

1 Journal submission: Journal of Neurophysiology
2 Manuscript type: Innovative Methodology

3

4

5 **Quantifying the signals contained in heterogeneous neural responses**
6 **and determining their relationships with task performance**

7

8 **Authors:** Marino Pagan and Nicole C. Rust

9 **Affiliation:** Department of Psychology, University of Pennsylvania, Philadelphia, PA 19104

10

11

12

13 **Running title:** Quantifying heterogeneous neural responses

14

15

16

17

18

19

20 **Corresponding author:**

21 Nicole Rust

22 Department of Psychology

23 University of Pennsylvania

24 3401 Walnut St. Room 317C

25 Philadelphia, PA 19104

26 Phone: (215) 898-4587

27 Email: nrust@psych.upenn.edu

28

29

30 **Acknowledgements:** This work was supported by National Eye Institute Grant R01EY020851
31 and a McKnight Foundation award to NCR

32

33

34

35 **Abstract.** The responses of high-level neurons tend to be mixtures of many different types of
36 signals. While this diversity is thought to allow for flexible neural processing, it presents a
37 challenge for understanding how neural responses relate to task performance and to neural
38 computation. To address these challenges, we have developed a new method to parse the
39 responses of individual neurons into weighted sums of intuitive signal components. Our method
40 computes the weights by projecting a neuron's responses onto a pre-defined orthonormal basis.
41 Once determined, these weights can be combined into measures of signal modulation, however,
42 in their raw form, these signal modulation measures are biased by noise. Here we introduce and
43 evaluate two methods for correcting this bias, and we report that an analytically derived
44 approach produces performance that is robust and superior to a bootstrap procedure. Using
45 neural data recorded from IT and perirhinal cortex as monkeys performed a delayed-match-to-
46 sample target search task, we demonstrate how the method can be used to quantify the
47 amounts of task-relevant signals in heterogeneous neural populations. We also demonstrate
48 how these intuitive quantifications of signal modulation can be related to single-neuron
49 measures of task performance (d').

50

51

52 **Introduction**

53
54 The responses of neurons at higher stages of neural processing in the brain tend to
55 reflect heterogeneous mixtures of many different types of task-relevant signals (e.g. Miller and
56 Desimone 1994, Brody et al. 2003, Buckley et al. 2009, Bennur and Gold 2011, Rigotti et al.
57 2013). This diversity is thought to be advantageous insofar as a population that contains a
58 diversity of neural responses is capable of performing a diversity of tasks (Rigotti et al. 2013).
59 However, response heterogeneity also makes these high-level brain areas difficult to
60 understand using classical single-neuron approaches, which inherently rely on identifying
61 regularities in the response properties of individual neurons across a population (e.g.
62 discovering that the majority of V1 neurons are tuned for orientation).

63 Here we present a method to deconstruct the responses of heterogeneous neurons as
64 weighted sums of intuitive signals. Our method is useful when applied to experimental designs
65 that involve changing multiple experimental parameters, which is of course a prerequisite for
66 uncovering signal “mixtures”. Examples include tasks that require finding a “match” to a target,
67 which involves changing the identities of the “stimuli” and the “target” (e.g. Maunsell et al. 1991,
68 Miller and Desimone 1994, Pagan et al. 2013). Likewise, tasks that require flexible rule-based
69 mappings of sensory stimuli onto behavioral responses involve manipulating the sensory
70 stimulus and the rule (e.g. Mansouri et al. 2007, Bennur and Gold 2011). A slightly less obvious
71 example is a task that requires a subject to remember the specific sequence with which objects
72 appear; the different conditions in such a task can be envisioned as combinations of object
73 identity and time (Naya and Suzuki 2011).

74 To address the challenges associated with understanding how the responses of a
75 heterogeneous neural population reflects different task-relevant components, we have
76 developed a method to parse the responses of individual neurons into weighted sums of
77 intuitive components. Our method computes the weights by projecting a neuron’s responses
78 onto a pre-defined orthonormal basis. Once determined, these weights can then be combined

79 to quantify different types of signal modulation in a manner that does not depend on sign (e.g.
80 firing rate increases or decreases). From a neural coding perspective, both firing rate increases
81 and decreases convey information and thus unsigned modulation measures more accurately
82 reflect signal magnitude. Additionally, because firing rate increases and decreases tend to be
83 balanced in many high-level brain areas (e.g. Maunsell et al. 1991, Miller and Desimone 1994,
84 Romo et al. 1999), the “average” signed modulation across a population is not a useful quantity
85 (i.e. because it takes on a value near zero) whereas the “average” unsigned (absolute valued or
86 squared) modulation is meaningful.

87 As we describe in detail below, our method is related to other approaches, including the
88 analysis-of-variance (ANOVA), the multiple linear regression (MLR), the principal components
89 analysis (PCA) and a recent PCA extension (de-mixed PCA; Machens 2010). While these
90 methods have advantages over our method for some applications, one advantage of our
91 method over the others is that it produces unsigned and unbiased estimates of signal
92 modulation magnitudes. Unbiased signal estimates are important when one wants to compare
93 signals across brain areas, across different points in time, or across different types of signals.
94 However, we note that our method is not ideal for describing exactly “how” neurons are tuned
95 for a particular parameter (e.g. for describing tuning curves).

96 In addition to introducing a new way to measure neural signals, we demonstrate how
97 these measures can be related to task performance. Quantifying task performance for individual
98 neurons by performing a Receiver Operating Characteristic (ROC) analysis or by calculating the
99 related discriminability measure d' is a common way to compare neural signals – between
100 different brain areas, between different points in time within the same brain area, or with
101 behavior (e.g. Newsome et al. 1989, Bennur and Gold 2011, Liebe et al. 2011, Adret et al. 2012,
102 Gu et al. 2012, Swaminathan and Freedman 2012). Understanding the underlying sources of
103 neural task performance differences (e.g. overall firing rate changes versus changes in different
104 types of tuning modulation) is crucial for accurate interpretation of what these differences mean

105 for neural coding. Here we show how our method can be used to derive a precise
106 understanding of how task performance depends on different types of signal modulation.

107
108

109 **Methods**

110
111

112 The data we use to describe our method has been reported previously (Pagan et al.
113 2013). All procedures were performed in accordance with the guidelines of the University of
114 Pennsylvania Institutional Animal Care and Use Committee. Briefly, we recorded neural
115 responses in IT and PRH as monkeys performed a delayed-match-to-sample, sequential target
116 search task that required treating the same images as targets and as distractors on different
117 trials (Fig 1a). Monkeys initiated a trial by fixating a small dot and after a short delay, a cue
118 indicating the target for that trial was presented, followed by a random number (0-3) of
119 distractors, and then the target match. Monkeys indicated the presence of the target match by
120 making a saccade to a specific location on the screen before the onset of the next stimulus and
121 were rewarded for correct responses. Altogether, four images were presented in all possible
122 combinations as a visual stimulus (“looking at”), and as a target (“looking for”), resulting in a
123 four-by-four matrix and at least 20 repeated trials of each condition were collected (Fig 1b).

124 Most of our methods are described in the Results section. Here we describe the
125 statistical procedures we used to evaluate the statistical significance of the observed differences
126 in the mean values of various indices between IT and PRH (Fig 5). Because many of these
127 measures were not normally distributed, we calculated these p-values via a bootstrap procedure.
128 On each iteration of the bootstrap, we randomly sampled the true values from each population,
129 with replacement, and we computed the difference between the means of the two newly created
130 populations. We computed the p-value as the fraction of 1000 iterations on which the difference
131 was flipped in sign relative to the actual difference between the means of the full dataset (e.g. if
132 the mean for PRH was larger than the mean for IT, the fraction of bootstrap iterations in which
the IT mean was larger than the PRH mean; Efron and Tibshirani 1994).

133

134

135 **Results**

136 The methods we describe here are useful for analyzing the neural data from
137 experiments in which experimental conditions are combinations of multiple stimulus parameters
138 (e.g. sensory stimuli combined with different task instructions). Additionally, they can be applied
139 to both parametric variation (e.g. systematic changes in motion direction) as well as non-
140 parametric variation (e.g. changes in object identity where the relationships between different
141 identities are not well-defined). The ultimate goal of our method is to measure the magnitude by
142 which a neuron's responses are modulated by different experimental parameters and below we
143 refer these modulation magnitudes as "signals". Our method involves parsing a neuron's firing
144 responses to N different combinations of the stimulus parameters (i.e. experimental conditions),
145 which we refer to as a "response matrix", into a weighted sum of N intuitively defined signals.
146 This process begins by constructing an orthonormal basis of N vectors. "Ortho" refers to the fact
147 that the vectors are "orthogonal", and this allows the original matrix to be deconstructed into a
148 weighted sum (i.e. none of the neural responses are counted twice). "Norm" refers to the fact
149 that all the vectors have the same length (i.e. the "norm" of each vector, computed as the
150 square root of the summed squared values, is equal to 1). "Basis" refers to the fact that
151 together, the vectors capture all possible types of response modulation that could occur given
152 the specific experimental design. As described in more detail below, once the orthonormal
153 basis is determined, the weights are calculated for each neuron by taking the projection (i.e. the
154 dot product) of the neuron's average firing rate responses and each basis vector and the
155 "signals" are determined by combining weights of the same type.

156

157 *Constructing an orthonormal basis*

158 To construct the basis, we begin by constructing a set of N vectors that capture the
159 types of modulation we are interested in. Next we apply the Gram-Schmidt process to convert
160 the set of vectors into an orthonormal basis. To describe the method, we apply this procedure to
161 an example experimental design taken from our previous work: a delayed-match-to-sample
162 (DMS) target search task (Pagan et al. 2013). In these experiments, monkeys viewed a series
163 of sequentially presented images and indicated when a “target match” appeared within a
164 sequence of “distractors” (Fig 1a). Altogether, monkeys viewed each of four visual images in
165 the context of each image as a target, resulting in a four-by-four matrix of experimental
166 conditions (Fig 1b). In this matrix, target matches fall along the diagonal and distractors fall off
167 the diagonal. This “response matrix” \mathbf{R} is computed as the average spike count response
168 across 20 repeated trials for each of the 16 experimental conditions. Below, we treat \mathbf{R} as a 16-
169 entry vector to perform our calculations.

170 To design an orthonormal basis for this task, we began by constructing a first vector
171 which corresponds the grand mean spike count response across all conditions; all entries in this
172 vector take on the same, constant value (e.g. $1/16$; Fig 1c). The remaining vectors are
173 designed to capture the types of modulation that neural responses might reflect, which follow
174 from the task design. In the case of our experiment, this included three vectors to describe the
175 visual modulation, reflected by columns in the response matrix (Fig 1c). Notably, while there
176 are four different visual images, only three are required to capture the visual modulation once
177 the mean firing rate response has also been defined (i.e. degrees of freedom for the visual
178 conditions = $4 - 1$). The second type of modulation is reflected by rows in this matrix, and
179 corresponds to response modulations that can be attributed to changing the identity of the
180 target; because target identity must be held in working memory during this task, we refer to this
181 as “working memory” modulation. The third type of modulation differentiates whether a
182 condition was a target match or a distractor, and this corresponds to modulation along the
183 diagonal. The final type of modulation is that which is required to describe responses that are

184 “peppered” across the matrix, such as differential responses to the same visual image under
 185 two different distractor conditions, and we refer to this modulation as “residual”. More
 186 technically, residual modulations reflect all nonlinear combinations of visual and working
 187 memory signals that are not diagonal.

