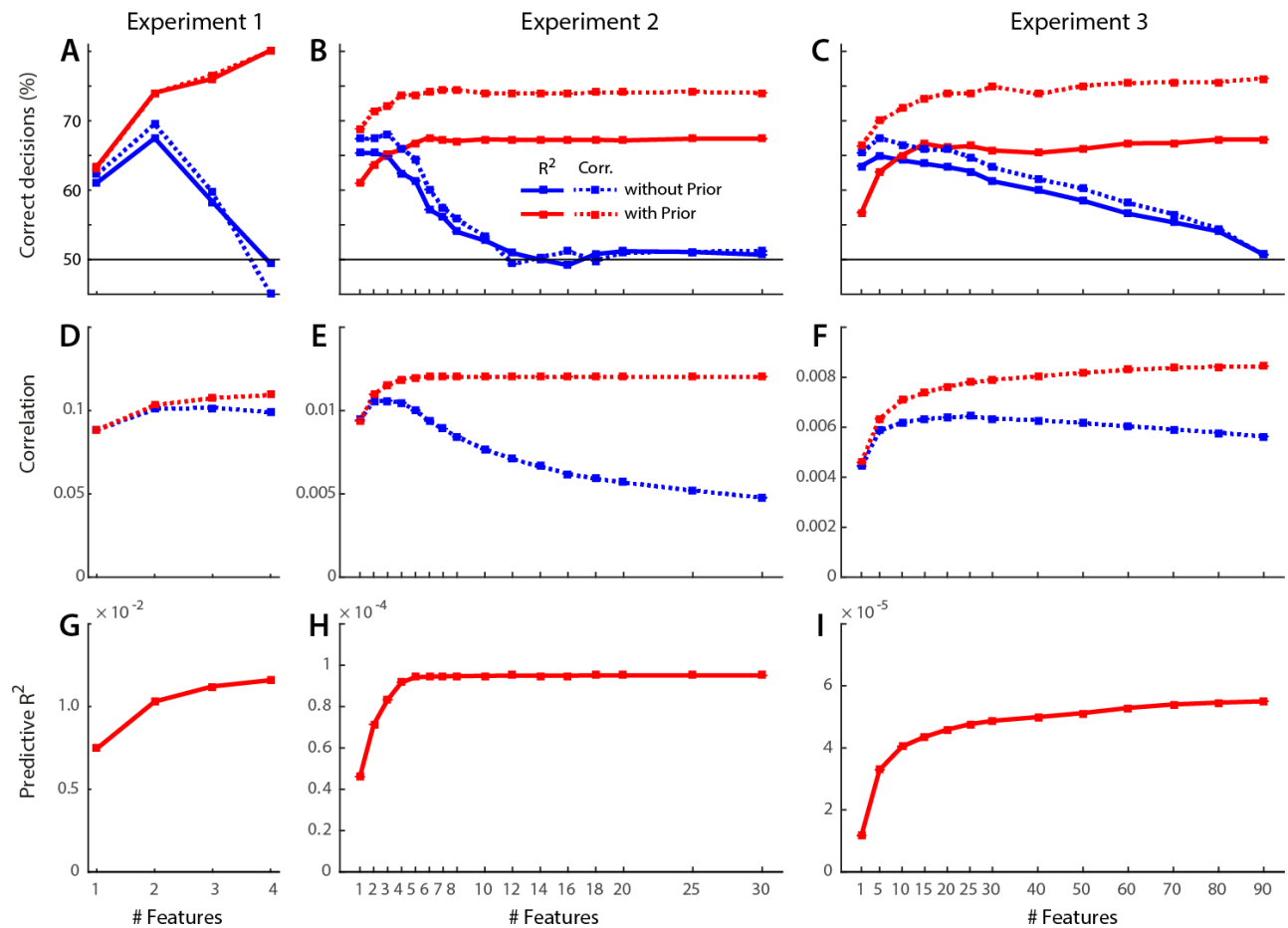
655 in favor of the correct model. Thus, even though there were 7 alternative models in Experiment 3,
656 chance performance for the pairwise comparisons was always 50%. The percentage of correct
657 decisions over all possible model pairs and simulation was used as a measure of model-selection
658 accuracy.

659 Results

660 *Encoding analysis without regularization*

661 When evaluating encoding models without using regularization, one compares the subspaces
662 spanned by the respective model features. To make different models distinguishable, one typically
663 needs to reduce the dimensionality of the model matrix \mathbf{M} , for example by using only the
664 eigenvectors with the n highest eigenvalues of the predicted second-moment matrix. The decision to
665 use a given number regressors is somewhat arbitrary: For example, Leo et al. [17] used 5
666 “synergies” (i.e. principal components of the kinematic data of 20 movements), as these explained
667 90% of the variance of the behavioral data.

668 Here we explore systematically how the number of principal components influences model
669 selection. For each experiment, we simulated data sets with a fixed signal-to-noise ratio (Exp. 1 and
670 Exp. 3: $s = 0.3$, Exp 2: $s = 0.1$; $\sigma_{\epsilon}^2 = 1$), and compared model selection accuracies using a number
671 of principal components ranging between one and the maximum number. We used both cross-
672 validated performance measures, R_{cv}^2 (Eq. 10) and r (the correlation between predicted and
673 observed values; Eq. 11) to perform model selection.



674

675 **Figure 5. Dependence of encoding model analysis on regularization and the number of**
 676 **included model features. (A-C) Percent correct model selections using either R_{cv}^2 (solid line)**
 677 **or correlation (dashed line) for encoding models without a prior (blue lines) and with a prior**
 678 **(red line). (D-F) Correlation between predicted and observed patterns. (G-I) Predictive R_{cv}^2**
 679 **for the encoding models with prior. All R_{cv}^2 values for models without prior are negative, and**
 680 **therefore not visible.**

681 Fig. 5A-C shows the percentage of correct model selections for Experiments 1-3. Results for
 682 encoding analysis without regularization are shown in blue. The dimensionality that differentiated
 683 best between competing models was 2, 3, and 5 features, respectively. As more features were
 684 included, the number of correct model selections declined. When the number of features was the
 685 same as the number of conditions minus 1 (due to the mean subtraction), i.e. the models became

686 saturated, model selection accuracy fell to chance. This is expected, as two saturated models span
687 exactly the same subspace and hence make identical predictions (Fig. 3D).

688 Using correlations as selection criterion led to more accurate decisions than using R_{cv}^2 .
689 Correlations (Fig. 5D-F, blue lines) were generally positive and peaked at a number of features that
690 was slightly higher than the optimal dimensionality for model selection. R_{cv}^2 values for encoding
691 without a prior were all negative (and are therefore not visible), because the approach does not
692 account for the noise in the data and hence leads to predictions that are too extreme – i.e. the
693 approach over-predicts the scale of the data. Correlations are insensitive to this problem as they
694 allow for arbitrary scaling between predicted and observed values.

