

1
2
3
4
5
6
7
8

Inferotemporal cortex multiplexes behaviorally-relevant target match signals and visual representations in a manner that minimizes their interference

9
10 Noam Roth and Nicole C. Rust*

11
12 Department of Psychology, University of Pennsylvania, Philadelphia, PA 19104

13
14 * Corresponding author
15 Email: nrust@psych.upenn.edu

16
17
18

19
20 Short title: Inferotemporal cortex during invariant object search

21
22
23

24 **Abstract:**

25

26 Finding a sought visual target object requires combining visual information about a scene with a
27 remembered representation of the target to create a “target match” signal that indicates when a
28 target is in view. Target match signals have been reported to exist within high-level visual brain
29 areas including inferotemporal cortex (IT), where they are mixed with representations of image
30 and object identity. However, these signals are not well understood, particularly in the context of
31 the real-world challenge that the objects we search for typically appear at different positions,
32 sizes, and within different background contexts. To investigate these signals, we recorded
33 neural responses in IT as two rhesus monkeys performed a delayed-match-to-sample object
34 search task in which target objects could appear at a variety of identity-preserving
35 transformations. Consistent with the existence of behaviorally-relevant target match signals in
36 IT, we found that IT contained a linearly separable target match representation that reflected
37 behavioral confusions on trials in which the monkeys made errors. Additionally, target match
38 signals were highly distributed across the IT population, and while a small fraction of units
39 reflected target match signals as target match suppression, most units reflected target match
40 signals as target match enhancement. Finally, we found that the potentially detrimental impact
41 of target match signals on visual representations was mitigated by target match modulation that
42 was approximately (albeit imperfectly) multiplicative. Together, these results support the
43 existence of a robust, behaviorally-relevant target match representation in IT that is configured
44 to minimally interfere with IT visual representations.

45

46

47 **Introduction:**

48

49 Finding a sought visual target object requires combining incoming visual information about the
50 identities of the objects in view with a remembered representation of a sought target object to
51 create a “target match” signal that indicates when a target has been found. During visual target
52 search, target match signals have been reported to emerge in the brain as early as visual areas
53 V4 (Bichot et al., 2005; Chelazzi et al., 2001; Haenny et al., 1988; Kosai et al., 2014; Maunsell
54 et al., 1991) and IT (Chelazzi et al., 1998; Chelazzi et al., 1993; Eskandar et al., 1992; Gibson
55 and Maunsell, 1997; Leuschow et al., 1994; Mruczek and Sheinberg, 2007; Pagan et al., 2013;
56 Woloszyn and Sheinberg, 2009). However, we understand very little about the nature of target
57 match signals, their behavioral relevance, and how these signals are mixed with visual
58 representations.

59

60 The nature of the target match signal has been investigated most extensively with traditional
61 versions of the delayed-match-to-sample (DMS) paradigm, which involves the presentation of a
62 cue image indicating a target’s identity, followed by the presentation of a random number of
63 distractors and then a target match (e.g. Haenny et al., 1988; Miller and Desimone, 1994; Pagan
64 et al., 2013). During classic DMS tasks in which the cue is presented at the beginning of each
65 trial (and the match is thus a repeat later on), IT has been reported to reflect target match
66 information with approximately equal numbers of neurons preferring target matches versus

67 those preferring distractors (i.e. “target match enhancement” and “target match suppression”,
68 respectively; Miller and Desimone, 1994; Pagan et al., 2013). Upon observing that target match
69 suppression also follows from the repetition of distractor images within a trial, and thus cannot
70 account for a signal that corresponds to a “target match” behavioral report, some have
71 speculated that target match enhancement alone reflects the signal used to make behavioral
72 judgments about whether a target match is present (Miller and Desimone, 1994). Others have
73 proposed that the responses of both target match enhanced and suppressed subpopulations
74 are incorporated to make behavioral judgments, particularly when a task requires
75 disambiguating changes in firing rate due to the presence of a target match from other factors
76 that impact overall firing rate, such as stimulus contrast (Engel and Wang, 2011). Notably, no
77 study to date has produced compelling evidence that either IT target match enhancement or
78 suppression accounts for (or correlates with) behavioral reports (e.g. on error trials).

79
80 Another limitation of the traditional DMS paradigm is that the cue image tends to be an exact
81 copy of the target match, whereas real-world object search involves searching for an object that
82 can appear at different positions, sizes and background contexts. One DMS study examined IT
83 neural responses during this type of object variation and reported the existence of target match
84 signals under these conditions (Leuschow et al., 1994). However, we still do not understand
85 how IT target match signals are intermingled with IT invariant object representations of the
86 currently-viewed scene. One intriguing proposal (Fig 1) suggests how visual and target match
87 signals might be multiplexed to minimize the interference between them. That is, insofar as
88 visual representations of different images are reflected as distinct patterns of spikes across the
89 IT population (reviewed by DiCarlo et al., 2012), this translates into a population representation
90 in which visual information is reflected by the population vector angle (Fig 1, ‘Visual
91 modulation’). If the introduction of target match modulation also changes population vector
92 angles in IT, this could result in perceptual confusions about the visual scene. However, if target
93 match modulation amounts to multiplicative rescaling of population response vector lengths, this
94 would minimize interference when superimposing visual memories and target match
95 representations within the same network (Fig 1, ‘Target match modulation’). The degree to
96 which the target match signal acts in this way remains unknown.

97
98
99 **Figure 1. Multiplexing visual and target match representations.** Shown are the hypothetical
100 population responses to two images, each viewed (at different times) as target matches versus
101 as distractors, plotted as the spike count response of neuron 1 versus neuron 2. In this
102 scenario, visual information (e.g. image or object identity) is reflected by the population
103 response pattern, or equivalently, the angle that each population response vector points. In
104 contrast, target match information is reflected by changes in population vector length (e.g.
105 multiplicative rescaling). Because target match information does not impact vector angle in this
106 hypothetical scenario, superimposing target match information in this way would mitigate the
107 impact of intermingling target match signals within underlying perceptual representations.

108
109
110 To investigate the nature of the IT target match signal, its behavioral relevance, and how it
111 intermingles with IT visual representations, we recorded neural signals in IT as monkeys
112 performed a modified delayed-match-to-sample task in which they were rewarded for indicating
113 when a target object appeared across changes in the objects’ position, size and background
114 context.

115 Results:

116

117 *The invariant delayed-match-to-sample task (IDMS)*

118

119 To investigate the target match signal, we trained two monkeys to perform an “invariant
120 delayed-match-to-sample” (IDMS) task that required them to report when target objects
121 appeared across variation in the objects’ positions, sizes and background contexts. In this task,
122 the target object was held fixed for short blocks of trials (~3 minutes on average) and each block
123 began with a cue trial indicating the target for that block (Fig 2a, “Cue trial”). Subsequent test
124 trials always began with the presentation of a distractor and on most trials this was followed by
125 additional distractors and then an image containing the target match (Fig 2a, “Test trial”). The
126 monkeys’ task required them to fixate during the presentation of distractors and make a saccade
127 to a response dot on the screen following target match onset to receive a reward. In cases
128 where the target match was presented for 400 ms and the monkey had still not broken fixation,
129 a distractor stimulus was immediately presented. To minimize the predictability of the match
130 appearing as a trial progressed, on a small subset of the trials the match did not appear and the
131 monkey was rewarded for maintaining fixation. Our experimental design differs from other
132 classic DMS tasks (e.g. Miller and Desimone, 1994; Pagan et al., 2013) in that it does not
133 incorporate a cue at the beginning of each test trial, to better mimic real-world object search
134 conditions in which target matches are not repeats of the same image presented shortly before.

135

136

137 **Figure 2. The invariant delayed-match-to-sample task. a)** Each block began with a cue trial
138 indicating the target object for that block. On subsequent trials, no cue was presented and
139 monkeys were required to maintain fixation throughout the presentation of distractors and make
140 a saccade to a response dot following the onset of the target match to receive a reward. **b)** The
141 experiment included 4 objects presented at each of 5 identity-preserving transformations (“up”,
142 “left”, “right”, “big”, “small”), for 20 images in total. In any given block, 5 of the images were
143 presented as target matches and 15 were distractors. **c)** The complete experimental design
144 included looking “at” each of 4 objects, each presented at 5 identity-preserving transformations
145 (for 20 images in total), viewed in the context of looking “for” each object as a target. In this
146 design, target matches (highlighted in gray) fall along the diagonal of each “looking at” / “looking
147 for” transformation slice. **d)** Percent correct for each monkey, calculated based on both misses
148 and false alarms (but disregarding fixation breaks). Mean percent correct is plotted as a function
149 of the position of the target match in the trial. Error bars (SEM) reflect variation across the 20
150 experimental sessions. **e)** Histograms of reaction times during correct trials (ms after stimulus
151 onset) during the IDMS task for each monkey, with means indicated by arrows and labeled.

152

153

154

155

156 Our experiment included a fixed set of 20 images, including 4 objects presented at each of 5
157 transformations (Fig 2b). Our goal in selecting these specific images was to make the task of
158 classifying object identity challenging for the IT population and these specific transformations
159 were built on findings from our previous work (Rust and DiCarlo, 2010). In any given block (e.g.
160 a squirrel target block), a subset of 5 of the images would be considered target matches and the
161 remaining 15 would be distractors (Fig 2b). Our full experimental design amounted to 20 images
162 (4 objects presented at 5 identity-preserving transformations), all viewed in the context of each

163 of the 4 objects as a target, resulting in 80 experimental conditions (Fig 2c). In this design,
164 “target matches” fall along the diagonals of each looking at / looking for matrix slice (where
165 “slice” refers to a fixed transformation; Fig 2c, gray). For each condition, we collected at least 20
166 repeats on correct trials. Monkeys generally performed well on this task (Fig 2d; mean percent
167 correct monkey 1 = 96%; monkey 2 = 87%). Their mean reaction times (computed as the time
168 their eyes left the fixation window relative to the target match stimulus onset) were 332 ms and
169 364 ms (Fig 2e).