188 Once this initial set of vectors is defined, we apply the Gram-Schmidt procedure to
 189 convert it into an orthonormal basis. Specifically, we define each of the N original vectors as \mathbf{v}_i ,
 190 and each of the vectors of the resulting orthonormal basis as \mathbf{b}_i . The Gram-Schmidt process is
 191 applied iteratively to each initially defined vector, and consists of two stages: first the vector is
 192 orthogonalized relative to all the vectors already incorporated into the final, orthonormal basis,
 193 and second, the resulting vector is normalized by its norm $\|\mathbf{b}_i\|$:

$$194 \quad \mathbf{b}_i = \mathbf{v}_i - (\mathbf{v}_i^T \cdot \mathbf{b}_1) \cdot \mathbf{b}_1 - (\mathbf{v}_i^T \cdot \mathbf{b}_2) \cdot \mathbf{b}_2 - \dots - (\mathbf{v}_i^T \cdot \mathbf{b}_{i-1}) \cdot \mathbf{b}_{i-1} \quad (1)$$

$$195 \quad \mathbf{b}_i = \frac{\mathbf{b}_i}{\|\mathbf{b}_i\|} \quad ; \quad \|\mathbf{b}_i\| = \sqrt{\sum_j b_{ij}^2} \quad (2)$$

196 where b_{ij} indicates the j-th element of the i-th vector \mathbf{b}_i .

197 The final orthonormal basis obtained for our experiment is shown in Figure 1d. A crucial
 198 requirement is that the originally defined vectors $\mathbf{v}_1 \dots \mathbf{v}_N$ span the full space; if this is not the
 199 case, the Gram-Schmidt process will fail to produce a valid orthonormal basis. It is possible to
 200 verify this simply by measuring the rank of the matrix obtained by juxtaposing the original
 201 vectors $[\mathbf{v}_1 \dots \mathbf{v}_N]$ and checking that it is equal to N.

202 There is no unique way to parse a set of N vectors into an orthonormal basis. For
 203 example, one might consider the “standard basis” as the set of vectors that define each
 204 experimental condition (e.g. 10000, 01000, 00100, etc). While this basis is orthonormal, it is not
 205 very useful because a projection of a neuron’s responses \mathbf{R} onto this basis would simply return
 206 the mean firing rate response to each experimental condition (i.e. each entry in \mathbf{R}). Decisions

207 about how to create the initial vectors when designing the basis depend on what one is trying to
208 achieve. We often find it useful to begin by considering the task “inputs” and whether the task
209 “output” (i.e. the solution) can be expressed as a linear or nonlinear combination of the inputs
210 because this approach formalizes the mapping between the computational goals of the task and
211 the neural signals. For the DMS task described above, the task inputs include “visual” and
212 “working memory” signals (i.e. the monkey is presented with the identity of the target, which he
213 holds in working memory, and the identity of the visual image). These are equivalent to the
214 “linear terms” of a two-factor ANOVA analysis. The solution for this task – differentiating
215 whether each condition is a target match or a distractor (i.e. the diagonal matrix) – cannot be
216 expressed by any linear combination of inputs but instead requires a nonlinear computation.
217 However, it is only one of many possible nonlinear vectors and it is thus essential to parse it
218 from the “residual” vectors, which also reflect nonlinear combinations of visual and memory
219 signals. We note that “diagonal” and “residual” signals would be combined into a single
220 “nonlinear interaction term” in a two-factor ANOVA (for a more extensive description of the
221 relationship between the orthonormal basis and the ANOVA, see the Discussion).

222 Not all experimental designs allow for orthogonalization, or equivalently, not all
223 experimental parameters can be orthogonalized. For example, Figure 1e depicts a modified
224 experimental design in which some visual images are always presented as distractors and
225 never as targets. In this case, there is no way to produce a component that captures “target
226 match” signals (e.g. one that reflects tuning for whether a condition is a target match or a
227 distractor) that can be orthogonalized with the “visual” components. This is because the
228 experimental design introduces a correlation between image identity and whether the condition
229 is a target match: once you know that the identity of an image is 5, 6, 7 or 8, you know with
230 certainty that the image is a distractor. Stated differently, in this experimental design “target
231 match” and “visual” signals are confounded. Thus an additional advantage of our method is that

232 it introduces a means to evaluate and improve a candidate experimental design through the
 233 attempted construction of a useful orthonormal basis.

234

235 *Computing and interpreting signal modulation magnitudes*

236 Once the orthonormal basis has been defined, we can compute the corresponding signal
 237 modulation magnitudes. A neuron's response matrix \mathbf{R} can always be decomposed into a
 238 weighted sum of the orthonormal components:

239
$$\mathbf{R} = \sum_{i=1}^{16} w_i \cdot \mathbf{b}_i \quad (3)$$

240 where \mathbf{b}_i indicates the i -th component, and w_i indicates the weight associated with the i -th
 241 component. The weights w_i are thus determined by computing the projection (i.e. the dot
 242 product) of the vector \mathbf{R} and each basis component \mathbf{b}_i :

243
$$w_i = \mathbf{R} \cdot \mathbf{b}_i^T \quad (4)$$

244 Ultimately, we are interested in quantifying how much of a neuron's firing rate modulation can
 245 be attributed to changes in specific type of experimental manipulation (e.g. the amount of firing
 246 rate modulation that can be attributed to changes in the visual stimulus), and thus we need to
 247 group together weights that correspond to the same type (e.g. the three visual weights). When
 248 doing so, it is important to consider that these weights can be negative as well as positive.
 249 Positive weights correspond to neural responses that directly resemble an orthonormal basis
 250 component whereas negative weights correspond to neural responses that are simply flipped in
 251 sign and thus they also reflect relevant firing rate modulations. To convert a set of weights into a
 252 measure of a particular type of modulation, we square the weights, sum across the set, and
 253 then take the square root. For the DMS task:

254

$$M_{vis} = \sqrt{\sum_{i \in vis} w_i^2} \quad ; \quad M_{wm} = \sqrt{\sum_{i \in wm} w_i^2} \quad ; \quad M_{diag} = |w_{diag}| \quad ; \quad M_{residual} = \sqrt{\sum_{i \in residual} w_i^2} \quad (5)$$

256 where M_{vis} is the amount of visual modulation, M_{wm} is working memory modulation, M_{diag} is
 257 diagonal modulation, and $M_{residual}$ is residual modulation. When computed this way, each type
 258 of signal modulation measures the standard deviation (i.e. the spread) of the responses,
 259 averaged across repeated trials, and has units of spike count. For example, a “visual signal
 260 equal to 2” means that the trial-averaged spike count was spread two standard deviations
 261 around the grand mean firing rate as a result of changes in the visual stimulus.

262 Next we introduce three different ways to normalize these signal modulation magnitudes,
 263 each designed to highlight a different aspect of signal modulation. First, one might wish to
 264 produce signal modulation measures that are not “raw” (Equation 5) but instead are compared
 265 to the amount of noise. This type of “signal-to-noise” modulation measure can be obtained by
 266 simply normalizing by the average trial-by-trial variability of a neuron:

$$267 \quad M'_{vis} = M_{vis} / \bar{\sigma}_{noise} ; M'_{wm} = M_{wm} / \bar{\sigma}_{noise} ; M'_{diag} = M_{diag} / \bar{\sigma}_{noise} ; M'_{residual} = M_{residual} / \bar{\sigma}_{noise} \quad (6)$$

268 where $\bar{\sigma}_{noise}$ is computed as:

$$269 \quad \bar{\sigma}_{noise} = \sqrt{\frac{1}{16} \cdot \sum_{i=1}^{16} \sigma_{i,noise}^2} \quad (7)$$

270 and $\sigma_{i,noise}^2$ indicates the trial-by-trial variability (variance) associated with the i-th condition. In
 271 this formulation, modulations are unitless and they measure the ratio between the signal and
 272 noise modulations. For example, a “visual signal equal to 2” now means that changes in the
 273 visual signal produced a spread in the trial average spike counts with a standard deviation two-
 274 fold larger than the standard deviation of the noise. To anticipate and prevent confusion, we
 275 note that the issue of whether a signal modulation estimate is biased by noise is distinct from
 276 the issue of normalizing the size of the signal relative to the size of the noise; the former is
 277 related to the issue of getting an accurate estimate of signal size (discussed below) whereas the
 278 latter informs how much a given amount of signal will be actually “useful” at conveying

279 information. In other words, a fixed amount of signal can provide perfect information in the
 280 absence of noise, or it can be almost impossible to detect within a very large amount of noise.

281 As a second consideration, we note that in some situations, including the DMS task,
 282 different types of signals have different numbers of components and it may be desirable to
 283 normalize by the number of components to arrive at a measure of modulation “per degree of
 284 freedom”:

$$285 \quad M_{vis} = \sqrt{\frac{1}{N_{vis}} \cdot \sum_{i \in vis} w_i^2} ; M_{wm} = \sqrt{\frac{1}{N_{wm}} \cdot \sum_{i \in wm} w_i^2} ; M_{diag} = |w_{diag}| ; M_{res} = \sqrt{\frac{1}{N_{res}} \cdot \sum_{i \in res} w_i^2} \quad (8)$$

286 where N_{vis} indicates the number of visual components (=3), N_{wm} indicates the number of
 287 working memory components (=3), and N_{res} indicates the number of residual components (=8).
 288 For example, a “visual signal equal to 2” now means that each visual component was (on
 289 average) responsible for spreading the trial-averaged spike count two standard deviations
 290 around the grand mean. This normalization can also be combined with the noise normalization
 291 in Equation 6 to produce a measurement of the signal-to-noise ratio per component.

292 Finally, one might wish to produce a measure of signal modulation that is affected by
 293 changes in the “pattern” of the response matrix but not by an overall rescaling of the firing rates
 294 whereas in their raw form, signal modulations (Equation 5) are directly proportional to the overall
 295 grand mean firing rate. Scale-invariant modulation measures can be computed by normalizing
 296 each type of signal modulation by the grand mean response to produce quantities that we refer
 297 to as “signal strengths”. This normalization is described in more detail in the section “Relating
 298 signal modulations and task performance”.

299 To illustrate an example of signal modulations, Fig 2 shows the result of our method
 300 applied to six neurons collected during the DMS task, including three neurons whose responses
 301 reflect relatively pure selectivity for signals of a single type (Fig 2, top), and three neurons
 302 whose responses reflect mixtures of different types of signals (Fig 2, bottom). Shown are the

303 response matrices for each neuron (Fig 2, left column), and the top five orthonormal
304 components rescaled by their weights. As described above, weights can be positive or negative
305 and negative weights invert the polarity of the orthonormal component (e.g. compare the
306 diagonal matrices in the 3rd and 6th rows of Fig 2). Also shown are the signal modulations
307 computed as the square root of the summed, squared weights for each type of signal,
308 normalized by the number of components for each signal type (Fig 2, right column; Equation 8).
309 To produce these plots, response matrices were computed by counting spikes in 25 ms
310 windows systematically shifted relative to response onset. Signal modulations are computed by
311 squaring the weights, so that both positive and negative weights contribute equally to measured
312 modulations. Signal modulations thus provide an intuitive quantification of “how much” of a
313 particular type of signal is reflected in the responses of a particular neuron, regardless of the
314 “sign” of that weight (i.e. responses increases or decreases) and regardless of “how” that
315 modulation is distributed across the different components (i.e. tuning). Importantly, computing
316 modulations in this way is biased (Fig 3-4), and this bias must be corrected for to get an
317 accurate measure of modulation (as described below).

318 As highlighted above, an orthonormal basis is not uniquely defined for a given
319 experimental design. This statement also applies to subsets of different types of components –
320 for example, one could define an orthonormal basis with “visual” vectors that are different from
321 those presented in Figure 1d but capture the visual modulations equally well because the three
322 new visual vectors will define the same linear subspace as the original ones. Thus the
323 combined projection of a neuron’s response vector onto the three visual components uniquely
324 captures the amount of modulation that can be attributed to changes in the identity of the visual
325 stimulus even if the specific visual vectors themselves are not uniquely defined. Under what
326 situations is a particular type of signal modulation uniquely defined? In our experiment, the
327 uniquely defined parsing of different signal types follows from the two-dimensional “looking at” /
328 “looking for” matrix structure this task, in which the “visual” and “working memory” conditions are

329 presented in all possible combinations and are thus independent from one another (Fig 1b).
330 Similarly, because the diagonal matrix is a single dimension, it is also uniquely defined. Finally,
331 the “residual” subspace is uniquely defined because it describes everything that remains after
332 the other uniquely defined subspaces have been accounted for. In contrast, if we were to, for
333 example, combine the first visual, the first working memory and the first residual dimension into
334 a measure of signal modulation, we would obtain a subspace that is strictly dependent on the
335 particular choice of basis, i.e. a different orthonormal basis would produce a different linear
336 subspace when the same three components are considered.