695 ***Encoding approaches with regularization***

696 From a Bayesian perspective, employing regularization (Eq. 13) is equivalent to adding a prior to
697 the feature weights. Note that this changes the representational hypotheses tested. For example, the
698 models for Experiment 3, based on the neural network representations, now predicted not only that
699 some weighted combination of the neural network features can account for the data, but more
700 specifically that the weights should be small and of similar magnitude, preserving the
701 representational geometry of the original neural network representations. In the model matrix, we
702 scaled each principal component of \mathbf{G} with the square root of the eigenvalue (Eq. 15), such that we
703 could employ ridge regression to obtain the best linear unbiased predictor for the held-out data
704 patterns.

705 For encoding models with a prior, model selection performance increased with increasing
706 number of features (red lines, Fig. 5A-C). Thus, dimensionality reduction of the model is not
707 necessary here. Furthermore, model selection was always more powerful with than without a prior
708 when correlation was used for model selection. This reflects the fact that the prior provides
709 additional information about the models to be compared. It enables us to compare well-defined
710 distributions of activity profiles instead of just subspaces.

711 For Experiments 2 and 3, the R_{cv}^2 criterion performed substantially worse than the
 712 correlation between predicted and observed activity patterns. The difference between the two
 713 criteria arises from the fact that correlations allow for an arbitrary scaling between predicted and
 714 observed activity patterns, whereas R_{cv}^2 penalizes deviation in scale. The scaling of the prediction in
 715 turn strongly depends on the choice of the scalar regularization coefficient. This fact is illustrated in
 716 Fig. 6, where we simulated data from Exp. 2 with a fixed noise and signal strength, and varied the
 717 regularization coefficient systematically. While R_{cv}^2 is highly sensitive to the choice of the
 718 regularization coefficient, the correlation criterion is not. Because the regularization coefficient is
 719 determined separately for each cross-validation fold and model, deviations from the optimal ridge
 720 will decrease model selection accuracy for the cross-validated R_{cv}^2 criterion, but not for the
 721 correlation criterion.

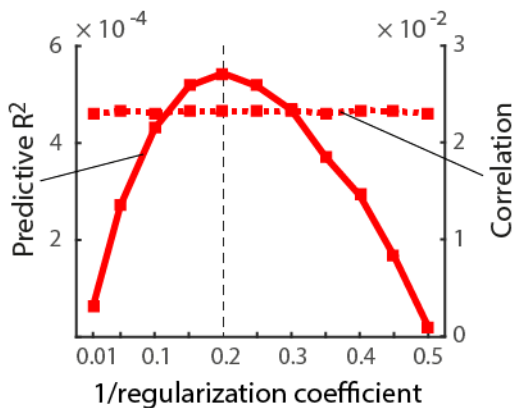


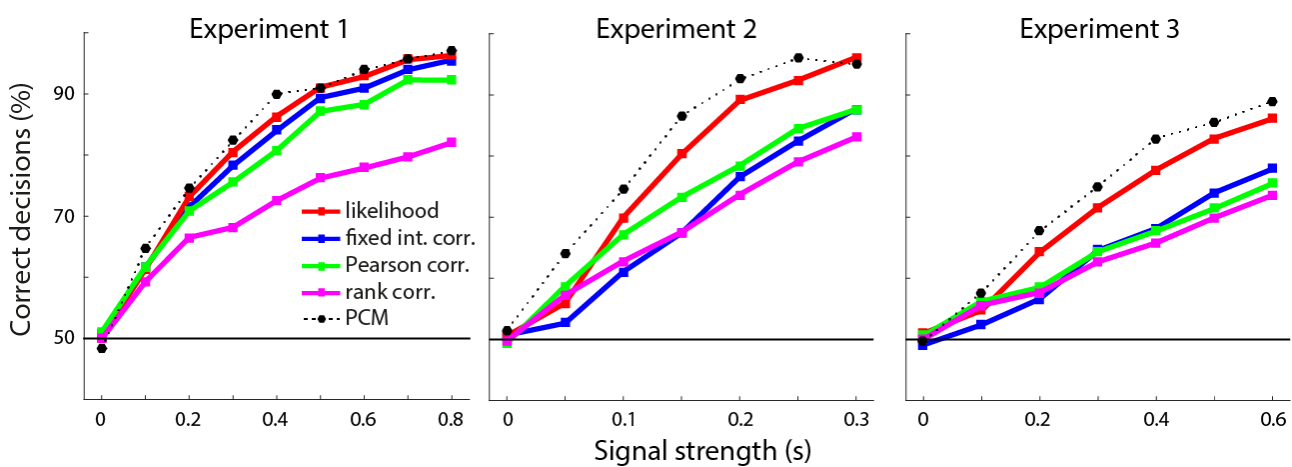
Figure 6. Sensitivity of the R_{cv}^2 (solid line) and correlation (dashed line) to the choice of the regularization coefficient. Simulations come from Experiment 2 with a true signal strength of $s=0.2$ and a noise variance of 1. For this combination, the optimal regularization coefficient is $s^{-1}\sigma_\epsilon^2$ (dashed

728 vertical line). The correlation criterion is generally robust against non-optimal choice of
 729 regularization coefficient.

730 In sum, using regularization improves model selection performance, even if the encoding
 731 model has fewer features than conditions or measurements. Rather than just comparing subspaces,
 732 the implicit prior on the weights means that a more specific hypothesis is being tested. From this
 733 perspective, it is unsurprising that we can adjudicate between these hypotheses with greater
 734 accuracy. Furthermore, the use of correlation instead of the predictive R_{cv}^2 makes model selection
 735 more robust against variations in the regularization coefficient.

736 **Representational similarity analysis**

737 When evaluating models with RSA, there is no need to restrict the model to a specific number of
 738 features – the second-moment matrix from all features can determine the predicted distances. As an
 739 empirical dissimilarity measure, we used the crossnobis estimator [29] and compared the predicted
 740 to the measured RDM. To select the winning model, we used rank-based correlation of
 741 dissimilarities [24], Pearson correlation, correlation with a fixed intercept (Eq. 24), and the
 742 likelihood of the observed distances under the normal approximation (Eq. 26) using the full
 743 variance-covariance matrix of the estimated dissimilarities.



744

745 **Figure 7. RSA model selection accuracies for different criteria of RDM fit.** Data sets for all
 746 three experiments were generated with varying signal strength (horizontal axis). The
 747 percentage of correct decisions using different criteria is shown (dotted line). Models were
 748 selected based on the Spearman rank correlation (purple), Pearson correlation (green), fixed
 749 intercept correlation (blue) or likelihood under the multinormal approximation (red). For
 750 comparison, the model selection accuracy for PCM is shown in the dotted line.