170
171 As two monkeys performed this task, we recorded neural activity in IT using 24-channel probes.
172 We performed two types of analyses on these data. The first type of analysis was performed on
173 the data recorded simultaneously across units within a single recording session (n=20 sessions,
174 including 10 sessions from each monkey). The second type of analysis was performed on data
175 that was concatenated across different sessions to create a pseudopopulation after screening
176 for units based on their stability, isolation, and task modulation (see Methods; n=204 units in
177 total, including 108 units from monkey 1 and 96 units from monkey 2; S1 Dataset). For all but
178 four of our analyses (Fig 4b, 4d, 8, 9), we counted spikes in a window that started 80 ms
179 following stimulus onset (to allow stimulus-evoked responses time to reach IT) and ended at 250
180 ms, which was always before the monkeys’ reaction times on these trials. For all but two of our
181 analyses (Fig 6, 7d), the data are extracted from trials with correct responses.

182
183
184
185
186
187

188 **Target match signals were reflected in IT during the IDMS task**

189
190 Distributions of stimulus-evoked firing rates for the 204 units recorded in our experiment are
191 shown in Figure 3. As is typical of IT and other high-level brain areas, we encountered a
192 heterogeneous diversity of units with regard to their tuning to different aspects of the IDMS task.
193 Figure 4a depicts the responses of four example units, plotted as five slices through our
194 experimental design matrix (Fig 2c), where each slice corresponds to viewing each of the four
195 objects at a fixed transformation (“Looking AT”) in the context of searching for each of the four
196 objects as a target (“Looking FOR”). Different types of task modulation produce distinct structure
197 in these response matrices. Visual modulation translates to vertical structure, (e.g. looking at the
198 same image while looking for different things) whereas target modulation translates to horizontal
199 structure (e.g. looking for the same object while looking at different things). In contrast, target
200 match modulation is reflected as a differential response to the same images presented as target
201 matches (which fall along the diagonal of each slice) versus distractors (which fall off the
202 diagonal), and thus manifests as diagonal structure in each slice.

203
204
205
206
207

Figure 3. Firing rate distributions. The firing rate response to each stimulus was calculated as
the mean across 20 trials in a window 80 - 250 ms following stimulus onset. **a)** Grand mean
firing rate across all 80 conditions. **b)** Maximum firing rates across the 80 conditions. Arrows
indicate the means (n=204 units).

208 **Figure 4. Quantifying modulation in IT during the IDMS task. a)** The response matrices
209 corresponding to four example IT units, plotted as the average response to five slices through
210 the experimental design, where each slice (a 4x4 matrix) corresponds to viewing each of four

211 objects ('Looking AT') in the context of each of four objects as a target ('Looking FOR'), at one
212 transformation ('Big', 'Left', 'Right', 'Small', 'Up'). To compute these responses, spikes were
213 counted in a window 80 –250 ms following stimulus onset, averaged across 20 repeated trials,
214 and rescaled from the minimum (black) to maximum (white) response across all 80 conditions.
215 **b)** Firing rate modulations were parsed into constituent types, where modulation was quantified
216 in units of standard deviation around each unit's grand mean spike count (see Results). The
217 evolution of average modulation magnitudes (across all the units for each animal; monkey 1: n =
218 108, monkey 2: n = 96), shown as a function of time relative to stimulus onset. The shaded area
219 indicates the spike count window used for subsequent analyses. **c)** Average modulation
220 magnitudes computed using the spike count window depicted in panel b. **d)** The average
221 temporal evolution of visual modulation plotted against the average temporal evolution of target
222 match modulation for groups of units organized into quantiles. Units with either target match or
223 visual modulation (n=203 of 204 units) were sorted by their ratios of target match over visual
224 modulation, computed in a window 80-250 ms following stimulus onset. The temporal evolution
225 of the mean across the population (black dotted line) is plotted among the temporal evolution of
226 each 25% quartile of the data, as well as the 95-100% quantile (labeled). Start times of each
227 trajectory (0 ms after stimulus onset) are indicated by a blue dot whereas end times of each
228 trajectory (250 ms after stimulus onset) are indicated by a red dot.

229
230 The first example unit (Fig 4a, 'visual, selective') only responded to one image (object 3
231 presented in the "big" transformation) and was unaffected by target identity. In contrast, the
232 second example unit ('Fig 4a, 'visual, invariant') responded fairly exclusively to one object, but
233 did so across four of the five transformations (all but "up"). This unit also had modest target
234 match modulation, reflected as a larger response to its preferred object (object 1) when it was a
235 distractor (i.e. when searching for targets 2-4) as compared to when it was a target (i.e. when
236 searching for target 1). In other words, this unit exhibited target match suppression. The third
237 example unit ("Fig 4a, 'one-object target match detector') consistently responded with a high
238 firing rate to object 3 presented as a target match across all transformations, but not to other
239 objects presented as target matches. This unit thus exhibited a form of target match
240 enhancement that was selective for object identity. The fourth example unit ("Fig 4a, 'four-object
241 target match detector') responded in a compelling way with a higher firing rate response to
242 nearly any image (any object at any transformation) presented as a target match as compared
243 to as a distractor, or equivalently target match enhancement that was invariant to object identity.
244 Given that the IDMS task requires an eye movement in response to images presented as target
245 matches and fixation to the same images presented as distractors, this unit reflects something
246 akin to the solution to the monkeys' task.

247
248 To quantify the amounts of these different types of modulations across the IT population, we
249 applied a procedure that quantified different types of modulation in terms of the number of
250 standard deviations around each unit's grand mean spike count (Pagan and Rust, 2014b). Our
251 procedure includes a bias-correction to ensure that modulations are not over-estimated by trial
252 variability and it is similar to a multi-way ANOVA, with important extensions (see Methods).
253 Modulation magnitudes were computed for the types described above, including visual, target
254 identity, and target match modulation, as well as "residual" modulations that are reflected as
255 nonlinear interactions between the visual stimulus and the target identity that are not captured
256 by target match modulation (e.g. specific distractor conditions). Notably, this analysis defines
257 target match modulation as a differential response to the same images presented as target
258 matches versus distractors, or equivalently, diagonal structure in the transformation slices

259 presented in Fig 4a. Consequently, units both like the “one-object target match detector” as well
260 as the “four-object target match detector” reflect target match modulation, as both units have a
261 diagonal component to their responses. What differentiates these two units is that the “one-
262 object target match detector” also reflects selectivity for image and target identity, reflected in
263 this analysis as a mixture of target match, visual, and target identity modulation.

264
265 Figure 4b illustrates these modulations computed in a sliding window relative to stimulus onset
266 and averaged across all units recorded in each monkey. As expected from a visual brain area,
267 we found that visual modulation was robust and delayed relative to stimulus onset (Fig 4b,
268 black). Visual modulation was considerably larger in monkey 1 as compared to monkey 2.
269 Target match modulation (Fig 4b, red) was also (as expected) delayed relative to stimulus onset
270 and was smaller than visual modulation, but it was well above the level expected by noise (i.e.
271 zero) and was similar in magnitude in both animals. In contrast, a robust signal reflecting
272 information about the target identity (Fig 4b, green) appeared before stimulus onset in monkey 1
273 and was weaker but also present in monkey 2, consistent with a persistent working memory
274 representation. Note that because the IDMS task was run in blocks with a fixed target, target
275 identity information was in fact present before the onset of each stimulus. Lastly, we found that
276 residual modulation was also present but was smaller than target match modulation in both
277 animals (Fig 4b, cyan). In sum, for a brief period following stimulus onset, visual and target
278 signals were present, but target match signals were not. After a short delay, target match
279 signals appeared as well. When quantified in a window positioned 80 to 250 ms following
280 stimulus onset and computed relative to the size of the target match signal (Fig 4c), visual
281 modulation was considerably larger than target match modulation (monkey 1: 2.9x, monkey 2:
282 2.0x; Fig 4c, gray), whereas the other types of modulations were smaller than target match
283 modulation (target modulation, monkey 1: 0.9x, monkey 2: 0.5x, Fig 4c green; residual
284 modulation, monkey 1: 0.6x, monkey 2: 0.9x Fig 4c, cyan).

285
286 To what degree do these population average traces (Fig 4b) reflect the evolution of visual and
287 target match signals in the same units as opposed to different units? To address this question,
288 we ranked units by their ratios of target match and visual modulation, and grouped them into
289 quantiles of neighboring ranks. Fig 4d shows a plot of the temporal evolution of visual
290 modulation plotted against the evolution of target match modulation for each 25% quartile. The
291 lowest-ranked quartile (Fig 4d, red) largely traversed and then returned along the y-axis,
292 consistent with units that were nearly completely visually modulated. Of interest was whether
293 quartiles with higher ratios of target match modulation would traverse the x-axis in an analogous
294 fashion (reflecting pure target match modulation) or whether these units would begin as visually
295 modulated and become target match modulated at later times. The trajectories for all three
296 higher quartiles (Fig 4d, orange, green, blue) reflected the latter scenario, as they all began with
297 a visually dominated component positioned above the unity line (Fig 4d, gray dashed). Later, the
298 trajectories become more horizontal, indicative of the emergence of target match modulation.
299 Similarly, the trajectory confined to just the top 5% (n=8) units (Fig 4d, purple dashed) began
300 with a visually dominated component that later evolved into strong target match modulation.
301 These results suggest that the evolution of visual to target match modulation is not happening
302 within distinct subpopulations, but rather is reflected within individual units.

303
304 To summarize, the results presented thus far verify the existence of a target match signal in IT
305 that is on average ~40% of the size of the visual modulation. Additionally, while the arrival of
306 target match modulation was delayed relative to the arrival of visual modulation, both types of
307 modulation tend to be reflected in the same units (as opposed to distinct subpopulations).

308
309
310
311
312

IT target match information was reflected as a highly distributed, linearly separable representation

313 The IDMS task required monkeys to determine whether each condition (an image viewed in the
314 context of a particular target block) was a target match or a distractor. This task ultimately maps
315 all the target match conditions onto one behavioral response (a saccade) and all the distractor
316 conditions onto another (maintain fixation), and as such, this task can be envisioned as a two-
317 way classification across changes in other parameters, including changes in target and image
318 identity (Fig 5a). One question of interest is the degree to which the target match versus
319 distractor classification can be made with a linear decision boundary (or equivalently a linear
320 decoder) applied to the IT neural data, as opposed to requiring a nonlinear decoding scheme. In
321 a previous study, we assessed the format of IT target match information in the context of the
322 classic DMS task design (Pagan et al., 2016; Pagan et al., 2013) and found that while a large
323 component was linear, a considerable nonlinear (quadratic) component existed as well.