337

338 *Bias and bias correction*

339 When estimating the amount of modulation in a signal, noise and limiting sampling size
340 are known to introduce a positive bias (Panzeri et al. 2007). For example, consider a
341 hypothetical neuron that responds with the exact same average firing rate response to each of a
342 set of experimental conditions. Because neurons are noisy, if we were to estimate these mean
343 rates based on a limited number of repeated trials, we would get different values for different
344 conditions and this could lead to the erroneous impression that the neural responses are in fact
345 modulated by the stimuli. Similarly, applying the orthonormal basis method to this data would
346 produce weights shifted away from zero as a result of noise for at least a subset of the basis
347 vectors. While the mean of the weights themselves would be unbiased (because noise would
348 shift the weights to both more positive and more negative values), the process of converting the
349 weights into signal modulation magnitudes by squaring (Equation 5) would result in a positive
350 mean bias.

351 To illustrate this bias, Fig 3 includes plots of summed raw (dotted) and bias-corrected
352 (solid) signal modulation magnitudes plotted as a function of time for the IT and PRH
353 populations (using the “closed form” bias correction described below). These results reveal
354 that under physiologically relevant conditions, these biases can be considerable when signals

355 are small or absent (e.g. at stimulus onset, biased estimates of visual modulation are ~1.7
 356 standard deviations in IT and PRH as compared to bias-corrected measures of ~0) and that
 357 these biases become smaller when signals are larger (e.g. at the peak of the visual signal, the
 358 bias is ~0.25 standard deviations in IT and PRH, which is only ~3% of the bias-corrected value).
 359 This is because the bias is additive in the domain of the squared weights but to compute signal
 360 modulations, we take the square root. The square root operation has the effect of enhancing
 361 the effect of the bias when the modulation is small and shrinking it when the modulation is larger.
 362 The reason why we prefer to take the square root rather than operating on the squared
 363 modulations is that we find that measures of signal modulations in units of “spike counts” are
 364 preferable to units of “squared spike counts” in that they more clearly map onto our intuitive
 365 definitions of signals (e.g. signals double when firing rates double).

366 To estimate bias, we compared two methods: an analytical solution and a bootstrap
 367 technique. Under the assumption that trial-by-trial variability is Gaussian distributed, which is a
 368 reasonable approximation of Poisson distributions when the mean spike counts are sufficiently
 369 large, the amount of bias can be derived and unbiased measures of the squared weights \hat{w}_i^2
 370 can be computed as (see Appendix):

$$371 \quad Bias_{closed\ form} = \frac{\mathbf{R} \cdot (\mathbf{b}_i^T)^2}{T} \quad ; \quad \hat{w}_i^2 = (\mathbf{R} \cdot \mathbf{b}_i^T)^2 - \frac{\mathbf{R} \cdot (\mathbf{b}_i^T)^2}{T} \quad (9)$$

372 where T equals the number of repeated trials for each experimental condition.

373 Because the analytical solution assumes that spike counts are Gaussian distributed
 374 whereas spike count distributions are known to deviate from this assumption, particularly at low
 375 firing rates, we also introduce a bootstrap procedure. The first step in estimating the bias for a
 376 given weight involves resampling with replacement T responses to each condition and
 377 recomputing the squared weight for these bootstrapped responses \tilde{w}_i^2 using Equation 4. Next,
 378 the bias can be estimated by subtracting the modulation computed from the actual responses

379 from the bootstrapped modulation estimates, and finally a corrected estimate for each squared
380 weight \hat{w}_i^2 can be computed simply by subtracting the bias (Efron and Tibshirani, 1994):

$$381 \quad \text{Bias}_{bootstrap} = \tilde{w}_i^2 - w_i^2 \quad ; \quad \hat{w}_i^2 = w_i^2 - (\tilde{w}_i^2 - w_i^2) \quad (10)$$

382 In practice, we find that the bias estimated on any one resampling can be noisy and thus we find
383 it useful to calculate the bias a number of times (e.g. 100) and average the bias across those
384 calculations.

385 To test our bias correction procedures, we performed simulations in which we created
386 “ground truth” neurons with known amounts of underlying modulation, simulated their trial-by-
387 trial variability as Poisson, and compared the ground truth and estimated modulation
388 magnitudes. To test the bias correction in a relevant regime, we performed these simulations
389 by creating a population of 150 “ground truth” neurons that were inspired by actual neurons we
390 recorded in IT and PRH (examples include the neurons shown in Fig 2). Specifically, we
391 computed each simulated neuron’s underlying responses by applying a bias correction to 150
392 randomly selected raw response matrices measured in our experiments, and we then used
393 these mean values to generate N Poisson simulated trials. Figure 4a shows the fractional bias
394 (total bias / total signal), computed for a population of 150 neurons, averaged over 100
395 simulated experiments. This plot reveals that, as expected, total bias decreases as a function of
396 the number of repeated trials (Fig 4a, black). With only 2 trials, the magnitude of the bias
397 exceeded the magnitude of the signal (fractional bias ~1.5), and fractional bias dropped to ~0.15
398 for 20 trials and ~0.025 for 100 trials. However, at all numbers of trials, the closed-form bias
399 correction did a very good job at correcting bias (maximal fractional bias remaining after
400 correction = 0.01 for 2 trials; Fig 4a, red). In contrast, fractional bias remained high after the
401 bootstrap correction for small numbers of trials (~0.7 for 2 trials), but converged to the closed-
402 form correction for more than 25 trials (Fig 4a, cyan). Poor bootstrap performance with small
403 sample size is a well-known phenomenon (Chernick 2007).

404 For a closer look at the closed-form bias correction, Figure 4b displays the distribution of
405 fractional error (total error / total signal) after correction across the 100 simulated experiments
406 when 20 trials for each condition were collected. This distribution is centered around 0, thus
407 confirming that the bias has been successfully removed, and it shows that the remaining
408 average fractional error is small (<0.012 in magnitude) for individual experiments. These results
409 support the validity of our procedure for estimating signal modulation magnitudes averaged
410 across a population. However, we caution the reader that while the average signal modulation
411 estimates are very accurate, no method can correct for the specific “noise” patterns within the
412 data for a particular neuron. To illustrate this, Figure 4c (left) displays the distribution of
413 fractional error remaining after correction for a representative simulated neuron across the 100
414 simulated experiments. On average, the error was zero (thus showing no bias), however, on
415 individual simulated experiments, the fractional error ranged from -0.45 to 0.46. Figure 4c
416 (right) shows the “ground truth” response matrix for this neuron as well as one response matrix
417 collected during a simulated experiment. As depicted by the “ground truth” matrix, this neuron
418 was largely visual modulated and responsive to the visual presentation of image 4 but the
419 response matrices collected in this simulated experiment reflects other types of modulation as a
420 result of trial-by-trial variability; even after bias correction, these translate to signal modulation
421 estimates that deviate from the true underlying value. These results demonstrate that the signal
422 modulation magnitudes computed for any individual neuron need to be interpreted cautiously.

423 The simulations reported above were performed using spike counting windows of 50 ms
424 and with that size counting window, we found that the closed-form bias correction was better at
425 estimating bias than the bootstrap bias correction (Fig 4a). We wondered whether the bootstrap
426 might perform better than the closed-form correction for smaller counting windows where spike
427 count distributions deviate more from the Gaussian assumption (e.g. are Poisson) and thus we
428 compared both types of correction for spike count windows of 2 ms. In these narrow windows,
429 the total fractional bias increased dramatically relative to the broader windows (bias was 16-fold

430 larger than signal for 2 trials; Fig 4d black) but the closed-form bias correction continued to
431 perform well (maximal fractional bias remaining after correction = -0.04 for 2 trials; Fig 4d, red).
432 In contrast, the bootstrap correction performed considerably worse at all numbers of trials, with
433 the most discrepant differences for small numbers of trials; even with 10 trials, the average
434 fractional bias remaining after bootstrap correction was 47% the magnitude of the signal (Fig 4d,
435 cyan). These results suggest that for small spike count windows and the numbers of trials
436 typically collected in these types of experiments (n=5-20), the bootstrap correction is highly
437 inaccurate. In contrast, the closed-form bias correction is highly accurate within this regime
438 despite its assumption of Gaussian distributed trial-by-trial variability.

439

440 *Relating signal modulations and task performance*

441

442 Quantifying the performance of individual neurons on a task by calculating the
443 discriminability measure d' is a commonly used approach to compare neurons within or between
444 brain areas. For tasks that involve multiple experimental parameters or require the combination
445 of multiple information sources to compute a solution, arriving at a quantitative understanding of
446 how different signal types relate to task performance can be challenging. Here we derive this
447 relationship for the DMS task (Fig 1a). We then go on to demonstrate how this type of
448 quantitative understanding can be used to (e.g.) determine which of many possible accounts
449 can explain why two populations have different average d' , by applying the analysis to data
450 collected in IT and PRH.

451 The delayed-match-to-sample task described in Figure 1a requires a subject to
452 determine whether each test image is a target match or a distractor, and thus can be envisioned
453 as a two-way discrimination between the set of all target matches versus the set of all
454 distractors. Because target matches and distractors correspond to conditions on versus off the
455 diagonal of the response matrix, respectively, we refer to this as “diagonal d' ”. Diagonal d' is

456 calculated as the absolute value of the difference between the mean response to all target
 457 matches and the mean response to all distractors, divided by their pooled standard deviation:

$$458 \quad |d| = \frac{|\mu_{Match} - \mu_{Distractor}|}{\sigma_{pooled}}, \quad \text{where } \sigma_{pooled} = \sqrt{\frac{4 \cdot \sigma_{Match}^2 + 12 \cdot \sigma_{Distractor}^2}{16}} \quad (11)$$

459 Because target match modulations in IT and PRH result from both increases and decreases in
 460 the firing rates (e.g. Fig 2, 3rd versus 6th rows), the absolute value of diagonal d' best quantifies
 461 the linearly discriminable match/distractor information in each neuron. Similar to the signal
 462 modulation bias described above, merely taking the absolute value of d' produces a biased
 463 estimate of performance in which any modulations, including noise, translate into positive d' . In
 464 particular, note that this bias is directly dependent on the numerator of the d' (i.e. the estimated
 465 absolute difference can be larger than 0 even if the true difference was 0), while the
 466 denominator corresponds to the classic estimator of the standard deviation and it does not have
 467 a direct impact on the bias (see Appendix). To correct for the bias of the d' we can thus focus on
 468 correcting for the bias of the numerator, which requires a calculation analogous to the one
 469 described above for the case of signal modulations (Equation 9). In particular, it is possible to
 470 show (see Appendix) that the bias of the squared numerator is equal to:

$$471 \quad \frac{\sum_{i=1}^4 \frac{1}{16} \cdot m_i + \sum_{i=1}^{12} \frac{1}{144} \cdot d_i}{T} \quad (12)$$

472 where m_i indicates the response to the i -th match and d_i indicates the response to the i -th
 473 distractor, and T indicates the number of trials. Therefore, a corrected estimate of the absolute
 474 d' can be obtained as:

$$475 \quad |\hat{d}| = \sqrt{\frac{(\mu_{Match} - \mu_{Distractor})^2 - \frac{\sum_{i=1}^4 \frac{1}{16} \cdot m_i + \sum_{i=1}^{12} \frac{1}{144} \cdot d_i}{T}}{\sigma_{pooled}^2}} \quad (13)$$

476 where $|\hat{d}|$ is set to 0 if the numerator takes on a negative value. Below, we also apply the bias-
 477 correction to estimate the orthonormal weights (Equation 9).