751 For Experiment 1 (Fig. 7), rank-based correlation performed substantially worse than the other
 752 criteria. The lower performance of rank correlation may have been exacerbated here by the fact that
 753 the two models predict relatively similar dissimilarity ranks. However, we expect lower
 754 performance for rank correlation in general, because this approach does not use all the information
 755 in the measured RDMs. It forgoes the assumption of a linear relationship between predicted and

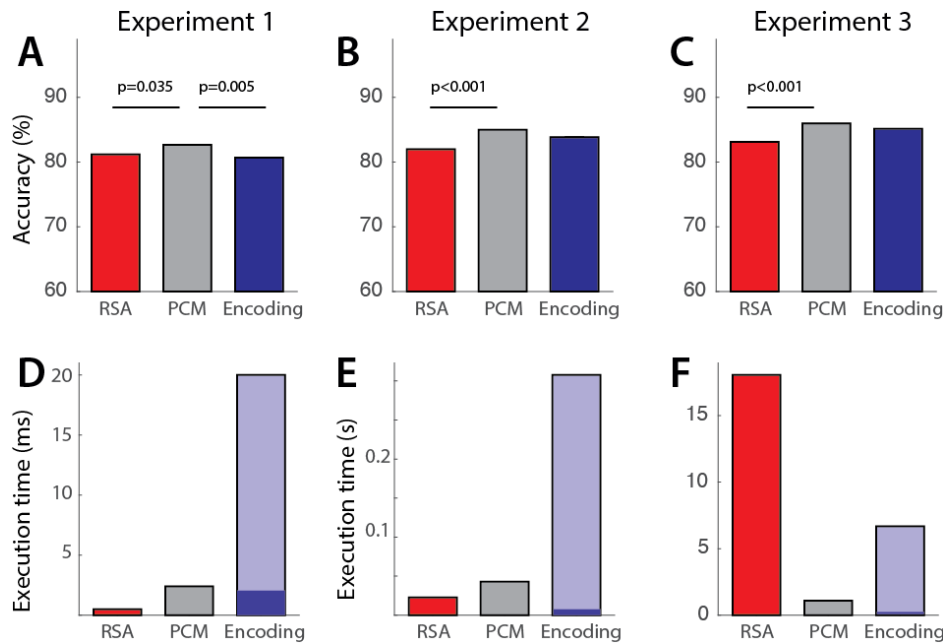
756 measured dissimilarities and therefore does not exploit the information in the continuous
757 magnitudes of the dissimilarities. Likelihood-based RSA yielded the best decisions; slightly better
758 than Pearson correlation and fixed-intercept correlation.

759 The advantage of the likelihood-based approach was clearer for Exp. 2 and 3. Here, it led to
760 about 10 percentage points greater accuracy of the decisions than the next-best RSA approach. This
761 advantage is likely due to the fact that Pearson correlations and especially fixed-intercept
762 correlations (Eq. 24) are sensitive to the observed value for the largest predicted dissimilarities, as
763 these data points have a large leverage on the estimated regression line. Indeed, some of the models
764 for Exp. 2 and 3 contain a few especially large dissimilarities, which will influence the model fit
765 strongly. The likelihood-based approach incorporates the knowledge that large dissimilarities are
766 measured with substantially larger variability [26], and hence discounts their influence. Notably,
767 rank-based correlation performed relatively well on these models as compared to Pearson
768 correlation, likely because rank correlation is robust to outliers and less dominated by the large
769 predicted distances.

770 In sum, these simulations show that it is advantageous to take the covariance structure of the
771 measured dissimilarities into account whenever the additional assumptions this requires are
772 justified.

773 ***Pattern component modeling***

774 In the same simulations, we also applied the direct likelihood-ratio test, as implemented by PCM.
775 As all the assumptions of the generative model behind PCM are met in the simulation, we would
776 expect, by the Neyman-Pearson lemma [20], that this method should provide us with highest
777 achievable model selection accuracy. Model selection performance (dotted line in Fig. 7) was
778 indeed systematically higher than for the best RSA-based method. For direct comparison of the so
779 far best methods – PCM, likelihood-based RSA, and encoding analysis with regularization (using
780 correlations as a model selection criterion) – we simulated the three Experiments at a single signal
781 strength (Fig. 8).



782

783

784

785

786

787

788

789

790

791

792

Figure 8. Model selection accuracy and execution time for likelihood-based RSA, PCM, and encoding analysis with regularization. (A-C) Model-selection accuracy was inferentially compared between the three techniques on the basis of $N=3,000$ simulations, using a likelihood-ratio test of counts of correct model decisions [49]. The signal-strength parameter for the simulation was set to $s = 0.3$ for Exp. 1, $s = 0.15$ for Exp. 2, and $s = 0.5$ for Exp. 3. All resulting significant differences (two-tailed, $p < 0.01$, uncorrected) are indicated by a horizontal line above the bars. (D-F) Execution times for the evaluation of a single data set under a single model. For encoding, the time is split into the time required to estimate regression coefficients (dark blue) and the time to determine the regularization constant (light blue).

793

794

795

796

797

In this simulation, PCM resulted in 1.48, 3.01 and 2.86 percentage points (for Exp. 1-3, respectively) better model selection accuracy than likelihood-based RSA, and 1.98, 1.17 and 0.85 percentage points higher model selection accuracies than an encoding analysis using correlations. PCM never performed worse than another method and performed significantly better than the other two approaches in 4 of 6 total comparisons across the three experiments (Fig. 8). There were no

798 significant performance differences between RSA and encoding analysis. Overall, however, all
799 three methods were very close in performance.

800 *Computational cost*

801 A practical concern is the speed at which the model comparison can be performed. This is usually
802 not important when evaluating the model fit on a small number of participants or ROIs. However, if
803 a larger number of models is evaluated continuously over the cortical surface using a searchlight
804 approach [50, 51], or in data sets with large numbers of participants, computational cost becomes a
805 practical issue. While we cannot treat this issue exhaustively, we provide here a brief overview over
806 the computation time required for the three methods for our specific examples and implementation.
807 In general, the computation time will of course depend on the number of conditions, the number of
808 channels, the exact variant of each technique. Our goal here is simply to give the reader a starting
809 point for making a choice for a particular application, trading off computational and statistical
810 efficiency.

811 Both RSA and PCM operate on the inner product matrix of the activity estimates, thus the
812 computational costs for these approaches is virtually independent of the number of voxels. PCM
813 works on the $MK \times MK$ inner product matrix of the activity estimates, whereas RSA operates on a K
814 $\times K$ matrix of distances between conditions. For a small number of conditions, this explains the
815 favorable computational cost of RSA. However, when using likelihood-based RSA, the covariance
816 matrix of the distances needs to be calculated and inverted. The size of this matrix is $(K(K-1)/2)^2$
817 and it therefore grows with the 4th power of the number of conditions K . For Exp. 3 (Fig. 8F) with K
818 = 96, this is computationally costly, whereas PCM still only needs to invert matrices of size $(MK)^2$.
819 Using RDM-correlation-based model selection (whether rank, Pearson or fixed-intercept), RSA is
820 much more computationally efficient (not shown).