324
325 To quantify the amount and format of target match information within IT, we began by
326 quantifying cross-validated performance for a two-way target match versus distractor
327 classification with a weighted linear population decoder (a Fisher Linear Discriminant, FLD).
328 Linear decoder performance began near chance and grew as a function of population size,
329 consistent with a robust IT target match representation (Fig 5b, white). To determine the degree
330 to which a component of IT target match information was present in a nonlinear format that
331 could not be accessed by a linear decoder, we measured the performance of a maximum
332 likelihood decoder designed to extract target match information regardless of its format
333 (combined linear and nonlinear, Pagan et al., 2016; Pagan et al., 2013, see Methods).
334 Performance of this nonlinear decoder (Fig 5b, gray) was slightly higher than the linear decoder
335 for the pooled data ($p = 0.022$), and was not consistently higher in both animals (monkey 1 $p =$
336 0.081 ; monkey 2 $p = 0.647$). These results suggest that under the conditions of our
337 measurements (e.g. the population sizes we recorded and the specific images used), IT target
338 match information is reflected almost exclusively in a linearly separable format during the IDMS
339 task. These results are at apparent odds with our previous reports of how IT target match
340 information is reflected during a classic DMS task (see Discussion).

341
342
343 **Figure 5.** *IT target match information is reflected via weighted linear scheme. a)* The target
344 search task can be envisioned as a two-way classification of the same images presented as
345 target matches versus as distractors. Shown are cartoon depictions where each point depicts a
346 hypothetical response of a population of two neurons on a single trial, and clusters of points
347 depict the dispersion of responses across repeated trials for the same condition. Included are
348 responses to the same images presented as target matches and as distractors. Here only 6
349 images are depicted but 20 images were used in the actual analysis. The dotted line depicts a
350 hypothetical linear decision boundary. **b)** Linear (FLD) and nonlinear (maximum likelihood)
351 decoder performance as a function of population size for a pseudopopulation of 204 units
352 combined across both animals, as well as for the data recorded in each monkey individually
353 (monkey 1: $n = 108$ units; monkey 2: $n = 96$ units.) Error bars (SEM) reflect the variability that
354 can be attributed to the random selection of units (for populations smaller than the full dataset)
355 and the random assignment of training and testing trials in cross-validation. **c)** Linear (FLD)
356 decoder performance as a function of the number of top-ranked units removed. Shaded error

357 (SEM) reflects the variability that can be attributed to the random assignment of training and
358 testing trials in cross-validation.

359
360

361 Next, we wanted to better understand how target match information was distributed across the
362 IT population. We thus performed an analysis targeted at the impact of excluding the N “best”
363 target match units for different values of N, with the rationale that if it were the case that the
364 majority of target match information was carried by a small subpopulation of units, performance
365 should fall quickly when those units are excluded. For this analysis, we considered the
366 magnitude but not the sign of the target match modulation (whereas we address questions
367 related to parsing target match modulation by sign, or equivalently target match enhancement
368 versus suppression, below in Figure 7). To perform this analysis, we excluded the top-ranked IT
369 units via a cross-validated procedure (i.e. based on the training data; see Methods). Consistent
370 with a few units that carry target match signals that are considerably stronger than the rest of
371 the population, we found that the slope of the performance drop following the exclusion of the
372 best units was steepest for the top 8% (n=16) ranked units, and that these units accounted for
373 ~25% of total population performance (Fig 5c). However, it was also the case that population
374 performance continued to decline steadily as additional units were excluded, and consequently,
375 population performance could not be attributed to a small fraction of top-ranked units alone (Fig
376 5c). For example, a 50% decrement in performance required removing 27% (n=55/204) of the
377 best-ranked IT population, and mean +/- SEM performance remained above chance up to the
378 elimination of 78% (n=160/204) of top-ranked units. These results are consistent with target
379 match signals that are strongly reflected in a few units (such as Fig 4a example unit 4), and are
380 more modestly distributed across a large fraction of the IT population (such as Fig 4a example
381 unit 2).

382

383 Taken together, these results suggest that IT target match information is reflected by a weighted
384 linear scheme and that target match performance depends on signals that are broadly
385 distributed across most of the IT population.

386

387

388 ***Projections along the IT linear decoding axis reflected behavioral confusions***

389

390 Upon establishing that the format of IT target match information during the IDMS task was linear
391 (on correct trials), we were interested in determining the degree to which behavioral confusions
392 were reflected in the IT neural data. To measure this, we focused on the data recorded
393 simultaneously across multiple units within each session, where all units observed the same
394 errors. With this data, we trained the linear decoder to perform the same target match versus
395 distractor classification described for Fig 5 using data from correct trials, and we measured
396 cross-validated performance on pairs of condition-matched trials: one for which the monkey
397 answered correctly, and the other for which the monkey made an error. On correct trials, target
398 match performance grew with population size and reached above chance levels in populations
399 of 24 units (Fig 6, black). On error trials, mean +/- SE of decoder performance fell below
400 chance, and these results replicated across each monkey individually (Fig 6, white). These
401 results establish that IT reflects behaviorally-relevant target match information insofar as
402 projections of the IT population response along the FLD decoding axis co-vary with the
403 monkeys' behavior.

404

405

406 **Figure 6.** *The IT FLD linear decoder axis reflects behavioral confusions.* Linear decoder
407 performance, applied to the simultaneously recorded data for each session, after training on
408 correct trials and cross-validating on pairs of correct and error trials matched for condition. Error
409 bars (SEM) reflect the variability that can be attributed to the random selection of units (for
410 populations smaller than the full dataset) and the random assignment of training and testing
411 trials in cross-validation. Results are shown for the data pooled across all sessions (main plot,
412 $n = 20$ sessions) as well as when the sessions are parsed by those collected from each animal
413 (monkey 1, $n = 10$ sessions; monkey 2, $n = 10$ sessions).
414

415
416 **Behaviorally-relevant target match signals were reflected as combinations of target**
417 **match enhancement and suppression:**
418

419 As described in the introduction, the IT target match signal has largely been studied via the
420 classic DMS paradigm (which includes the presentation of the cue at the beginning of the trial)
421 and previous results have reported approximately balanced mixtures of target match
422 enhancement and suppression (Miller and Desimone, 1994; Pagan et al., 2013). While some
423 have speculated that target match enhancement alone reflects the behaviorally-relevant target
424 match signal (Miller and Desimone, 1994), others have argued that enhancement and
425 suppression are both behaviorally-relevant (Engel and Wang, 2011). The results presented
426 above demonstrate that during the IDMS task, the representation of target match information is
427 largely linear, and projections along the FLD weighted linear axis reflect behavioral confusions.
428 To what degree does IT target match information, including the reflection of behavioral
429 confusions, follow from units that reflect target match information with target enhancement
430 (positive weights) as compared to target suppression (negative weights)? In our study, this
431 question is of particular interest in light of the fact that our experimental design does not include
432 the presentation of a cue at the beginning of each trial, and thus minimizes the degree to which
433 target match suppression follows passively from stimulus repetition.
434

435 To investigate this question, we computed a target match modulation index for each unit as the
436 average difference between the responses to the same images presented as target matches
437 versus as distractors, divided by the sum of those two quantities. This index takes on positive
438 values for target match enhancement and negative values for target match suppression. In both
439 monkeys, this index was significantly shifted toward positive values (Fig 7a; Wilcoxon sign rank
440 test, monkey 1: mean = 0.063 $p = 8.44e^{-6}$; monkey 2: mean = 0.071, $p = 2.11e^{-7}$). Notably, while
441 these distributions were dominated by units that reflected target match enhancement, a small
442 fraction of IT units in both monkeys reflected statistically reliable target match suppression as
443 well (fraction of units that were significantly target match enhanced and suppressed,
444 respectively, monkey 1: 49.1%, 17.6%; monkey 2: 41.7%, 8.3%; bootstrap significance test,
445 $p < 0.01$).
446

447
448 **Figure 7.** *Target match signals are reflected as mixtures of enhancement and suppression.* **a)** A
449 target match modulation index, computed for each unit by calculating the mean spike count
450 response to target matches and to distractors, and computing the ratio of the difference and the
451 sum of these two values. Dark bars in each histogram indicate the proportions for all units
452 (Monkey 1: $n = 108$; monkey 2: $n = 96$) whereas gray bars indicate the fractions of units whose
453 responses to target matches versus distractors were statistically distinguishable (bootstrap

454 significance test, $p < 0.01$). Arrows indicate the distribution means. **b)** Target match modulation
455 index, computed and plotted as in (a), but after excluding responses to repeated presentation of
456 the same object within a trial. Included are units in which there were at least 10 repeated trials
457 for each condition ($n = 176$ of 204 possible units). **c)** Performance of the FLD classifier for the
458 combined population ($n=204$ units), computed for all units (as described for Fig 5b), target
459 match enhanced units (“E units”) or target match suppressed units (“S units”). **d)** Performance of
460 the FLD classifier for populations of size 24 recorded in each session when trained on correct
461 trials and tested on condition-matched pairs of correct (“Corr.”) and error (“Err.”) trials (as
462 described for Fig 6), computed for all units, E units, and S units.

463
464

465 In our experiment, the same images were not repeated within a trial but the same objects,
466 presented under different transformations, could be. To what degree did the net target match
467 enhancement that we observed follow from distractor suppression as a consequence of
468 adaptation to object repetitions? To assess this, we recomputed target match modulation
469 indices in a manner that only incorporated the responses to the first presentation of each object
470 in a trial. Because this sub-selection reduced the number of distractor trials available for each
471 condition, we equated these with equal numbers of (randomly selected) target match trials. A
472 unit was only incorporated in the analysis if it had at least 10 trials per condition, yielding a
473 subpopulation of 176 (of 204 possible) units. In the absence of distractor object repetitions,
474 target match indices remained shifted toward net enhancement (Fig 7b; Wilcoxon sign rank test,
475 mean = 0.078 $p = 8.09e^{-11}$; fraction of units that were significantly target match enhanced and
476 suppressed, respectively: 30.0%, 6.3%, bootstrap significance test, $p < 0.01$), and the target
477 match indices computed without repeated distractors were not statistically distinguishable from
478 target match indices computed for the full dataset equated for numbers of trials, randomly
479 selected (not shown; mean = 0.067, $p = 0.33$). We thus conclude that the dominance of target
480 match enhancement in our population was not a consequence of distractor suppression that
481 follows from object repetitions within a trial.