478 To derive the relationship between the signal modulations of a neuron (Fig 2) and its
 479 diagonal d' , we begin by computing the orthonormal basis weights (Equation 4). As described
 480 above, rescaling a neuron's firing rate responses (e.g. multiplying a neuron's response matrix by
 481 two) will result in a rescaling (i.e. a doubling) of all its weights, and consequently, its signal
 482 modulations will also increase. To describe the signal modulations in a manner that does not
 483 depend on the overall scaling of firing, we normalized the weights by the grand mean spike
 484 count \overline{SC} . We then considered the grand mean spike count as a separate term.

485 Using the orthonormal basis, the diagonal d' of a neuron can be deconstructed as a
 486 function of three intuitive "signal strengths" (see the Appendix for the derivation):

$$487 \quad |d'| = \sqrt{\frac{D}{ND + 1/\overline{SC}}} \quad (14)$$

488 The first signal strength, "D", which we call the "Diagonal strength" (Fig 5, red) is computed as:

$$489 \quad D = \frac{1}{3} \cdot \left(\frac{w_{diag}}{\overline{SC}} \right)^2 \quad (15)$$

490 where w_{diag} corresponds to the weight applied to the diagonal basis component (Fig 1d). This
 491 signal strength determines the distance between the average of the diagonal responses (the
 492 target matches), and the average of the off-diagonal responses (the distractors), averaged
 493 across all images and all trials (Fig 5b, red line), and this term is proportional to diagonal d' (Fig
 494 5c).

495 The second signal strength, "ND", which we call the "Non-diagonal strength" (Fig 5,
 496 cyan) is computed as:

$$497 \quad ND = \frac{1}{16} \cdot \sum_{\substack{i \neq diag \\ i \neq mean}} \left(\frac{w_i}{\overline{SC}} \right)^2 \quad (16)$$

498 where the weights used are those corresponding to the visual, working memory and
499 residual components (Fig 1d). This term determines the spread of the firing rate responses
500 within the target matches and within the distractors (Fig 5b; cyan line) and it is inversely related
501 to diagonal d' (Fig 5c).

502 The final term $1/\overline{SC}$ (Fig 5, lavender) is designed to capture the trial-by-trial variability of
503 a neuron. When trial-by-trial variability is generated by a Poisson process, the grand mean spike
504 count can be used as a good approximation of the variance across trials within each condition,
505 averaged across the 16 conditions and this term can be described by the inverse of the grand
506 mean spike count (see Appendix). This term is also inversely related to diagonal d' , as an
507 increase in the spread within each condition will produce an overall increase in the spread
508 across the set of all target matches and the set of all distractors.

509 We now demonstrate how understanding the relationship between different signal types
510 and single-neuron task performance can be used to gain insight into neural processing by
511 applying these analyses to our data from IT and PRH. We begin with the observation that
512 diagonal d' was significantly higher in PRH as compared to IT (mean IT=0.11, PRH=0.19,
513 $p<0.001$, Fig 5a). We can use the derivation of diagonal d' presented above to discriminate
514 between different possible explanations of why diagonal d' is higher in PRH. Our decomposition
515 suggests three possible factors that might account for this result (that are not mutually
516 exclusive): 1) the diagonal strength (Fig 5c, red) could be higher in PRH than in IT, 2) the non-
517 diagonal strength (Fig 5c, cyan) could be lower in PRH than in IT and/or 3) grand mean firing
518 responses (Fig 5c, lavender) could be higher in PRH than in IT. First and foremost, the diagonal
519 strength was significantly higher in PRH than in IT (Fig 5d), suggesting that this factor
520 contributed to higher average neuron diagonal d' in PRH. Second, the non-diagonal strength
521 was not significantly different between IT and PRH (Fig 5e), suggesting that this factor could not
522 account for the difference in neuron diagonal d' . Finally, the grand mean firing rates were

523 slightly lower in PRH as compared to IT but not significantly so (Fig 5f), and notably, lower firing
524 rates in PRH are the opposite of what would be required to account for higher average PRH
525 neuron diagonal d' (Fig 5c). Taken together, these results suggest that higher neuron diagonal d'
526 results from a two-fold increase in diagonal structure within the response matrices of PRH
527 neurons as compared to IT neurons, as opposed to alternative explanations (such as increases
528 in firing rate in PRH or more non-diagonal modulation in IT).

529
530

531

532 **Discussion**

533 In our own work, we have found this method of estimating signal modulation
534 magnitudes to be useful for a variety of applications. For example, we have used these signal
535 quantifications as a benchmark to assess model performance (Pagan et al. 2013). We have
536 also used these methods to compare the latencies with which specific types of signals arrive in
537 different brain areas to infer the direction of information flow between them (Pagan et al. 2013).
538 As described above, these methods can also be used to uncover the underlying source of
539 differences in single-neuron performance measures between brain areas to gain insights into
540 neural coding. These are but a few examples of the potential uses of this method.

541

542 *Relationship to other analyses*

543 The method we describe here is similar to a multi-way ANOVA, but it incorporates two
544 important extensions: it parses the signal into more terms and it produces a bias-corrected
545 estimate of signal modulation. For the DMS task described above, a two-way ANOVA would
546 parse the total response variance into two linear terms, a nonlinear interaction term, and an
547 error term. The two ANOVA linear terms map directly onto the summed squared projections
548 onto the visual and working memory orthonormal basis vectors (e.g. in Equation 5, the

549 computation of M_{vis} and M_{wm} before taking the square root). Similarly, the ANOVA nonlinear
550 interaction term maps onto the summed squared projections onto the “diagonal” and “residual”
551 terms in our analysis. We note that parsing the diagonal signals from the other nonlinear terms
552 is crucial in our analysis because this signal reflects the task solution, whereas the other types
553 of nonlinear terms do not. The final term in the ANOVA analysis, the error term, is equal to the
554 square of the $\bar{\sigma}_{noise}$ term described in Equation 6. We remind the reader that in this raw form,
555 the values of the orthonormal basis, as well as the ANOVA, are biased due to trial-by-trial
556 variability (i.e. response matrix structure that arises from noise). The ANOVA deals with this
557 bias by computing the probability (the p-value) that each term is significantly higher than
558 expected by chance, given the trial-by-trial variability, by considering the ratio between each
559 term and the error term (the “F statistic”), based on the assumption that the noise is Gaussian-
560 distributed. However, the ANOVA does not produce bias-corrected estimates of signal
561 modulation whereas here we describe two ways to estimate and correct for this bias.

562 Our method also has similarities with an approach related to the ANOVA, multiple linear
563 regression (MLR). Similar to our procedure, MLR seeks to describe a neuron’s responses as a
564 weighted sum of multiple terms. In practice, it is most often applied to continuous variables (e.g.
565 motion direction or color), and often in cases in which one has an specific underlying model of
566 how different stimulus parameters combine to determine a neuron’s response (e.g. knowledge
567 that neurons have Gaussian shaped tuning functions for motion direction). MLR can also be
568 applied in non-parametric cases and when used in this way, multiple terms are required to
569 capture response modulation of a single variable type (e.g. for object identity, response =
570 baseline + weight_1*identity_1 + weight_2*identity_2 + ...) and in fact, our method could be
571 described as an MLR with the regressors specified by an intuitive orthonormal basis. When
572 viewed from this perspective, our method can provide multiple insights into those wishing to
573 perform this type of MLR. First, a crucial consideration with MLR is the degree to which the

574 different regressors are correlated with one another, because the values of the weights (i.e. the
575 “beta coefficients”) can be misleading in case of strongly correlated regressors. One solution to
576 this problem is to orthogonalize the variables of interest although we note that for some data
577 sets, the experimental variables simply cannot be orthogonalized (e.g. Fig 1e). Our method
578 provides a straightforward way to evaluate the degree to which different candidate experimental
579 designs can be orthogonalized for MLR. Second, determining the weights for a complete
580 orthonormal basis guarantees a full account of a neuron’s spike count modulation whereas an
581 MLR against a few (e.g. linear) terms might provide only a partial account. Finally, if one
582 desires to convert MLR “beta coefficients” into positive-valued measures of modulation, these
583 measures will be biased in the exact same manner we describe above and here we introduce a
584 way to correct for that bias.

585 Our method also has similarities with principal components analysis (PCA) and related
586 techniques (Machens 2010). A PCA applied to the mean firing rate responses of a population of
587 neurons to N experimental conditions returns an orthonormal basis of N “eigenvectors” and
588 each neuron’s mean firing rate response to the N conditions can be decomposed into a
589 weighted sum of these vectors by projecting the neuron’s responses onto the basis, as
590 described above. PCA differs from our method in that the eigenvectors are produced via a
591 procedure that iteratively determines the stimulus dimensions that account for the most
592 response variance across the population with the constraint that each successive vector must
593 be orthogonal to all the others. Consequently, PCA dimensions are not guaranteed to be
594 intuitive. As an illustration, Fig 6a shows the results of a PCA applied to our IT and PRH
595 populations. While the two largest eigenvectors for each population are primarily visual, they
596 are not purely so, and eigenvectors of rank three and lower capture mixtures of different types
597 of modulation. Thus PCA is not very useful in providing an intuitive description of the types of
598 signals reflected in these populations. Rather, PCA is most often used as a “dimensionality
599 reduction” technique. For example, in the case of the reverse correlation method “spike-

600 triggered covariance” (Schwartz et al. 2006) one applies a PCA to the spike-triggered stimuli in
601 an attempt to find a small number of stimulus dimensions that can account for individual
602 neuron’s responses within a linear-nonlinear model framework.

603 One extension of the PCA framework, demixed PCA (dPCA; Machens 2010, Brendel et
604 al. 2011) has recently been introduced as a solution to the “mixing” issues described above for
605 PCA. dPCA allows one to specify the experimental parameters that should not be mixed and
606 thus to perform dimensionality reduction within specific linear subspaces. It is advantageous
607 over our method in scenarios in which (e.g.) one wants to determine whether the responses to a
608 specific type of stimulus parameter can be captured with a simple (i.e. low-dimensional)
609 description, or equivalently to uncover specific types of “tuning”. For example, dPCA has
610 provided important insights into how the memory delay period activity of neurons in prefrontal
611 cortex depends on time (Machens et al. 2010). The results of a dPCA applied to our data in IT
612 and PRH are shown in Fig 6b. The input to dPCA included the neural responses to all conditions
613 as well as the task parameters (i.e. the visual and target identities) associated with each
614 condition. This information, which is not provided to traditional PCA, allows dPCA to search for a
615 set of components that capture most of the modulation while avoiding mixing different types of
616 signals (e.g. visual and working memory). In contrast to a regular PCA (Fig 6a), the first three
617 components in each area are almost exclusively visual, and the fourth component for PRH
618 corresponds to the “diagonal” component of the orthonormal basis. However, one can also see
619 from this analysis that if the desired outcome is a characterization of “how much” of specific,
620 pre-defined signal types are present in a population, the orthonormal basis provides a better
621 approach for two reasons: 1) the components retrieved by dPCA still present some degree of
622 “noise” and thus if the relevant axes are known in advance it is better to measure their
623 modulations directly and 2) in situations in which one wants to make a quantitative comparison
624 between two populations, some compromise has to be established when different dPCA
625 components are retrieved for each populations (e.g. compare IT and PRH in Fig. 6b).

626 Finally, a complementary approach for quantifying signals is to measure single neuron
627 performance either by a Receiver Operating Characteristic analysis (e.g. Newsome et al. 1989,
628 Bennur and Gold 2011, Swaminathan and Freedman 2012) or by the related (boundless)
629 discriminability measure d' (e.g. Liebe et al. 2011, Adret et al. 2012, Gu et al. 2012). Under the
630 assumption that trial-by-trial variability is Gaussian distributed, one can convert between the two
631 measures with a simple nonlinear function (i.e. the complementary error function; Dayan and
632 Abbott 2001). As our results show, in a multi-parameter task like DMS, single-neuron task
633 performance does not necessarily depend on a single type of signal but instead can reflect the
634 combination of multiple signal types. Additionally, it is important to note that if one wishes to
635 compute a measure of task performance that is unsigned (i.e. by taking the absolute value or
636 squaring), these task performance measures will be biased. However, this bias can be
637 estimated and corrected using the approaches we describe here.