821 For encoding models, conducting the actual ridge regression for each cross-validation fold
822 (dark blue area) is extremely fast and efficient. The main cost of the technique lies in the
823 determination of the optimal ridge coefficient (light blue area). In our simulations, we use restricted

824 maximum likelihood estimation (Eq. 18) to do so – therefore this cost is always M times higher than
825 for PCM alone. Depending on the implementation, generalized cross-validation [43] may offer a
826 considerable speed-up. If very high speeds are required, one could use a constant ridge coefficient
827 and accept the possible loss in model selection accuracy. In sum, while PCM is computationally
828 feasible across the three experiments, encoding models were less efficient in the present
829 implementation and likelihood-based RSA was less efficient than PCM for the condition-rich
830 scenario of Experiment 3. Alternative variants of encoding models (with fixed ridge coefficient)
831 and RSA (with correlation-based model selection) are less statistically efficient, but beat PCM in
832 terms of computational efficiency.

833 Discussion

834 In this paper, we defined representational models as formal hypotheses about the distribution of the
835 activity profiles in the space defined by the experimental conditions. That is, a representational
836 model specifies, which features are represented in a brain region, and how strongly they are
837 represented. The “strength” of representation of a feature has two aspects: the number of responses
838 (e.g. neurons) dedicated to a feature and the scaling of their activity profiles relative to the noise.
839 The second-moment matrix of the activity profiles captures the combined effect of these aspects of
840 feature strength. Two distinct representations with identical second-moment matrices therefore
841 support linear decoding of any given feature at the same signal-to-noise ratio. This holds
842 independent of the question whether the distribution of activity profiles is Gaussian. It motivates
843 using the second moment as a summary statistic for characterizing representations. RSA, PCM and
844 encoding models offer different tests of representational models, but all three depend, explicitly or
845 implicitly, on the second-moment matrix to characterize each representational hypothesis. Thus,
846 these methods are deeply related and should be understood as part of the same multivariate toolbox.
847 The main characteristics of the three methods are summarized in Table 1.

848 ***Encoding models without prior define subspaces, not distributions of activity***
849 ***profiles***

850 There is a fundamental difference between encoding models with and without weight priors.
851 Without a prior on the feature weights, encoding models test how well the subspace spanned by the
852 model features captures the observed activity profiles. For models to be discriminable, the
853 dimensionality (i.e. the number of features) of each model must be substantially lower than the
854 number of experimental conditions. As the number of model dimensions increases, the subspaces of
855 competing models increasingly overlap. Once the number of features matches the number of
856 experimental conditions, their subspaces comprise the entire space of activity profiles, each
857 perfectly fits the training data, and their predictions for unseen data become identical.

858 A subspace specifies what activity profiles are possible and what activity profiles are
859 impossible (though they might still arise as estimates because of the noise). A subspace might be
860 conceptualized as an infinite flat distribution over the subspace dimensions, with 0 probability
861 outside the subspace. However, a uniform distribution on an infinite interval has an infinite second
862 moment and hence does not specify the neural representation uniquely.

863 L2-norm regularization (i.e. ridge regression) is equivalent to imposing a Gaussian prior on
864 the regression weights. With such a prior, the representational model specifies a probability
865 distribution with a finite second moment. When changing the form of regularization, one also
866 changes the implicit prior, and hence the representational model that is being tested. Thus,
867 regularization is not simply a trick for stabilizing the fit. Instead, the weight prior forms an integral
868 part of the model, which determines the strength with which each feature is encoded according to
869 the model. Choosing a specific form of regularisation therefore constitutes a decision about the
870 neuroscientific hypothesis to be tested rather than a methodological consideration.

871 ***Encoding models tests hypotheses about activity profile distributions, not features***
872 ***sets***

873 Encoding models do not support inferences about the particular feature set generating a
874 representation, because infinitely many feature sets can span the same space. Even when using a
875 prior, the feature set that characterizes a given representational model is not unique. Features should
876 not in general be constrained to be orthogonal in the space of experimental conditions, because the
877 structure of the model is not usually meant to depend on the experimental conditions chosen.
878 Whether the features chosen are orthogonal or not, there is an infinite number of basis sets of
879 features that express the same representational model (inducing the same second moment of activity
880 profiles, Eq. 3). For example, two equally long correlated feature vectors can equally well describe
881 a distribution with elliptical isoprobability-density contours (Fig. 3A) as two orthogonal features,
882 with one vector longer than the other. Thus, when one representational model is shown to be
883 superior to others, it does not imply anything special about the feature set chosen to express that
884 model. These complications need to be kept in mind in the interpretation of the results of encoding
885 model analyses. It is very tempting to attribute meaning to the particular feature basis chosen,
886 especially when they are mapped onto the cortical surface [13, 17]. When interpreting these maps,
887 one needs to remember that a feature set only describes a distribution of activity profiles, and that
888 very different maps can emerge when the same distribution is described by a rotated set. In PCM
889 and RSA, the equivalence of different feature sets is made explicit, as they lead to the same second-
890 moment and representational dissimilarity matrices.

891 ***Likelihood-based RSA is more sensitive than correlation-based RSA***

892 When using RSA to test representational models, the crossnobis estimator provides a highly
893 reliable measure of dissimilarity with the added advantage of having an interpretable zero-point
894 [29]. Rank-based, Pearson, and fixed-intercept correlation provide fast and straightforward ways of
895 measuring the correspondence between predicted and observed distances, so as to select the
896 representational model most consistent with the data. However, using simple correlations ignores

897 the dependence of the distance estimates, as well as their unequal variances. In other words, the
898 sampling distribution of the estimated RDM in the space spanned by the dissimilarities (one
899 dimension per pair of conditions) is not isotropic. This problem is addressed in likelihood-based
900 RSA, which uses a multivariate-normal approximation to the sampling distribution of the
901 crossnobis RDM estimate [26]. The approximation provides an analytical expression for the
902 statistical dependency of distance estimates, as well as their signal-dependent variances. In the
903 simulations, likelihood-based RSA was shown to be more powerful than correlation-based RSA. Its
904 model-selection accuracy was only slightly below the theoretical upper bound, as established by
905 PCM. Likelihood-based RSA might therefore become the approach of choice when comparing
906 representational models using crossnobis estimates.

907 There are situations, however, in which the models are not specific enough to support ratio-
908 scale predictions of representational dissimilarities. Moreover, for measurement modalities like
909 fMRI, it might be undesirable to assume a linear relationship between predicted and measured
910 representational dissimilarities. Rank-correlation-based RSA [22, 24] provides a robust method that
911 is not dependent on the assumption of a linear reflection of the underlying neural dissimilarities in
912 the data RDM. It is also more computationally efficient in the context of condition-rich designs.
913 Likelihood-based RSA becomes computationally expensive as the number of conditions increases.
914 A practical compromise might be to only use the diagonal of the variance-covariance matrix, which
915 would dramatically reduce computational complexity at the expense of neglecting dependencies
916 among dissimilarity estimates.