482

483 To determine the degree to which target match enhanced versus target match suppressed units
484 contributed to population target match classification performance, we computed performance of
485 the FLD linear decoder when isolated to the target match enhanced or target match suppressed
486 subpopulations. More specifically, we focused on the combined data across the two monkeys
487 (to maximize the numbers of units, particularly given small fraction that were target match
488 suppressed), and we computed performance for variants of the FLD classifier in which the sign
489 of modulation was computed for each unit based on the training data. Cross-validated
490 performance was determined for either the subset of target match enhanced units or the subset
491 of target match suppressed units with the goal of determining their respective contributions to
492 overall population performance (while accounting for the fact that their proportions were not
493 equal). When the analysis was isolated to target match enhanced units (“E units”), performance
494 was virtually identical to the intact population (Fig 7c, mean \pm SEM performance for all units =
495 90.9 \pm 0.02% vs. E units = 90.6 \pm 0.02%), consistent with target match enhancement as the
496 primary type of modulation driving population performance. When the analysis was isolated to
497 target match suppressed units (“S units”), performance on correct trials was lower than that of
498 the intact population but still well above chance (Fig 7c, performance for S units = 64.4 \pm
499 0.03%). This suggests that while target match suppressed units are smaller in number, the
500 target match suppressed units that do exist do in fact carry reliable target match signals.

501

502 What were the relative contributions of E units versus S units to error trial confusions? To
503 address this question, we repeated the error trial analysis described above for Figure 6, but
504 isolated to E or S units. Specifically, we repeated the analysis presented in Figure 6 where we
505 considered the simultaneously recorded data collected across 24 units for each session, but
506 isolated to the E or S units as described for Figure 7c (based on the training data), and we
507 compared cross-validated performance on condition-matched correct versus error trials. E units
508 classified correct trials above chance and misclassified error trials below chance at rates similar
509 to the entire population (Fig 7d, “All units” vs. “E units”), consistent with a larger overall
510 proportion of E units. In contrast, performance of the S units on correct trials was weaker and
511 mean \pm SEM performance was not above chance (53.0 \pm 0.04%; Fig 7d “S units, Corr.”),
512 consistent with smaller numbers of these units in IT. Similarly, performance of S units on
513 correct trials was slightly but not significantly higher than performance on error trials (mean \pm
514 SE performance on error trials = 46.6 \pm 0.02%; $p = 0.090$, Fig 7d, “S units, Err.”). These results
515 indicate that the reflection of behavioral confusions in the IT neural data arises primarily from the
516 activity of E units, but suggest that behavioral confusions may also be weakly reflected in S
517 units.

518
519 As a complementary analysis of behavioral relevance, we also examined the degree to which
520 the responses to target matches reflected pre-saccadic activity by comparing the same
521 responses time-locked to stimulus onset versus saccade onset (Fig 8). The saccade-aligned
522 response was smaller and more diffuse than the stimulus-aligned response and saccade-
523 aligned responses peaked well before saccade onset (\sim 200 ms), suggesting that on average, IT
524 responses to target matches do not reflect characteristic pre-saccadic activity.

525
526 **Figure 8.** *Comparison of stimulus-aligned versus reaction time-aligned responses to target*
527 *matches. a)* Grand mean PSTH across all units ($n=204$) for all target match stimuli, aligned to
528 stimulus onset. *b)* Grand mean PSTH across all units ($n=204$) for all target match stimuli,
529 aligned to behavioral reaction time.

530
531 Together, these results suggest that in the IDMS experiment, target match signals were
532 dominated by target match enhancement, but a smaller, target match suppressed subpopulation
533 exists as well. Additionally, they suggest that the reflection of behavioral confusions in IT neural
534 responses could largely be attributed to units that are target match enhanced, but behavioral
535 confusions were weakly reflected in units that are target match suppressed. Finally, while IT
536 responses reflect behavioral confusions, they were not well-aligned to reaction times.

537
538
539 **The IT target match representation was configured to minimize interference with IT visual**
540 **representations:**

541
542 As a final topic of interest, we wanted to understand how the representation of target match
543 information was multiplexed with visual representations in IT and more specifically, whether IT
544 had a means of minimizing the potentially detrimental impact of mixing these two types of
545 signals. One possible way to achieve this is multiplicative rescaling, as described in Figure 1. To
546 what degree is this happening in IT? As a first step toward addressing this question, we
547 quantified the impact of target match modulation as the representational similarity between the
548 IT population response vectors corresponding to the same images presented as target matches
549 versus as distractors, using a scale-invariant measure of similarity (the Pearson correlation,

550 reviewed by Kriegeskorte and Kievit, 2013). More specifically, we measured the Pearson
551 correlation between pairs of population response vectors via a split-halves procedure (see
552 Methods), and we compared the representational similarity for the same images presented as
553 target matches versus as distractors with other benchmarks in our experiment, including: within
554 the same experimental condition (i.e. random splits across repeated trials); between images
555 containing different transformations of the same object; and between images containing different
556 objects.

557
558 Shown in Figure 9a is the representational similarity matrix corresponding to all possible
559 pairwise combinations of the 20 images used in this experiment, averaged across the matrices
560 computed when the pairs of response vectors under consideration were target matches and
561 when they were distractors, computed with spike count windows 80-250 ms relative to stimulus
562 onset (see Methods). The matrix is organized such that the five transformations corresponding
563 to each object are grouped together. Figure 9b reorganizes the data into plots of the mean and
564 standard error of representational similarity computed for different pairwise comparisons. As
565 expected, we found that the representational similarity was the highest for random splits of the
566 trials corresponding to the same images, presented under the same conditions (Fig 9b, “Same
567 image & condition”, mean = 0.43), which can be regarded as the noise ceiling in our data. In
568 comparison, the representational similarity was significantly lower for different transformations of
569 the same object (Fig 9b, “Different transforms.”; mean = 0.14; $p = 1.14e^{-8}$) as well as for different
570 objects (Fig 9b, “Different objects”; mean = -0.02; $p = 1.92e^{-29}$). We note that a representational
571 similarity value of zero reflects the benchmark of IT population responses that are orthogonal,
572 and this was the case for the representation of different objects in IT. It was also the case that
573 representational similarity was significantly lower for different objects as compared to different
574 transformations of the same object ($p=1.43e^{-7}$), consistent with an IT representation that was
575 tolerant to changes in identity-preserving transformations. With these benchmarks established,
576 what impact did target match modulation have on IT visual representations? The average
577 representational similarity for the same images presented as target matches as compared to
578 distractors was significantly lower than the noise ceiling (Fig 9b, “Matches versus distractors”;
579 mean = 0.28; $p = 2.09e^{-7}$) but was significantly higher than presenting the same object under a
580 new transformation (Fig 9b, $p = 0.0016$) or presenting a different object (Fig 9b, $p = 3.057e^{-20}$).
581 These results suggest that the multiplexing of IT target match signals was not perfect, but also
582 had a smaller impact on the population response than changing either the transformation in
583 which an object was viewed in or the object in view. These results, computed for broad spike
584 count windows (80-250 ms), were qualitatively replicated in narrower windows positioned early
585 (80-130 ms), midway (140-190 ms) and late (200-250 ms) relative to stimulus onset (Fig 9c).
586 Most notably, representational similarity for matches and distractors remained significantly
587 higher than representational similarity for different transformations of the same object in all
588 epochs (Fig 9c, “Mtch. v. Dstr.” vs. “Diff. trans.”, early $p = 0.0023$, mid $p = 0.0081$, late $p =$
589 0.0092). These results confirm that the impact of target match modulation on IT population
590 representational similarity remains modest throughout the stimulus-evoked response period.

591
592 **Figure 9. Target match signaling has minimal impact on the IT visual population response. a)**
593 The representational similarity matrix, computed as the average Pearson correlation between
594 the population response vectors computed for all possible pairs of images. Before computing
595 the correlations between pairs of population response vectors, the responses of each unit were
596 z-normalized to ensure that correlation values were not impacted by differences in overall firing
597 rates across units (see Methods). Correlations were computed based on a split halves

598 procedure. Shown are the average correlations, computed between images with a fixed target
599 and averaged across all possible targets, as well as averaged across 1000 random splits. The
600 matrix is organized such that different transformations of the same object are grouped together,
601 in the same order as depicted in Fig 2. **b)** The average representational similarity, computed
602 across: “Same image and condition”: different random splits of the 20 trials into two sets of 10
603 trials each; “Different transforms.”: images containing different transformations of the same
604 object, computed with a fixed target identity; “Different objects”: images containing different
605 objects, computed with a fixed target identity; “Match versus distractor”: the same image viewed
606 as a target match as compared to as a distractor, averaged across all 9 possible distractor
607 combinations (see Methods). **c)** The analysis described for panel b applied to different time
608 epochs. Error bars (SEM) reflect variability across the 20 images.

609
610 To what degree does the modest impact of target match modulation follow from the
611 multiplicative mechanism highlighted in Figure 1? One requirement for multiplicative population
612 responses are individual units whose responses are themselves multiplicatively rescaled. To
613 determine the degree to which our recorded IT units were multiplicative, we computed the
614 impact of target match modulation as a function of stimulus rank and compared it to the
615 benchmarks expected for multiplicative rescaling as well as other alternatives (including
616 subtraction and sharpening; Fig 10a,c). Specifically, we ranked the responses of each unit to the
617 20 images separately (after averaging across target matches and distractors), and we then
618 computed the average across all units at each rank for target matches and distractors
619 separately. Average IT target match modulation was much better described as multiplicative
620 than as subtractive or sharpening (Fig 10b,d).

621
622 **Figure 10.** *The impact of target match modulation on the visual responses of individual units.*
623 **a)** Cartoon depiction of the impact of different types of target match modulation on the rank-
624 ordered responses to different images. **b)** Mean and SEM of the rank-order responses across
625 units, after ranking the responses for each unit separately (based on the averaged response to
626 target matches and distractors). **c)** The cartoons in panel a, replotted as the difference between
627 target matches and distractors at each rank to visualize the differences between them. **d)** The
628 analysis described in panel c, applied to the data in panel b, reveals that the impact of target
629 match modulation is better described as multiplicative than as subtractive or as sharpening.