638

639

640

641 **Appendix**

642
643 *Derivation of the bias correction for signal modulations*

644 When estimating the amount of modulation (or information) in a signal, noise and limited
645 sample size are known to introduce a positive bias (e.g. Treves and Panzeri 1995). Here we
646 quantify the magnitude of this bias for the weights associated to the different signal components
647 (Equation 4), which are used to produce the estimated modulation components (Equation 5) of
648 a single neuron.

649 We begin by making the simplifying assumption that responses to each condition j are
650 normally distributed with mean μ_j and variance σ_j^2 , and that for each condition, μ_j is
651 approximately equal to σ_j^2 . We indicate the estimate of the mean response μ_j to each
652 condition as r_j defined as the average of the responses sampled over T trials. The value of the
653 estimate r_j will itself be normally distributed, with mean equal to the true mean μ_j and variance
654 equal to σ_j^2 / T .

655 By expanding Equation 4, the estimated weight associated to each component i can
656 also be written as:

657
$$w_i = \mathbf{R} \cdot \mathbf{b}_i^T = \sum_{j=1}^{16} r_j \cdot b_{ij} \quad (17)$$

658 where r_j indicates the neuron's average response to the j -th condition, and b_{ij} indicates the j -th
659 entry of the i -th basis component. Since w_i is a linear combination of normally distributed
660 variables, it will also be normally distributed. The mean of w_i will thus be equal to the linear
661 combination of the means of the estimated r_j (i.e. the true mean responses μ_j) and the entries
662 b_{ij} , while the variance will be equal to the linear combination of the variances of r_j (i.e. σ_j^2 / T)
663 and the squared entries b_{ij}^2 :

664
$$\text{Mean}(w_i) = \sum_{j=1}^{16} \mu_j \cdot b_{ij} \quad ; \quad \text{Variance}(w_i) = \frac{\sum_{j=1}^{16} \sigma_j^2 \cdot b_{ij}^2}{T} \quad (18)$$

665 Note that $\sum_{j=1}^{16} \mu_j \cdot b_{ij}$ is the value of the “true” weight, and this estimate of w_i is unbiased.

666 However, a bias is introduced by squaring the estimated weight w_i . The square of a normally
 667 distributed variable with non-zero mean takes the form of a non-central chi-squared distribution,
 668 whose mean is equal to the sum of the squared mean of the original normally distributed
 669 variable plus its variance. In our case:

670
$$\text{Mean}(w_i^2) = [\text{Mean}(w_i)]^2 + \text{Variance}(w_i) = \left(\sum_{j=1}^{16} \mu_j \cdot b_{ij} \right)^2 + \frac{\sum_{j=1}^{16} \sigma_j^2 \cdot b_{ij}^2}{T} \quad (19)$$

671 Under the assumption that the variance for each condition is equal to its mean:

672
$$\text{Mean}(w_i^2) = \left(\sum_{j=1}^{16} \mu_j \cdot b_{ij} \right)^2 + \frac{\sum_{j=1}^{16} \sigma_j^2 \cdot b_{ij}^2}{T} = \left(\sum_{j=1}^{16} \mu_j \cdot b_{ij} \right)^2 + \frac{\sum_{j=1}^{16} \mu_j \cdot b_{ij}^2}{T} \quad (20)$$

673 where the first term corresponds to the “true” squared weight, and the second term represents
 674 the additive bias. If we substitute the true mean responses with their estimates, we obtain an
 675 estimator of the bias:

676
$$\text{Bias} = \frac{\sum_{j=1}^{16} \mu_j \cdot b_{ij}^2}{T} \quad \text{Bias estimator} = \frac{\sum_{j=1}^{16} r_j \cdot b_{ij}^2}{T} = \frac{\mathbf{R} \cdot (\mathbf{b}_i^T)^2}{T} \quad (21)$$

677 where \mathbf{R} indicates the neuron’s response matrix (“flattened” into a vector), and $(\mathbf{b}_i^T)^2$ indicates
 678 the i -th basis function, squared element-by-element. Finally, an unbiased estimator of w_i^2 is
 679 given by:

680
$$\hat{w}_i^2 = (\mathbf{R} \cdot \mathbf{b}_i^T)^2 - \frac{\mathbf{R} \cdot (\mathbf{b}_i^T)^2}{T} \quad (22)$$

681

682 *Derivation of the bias correction for the diagonal d'*

683 The equation for the absolute value of the diagonal d' is presented in Equation 11. As in
 684 the case of signal modulations, for simplicity we proceed by estimating the bias for the squared
 685 diagonal d'. As above, we make the simplifying assumption that responses to each condition j
 686 are normally distributed with mean μ_j and variance σ_j^2 , and that for each condition, μ_j is
 687 approximately equal to σ_j^2 .

688 The numerator of the squared diagonal d' is given by the square of the difference
 689 between the mean match response and the mean distractor response. Since the response to
 690 each condition is assumed to be normally distributed, the difference between mean match and
 691 mean distractor is a linear combination of normal random variables and is also normally
 692 distributed. The numerator is equal to the square of this value, and it thus follows a non-central
 693 chi-squared distribution, whose mean is equal to the sum of the squared mean of the original
 694 normally distributed variable plus its variance:

$$\begin{aligned}
 695 \quad \text{Mean} \left[\left(\mu_{Match} - \mu_{Distractor} \right)^2 \right] &= \left[\text{Mean} \left(\mu_{Match} - \mu_{Distractor} \right) \right]^2 + \text{Variance} \left(\mu_{Match} - \mu_{Distractor} \right) = \dots \\
 696 \quad \dots &= \left(\sum_{i=1}^4 \frac{1}{4} \cdot m_i - \sum_{i=1}^{12} \frac{1}{12} \cdot d_i \right)^2 + \frac{\sum_{i=1}^4 \frac{1}{16} \cdot \sigma_{m_i, noise}^2 + \sum_{i=1}^{12} \frac{1}{144} \cdot \sigma_{d_i, noise}^2}{T} \\
 697 \quad \dots &= \left(\sum_{i=1}^4 \frac{1}{4} \cdot m_i - \sum_{i=1}^{12} \frac{1}{12} \cdot d_i \right)^2 + \frac{\sum_{i=1}^4 \frac{1}{16} \cdot m_i + \sum_{i=1}^{12} \frac{1}{144} \cdot d_i}{T} \quad (23)
 \end{aligned}$$

698 where m_i indicates the mean response to the i-th match and $\sigma_{m_i, noise}^2$ is its corresponding trial-
 699 by-trial variance, d_i indicates the mean response to the i-th distractor and $\sigma_{d_i, noise}^2$ is its
 700 corresponding trial-by-trial variance, and we used the assumption that the trial-by-trial variance

701 for a given condition is equal to its corresponding mean response. The bias of the numerator of
 702 the squared d' is then equal to:

$$703 \quad \text{Bias} \left[\left(\mu_{Match} - \mu_{Distractor} \right)^2 \right] = \frac{\sum_{i=1}^4 \frac{1}{16} \cdot m_i + \sum_{i=1}^{12} \frac{1}{144} \cdot d_i}{T} \quad (24)$$

704 The denominator of the squared diagonal d' is equal to the pooled variance of the noise
 705 across matches and distractors. In the general case in which different conditions elicit different
 706 amounts of trial-by-trial variability, the denominator results in a linear combination of chi-squared
 707 variables, and its parameters can only be estimated in an approximated form (Satterthwaite
 708 1946). However, one can note that the estimate of each individual trial-by-trial variance is
 709 unbiased, and therefore a linear combination of unbiased quantities is unbiased itself, thus no
 710 bias is introduced by the denominator. Consequently, the bias of the diagonal d' can be
 711 corrected by subtracting the bias of the squared numerator according to Equation 24, dividing it
 712 by the estimate of the pooled variance, and taking the square root (Equation 13).

713

714 *Derivation of diagonal d' as a function of the orthonormal basis:*

715 Here we demonstrate that a neuron's diagonal d' can be deconstructed into a function of
 716 three "signal strengths" defined in terms of the orthonormal basis presented in Fig 1d. Diagonal
 717 d' is defined as the absolute value of the difference between the mean response to all target
 718 matches and the mean response to all distractors, divided by the pooled standard deviation of
 719 the noise (Equation 11).

720 The numerator of diagonal d' can thus be expressed as the absolute value of the dot
 721 product between the flattened response matrix \mathbf{R} and a similarly formatted vector \mathbf{c} , in which the
 722 target matches are scaled by $1/4$ and the distractors are scaled by $-1/12$:

723
$$|\mu_{Match} - \mu_{Distractor}| = \left| \sum_{i=1}^4 \frac{1}{4} \cdot m_i - \sum_{i=1}^{12} \frac{1}{12} \cdot d_i \right| = |\mathbf{R} \cdot \mathbf{c}| \quad (25)$$

724 where m_i denotes the mean response to the i -th match and d_i denotes the mean response to
 725 the i -th distractor. The orthonormal basis function corresponding to the diagonal modulation
 726 \mathbf{b}_{diag} is equal to \mathbf{c} multiplied by $\sqrt{3}$ to impose unitary norm. As a result, the numerator of a
 727 neuron's diagonal d' can be rewritten as:

728
$$|\mu_{Match} - \mu_{Distractor}| = |\mathbf{R} \cdot \mathbf{c}| = \left| \mathbf{R} \cdot \frac{\mathbf{b}_{diag}}{\sqrt{3}} \right| = \sqrt{\frac{1}{3} \cdot \mathbf{w}_{diag}^2} \quad (26)$$

729 The denominator of the diagonal d' is equal to the pooled standard deviation, i.e. the
 730 square root of the pooled variance (Equation 11). Our goal is to arrive at a formulation of the
 731 pooled standard deviation as a function of the orthonormal basis weights.

732 We begin by expanding the terms for the variance of spike count responses to target
 733 matches σ_{Match}^2 and to distractors $\sigma_{Distractor}^2$. If we indicate with m_{it} the response to the i -th match
 734 on the t -th trial, σ_{Match}^2 can be rewritten as:

735
$$\sigma_{Match}^2 = \frac{1}{80} \cdot \sum_{i=1}^4 \sum_{t=1}^{20} (m_{it} - \mu_{Match})^2 = \frac{1}{80} \cdot \sum_{i=1}^4 \sum_{t=1}^{20} (m_{it} - m_i + m_i - \mu_{Match})^2 = \dots$$

736
$$= \frac{1}{4} \cdot \sum_{i=1}^4 (m_i - \mu_{Match})^2 + \frac{1}{4} \cdot \sum_{i=1}^4 \frac{1}{20} \cdot \sum_{t=1}^{20} (m_{it} - m_i)^2 = \frac{1}{4} \cdot \sum_{i=1}^4 (m_i - \mu_{Match})^2 + \overline{\sigma}_{noise, Match}^2 \quad (27)$$

737 where $\overline{\sigma}_{noise, Match}^2$ indicates the average trial-by-trial variability across the 4 matches, and m_i
 738 denotes the mean response to the i -th match. Similarly, if we indicate with d_{it} the response to
 739 the i -th distractor on the t -th trial, $\sigma_{Distractor}^2$ can be written as:

740
$$\sigma_{Distractor}^2 = \frac{1}{240} \cdot \sum_{i=1}^{12} \sum_{t=1}^{20} (d_{it} - \mu_{Distractor})^2 = \frac{1}{12} \cdot \sum_{i=1}^{12} (d_i - \mu_{Distractor})^2 + \overline{\sigma}_{noise, Distractor}^2 \quad (28)$$