917 ***Which method is best?***

918 For all simulations, model selection using PCM [18] was better than competing methods. This is
919 not surprising, as the data were simulated exactly according to the generative model underlying this
920 approach (Gaussian distribution of noise and signal, independence across voxels). In this case, PCM
921 implements the likelihood-ratio test, which by the Neyman-Pearson lemma [20] is the most
922 powerful test. Beyond confirming what we know from theory, the simulations were important

923 because they revealed how close the other two techniques come to the theoretical upper bound
924 established by PCM. Results showed that encoding models with a prior and likelihood-based RSA
925 perform near-optimally. In practice, we therefore expect these three techniques to provide similar
926 answers. When its assumptions hold, PCM has clear advantages for model comparison, providing
927 optimal power at reasonable computational cost. However, the other two techniques have other
928 advantages that make them attractive for specific applications.

929 RSA using RDM correlation for model selection gives up statistical efficiency for
930 computational efficiency, and beats PCM at the latter. When rank correlation is used to compare
931 RDMs, the inference does not rely on a linear relationship between the true dissimilarities and the
932 estimated dissimilarities, an assumption that might be violated in many contexts. RSA also provides
933 readily interpretable intermediate statistics (cross-validated distances), which are closely related to
934 linear decoders for all pairs of stimuli. These statistics can be used to test whether two conditions
935 have different activity patterns [24, 26], or whether the dissimilarity is larger for one pair than for
936 another pair of conditions. Multidimensional scaling of the stimuli on the basis of their
937 representational dissimilarities also provides an intuitive visualization of the representational
938 structure [22], which can be very helpful in the generation of novel representational hypotheses.

939 In contrast, PCM sometimes demands complicated approaches to answer simple questions:
940 For example, to test the hypothesis that a difference between two conditions is encoded, one would
941 need to fit one model that allows for separate patterns and one model that does not – and then
942 compare the marginal likelihood of these models. Furthermore, PCM requires the noise to be
943 explicitly modeled, whereas RSA removes the bias arising from noise through cross-validated
944 distances.

945 Encoding analysis explicitly estimates the first-level parameters that describe the response
946 for each individual voxel. This enables the mapping of the estimated features onto the cortical
947 surface to study their spatial distribution [13, 17].

948 In sum, the three methods are deeply related in that they test hypotheses about the second
949 moment of the activity profiles. However, each method constitutes a unique perspective on the data
950 and supports different kinds of exploratory analyses. We view the methods as complementary tools
951 that are part of a single coherent toolkit for analyzing representations.

952 ***Single-voxel vs. multi-voxel inference***

953 An important issue, which we have not touched upon so far, is whether to perform model
954 comparison on single or multiple voxels. While RSA and PCM are usually applied to groups of
955 voxels (such as for ROIs or searchlights), encoding models are often compared on the single-voxel
956 level. This tendency, however, is not strictly inherent in methodological constraints: The searchlight
957 approach for RSA and PCM can be reduced to single voxels, and encoding models can be combined
958 with multi-voxel searchlights. Analyses with coarser granularity give up some spatial precision of
959 the map in exchange for greater statistical power. Searchlight mapping boosts power (1) by locally
960 combining the evidence, (2) by enabling the use of a multivariate noise normalization, and (3) by
961 reducing the effective number of multiple comparisons [52]. There is no reason to assume that a
962 single-voxel searchlight is always the optimal choice when balancing spatial precision and power.
963 Based on our previous results [29], we expect that ignoring voxel dependencies will entail a loss of
964 sensitivity when making inferences on representational models for regions of interest comprising
965 multiple responses.

966 ***Testing models without overfitting to the noise and to the sample of experimental*** 967 ***conditions***

968 Whenever a model is fitted using experimental data, its parameters will necessarily be
969 overfitted to the data to some extent. Assessing the performance of a fitted model therefore requires
970 independent test data. An important question is whether the test data should consist in independent
971 measurements for the same experimental conditions or in measurements for a fresh sample of
972 experimental conditions (e.g. a different sample of visual images). The simple answer is that it
973 depends on the inference we would like to make. If our hypothesis is restricted to the present set of

974 conditions (e.g. five finger movements), we need only account for overfitting to the noise in the
975 data and require different measurements for the same conditions. If our hypothesis is about a
976 population of conditions (e.g. all natural images), we need to account for overfitting to the condition
977 sample and require measurements for an independent random sample of conditions from the
978 population of conditions covered by our hypothesis.

979 However, overfitting only needs to be accounted for when the model being tested had
980 parameters fitted in the first place. Encoding models always require independent test sets to account
981 for the over-fitting of the first-level parameters of the representational model (feature weights).
982 RSA and PCM, by contrast, rely explicitly on summary statistics of the responses. Therefore, only
983 second-level parameters related to the strength of the signal and noise need to be fit (see Table 1).
984 The representational models considered here had the same number of such second-level parameters
985 and could therefore be compared directly. Even in this context, however, cross-validation across
986 stimuli can provide valuable information about the generalizability of the model.

987 ***What about decoding approaches?***

988 Linear decodeability is central to our definition of a neural representation. Correspondingly,
989 decoding is widely used in multivariate analysis of brain imaging data [7-9]. Can it serve us also as
990 a tool for comparing representational models? While one can use standard decoding approaches to
991 determine whether specific features are represented in an area or not, it does not lend itself to the
992 comparison of full representational models (as defined here). Representational models determine
993 (via the second moment matrix) the decodability of any linear feature, and therefore constitute a
994 hypothesis-driven generalization of the concept of decoding. This is most obvious in RSA, where
995 the RDM assembles all pairwise condition discriminabilities. It is of course possible to use
996 decoding in the context of the methods considered here. For example, some studies have used
997 encoding models to decode stimuli [11, 12, 17]. Decoding accuracy can then serve, instead of
998 correlation or R_{cv}^2 , to evaluate the performance of an encoding model on held-out data. While this
999 approach is motivated by the intuitive demonstration of mind reading it does not provide a

particularly natural or powerful approach to adjudicating between representational models. Alternatively, we could use classification accuracy as a measure of dissimilarity between two conditions in the context of RSA [53]. However, classification essentially converts a continuous measure of dissimilarity into a binary label of correct / incorrect. It is therefore expected to be less informative than the underlying continuous measure, and we have shown previously that this entails a loss of sensitivity in practice [29]. In sum, decoding is not particularly useful for the evaluation of representational models [10, 19] and should therefore be limited to situations, in which the quality of the decoding itself is the measure of interest.