630
631 A second requirement for multiplicative population response vectors is homogeneity in target
632 match modulation across units (Fig 11a, cyan). Variation across units in terms of the
633 magnitudes of target match modulation (Fig 11a, left, red), and/or variation that includes
634 mixtures of target match enhancement and suppression (Fig 11a, right, red) can produce
635 changes in population response vector positions that could be confounded with changes in the
636 visual identity, if the variations were sufficiently large. Where does the amount of target match
637 modulation heterogeneity that we observed (e.g. Fig 7a) fall relative to the benchmarks of the
638 best versus worst format that it could possibly take? To investigate this question, we performed
639 a series of data-based simulations targeted at benchmarking our results relative to “best case”
640 and “worse case” scenarios for multiplexing given the magnitudes of target match modulation in
641 our data. As a first “replication” simulation, we replicated the responses recorded for each unit
642 by preserving the magnitudes and types of signals as well as each unit’s grand mean spike
643 count and we simulated trial variability with an independent, Poisson process (see Methods).
644 The pattern of representational similarities reflected in the raw data (Fig 9b) were approximated
645 in simulation (Fig 11b), suggesting that this simulation procedure was effective at capturing
646 important elements of the data. In the other simulations described below, we began in the same

647 way: by preserving the amounts and types of visual, target and residual modulation recorded in
648 each unit, as well as each unit's grand mean firing rate. What differed between the simulations
649 was how that target match modulation was distributed across units.

650
651 To simulate the “best case scenario” in our data, we approximated multiplicative rescaling by
652 distributing the total target match modulation across units in equal proportions relative to their
653 magnitudes of visual modulation. In this simulation, target match modulation was introduced
654 with the same sign (target match enhancement) across all units, consistent with the average
655 sign reflected in the raw data (Fig 7a). Representational similarity between target matches and
656 distractors in this multiplicative, same-sign simulation was statistically indistinguishable from the
657 noise ceiling (Fig 11c, $p = 0.395$), confirming intuitions that a population can (in principle)
658 multiplex target match signals in a multiplicative manner that has minimal interference with
659 visual representations. To simulate a “worse case scenario” for our data, we increased the
660 amount of target match modulation heterogeneity across units by both distributing target match
661 modulation uniformly (as opposed to proportionally) across units as well as preserving the
662 original sign of each unit's target match modulation (i.e. target match enhancement or
663 suppression). Representational similarity between target matches and distractors in this
664 uniform, mixed-sign simulation fell to levels measured for different transformations of the same
665 object (Fig 11c), confirming that our data do not reflect a “worst case scenario” given the
666 magnitudes of target match modulation that we observed. Together, these results suggest that
667 in line with Fig 1, the impact of target match modulation on IT visual representations is modest
668 (Fig 9) as a consequence of modulation that is approximately (albeit imperfectly) multiplicative,
669 due both to individual units with target match modulation that is multiplicative on average, as
670 well as target match modulation that is approximately (albeit imperfectly) functionally
671 homogenous.

672
673 **Figure 11. Benchmarking the impact of target match modulation heterogeneity across units. a)**
674 Cartoon depiction of how heterogeneity across units in target match modulation magnitudes
675 (left) and modulation signs (right) can lead to changes in the population response to the same
676 images presented as target matches versus distractors. **b)** Three simulated variants of the
677 recorded data (see Results), including target match modulation for each unit that was:
678 replicated; enforced to be multiplicative and reflected with the same-sign across all units (i.e.
679 target match enhancement); enforced to be uniform and reflected with mixed-signs across units
680 (i.e. target match enhancement or suppression, as determined by the original data).

681
682

683 Discussion:

684

685 Successfully finding a sought target object, such as your car keys, requires your brain to
686 compute a target match signal that reports when a target is in view. Target match signals have
687 been reported to exist in IT, but these signals are not well understood, particularly in the context
688 of the real-world problem of searching for an object that can appear at different identity-
689 preserving transformations. We recorded responses in IT as two monkeys performed a
690 delayed-match-to-sample task in which a target object could appear at different positions, sizes,
691 and background contexts. We found that the IT population reflected a target match
692 representation that was largely linear, and that it reflected behavioral confusions on trials in
693 which the monkeys made errors. IT target match signals were broadly distributed across most
694 IT units, and while they were dominated by target match enhancement, we also found evidence
695 for reliable target match suppression. Finally, we found that IT target match modulation was
696 configured in such a manner as to minimally impact IT visual representations. Together, these
697 results support the existence of a robust, behaviorally-relevant target match representation in IT
698 that is multiplexed with IT visual representations.

699

700 Our results support the existence of a robust target match representation in IT during this task
701 that reflects confusions on trials in which the monkeys make errors (Fig 6); this result has not
702 been reported previously. One earlier study also explored the responses of IT neurons in the
703 context of a DMS task in which, like ours, the objects could appear at different identity-
704 preserving transformations (Leuschow et al., 1994), but this study did not sort neural responses
705 based on behavior. Another study examined IT neural responses as monkeys performed a
706 visual target search task that involved free viewing as well as image manipulation during the
707 time of the saccade (Mruczek and Sheinberg, 2007). They reported higher firing rates in IT
708 neurons during trial sequence that normally led to a reward (an association between a target
709 object and a saccade to a response target) versus swap trials in which this sequence was
710 disrupted. Another study (from our lab) used a classic DMS design reported that IT population
711 classifications on error trials fell to chance (Pagan et al., 2013), but this study did not find
712 evidence for significant error trial misclassifications.

713

714 IT target match signals have been investigated most extensively in IT via a classic version of the
715 delayed-match-to-sample (DMS) paradigm where each trial begins with a visual cue indicating
716 the identity of the target object, and this cue is often the same image as the target match
717 (Eskandar et al., 1992; Miller and Desimone, 1994; Pagan et al., 2013). In this paradigm,
718 approximately half of all IT neurons that differentiate target matches from distractors do so with
719 enhanced responses to matches whereas the other half are match suppressed (Miller and
720 Desimone, 1994; Pagan et al., 2013). Because match suppressed responses also follow from
721 the repetition of distractors within a trial, some have speculated that the match enhanced
722 neurons alone carry behaviorally-relevant target match information (Miller and Desimone, 1994).
723 In general agreement with those notions, the target match signal is dominated by target match
724 enhancement in situations where the cue and target match are presented at different locations
725 (Chelazzi et al., 1993). Conversely, others have argued that a representation comprised
726 exclusively of match enhanced neurons would confuse the presence of a match with
727 modulations that evoke changes in overall firing rate, such as changes in stimulus contrast
728 (Engel and Wang, 2011). Additionally, these authors proposed that match suppressed neurons
729 could be used in these cases to disambiguate target match versus stimulus-induced modulation.
730 In our experiment, the IDMS task was run in blocks containing a fixed target to minimize the

731 impact of passive stimulus repetition of the target match. We found evidence for net target
732 match enhancement in our data (Fig 7a), and that this in turn translated into a type of
733 homogeneity that minimized the potentially detrimental impact of target match modulation on
734 visual representations (Fig 11). However, we also found evidence for a smaller subpopulation
735 of units that reflected reliable target match suppression. Whether the amount of target match
736 suppression that we observed is sufficient for the disambiguation strategy proposed by Engel
737 and Wang (2011) is thus unclear - because our experiment did not include variation in
738 parameters that change overall firing rate (such as contrast), we cannot directly test it with our
739 data.

740
741 How does the target match signal arrive in IT? Computation of the target match signal requires a
742 comparison of the content of the currently-viewed scene with a remembered representation of
743 the sought target. The existence of target match signals in IT could reflect the implementation of
744 the comparison in IT itself or, alternatively, this comparison might be implemented in a higher-
745 order brain area (such as prefrontal cortex) and fed-back to IT. Examination of the timing of the
746 arrival of this signal in IT (which peaks at 150 ms; Fig 4b) relative to the monkeys' median
747 reaction times (~340 ms; Fig 2e), does not rule out the former scenario. The fact that neural
748 responses to target matches were more time locked to stimulus onset than they are to reaction
749 times suggests that this activity does not reflect classic signatures of motor preparation.
750 Additional insights into whether or not target match signals are computed in IT might be gained
751 through analyses of the responses on cue trials, particularly with regard to whether signatures of
752 the visually-evoked responses to cues persist throughout each block, however, our experimental
753 design included too few cue presentations for such analyses. Thus while our data are consistent
754 with target match computations within IT cortex, we cannot definitively distinguish this proposal
755 from alternative scenarios with this data. Additionally, in this study monkeys were trained
756 extensively on the images used in these experiments and future experiments will be required to
757 address the degree to which these results hold under more everyday conditions in which
758 monkeys are viewing images and objects for the first time.

759
760 In a previous series of reports, we investigated target match signals in the context of the classic
761 DMS design in which target matches were repeats of cues presented earlier in the trial and each
762 object was presented on a gray background (Pagan and Rust, 2014a; Pagan et al., 2016;
763 Pagan et al., 2013). One of our main findings from that work was that the IT target match
764 representation was reflected in a partially nonlinearly separable format, whereas an IT
765 downstream projection area, perirhinal cortex, contained the same amount of target match
766 information but in a format that was largely linearly separable. In the data we present here, we
767 did not find evidence for a nonlinear component of the IT target match representation, reflected
768 as consistently higher performance of a maximum likelihood as compared to linear decoder (Fig
769 5b). The source of these differences is unclear. They could arise from the fact that the IDMS
770 task requires an "invariant" visual representation of object identity, which first emerges in a
771 linearly separable format in the brain area that we are recording from (Rust and DiCarlo, 2010),
772 whereas in more classic forms of the DMS task, the integration of visual and target information
773 could happen in a different manner and/or a different brain area. Alternatively, these differences
774 could arise from the fact that during IDMS, images are not repeated within a trial, and the
775 stronger nonlinear component revealed in DMS may be produced by stimulus repetition. It may
776 also be the case that nonlinearly separable information is in fact present in IT during IDMS but
777 was not detectable under the specific conditions used in our experiments. For example, the
778 proportion of nonlinearly separable information grows as a function of population size, and it

779 may be the case that it is detectable during IDMS for larger sized populations. Our current data
780 cannot distinguish between these alternatives.