741 where $\bar{\sigma}_{noise, Distractor}^2$ is the average trial-by-trial variability across the 12 distractors and d_i is the
742 mean response to the i -th distractor. Now we substitute σ_{Match}^2 and $\sigma_{Distractor}^2$ from equations 27
743 and 28 into equation 11, and express the pooled standard deviation as:

$$\begin{aligned}
744 \quad \sigma_{pooled} &= \sqrt{\frac{1}{16} \cdot \left[\sum_{i=1}^4 (m_i - \mu_{Match})^2 + \sum_{i=1}^{12} (d_i - \mu_{Distractor})^2 \right] + \frac{4 \cdot \bar{\sigma}_{noise, Match}^2 + 12 \cdot \bar{\sigma}_{noise, Distractor}^2}{16}} = \dots \\
745 \quad &= \sqrt{\frac{1}{16} \cdot \left[\sum_{i=1}^4 (m_i - \mu_{Match})^2 + \sum_{i=1}^{12} (d_i - \mu_{Distractor})^2 \right] + \bar{\sigma}_{noise}^2} = \sqrt{\sigma_{MD}^2 + \bar{\sigma}_{noise}^2} \quad (29)
\end{aligned}$$

746 where $\bar{\sigma}_{noise}^2$ indicates the average trial-by-trial variability across all conditions (as defined in
747 Equation 7), and σ_{MD}^2 indicates the sum of the variance across matches and the variance
748 across distractors:

$$749 \quad \sigma_{MD}^2 = \frac{1}{16} \cdot \left[\sum_{i=1}^4 (m_i - \mu_{Match})^2 + \sum_{i=1}^{12} (d_i - \mu_{Distractor})^2 \right] \quad (30)$$

750 We now wish to express σ_{MD}^2 as a function of the orthonormal basis components. Here we
751 indicate the average response to the i -th condition as r_i and the grand mean spike count across
752 all conditions as \overline{SC} , and we derive an expansion of the sum of the squared responses by
753 substituting $\overline{SC} = \frac{1}{4} \cdot \mu_{Match} + \frac{3}{4} \cdot \mu_{Distractor}$:

$$\begin{aligned}
754 \quad \sum_{i=1}^{16} r_i^2 &= \sum_{i=1}^{16} (r_i - \overline{SC})^2 + 16 \cdot \overline{SC}^2 = \dots \\
755 \quad &= \sum_{i=1}^4 (m_i - \mu_{Match} + \mu_{Match} - \overline{SC})^2 + \sum_{i=1}^{12} (d_i - \mu_{Distractor} + \mu_{Distractor} - \overline{SC})^2 + 16 \cdot \overline{SC}^2 = \dots \\
756 \quad &= \sum_{i=1}^4 (m_i - \mu_{Match})^2 + \sum_{i=1}^{12} (d_i - \mu_{Distractor})^2 + 3 \cdot (\mu_{Match} - \mu_{Distractor})^2 + 16 \cdot \overline{SC}^2 = \dots \\
757 \quad &= 16 \cdot \sigma_{MD}^2 + 3 \cdot (\mu_{Match} - \mu_{Distractor})^2 + 16 \cdot \overline{SC}^2 \quad (31)
\end{aligned}$$

758 Equation 31 can be rearranged as:

759
$$\sigma_{MD}^2 = \frac{1}{16} \cdot \left[\sum_{i=1}^{16} r_i^2 - 16 \cdot \overline{SC}^2 - 3 \cdot (\mu_{Match} - \mu_{Distractor})^2 \right] \quad (32)$$

760 The diagonal basis function \mathbf{b}_{diag} and the grand mean basis function \mathbf{b}_{mean} are defined such that
 761 their weights w_{diag} and w_{mean} take the following values:

762
$$w_{diag}^2 = (\mathbf{R} \cdot \mathbf{b}_{diag}^T)^2 = 3 \cdot (\mu_{Match} - \mu_{Distractor})^2 \quad ; \quad w_{mean}^2 = (\mathbf{R} \cdot \mathbf{b}_{mean}^T)^2 = 16 \cdot \overline{SC}^2 \quad (33)$$

763 Because $\mathbf{b}_1 \dots \mathbf{b}_{16}$ form an orthonormal basis,

764
$$\sum_{i=1}^{16} w_i^2 = \sum_{i=1}^{16} (\mathbf{R} \cdot \mathbf{b}_i^T)^2 = \sum_{i=1}^{16} r_i^2 \quad (34)$$

765 Substituting equations 33 and 34 into equation 32 allows us to derive:

766
$$\sigma_{MD}^2 = \frac{1}{16} \cdot \left[\sum_{i=1}^{16} w_i^2 - w_{diag}^2 - w_{mean}^2 \right] = \frac{1}{16} \cdot \sum_{\substack{i \neq diag, \\ i \neq mean}} w_i^2 \quad (35)$$

767 We now substitute equation 35 into equation 29:

768
$$\sigma_{pooled} = \sqrt{\frac{1}{16} \cdot \sum_{\substack{i \neq diag, \\ i \neq mean}} w_i^2 + \overline{\sigma}^2_{noise}} \quad (36)$$

769 Diagonal d' can thus be written as:

770
$$|d'| = \frac{|\mu_{Match} - \mu_{Distractor}|}{\sigma_{pooled}} = \frac{\frac{1}{3} \cdot w_{diag}^2}{\sqrt{\frac{1}{16} \cdot \sum_{\substack{i \neq diag, \\ i \neq mean}} w_i^2 + \overline{\sigma}^2_{noise}}} \quad (37)$$

771 In the raw formulation of the weights, rescaling a neuron's firing rate responses (e.g. multiplying
 772 a neuron's response matrix by two) results in a rescaling (i.e. a doubling) of all its deconstructed
 773 matrix weights, and consequently, modulations due to changes in the pattern of responses
 774 within the matrix and overall firing rates are entangled. To capture matrix structure in a manner
 775 that does not depend on the overall scaling of firing, we compute the "normalized weights" by
 776 dividing each weight by the grand mean spike count \overline{SC} . Dividing both numerator and

777 denominator of Equation 37 by \overline{SC} allows us to express diagonal d' as a function of the
 778 normalized weights:

$$779 \quad |d| = \frac{\sqrt{\frac{1}{3} \cdot w_{diag}^2}}{\sqrt{\frac{1}{16} \cdot \sum_{\substack{i \neq diag, \\ i \neq mean}} w_i^2 + \bar{\sigma}_{noise}^2}} = \frac{\sqrt{\frac{1}{3} \cdot \left(\frac{w_{diag}}{\overline{SC}}\right)^2}}{\sqrt{\frac{1}{16} \cdot \sum_{\substack{i \neq diag, \\ i \neq mean}} \left(\frac{w_i}{\overline{SC}}\right)^2 + \frac{1}{\overline{SC}} \cdot \left(\frac{\bar{\sigma}_{noise}^2}{\overline{SC}}\right)}} \quad (38)$$

780 Finally, we express diagonal d' as a function of three components:

$$781 \quad |d| = \sqrt{\frac{D}{ND + \frac{1}{\overline{SC}}}} \quad (39)$$

782 where:

$$783 \quad D = \frac{1}{3} \cdot \left(\frac{w_{diag}}{\overline{SC}}\right)^2 \quad ND = \frac{1}{16} \cdot \sum_{\substack{i \neq diag, \\ i \neq mean}} \left(\frac{w_i}{\overline{SC}}\right)^2 \quad \frac{1}{\overline{SC}} = \frac{1}{\overline{SC}} \cdot \left(\frac{\bar{\sigma}_{noise}^2}{\overline{SC}}\right) \quad (40)$$

784 using the assumption that \overline{SC} is equal to $\bar{\sigma}_{noise}^2$.

785 **Figure legends**

786

787 **Figure 1.** *Constructing an orthonormal basis for a delayed-match-to-sample (DMS) task.* a)

788 Each trial of the DMS task began with the presentation of a cue indicating the target for that trial,

789 followed by the presentation of 0-3 distractors, and then the target match. Images were

790 presented for 400 ms followed by a 400 ms blank. Monkeys were required to maintain fixation

791 throughout the distractors and saccade to a response dot after the target match appeared

792 (within 800 ms) to receive a reward. b) The experimental design included four images each

793 presented as a visual stimulus (“looking at”) in the context of every other image as a target

794 (“looking for”), thus defining a 4x4 matrix. In this matrix, target matches fall along the diagonal

795 and distractors fall off the diagonal. c) The matrices produced by the first stage of the

796 orthonormal basis design process (see Text). d) The matrices produced by applying the Gram

797 Schmidt process to the matrices described in panel c. e) An experimental design in which the

798 “target match” and “visual” conditions cannot be orthogonalized (see Text).

799

800 **Figure 2.** *Example neurons.* Each row depicts a single example neuron, where the responses

801 of the top three neurons reflect relatively pure selectivity and the responses of the bottom three

802 neurons reflect mixtures of different selectivity types. The first column shows the mean spike

803 count responses computed within a window 50-250 ms following stimulus onset to each of the

804 16 conditions (the “response matrix”), averaged over 20 repeated trials, normalized to range

805 from the minimum (black) to the maximum (white). The next five columns show the orthonormal

806 components with the five largest weights, plotted as shown in Figure 1d but with intensity scaled

807 by the weight applied to each component. The response matrix can be reconstructed as a

808 weighted sum of these matrices (once the grand mean spike count is also factored in, which is

809 not shown). The rightmost column shows the temporal evolution of the closed-form bias

810 corrected signal modulation magnitudes for each type of signal, computed as the square root of

811 the sum of squares of the bias corrected weights normalized by the number of components for
812 each signal type (Equation 8). To perform this analysis, spikes were counted in 25 ms sliding
813 windows shifted 1 ms for each successive time bin. The example neuron depicted in the fourth
814 row was recorded in IT; the other neurons were recorded in PRH.

815

816 **Figure 3. Empirical demonstration of bias.** Raw (dotted) and closed-form bias-corrected (solid)
817 measures of signal modulation summed across the IT (top row) and PRH (bottom row)
818 populations plotted with the same conventions as Fig 2, right.

819

820 **Figure 4. Evaluation of bias-correction procedures.** To evaluate the accuracy of our signal
821 modulation measures, a population of 150 simulated neurons with known amounts of signal
822 modulation were created from measured responses in IT and PRH. a) Fractional bias,
823 calculated as the ratio of the total bias (summed across all signal modulations) divided by the
824 total signal (summed across all signal modulations), plotted for the uncorrected simulated
825 population (black), the closed-form bias correction (red), and the bootstrap bias correction
826 (cyan) as a function of the number of Poisson trials collected in each simulated experiment
827 when spikes were counted in 50 ms bins centered 125 ms after stimulus onset. Shown are the
828 averages over 100 simulated experiments. The plot on the right shows an enlargement of the
829 dotted region indicated on the left. b) A histogram of the average (across the 150 simulated
830 neurons) fractional error (total error / total signal) remaining after the closed-form correction for
831 each of 100 simulated experiments with 20 Poisson trials to show that the fractional bias
832 measured per experiment is always near zero. c) *Left*, A histogram of the fractional bias
833 remaining for one representative simulated neuron to show that fractional error measured per
834 neuron can be large. *Right*, the “ground truth” response matrix for this neuron plotted along with
835 one example matrix measured from a simulated experiment that produced an extreme fractional
836 error. d) The results of the same analysis presented in panel a, but performed from responses

837 counted in 2 ms windows. As in panel a, the plot on the right shows an enlargement of the
838 dotted region indicated on the left.

839

840 **Figure 5.** *Relating signal modulation magnitudes with task performance (d').* a) Diagonal d' ,
841 computed as the absolute value of d' , followed by bias correction (Equation 13). Histograms are
842 shown for 167 neurons recorded IT and 164 neurons recorded in perirhinal cortex. PRH
843 neurons with neuron diagonal d' of 1.1, 1.3, and 2.1 are included in the last bin. Arrows indicate
844 means. b) The diagonal d' calculation was based on the distributions of responses over trials to
845 the target matches (dashed) and to the distractors (solid lines). c) Diagonal d' can be
846 expressed as a function of three intuitive components. b-c) The diagonal strength (red) is
847 computed as a function of the normalized diagonal weight (Equation 15) and this component
848 determines the distance between the average response to all target matches and the average
849 response to all distractors (red line); diagonal d' is proportional to this component (see Text).
850 The non-diagonal strength (cyan) is computed as a function of the combined normalized non-
851 diagonal weights (Equation 16) and this component determines the spread within the target
852 matches and within the distractors (cyan line); diagonal d' is inversely related to this component
853 (see Text). The final component (lavender) captures the trial-by-trial variability of the neuron
854 (lavender line). Within a Poisson process and after normalization, this term can be estimated by
855 the inverse of the grand mean spike count across all conditions (Equation 40), and thus
856 diagonal d' increases monotonically with the grand mean spike count (see Text). d) Diagonal
857 strength, b) Non-diagonal strength and c) Grand mean firing rate, in IT (gray) and PRH (white).
858 Arrows indicate means. In subpanel a, the first bin includes neurons with a diagonal strength
859 less than 0.001 and the broken axis extends to a proportion of 0.46 in IT and 0.34 in PRH.