Flexible representational models

All models considered here were “fixed”, i.e., they did not include free parameters that would change the predicted second-moment matrix. In many applications, however, the relative importance of different features (for example encoding strength for orientation and color) are unknown. In this case, the predicted second moment can be expressed as the weighted sum of different pattern components, i.e. $\mathbf{G} = \sum_i \omega_i \mathbf{C}_i$ [18, 54-56], with the weights being free second-level parameters. In other situations, \mathbf{G} is a nonlinear function of free model parameters: For example, \mathbf{G} depends non-linearly on the spatial tuning width in population receptive field modeling [57]. Both RSA and PCM already provide a mechanism to estimate such parameters, as both approaches already need to estimate the signal strength parameters s by maximizing the respective likelihood function (Eq. 17, 28) – and the analytical derivatives of the likelihood (Eq. 17, 28) with respect to the parameters are easily obtained. In the context of encoding approaches using ridge regression, free model parameters that change the model structure would result in independent scaling of different features, rotations, or extensions of the model matrix \mathbf{M} . At the time of writing there are no published examples of such parameter optimization in the context of cross-validated encoding models that we know of.

The inclusion of free parameters into the model also enables the specification of measurement models. Representational models ideally test hypotheses about the distribution of

026 activation profiles of the core computational elements – i.e. neurons. When using indirect measures
027 of brain activity such as fMRI or MEG, the distribution of activity profiles across measurement
028 channels is also influenced by the measurement process, which samples and mixes neuronal activity
029 signals in the measurement channels [27, 58-61]. As the underlying brain computational models
030 become more specific and detailed, the corresponding measurement models will also have to be
031 improved.

032 ***Non-Gaussian representational models***

033 A second restriction of the present exploration is that we focused on approaches that
034 characterize the distribution of activity profiles by its second moment. If the true distribution of the
035 activity profiles is a multivariate Gaussian, then this is a fully sufficient approach. However, a
036 representational hypothesis may not only predict that the response for condition A is uncorrelated to
037 the response for condition B, but, for example, that channels either respond to A or B, but not to
038 both A and B. Such tuning is for example prevalent in primary visual cortex, where neurons (and
039 voxels) respond to a stimulus in a *one* specific part of the visual field, but less often two or more
040 disparate locations [57]. This would correspond to a non-Gaussian prior on the feature weights. In a
041 recent publication, Norman-Haignere and colleagues [62] suggested a likelihood-based method, in
042 which the Gaussian prior on the feature weights \mathbf{W} is replaced with a Gamma distribution,
043 essentially providing a non-Gaussian extension of PCM. It will be interesting to determine to what
044 degree such non-Gaussian distributions are present in fMRI or single-cell data, and what
045 computational function these may play.

046 It is important to stress that the approaches considered here are still appropriate when the
047 distribution of activity profiles is truly non-Gaussian. Even in the non-Gaussian case, the second
048 moment determines the linear decodability of all possible features, and hence what is explicitly
049 represented, and thus remains essential for characterizing the representation. Taking into account
050 higher moments of the activity profile distribution would enable us to distinguish between

051 representations that afford the same linear readout of features, but achieve this by distinct
052 population codes.

053 ***Conclusions***

054 If advances in brain-activity measurements are to yield theoretical insights into brain computation,
055 they need to be complemented by analytical methods to test computational models of information
056 processing [63]. The main purpose of this paper was to provide a clear definition of one important
057 class of models – representational models – and to compare three important approaches of testing
058 these. We have shown that PCM, RSA and encoding analysis are all closely related, testing
059 hypotheses about the distribution of activity profiles. Moreover, all three approaches, in their
060 dominant implementations, are sensitive only to distinctions between representations that are
061 reflected in the second moment of the activity profiles. Thus, these three methods are properly
062 understood as components of a single analytical framework. Each of the three methods has
063 particular advantages and disadvantages and preferred areas of application.

- 064 1. PCM provides an analytic expression for the marginal likelihood of the data under the
065 model, and therefore constitutes the most powerful test for adjudicating between
066 representational models if the assumptions hold. Its analytical tractability and relative
067 computational efficiency are further attractive features, especially when considering models
068 with increasing numbers of free parameters.
- 069 2. RSA provides highly interpretable intermediate statistics and is therefore ideally suited for
070 the visualization and exploratory analysis. Furthermore, simple models are often more easily
071 tested than with PCM. The normal approximation to the distribution of estimated distances
072 enables inference that is nearly as powerful as the likelihood-ratio test provided by PCM.
073 Finally, dissimilarity-rank-based RSA, though less sensitive, provides a means of inference
074 that does not rely on the assumption of a linear relationship between predicted and measured
075 dissimilarities and is computationally efficient even for condition-rich designs.

076 3. Encoding approaches enable the voxel-wise mapping of model features onto the cortical
077 surface. They therefore are the natural choice when the spatial distribution of features or the
078 voxel-wise comparison of representational models is the main interest.

079

080 We hope that the general framework presented here will enable researchers to combine these
081 approaches to make progress revealing the computational mechanisms of biological brains.