781
782 Our results also add to the growing literature that suggests the brain “mixes” the modulations for
783 different task-relevant parameters within individual neurons, even at the highest stages of
784 processing (Freedman and Assad, 2009; Kobak et al., 2016; Mante et al., 2013; Meister et al.,
785 2013; Raposo et al., 2014; Rigotti et al., 2013; Rishel et al., 2013; Zoccolan et al., 2007). A
786 number of explanations have been proposed to account for mixed selectivity. Some studies
787 have documented situations in which signal mixing is an inevitable consequence of the
788 computations required for certain tasks, such as identifying objects invariant to the view in which
789 they appear (Zoccolan et al., 2007). Others have suggested that mixed selectivity may be an
790 essential component of the substrate required to maintain a representation that can rapidly and
791 flexibly switch with changing task demands (Raposo et al., 2014; Rigotti et al., 2013). Still others
792 have maintained that broad tuning across different types of parameters is important for learning
793 new associations (Barak et al., 2013). Our results suggest that IT mixes visual and target match
794 information within individual units. This could reflect the fact that the comparison of visual and
795 target match information happens within IT itself, and multiplexing is simply a byproduct of that
796 computation. Alternatively, if the comparison is performed elsewhere, this would reflect its
797 feedback to IT for some unknown purpose. In either case, our results suggest that the
798 multiplexing happens in a manner that is largely but imperfectly multiplicative (Fig 10-11) and
799 thus configured to minimize interference of visual representations when also signaling target
800 match information.

801
802 **Acknowledgements:**

803
804 We thank Margot P. Wohl for her contributions to early phases of this work. This work was
805 supported by the National Eye Institute of the National Institutes of Health (award
806 R01EY020851), the Simons Foundation (through an award from the Simons Collaboration on
807 the Global Brain), and the McKnight Endowment for Neuroscience.

808

809 METHODS

810

811 Experiments were performed on two adult male rhesus macaque monkeys (*Macaca mulatta*)
812 with implanted head posts and recording chambers. All procedures were performed in
813 accordance with the guidelines of the University of Pennsylvania Institutional Animal Care and
814 Use Committee and this study was approved under protocol 804222.

815

816 **The invariant delayed-match-to-sample (IDMS) task:**

817

818 All behavioral training and testing was performed using standard operant conditioning (juice
819 reward), head stabilization, and high-accuracy, infrared video eye tracking. Stimuli were
820 presented on an LCD monitor with an 85 Hz refresh rate using customized software
821 (<http://mworks-project.org>).

822

823 As an overview, the monkeys' task required an eye movement response to a specific location
824 when a target object appeared within a sequence of distractor images (Fig 2a). Objects were
825 presented across variation in the objects' position, size and background context (Fig 2b).
826 Monkeys viewed a fixed set of 20 images across switches in the identity of 4 target objects,
827 each presented at 5 identity-preserving transformations (Fig 2c). Monkeys were trained
828 extensively on the set of 20 images shown in Fig 2b before testing. We ran the task in short
829 blocks (~3 min) with a fixed target before another target was pseudorandomly selected. Our
830 design included two types of trials: cue trials and test trials (Fig 2a). Only test trials were
831 analyzed for this report.

832

833 Trials were initiated by the monkey fixating on a red dot (0.15°) in the center of a gray screen,
834 within a square window of $\pm 1.5^\circ$, followed by a 250 ms delay before a stimulus appeared. Cue
835 trials, which indicated the current target object, were presented at the beginning of each block
836 and after three subsequent trials with incorrect responses. To minimize confusion, cue trials
837 were designed to be distinct from test trials and began with the presentation of an image of each
838 object that was distinct from the images used on test trials (a large version of the object
839 presented at the center of gaze on a gray background; Fig 2a). Test trials, which are the focus
840 of this report, always began with a distractor image, and neural responses to this image were
841 discarded to minimize non-stationarities such as stimulus onset effects. Distractors were drawn
842 randomly from a pool of 15 possible images within each block without replacement until each
843 distractor was presented once on a correct trial, and the images were then re-randomized. On
844 most trials, a random number of 1-6 distractors were presented, followed by a target match (Fig
845 2a). On a small fraction of trials, 7 distractors were shown, and the monkey was rewarded for
846 fixating through all distractors. Each stimulus was presented for 400 ms (or until the monkeys'
847 eyes left the fixation window) and was immediately followed by the presentation of the next
848 stimulus. Following the onset of a target match image, monkeys were rewarded for making a
849 saccade to a response target within a window of 75 – 600 ms to receive a juice reward. In
850 monkey 1 this target was positioned 10 degrees above fixation; in monkey 2 it was 10 degrees
851 below fixation. If 400 ms following target onset had elapsed and the monkey had not moved its
852 eyes, a distractor stimulus was immediately presented. If the monkey continued fixating beyond
853 the required reaction time, the trial was considered a "miss". False alarms were differentiated
854 from fixation breaks via a comparison of the monkeys' eye movements with the characteristic
855 pattern of eye movements on correct trials: false alarms were characterized by the eyes leaving
856 the fixation window via its top (monkey 1) or bottom (monkey 2) outside the allowable correct

857 response period and traveling more than 0.5 degrees whereas fixation breaks were
858 characterized by the eyes leaving the fixation window in any other way. Within each block, 4
859 repeated presentations of the 20 images were collected, and a new target object was then
860 pseudorandomly selected. Following the presentation of all 4 objects as targets, the targets
861 were re-randomized. At least 20 repeats of each condition were collected. Overall, monkeys
862 performed this task with high accuracy. Disregarding fixation breaks (monkey 1: 11% of trials,
863 monkey 2: 8% of trials), percent correct on the remaining trials was as follows: monkey 1: 96%
864 correct, 1% false alarms, and 3% misses; monkey 2: 87% correct, 3% false alarms, and 10%
865 misses.

866
867

868 **Neural recording:**

869

870 The activity of neurons in IT was recorded via a single recording chamber in each monkey.
871 Chamber placement was guided by anatomical magnetic resonance images in both monkeys,
872 and in one monkey, Brainsight neuronavigation (<https://www.rogue-research.com/>). The region
873 of IT recorded was located on the ventral surface of the brain, over an area that spanned 4 mm
874 lateral to the anterior middle temporal sulcus and 15-19 mm anterior to the ear canals. Neural
875 activity was largely recorded with 24-channel U probes (Plexon, Inc) with linearly arranged
876 recording sites spaced with 100 mm intervals, with a handful of units recorded with single
877 electrodes (Alpha Omega, glass-coated tungsten). Continuous, wideband neural signals were
878 amplified, digitized at 40 kHz and stored using the OmniPlex Data Acquisition System (Plexon).
879 Spike sorting was done manually offline (Plexon Offline Sorter). At least one candidate unit was
880 identified on each recording channel, and 2-3 units were occasionally identified on the same
881 channel. Spike sorting was performed blind to any experimental conditions to avoid bias. A
882 multi-channel recording session was included in the analysis if the animal performed the task
883 until the completion of 20 correct trials per stimulus condition, there was no external noise
884 source confounding the detection of spike waveforms, and the session included a threshold
885 number of task modulated units (>4 on 24 channels). The sample size (number of units
886 recorded) was chosen to approximately match our previous work (Pagan and Rust, 2014a;
887 Pagan et al., 2016; Pagan et al., 2013).

888

889 For all the analyses presented in this paper except Fig 4b,d, Fig 8, and Fig 9c, we measured
890 neural responses by counting spikes in a window that began 80 ms after stimulus onset and
891 ended at 250 ms. On 1.9% of all correct target match presentations, the monkeys had reaction
892 times faster than 250 ms, and those instances were excluded from analysis such that spikes
893 were only counted during periods of fixation. When combining the units recorded across
894 sessions into a larger pseudopopulation, we screened for units that met three criteria. First, units
895 had to be modulated by our task, as quantified by a one-way ANOVA applied to our neural
896 responses (80 conditions * 20 repeats) with $p < 0.01$. Second, we applied a loose criterion on
897 recording stability, as quantified by calculating the variance-to-mean for each unit (computed by
898 fitting the relationship between the mean and variance of spike count across the 80 conditions),
899 and eliminating units with a variance-to-mean ratio > 5 . Finally, we applied a loose criterion on
900 unit recording isolation, quantified by calculating the signal-to-noise ratio (SNR) of the waveform
901 (as the difference between the maximum and minimum points of the average waveform, divided
902 by twice the standard deviation across the differences between each waveform and the mean
903 waveform), and excluding (multi)units with an SNR < 2 . This yielded a pseudopopulation of 204
904 units (of 563 possible units), including 108 units from monkey 1 and 96 units from monkey 2.

905
906
907

Quantifying single-unit modulation magnitudes:

908 To quantify the degree to which individual units were modulated by different types of task
909 parameters (Fig 4b-d), we applied a bias-corrected procedure described in detail by (Pagan and
910 Rust, 2014b) and summarized here. Our measure of modulation is similar to a multi-way
911 ANOVA, with important extensions. Specifically, a two-way ANOVA applied to a unit's
912 responses (configured into a matrix of 4 targets * 20 images * 20 trials for each condition) would
913 parse the total response variance into two linear terms, a nonlinear interaction term, and an
914 error term. We make 3 extensions to the ANOVA analysis. First, an ANOVA returns measures
915 of variance (in units of spike counts squared) whereas we compute measures of standard
916 deviation (in units of spike count) such that our measures of modulation are intuitive (e.g.,
917 doubling firing rates causes signals to double as opposed to quadruple). Second, while the
918 linear terms of the ANOVA map onto our "visual" and "target identity" modulations (after
919 squaring), we split the ANOVA nonlinear interaction term into two terms, including target match
920 modulation (i.e. Fig 2c gray versus white) and all other nonlinear "residual" modulation. This
921 parsing is essential, as target match modulation corresponds to the signal for the IDMS task
922 whereas the other types of modulations are not. Finally, raw ANOVA values are biased by trial-
923 by-trial variability (which the ANOVA addresses by computing the probability that each term is
924 higher than chance given this noise) whereas our measures of modulation are bias-corrected to
925 provide an unbiased estimate of modulation magnitude.

926
927 The procedure begins by developing an orthonormal basis of 80 vectors designed to capture all
928 types of modulation with intuitive groupings. The number of each type is imposed by the
929 experimental design. This basis \mathbf{b} included vectors \mathbf{b}_i that reflected 1) the grand mean spike
930 count across all conditions (\mathbf{b}_1 , 1 dimension), 2) whether the object in view was a target or a
931 distractor (\mathbf{b}_2 , 1 dimension), 3) visual image identity ($\mathbf{b}_3 - \mathbf{b}_{21}$, 19 dimensions), 4) target object
932 identity ($\mathbf{b}_{22} - \mathbf{b}_{24}$, 3 dimensions), and 5) "residual", nonlinear interactions between target and
933 object identity not captured by target match modulation ($\mathbf{b}_{25} - \mathbf{b}_{80}$, 56 dimensions). A Gram-
934 Schmidt process was used to convert an initially designed set of vectors into an orthonormal
935 basis.