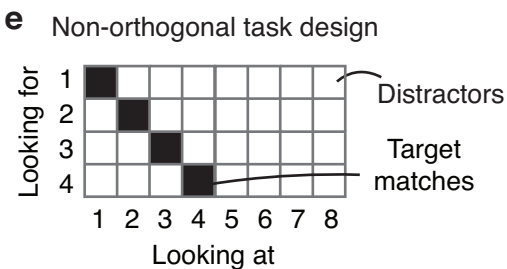
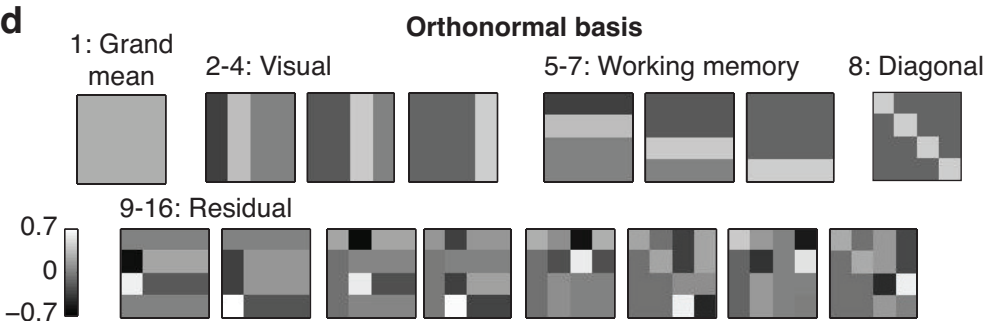
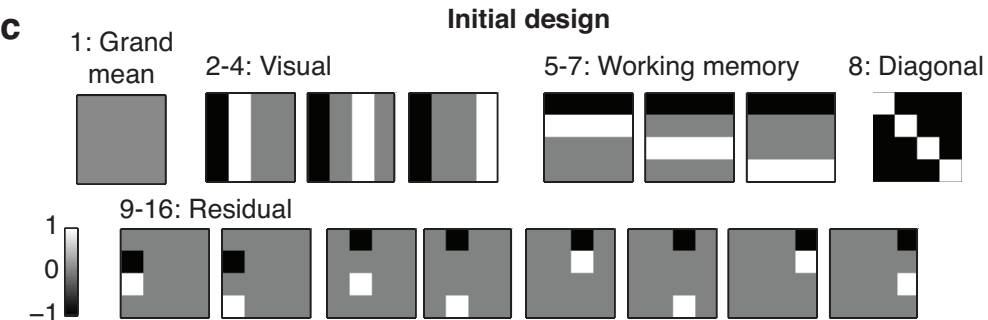
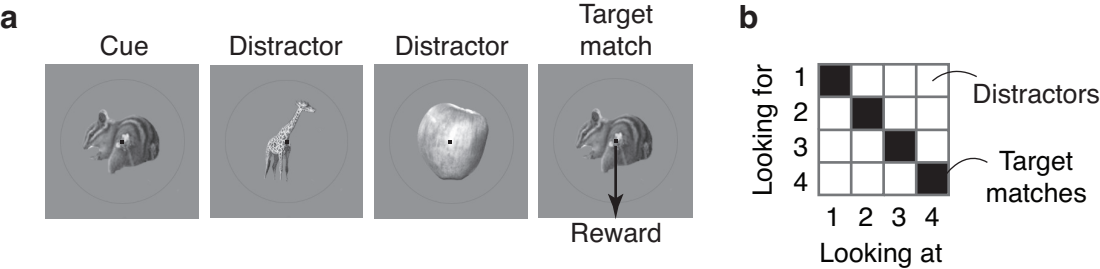
860

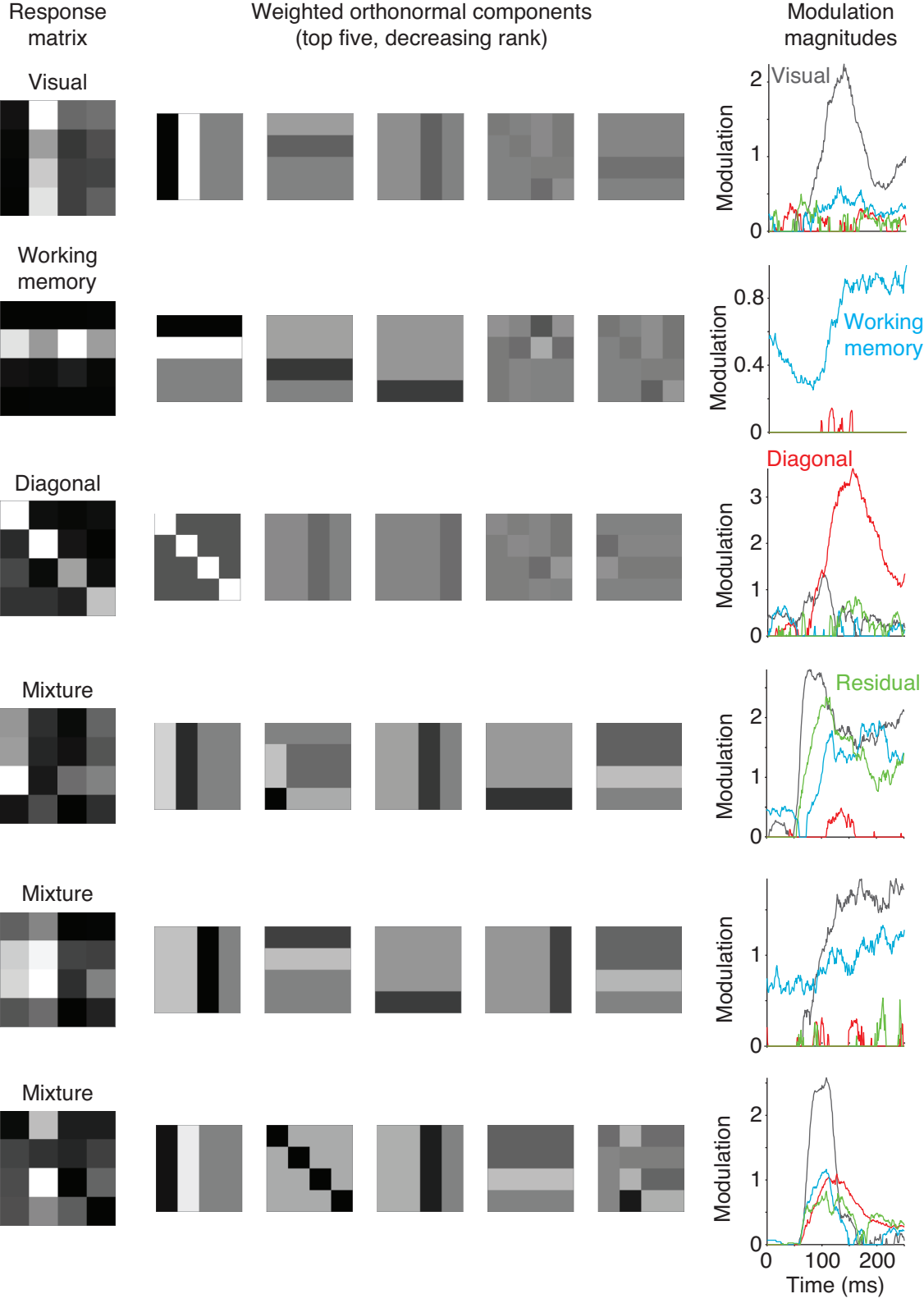
861 **Figure 6.** *Results of PCA and dPCA.* a) Illustration of the orthonormal components
862 corresponding to the eight largest eigenvectors obtained applying PCA to our IT and PRH data.
863 b) The eight largest orthonormal components resulting from the application of dPCA.
864