082 Bibliography

- 083 1. Stevenson IH, Kording KP. How advances in neural recording affect data analysis. *Nat*
084 *Neurosci.* 2011;14(2):139-42. doi: 10.1038/nn.2731.
- 085 2. DiCarlo JJ, Zoccolan D, Rust NC. How does the brain solve visual object recognition?
086 *Neuron.* 2012;73(3):415-34. doi: 10.1016/j.neuron.2012.01.010.
- 087 3. deCharms RC, Zador A. Neural representation and the cortical code. *Annu Rev Neurosci.*
088 2000;23:613-47. doi: 10.1146/annurev.neuro.23.1.613.
- 089 4. Kriegeskorte N. Pattern-information analysis: from stimulus decoding to computational-
090 model testing. *Neuroimage.* 2011;56(2):411-21. doi: 10.1016/j.neuroimage.2011.01.061.
- 091 5. DiCarlo JJ, Cox DD. Untangling invariant object recognition. *Trends Cogn Sci.*
092 2007;11(8):333-41. doi: 10.1016/j.tics.2007.06.010.
- 093 6. Quiroga RQ, Reddy L, Kreiman G, Koch C, Fried I. Invariant visual representation by
094 single neurons in the human brain. *Nature.* 2005;435(7045):1102-7. doi:
095 10.1038/nature03687.
- 096 7. Haxby JV, Gobbini MI, Furey ML, Ishai A, Schouten JL, Pietrini P. Distributed and
097 overlapping representations of faces and objects in ventral temporal cortex. *Science.*
098 2001;293(5539):2425-30. Epub 2001/09/29. doi: 10.1126/science.1063736.
- 099 8. Norman KA, Polyn SM, Detre GJ, Haxby JV. Beyond mind-reading: multi-voxel pattern
100 analysis of fMRI data. *Trends Cogn Sci.* 2006;10(9):424-30. Epub 2006/08/11. doi:
101 10.1016/j.tics.2006.07.005 [doi].
- 102 9. Pereira F, Mitchell T, Botvinick M. Machine learning classifiers and fMRI: a tutorial
103 overview. *Neuroimage.* 2009;45(1 Suppl):S199-209. Epub 2008/12/17. doi:
104 10.1016/j.neuroimage.2008.11.007.
- 105 10. Naselaris T, Kay KN, Nishimoto S, Gallant JL. Encoding and decoding in fMRI.
106 *Neuroimage.* 2011;56(2):400-10. Epub 2010/08/10. doi: 10.1016/j.neuroimage.2010.07.073.
- 107 11. Mitchell TM, Shinkareva SV, Carlson A, Chang KM, Malave VL, Mason RA, et al.
108 Predicting human brain activity associated with the meanings of nouns. *Science.*
109 2008;320(5880):1191-5. doi: 10.1126/science.1152876.
- 110 12. Kay KN, Naselaris T, Prenger RJ, Gallant JL. Identifying natural images from human brain
111 activity. *Nature.* 2008;452(7185):352-5. Epub 2008/03/07. doi: 10.1038/nature06713.
- 112 13. Huth AG, de Heer WA, Griffiths TL, Theunissen FE, Gallant JL. Natural speech reveals the
113 semantic maps that tile human cerebral cortex. *Nature.* 2016;532(7600):453-8. doi:
114 10.1038/nature17637.
- 115 14. Georgopoulos AP, Schwartz AB, Kettner RE. Neuronal population coding of movement
116 direction. *Science.* 1986;233(4771):1416-9.
- 117 15. Sergio LE, Hamel-Paquet C, Kalaska JF. Motor cortex neural correlates of output
118 kinematics and kinetics during isometric-force and arm-reaching tasks. *J Neurophysiol.*
119 2005;94(4):2353-78.

- 120 16. Sergio LE, Kalaska JF. Systematic changes in directional tuning of motor cortex cell activity
121 with hand location in the workspace during generation of static isometric forces in constant
122 spatial directions. *J Neurophysiol.* 1997;78(2):1170-4. Epub 1997/08/01.
- 123 17. Leo A, Handjaras G, Bianchi M, Marino H, Gabiccini M, Guidi A, et al. A synergy-based
124 hand control is encoded in human motor cortical areas. *Elife.* 2016;5. doi:
125 10.7554/eLife.13420.
- 126 18. Diedrichsen J, Ridgway GR, Friston KJ, Wiestler T. Comparing the similarity and spatial
127 structure of neural representations: a pattern-component model. *Neuroimage.*
128 2011;55(4):1665-78. Epub 2011/01/25. doi: 10.1016/j.neuroimage.2011.01.044.
- 129 19. Friston K, Chu C, Mourao-Miranda J, Hulme O, Rees G, Penny W, et al. Bayesian decoding
130 of brain images. *Neuroimage.* 2008;39(1):181-205. Epub 2007/10/09. doi:
131 10.1016/j.neuroimage.2007.08.013.
- 132 20. Neyman J, Pearson ES. On the problem of the most efficient test of statistical hypotheses.
133 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering*
134 *Sciences.* 1933;231:289–337. doi: doi:10.1098/rsta.1933.0009.
- 135 21. Kriegeskorte N, Mur M, Ruff DA, Kiani R, Bodurka J, Esteky H, et al. Matching categorical
136 object representations in inferior temporal cortex of man and monkey. *Neuron.*
137 2008;60(6):1126-41. Epub 2008/12/27. doi: 10.1016/j.neuron.2008.10.043 [doi].
- 138 22. Kriegeskorte N, Mur M, Bandettini P. Representational similarity analysis - connecting the
139 branches of systems neuroscience. *Front Syst Neurosci.* 2008;2:4. Epub 2008/12/24. doi:
140 10.3389/neuro.06.004.2008.
- 141 23. Kriegeskorte N, Kievit RA. Representational geometry: integrating cognition, computation,
142 and the brain. *Trends Cogn Sci.* 2013;17(8):401-12. Epub 2013/07/24. doi:
143 10.1016/j.tics.2013.06.007.
- 144 24. Nili H, Wingfield C, Walther A, Su L, Marslen-Wilson W, Kriegeskorte N. A toolbox for
145 representational similarity analysis. *PLoS Comput Biol.* 2014;10(4):e1003553. Epub
146 2014/04/20. doi: 10.1371/journal.pcbi.1003553.
- 147 25. Ejaz N, Hamada M, Diedrichsen J. Hand use predicts the structure of representations in
148 sensorimotor cortex. *Nat Neurosci.* 2015;18(7):1034-40. Epub 2015/06/02. doi:
149 10.1038/nn.4038.
- 150 26. Diedrichsen J, Zareamoghaddam H, Provost S. The distribution of crossvalidated
151 mahalanobis distances. *ArXiv.* 2016.
- 152 27. Kriegeskorte N, J. D. Inferring brain-computational mechanisms with models of activity
153 measurements. *Proceedings of the Royal Society.* 2016.
- 154 28. Cai MB, Schuck NW, Pillow J, Niv Y. A Bayesian method for reducing bias in neural
155 representational similarity analysis. *BioRxiv.* 2016. doi: 10.1101/073932.
- 156 29. Walther A, Nili H, Ejaz N, Alink A, Kriegeskorte N, Diedrichsen J. Reliability of
157 dissimilarity measures for multi-voxel pattern analysis. *Neuroimage.* 2016;137:188-200.
158 doi: 10.1016/j.neuroimage.2015.12.012.
- 159 30. Ledoit O, Wolf M. Improved estimation of the covariance matrix of stock returns with an
160 application to portfolio selection. *Journal of Empirical Finance.* 2003;10(5)(603–621).
- 161 31. Misaki M, Kim Y, Bandettini PA, Kriegeskorte N. Comparison of multivariate classifiers
162 and response normalizations for pattern-information fMRI. *Neuroimage.* 2010;53(1):103-18.
163 Epub 2010/06/29. doi: 10.1016/j.neuroimage.2010.05.051.
- 164 32. Fisher RA. The use of multiple measurements in taxonomic problems. *Annals of Eugenics.*
165 1936;7(2):179-88.
- 166 33. Friston KJ, Holmes AP, Poline JB, Grasby PJ, Williams SC, Frackowiak RS, et al. Analysis
167 of fMRI time-series revisited. *Neuroimage.* 1995;2(1):45-53.
- 168 34. Worsley KJ, Friston KJ. Analysis of fMRI time-series revisited--again. *Neuroimage.*
169 1995;2(3):173-81.
- 170 35. Laird NM, Ware JH. Random-effects models for longitudinal data. *Biometrics.*
171 1982;38(4):963-74.