936
937 Because this basis spans the space of all possible responses for our task, each trial-averaged
938 vector of spike count responses to the 80 experimental conditions \mathbf{R} can be re-expressed as a
939 weighted sum of these basis vectors. To quantify the amounts of each type of modulation
940 reflected by each unit, we began by computing the squared projection of each basis vector
941 \mathbf{b}_i and \mathbf{R} . An analytical bias correction, described and verified in (Pagan and Rust, 2014b), was
942 then subtracted from this value:

943

$$944 \quad (8) \quad w_i^2 = (\mathbf{R} \cdot \mathbf{b}_i^T)^2 - \frac{\sigma_i^2 \cdot (\mathbf{b}_i^T)^2}{m}$$

945

946 where σ_i^2 indicates the trial variance, averaged across conditions (n=80), and where m indicates
947 the number of trials (m=20). When more than one dimension existed for a type of modulation,
948 we summed values of the same type. Next, we applied a normalization factor (1/(n-1)) to convert
949 these summed values into variances. Finally, we computed the square root of these quantities
950 to convert them into modulation measures that reflected the number of spike count standard
951 deviations around each unit's grand mean spike count.

952

953 Target match modulation was thus computed as:

954

$$955 \quad (9) \quad \sigma_{TM} = \sqrt{\frac{1}{n-1} \cdot w_2^2}$$

956

957 visual modulation was computed as:

958

$$959 \quad (10) \quad \sigma_{Vis} = \sqrt{\frac{1}{n-1} \cdot \sum_{i=3}^{21} w_i^2}$$

960

961 target identity modulation was computed as:

962

$$963 \quad (11) \quad \sigma_{TI} = \sqrt{\frac{1}{n-1} \cdot \sum_{i=22}^{24} w_i^2}$$

964

965 and residual modulation was computed as:

966

$$967 \quad (12) \quad \sigma_{res} = \sqrt{\frac{1}{n-1} \cdot \sum_{i=25}^{80} w_i^2}$$

968

969

970 When estimating modulation population means (Fig 4b,c), the bias-corrected squared values
971 were averaged across units before taking the square root. Because these measures were not
972 normally distributed, standard error about the mean was computed via a bootstrap procedure.
973 On each iteration of the bootstrap (across 1000 iterations), we randomly sampled values from
974 the modulation values for each unit in the population, with replacement. Standard error was
975 computed as the standard deviation across the means of these newly created populations.

976

977 To quantify the sign of the modulation corresponding to whether an image was presented as a
978 target match versus as a distractor (Fig 7a,b), we calculated a target match modulation index for
979 each unit by computing its mean spike count response to target matches and to distractors, and
980 computing the ratio of their difference and their sum.

981

982

983 **Population performance:**

984

985 To determine the performance of the IT population at classifying target matches versus
986 distractors, we applied two types of decoders: a Fisher Linear Discriminant (a linear decoder)
987 and Maximum Likelihood decoder (a nonlinear decoder) using approaches that are described
988 previously in detail (Pagan et al., 2013) and are summarized here.

989

990 When applied to the pseudopopulation data (Fig 5b, Fig 7b), all decoders were cross-validated
991 with the same resampling procedure. On each iteration of the resampling, we randomly shuffled
992 the trials for each condition and for each unit, and (for numbers of units less than the full
993 population size) randomly selected units. On each iteration, 18 trials from each condition were
994 used for training the decoder, 1 trial was used to determine a value for regularization, and 1 trial
995 from each condition was used for cross-validated measurement of performance.

996

997 To ensure that decoder performance was not biased by unequal numbers of target matches and
998 distractors, on each iteration of the resampling we included 20 target match conditions and 20
999 (of 60 possible) distractor conditions. Each set of 20 distractors was selected to span all
1000 possible combinations of mismatched object and target identities (e.g. objects 1, 2, 3, 4 paired
1001 with targets 4, 3, 2, 1), of which there are 9 possible sets. To compute proportion correct, a
1002 mean performance value was computed on each resampling iteration by averaging binary
1003 performance outcomes across the 9 possible sets of target matches and distractors, each which
1004 contained 40 test trials. Mean and standard error of performance was computed as the mean
1005 and standard deviation of performance across 1000 resampling iterations. Standard error thus
1006 reflected the variability due to the specific trials assigned to training and testing and, for
1007 populations smaller than the full size, the specific units chosen.

1008
1009

1010 *Fisher Linear Discriminant:*

1011
1012

The general form of a linear decoding axis is:

1013

$$1014 (1) f(x) = \mathbf{w}^T x + b,$$

1015

1016 where \mathbf{w} is an N-dimensional vector (where N is the number of units) containing the linear
1017 weights applied to each unit, and b is a scalar value. We fit these parameters using a Fisher
1018 Linear Discriminant (FLD), where the vector of linear weights was calculated as:

1019

$$1020 (2) \mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2)$$

1021

and b was calculated as:

1022
1023

$$1024 (3) b = \mathbf{w} \cdot \frac{1}{2}(\mu_1 + \mu_2) = \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 - \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2$$

1025

1026 Here μ_1 and μ_2 are the means of the two classes (target matches and distractors, respectively)
1027 and the mean covariance matrix is calculated as:

1028

$$1029 (4) \Sigma = \frac{\Sigma_1 + \Sigma_2}{2}$$

1030

1031 where Σ_1 and Σ_2 are the regularized covariance matrices of the two classes. These covariance
1032 matrices were computed using a regularized estimate equal to a linear combination of the
1033 sample covariance and the identity matrix I (Pagan et al., 2016):

1034

$$1035 (5) \Sigma_i = \gamma \Sigma_i + (1 - \gamma) \cdot I$$

1036

1037 We determined γ by exploring a range of values from 0.01 to 0.99, and we selected the value
1038 that maximized average performance across all iterations, measured with the cross-validation
1039 “regularization” trials set aside for this purpose (see above). We then computed performance for
1040 that value of γ with separately measured “test” trials, to ensure a fully cross-validated measure.
1041 Because this calculation of the FLD parameters incorporates the off-diagonal terms of the
1042 covariance matrix, FLD weights are optimized for both the information conveyed by individual
1043 units as well as their pairwise interactions.

1044

1045 To compare FLD performance on correct versus error trials (Fig 6, 7d), we used the same
1046 methods described above with the following modifications. First, the analysis was applied to the
1047 simultaneously recorded data within each session, and the correlation structure on each trial
1048 was kept intact on each resampling iteration. Second, when more than 24 units were available,
1049 a subset of 24 units were selected as those with the most task modulation, quantified via the p-
1050 value of a one-way ANOVA applied to each unit's responses (80 conditions * 20 repeats).
1051 Finally, on each resampling iteration, each error trial was randomly paired with a correct trial of
1052 the same condition and cross-validated performance was performed exclusively for these pairs
1053 of correct and error responses. As was the case for the pseudopopulation analysis, training
1054 was performed exclusively on correct trials. A mean performance value was computed on each
1055 resampling iteration by averaging binary performance outcomes across all possible error trials
1056 and their condition-matched correct trial pairs, and averaging across different recording
1057 sessions. Mean and standard error of performance was computed as the mean and standard
1058 deviation of performance across 100 resampling iterations. Standard error thus reflected error in
1059 a manner similar to the pseudopopulation analysis - the variability due to the specific trials
1060 assigned to training and testing and, for populations smaller than the full size, the specific units
1061 chosen.

1062
1063 In the case of the ranked-FLD (Fig 5c), all units were considered on each resampling iteration,
1064 and weights were computed for each unit (with the training data) as described by Equation 2.
1065 Weights were then ranked by their magnitude (the absolute values of the signed quantities) and
1066 the top N units were selected for different population size N. Finally, both the weights and the
1067 threshold were recalculated before cross-validated testing with the training data. In the case of
1068 the signed versions of the FLD (which isolated the analysis to target matched enhanced or
1069 suppressed units, Fig 7c-d), the process was similar in that all units were considered on each
1070 resampling iteration and weights were computed for each unit (with the training data) as
1071 described by Equation 2. Weights were then isolated to all of those that were positive "E units"
1072 or all that were negative "S units". Finally, the weights and the threshold were recalculated
1073 before cross-validated testing with the training data.

1074
1075

1076 *Maximum likelihood decoder:*

1077

1078 As a measure of total IT target match information (combined linear and nonlinear), we
1079 implemented the maximum likelihood decoder (Fig 5b) introduced in our previous work (Pagan
1080 et al., 2016; Pagan et al., 2013). We began by using the set of training trials to compute the
1081 average response r_{uc} of each unit u to each of the 40 conditions c . We then computed the
1082 likelihood that a test response k was generated from a particular condition as a Poisson-
1083 distributed variable:

1084

1085
$$(7) \text{lik}_{u,c}(k) = \frac{(r_{uc})^k \cdot e^{-r_{uc}}}{k!}$$

1086

1087 The likelihood that a population response vector was generated in response to each condition
1088 was then computed as the product of the likelihoods of the individual units. Next, we computed
1089 the likelihood that each test vector arose from the category target match as compared to the
1090 category distractor as the product of the likelihoods across the conditions within each category.
1091 We assigned the population response to the category with the maximum likelihood, and we

1092 computed performance as the fraction of trials in which the classification was correct based on
1093 the true labels of the test data.

1094

1095

1096 **Representational similarity:**

1097

1098 Before computing representational similarity (Fig 9a), the responses of each unit were z-
1099 normalized to have a mean of zero and standard deviation of 1. To compute measures of the
1100 representational similarity between pairs of population response vectors, the 20 repeated trials
1101 for each (of 80) experimental conditions were randomly split into two sets of 10 trials, and the
1102 average population response vector was computed. To obtain measures of the noise ceiling,
1103 Pearson correlation was computed between many random splits of the data for each of the 80
1104 conditions. The mean across 1000 random splits was computed for each condition and the
1105 values were averaged across the splits as well as the 4 target conditions for each image,
1106 resulting in 20 correlations values (1 for each image). Fig 9b “Same image and condition”
1107 depicts the mean and standard deviation across the 20 images. Measures of the
1108 representational similarity between different conditions were computed in a comparable way, by
1109 also selecting 10 (of 20) trials before computing the mean population response vectors. To
1110 measure the representational similarity between the same objects presented at different
1111 transformations, Pearson correlation was computed for all possible pairs of the 5 images
1112 corresponding to each object under the conditions of a fixed target. A mean value was
1113 computed as the average across 1000 random splits and the pairwise comparison between 1
1114 image and other images containing the same object for each of 20 images, and Fig 9b “Different
1115 transforms.” depicts the mean and standard deviation across the 20 images. Fig 9b “Different
1116 objects” was computed in a similar manner, but for all possible pairs of one image and the
1117 images containing other objects. Finally, Fig 9b “Match versus distractor” was computed in a
1118 similar manner, but for all possible pairs of one image presented as a target match (viewing an
1119 image while looking for that object as a target) and the three distractor conditions (the three
1120 other targets). The same procedures were carried out for Fig 9c and 11b.