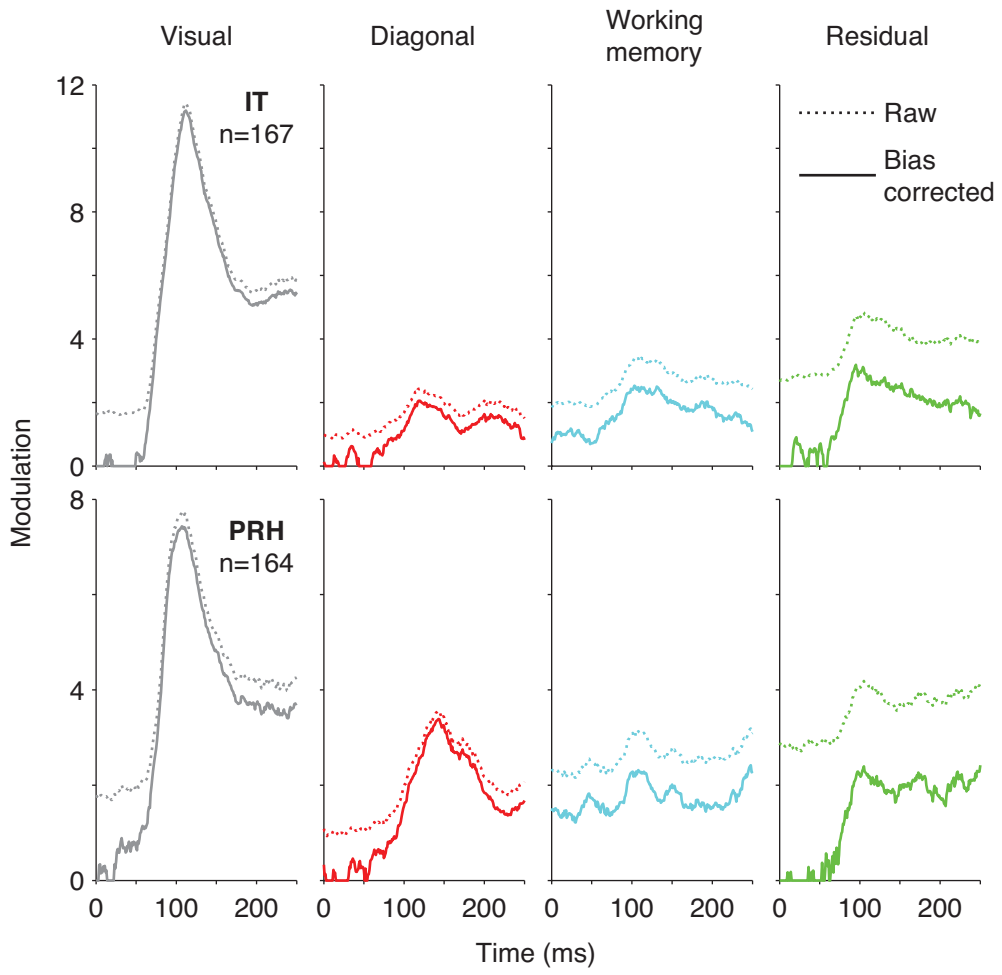
865 **References:**

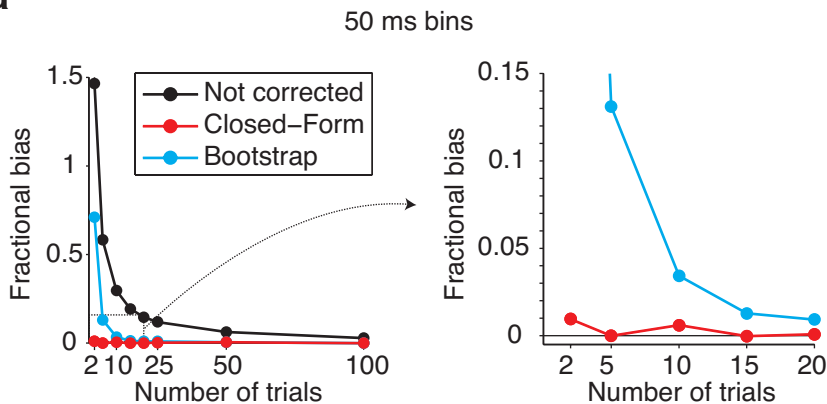
- 866 Adret, P., C. D. Meliza and D. Margoliash (2012). "Song tutoring in presinging zebra finch
867 juveniles biases a small population of higher-order song-selective neurons toward the tutor
868 song." *J Neurophysiol* 108(7): 1977-1987.
- 869 Bennur, S. and J. I. Gold (2011). "Distinct representations of a perceptual decision and the
870 associated oculomotor plan in the monkey lateral intraparietal area." *J Neurosci* 31(3): 913-921.
- 871 Brendel, W., R. Romo and C. K. Machens (2011). Demixed principal component analysis.
872 *Advances in Neural Information Processing Systems* 24.
- 873 Brody, C. D., A. Hernandez, A. Zainos and R. Romo (2003). "Timing and neural encoding of
874 somatosensory parametric working memory in macaque prefrontal cortex." *Cerebral Cortex*
875 13(11): 1196-1207.
- 876 Buckley, M. J., F. A. Mansouri, H. Hoda, M. Mahboubi, P. G. Browning, S. C. Kwok, A. Phillips
877 and K. Tanaka (2009). "Dissociable components of rule-guided behavior depend on distinct
878 medial and prefrontal regions." *Science* 325(5936): 52-58.
- 879 Chernick, M. R. (2007). *Bootstrap Methods: A Guide for Practitioners and Researchers*, 2nd
880 Edition, Wiley.
- 881 Dayan, P. and L. F. Abbott (2001). *Theoretical Neuroscience*, MIT Press.
- 882 Efron, B. and R. J. Tibshirani (1994). *An introduction to the bootstrap*. Boca Raton, CRC Press.
- 883 Gu, Y., G. C. Deangelis and D. E. Angelaki (2012). "Causal links between dorsal medial
884 superior temporal area neurons and multisensory heading perception." *J Neurosci* 32(7): 2299-
885 2313.
- 886 Liebe, S., N. K. Logothetis and G. Rainer (2011). "Dissociable effects of natural image structure
887 and color on LFP and spiking activity in the lateral prefrontal cortex and extrastriate visual area
888 V4." *J Neurosci* 31(28): 10215-10227.
- 889 Machens, C. K. (2010). "Demixing population activity in higher cortical areas." *Front Comput*
890 *Neurosci* 4: 126.
- 891 Machens, C. K., R. Romo and C. D. Brody (2010). "Functional, but not anatomical, separation of
892 "what" and "when" in prefrontal cortex." *J Neurosci* 30(1): 350-360.
- 893 Mansouri, F. A., M. J. Buckley and K. Tanaka (2007). "Mnemonic function of the dorsolateral
894 prefrontal cortex in conflict-induced behavioral adjustment." *Science* 318(5852): 987-990.

- 895 Maunsell, J. H. R., G. Sclar, T. A. Nealey and D. D. Depriest (1991). "Extraretinal
896 Representations in Area-V4 in the Macaque Monkey." *Visual Neuroscience* 7(6): 561-573.
- 897 Miller, E. K. and R. Desimone (1994). "Parallel Neuronal Mechanisms for Short-Term-Memory."
898 *Science* 263(5146): 520-522.
- 899 Naya, Y. and W. A. Suzuki (2011). "Integrating what and when across the primate medial
900 temporal lobe." *Science* 333(6043): 773-776.
- 901 Newsome, W. T., K. H. Britten and J. A. Movshon (1989). "Neuronal correlates of a perceptual
902 decision." *Nature* 341(6237): 52-54.
- 903 Pagan, M., U. L.S., W. M.P. and N. C. Rust (2013). "Signals in inferotemporal cortex and
904 perirhinal cortex suggest an untangling of visual target information." *Nature Neuroscience* 16:
905 1132-1139.
- 906 Panzeri, S., R. Senatore, M. A. Montemurro and R. S. Petersen (2007). "Correcting for the
907 sampling bias problem in spike train information measures." *J Neurophysiol* 98(3): 1064-1072.
- 908 Rigotti, M., O. Barak, M. R. Warden, X. J. Wang, N. D. Daw, E. K. Miller and S. Fusi (2013).
909 "The importance of mixed selectivity in complex cognitive tasks." *Nature* 497(7451): 585-590.
- 910 Romo, R., C. D. Brody, A. Hernandez and L. Lemus (1999). "Neuronal correlates of parametric
911 working memory in the prefrontal cortex." *Nature* 399(6735): 470-473.
- 912 Satterthwaite, F. E. (1946). "An approximate distribution of estimates of variance components."
913 *Biometrics Bulletin* 2: 110-114.
- 914 Schwartz, O., J. W. Pillow, N. C. Rust and E. P. Simoncelli (2006). "Spike-triggered neural
915 characterization." *J Vis* 6(4): 484-507.
- 916 Swaminathan, S. K. and D. J. Freedman (2012). "Preferential encoding of visual categories in
917 parietal cortex compared with prefrontal cortex." *Nat Neurosci* 15(2): 315-320.
- 918 Treves, A. and S. Panzeri (1995). "The upward bias in measures of information derived from
919 limited data samples." *Neural Computation* 7: 399 - 407.
920



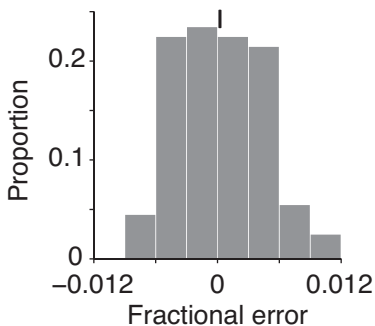




a**b**

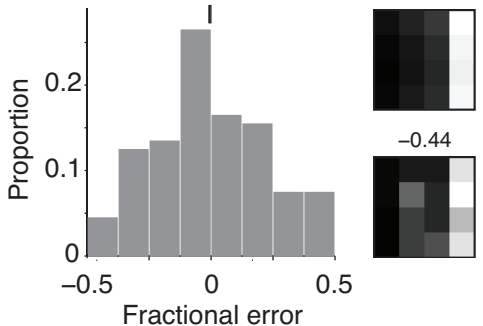
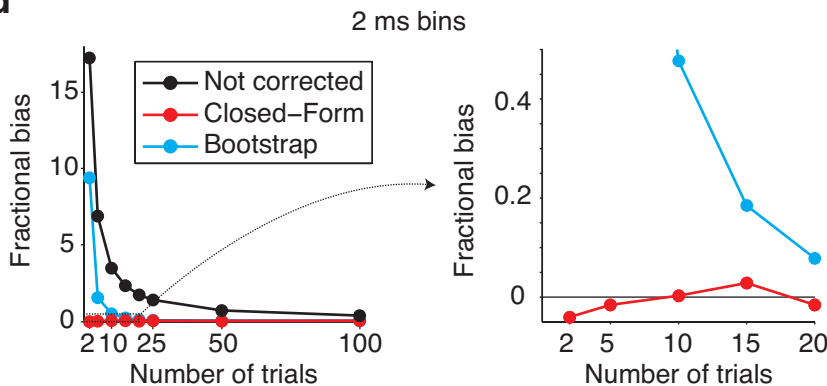
Closed-form correction
150 neurons, 20 trials

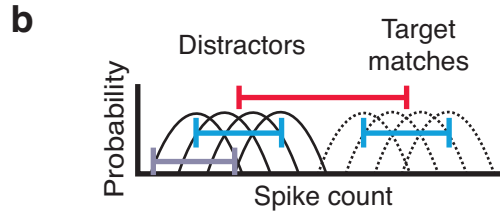
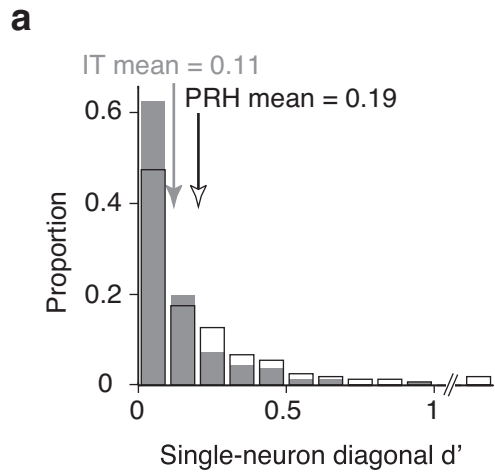
mean = 0.0001

**c**

Closed-form correction
Example neuron, 20 trials

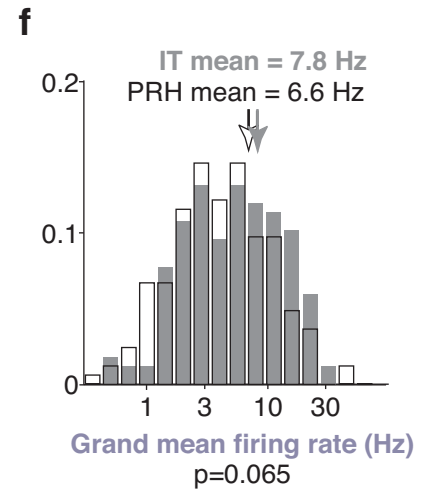
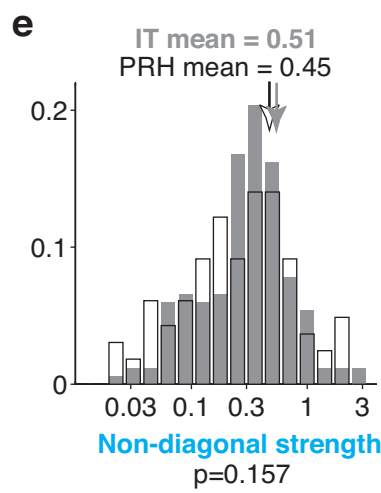
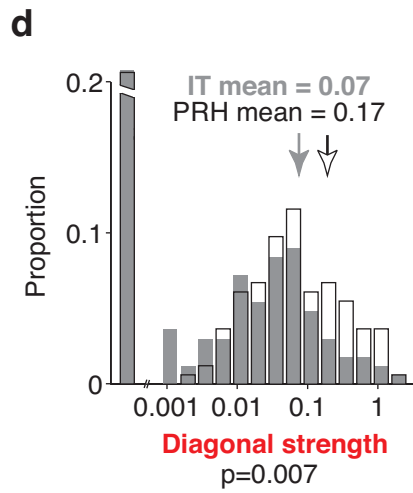
Mean = -0.0001

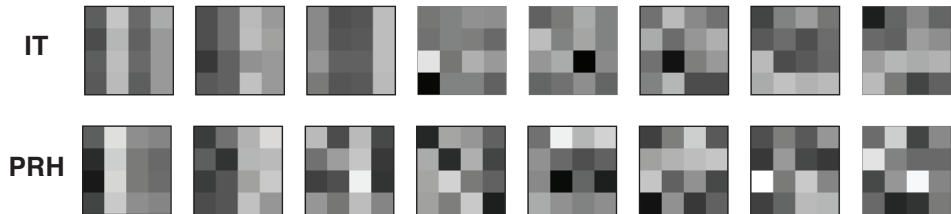
**d**



c

$$|d'| = f \left(\frac{\text{Diagonal Strength (D)}}{\text{Non-diagonal Strength (ND)} + \frac{1}{\text{Grand Mean Spike count (SC)}}} \right)$$



a**PCA: Eigenvectors (decreasing rank)****b****dPCA: Eigenvectors (decreasing rank)**