- 1172 36. Yu BM, Cunningham JP, Santhanam G, Ryu SI, Shenoy KV, Sahani M. Gaussian-process
1173 factor analysis for low-dimensional single-trial analysis of neural population activity. *J*
1174 *Neurophysiol.* 2009;102(1):614-35. doi: 10.1152/jn.90941.2008.
- 1175 37. Carlson T, Tovar DA, Alink A, Kriegeskorte N. Representational dynamics of object vision:
1176 the first 1000 ms. *J Vis.* 2013;13(10). doi: 10.1167/13.10.1.
- 1177 38. Wardle SG, Kriegeskorte N, Grootswagers T, Khaligh-Razavi SM, Carlson TA. Perceptual
1178 similarity of visual patterns predicts dynamic neural activation patterns measured with
1179 MEG. *Neuroimage.* 2016;132:59-70. doi: 10.1016/j.neuroimage.2016.02.019.
- 1180 39. Cichy RM, Pantazis D, Oliva A. Resolving human object recognition in space and time. *Nat*
1181 *Neurosci.* 2014;17(3):455-62. doi: 10.1038/nn.3635.
- 1182 40. Kobak D, Brendel W, Constantinidis C, Feierstein CE, Kepecs A, Mainen ZF, et al.
1183 Demixed principal component analysis of neural population data. *Elife.* 2016;5. doi:
1184 10.7554/eLife.10989.
- 1185 41. Robinson GK. That BLUP is a Good Thing: The Estimation of Random Effects. *Statistical*
1186 *Science.* 1991;6(1):15–32. doi: 10.1214/ss/1177011926.
- 1187 42. Murphy KP. *Machine Learning: A probabilistic perspective.* Cambridge, MA: MIT press;
1188 2012.
- 1189 43. Golub GH, Heath M, Wahba G. Generalized Cross-Validation as a Method for Choosing a
1190 Good Ridge Parameter. *Technometrics.* 1979;21(2):215-23. doi:
1191 10.1080/00401706.1979.10489751.
- 1192 44. Harville DA. Maximum Likelihood Approaches to Variance Component Estimation and to
1193 Related Problems. *Journal of the American Statistical Association.* 1977;72(358):320-38.
1194 doi: 10.1080/01621459.1977.10480998.
- 1195 45. Diedrichsen J, Yokoi A, Arbucl S. *Pattern component modeling toolbox.* 2016.
- 1196 46. Eisenhauer JG. Regression through the origin. *Teaching Statistics.* 2003;25(3):76-80.
- 1197 47. Ingram JN, Kording KP, Howard IS, Wolpert DM. The statistics of natural hand
1198 movements. *Exp Brain Res.* 2008;188(2):223-36. Epub 2008/03/29. doi: 10.1007/s00221-
1199 008-1355-3.
- 1200 48. Krizhevsky A, Sutskever I, Hinton GE. ImageNet Classification with Deep Convolutional
1201 Neural Networks. *NIPS; Lake Tahoe, Nevada*2012.
- 1202 49. Sokal RR, Rohlf FJ. *Biometry: the principles and practice of statistics in biological research.*
1203 2nd ed. San Fransisco: W. H. Freeman; 1981.
- 1204 50. Oosterhof NN, Wiestler T, Downing PE, Diedrichsen J. A comparison of volume-based and
1205 surface-based multi-voxel pattern analysis. *Neuroimage.* 2011;56(2):593-600. Epub
1206 2010/07/14. doi: 10.1016/j.neuroimage.2010.04.270.
- 1207 51. Kriegeskorte N, Goebel R, Bandettini P. Information-based functional brain mapping. *Proc*
1208 *Natl Acad Sci U S A.* 2006;103(10):3863-8.
- 1209 52. Kriegeskorte N, Bandettini P. Analyzing for information, not activation, to exploit high-
1210 resolution fMRI. *Neuroimage.* 2007;38(4):649-62. doi: 10.1016/j.neuroimage.2007.02.022.
- 1211 53. O'Toole AJ, Jiang F, Abdi H, Haxby JV. Partially distributed representations of objects and
1212 faces in ventral temporal cortex. *J Cogn Neurosci.* 2005;17(4):580-90. doi:
1213 10.1162/0898929053467550.
- 1214 54. Khaligh-Razavi SM, Kriegeskorte N. Deep supervised, but not unsupervised, models may
1215 explain IT cortical representation. *PLoS Comput Biol.* 2014;10(11):e1003915. doi:
1216 10.1371/journal.pcbi.1003915.
- 1217 55. Jozwik KM, Kriegeskorte N, Mur M. Visual features as stepping stones toward semantics:
1218 Explaining object similarity in IT and perception with non-negative least squares.
1219 *Neuropsychologia.* 2016;83:201-26. doi: 10.1016/j.neuropsychologia.2015.10.023.
- 1220 56. Khaligh-Razavi SM, Henriksson L, Kay K, Kriegeskorte N. Fixed versus mixed RSA:
1221 Explaining visual representations by fixed and mixed feature sets from shallow and deep
1222 computational models *BioRxiv.* 2016.

- 223 57. Dumoulin SO, Wandell BA. Population receptive field estimates in human visual cortex.
224 Neuroimage. 2008;39(2):647-60. doi: 10.1016/j.neuroimage.2007.09.034.
- 225 58. Kriegeskorte N, Cusack R, Bandettini P. How does an fMRI voxel sample the neuronal
226 activity pattern: compact-kernel or complex spatiotemporal filter? Neuroimage.
227 2010;49(3):1965-76. doi: 10.1016/j.neuroimage.2009.09.059.
- 228 59. Kamitani Y, Tong F. Decoding the visual and subjective contents of the human brain. Nat
229 Neurosci. 2005;8(5):679-85. Epub 2005/04/27. doi: 10.1038/nn1444.
- 230 60. Chaimow D, Yacoub E, Ugurbil K, Shmuel A. Modeling and analysis of mechanisms
231 underlying fMRI-based decoding of information conveyed in cortical columns. Neuroimage.
232 2011;56(2):627-42. doi: 10.1016/j.neuroimage.2010.09.037.
- 233 61. Ramirez FM, Cichy RM, Allefeld C, Haynes JD. The neural code for face orientation in the
234 human fusiform face area. J Neurosci. 2014;34(36):12155-67. doi:
235 10.1523/JNEUROSCI.3156-13.2014.
- 236 62. Norman-Haignere S, Kanwisher NG, McDermott JH. Distinct Cortical Pathways for Music
237 and Speech Revealed by Hypothesis-Free Voxel Decomposition. Neuron. 2015;88(6):1281-
238 96. Epub 2015/12/22. doi: 10.1016/j.neuron.2015.11.035.
- 239 63. Kording K, Jonas E. Could a neuroscientist understand a microprocessor? . BioRxiv. 2016.
240 doi: 10.1101/055624.
- 241