1121

1122 **Simulations:**

1123

1124 To better understand our results, we performed a number of data-based simulations (Fig 11b).
1125 Each simulation began by computing the bias-corrected weights for each unit as described
1126 above. For the “replication” simulation, we rectified bias-corrected modulations that fell below
1127 zero, recomputed the noise-corrected mean spike count responses for each condition, and
1128 generated trial variability with an independent Poisson process.

1129 For the “multiplicative, same-sign” simulation, we replaced the target match modulation for each
1130 unit with an amount that ensured the population total was distributed proportional to each unit’s
1131 total visual modulation (Equation 10), and always reflected as target match enhancement. For
1132 the “uniform, mixed-sign” simulation, we replaced each unit’s target match modulation with the
1133 same amount, reflected with the sign determined in the original data.

1134

1135 **Statistical tests:**

1136

1137 When comparing performance between the FLD and maximum likelihood classifier (Fig 5b), we
1138 reported *P* values as an evaluation of the probability that differences were due to chance. We

1139 calculated P values as the fraction of resampling iterations on which the difference was flipped
1140 in sign relative to the actual difference between the means of the full data set (for example, if the
1141 mean of decoding measure 1 was larger than the mean of decoding measure 2, the fraction of
1142 iterations in which the mean of measure 2 was larger than the mean of measure 1).

1143
1144 When evaluating whether each unit had a statistically different response to target matches as
1145 compared to distractors (Fig 7a-b, light bars), we recomputed each unit's modulation index by
1146 resampling trials with replacement on $n = 1000$ resampling iterations. A unit was considered as
1147 statistically significant if its resampled modulation indices were flipped in sign from the unit's
1148 actual modulation index less than 0.01% of the resampling iterations. When evaluating whether
1149 the single unit modulation indices (Fig 7a-b) were significantly different from zero, we reported P
1150 values as computed by a Wilcoxon sign rank test. When evaluating whether the single unit
1151 modulation indices computed without repeated distractors (Fig 7b) were significantly different
1152 from modulation indices computed with repeated distractors, we reported P values computed via
1153 a matched t test.

1154
1155 When comparing the representational similarity of different groupings of the IT population
1156 response (Fig 7b), we computed a mean Pearson correlation value for each of the 20 images
1157 (as described above), and reported P values as the probability that the observed differences in
1158 means across the 20 images were due to chance via a two-sample t -test.

1159

1160 **Animal husbandry, enrichment, and care:**

1161
1162 Monkeys received a nutritionally balanced diet of biscuits as well as daily supplements of fruit
1163 and nuts. Monkeys were housed in Allentown cages with space that exceeded the minimums
1164 described in the "Guide for Care and Use of Laboratory Animals". Additionally, monkeys had
1165 periodic access to larger playcages that included a variety of enrichment items, such as swings.
1166 Monkeys were also provided daily enrichment by social housing when possible and through the
1167 introduction of toys, games, and puzzles that involved manipulation to receive food treats. To
1168 maintain task motivation, access to water was regulated prior to experimental sessions.
1169 Monkeys received a minimum 20 mL/kg of water a day five days a week and a minimum of 40
1170 mL/kg on the other two days. When off study, animals were allowed unrestricted access to
1171 water. Animals on regulated access were monitored daily for health status and hydration. Daily
1172 hydration status was assessed by body weight, skin turgor, urine and fecal output, and overall
1173 demeanor. Following this study, both animals were used in one other neurophysiology study.
1174 Following the conclusion of the second study, both animals were euthanized in a manner
1175 consistent with the recommendations of the Panel on Euthanasia of the American Veterinary
1176 Medical Association, including sedation followed by the introduction of the euthanasia solution
1177 Euthasol.

1178

1179

1180 **S1 Dataset.** *IT neural data.* The data include the spike count responses recorded from each
1181 monkey, organized into 5-dimensional matrices as units (monkey 1: $n = 108$; monkey 2: $n = 96$)
1182 x targets ($n = 4$) x objects ($n = 4$) x transformations ($n = 5$) x trials ($n = 20$). Spikes were counted
1183 from 80 to 250 ms, and were extracted from trials with correct responses.

1184

1185

1186

1187

1188 **References:**

- 1189
- 1190 Barak, O., Rigotti, M., and Fusi, S. (2013). The sparseness of mixed selectivity neurons controls
1191 the generalization-discrimination trade-off. *J Neurosci* *33*, 3844-3856.
- 1192 Bichot, N.P., Rossi, A.F., and Desimone, R. (2005). Parallel and serial neural mechanisms for
1193 visual search in macaque area V4. *Science* *308*, 529-534.
- 1194 Chelazzi, L., Duncan, J., Miller, E.K., and Desimone, R. (1998). Responses of neurons in
1195 inferior temporal cortex during memory-guided visual search. *J Neurophysiology* *80*, 2918-2940.
- 1196 Chelazzi, L., Miller, E.K., Duncan, J., and Desimone, R. (1993). A neural basis for visual search
1197 in inferior temporal cortex. *Nature* *363*, 345-347.
- 1198 Chelazzi, L., Miller, E.K., Duncan, J., and Desimone, R. (2001). Responses of neurons in
1199 macaque area V4 during memory-guided visual search. *Cereb Cortex* *11*, 761-772.
- 1200 DiCarlo, J.J., Zoccolan, D., and Rust, N.C. (2012). How does the brain solve visual object
1201 recognition? *Neuron* *73*, 415-434.
- 1202 Engel, T.A., and Wang, X.J. (2011). Same or different? A neural circuit mechanism of similarity-
1203 based pattern match decision making. *J Neurosci* *31*, 6982-6996.
- 1204 Eskandar, E.N., Richmond, B.J., and Optican, L.M. (1992). Role of inferior temporal neurons in
1205 visual memory I. Temporal encoding of information about visual images, recalled images, and
1206 behavioral context. *Journal of Neurophysiology* *68*, 1277-1295.
- 1207 Freedman, D.J., and Assad, J.A. (2009). Distinct encoding of spatial and nonspatial visual
1208 information in parietal cortex. *J Neurosci* *29*, 5671-5680.
- 1209 Gibson, J.R., and Maunsell, J.H.R. (1997). The sensory modality specificity of neural activity
1210 related to memory in visual cortex. *Journal of Neurophysiology* *78*, 1263-1275.
- 1211 Haenny, P.E., Maunsell, J.H., and Schiller, P.H. (1988). State dependent activity in monkey
1212 visual cortex. II. Retinal and extraretinal factors in V4. *Exp Brain Res* *69*, 245-259.
- 1213 Kobak, D., Brendel, W., Constantinidis, C., Feierstein, C.E., Kepecs, A., Mainen, Z.F., Qi, X.L.,
1214 Romo, R., Uchida, N., and Machens, C.K. (2016). Demixed principal component analysis of
1215 neural population data. *Elife* *5*.
- 1216 Kosai, Y., El-Shamayleh, Y., Fyall, A.M., and Pasupathy, A. (2014). The role of visual area V4 in
1217 the discrimination of partially occluded shapes. *J Neurosci* *34*, 8570-8584.
- 1218 Kriegeskorte, N., and Kievit, R.A. (2013). Representational geometry: integrating cognition,
1219 computation, and the brain. *Trends Cogn Sci* *17*, 401-412.
- 1220 Leuschow, A., Miller, E.K., and Desimone, R. (1994). Inferior temporal mechanisms for invariant
1221 object recognition. *Cerebral Cortex* *5*, 523-531.

- 1222 Mante, V., Sussillo, D., Shenoy, K.V., and Newsome, W.T. (2013). Context-dependent
1223 computation by recurrent dynamics in prefrontal cortex. *Nature* *503*, 78-84.
- 1224 Maunsell, J.H., Sclar, G., Nealey, T.A., and DePriest, D.D. (1991). Extraretinal representations
1225 in area V4 in the macaque monkey. *Vis Neurosci* *7*, 561-573.
- 1226 Meister, M.L., Hennig, J.A., and Huk, A.C. (2013). Signal multiplexing and single-neuron
1227 computations in lateral intraparietal area during decision-making. *J Neurosci* *33*, 2254-2267.
- 1228 Miller, E.K., and Desimone, R. (1994). Parallel neuronal mechanisms for short-term memory.
1229 *Science* *263*, 520-522.
- 1230 Mruczek, R.E., and Sheinberg, D.L. (2007). Activity of inferior temporal cortical neurons predicts
1231 recognition choice behavior and recognition time during visual search. *J Neurosci* *27*, 2825-
1232 2836.
- 1233 Pagan, M., and Rust, N.C. (2014a). Dynamic target match signals in perirhinal cortex can be
1234 explained by instantaneous computations that act on dynamic input from inferotemporal cortex.
1235 *J Neurosci* *34*, 11067-11084.
- 1236 Pagan, M., and Rust, N.C. (2014b). Quantifying the signals contained in heterogeneous neural
1237 responses and determining their relationships with task performance. *J Neurophysiol* *112*, 1584-
1238 1598.
- 1239 Pagan, M., Simoncelli, E.P., and Rust, N.C. (2016). Neural Quadratic Discriminant Analysis:
1240 Nonlinear Decoding with V1-Like Computation. *Neural Comput*, 1-29.
- 1241 Pagan, M., Urban, L.S., Wohl, M.P., and Rust, N.C. (2013). Signals in inferotemporal and
1242 perirhinal cortex suggest an untangling of visual target information. *Nature neuroscience* *16*,
1243 1132-1139.
- 1244 Raposo, D., Kaufman, M.T., and Churchland, A.K. (2014). A category-free neural population
1245 supports evolving demands during decision-making. *Nature neuroscience* *17*, 1784-1792.
- 1246 Rigotti, M., Barak, O., Warden, M.R., Wang, X.J., Daw, N.D., Miller, E.K., and Fusi, S. (2013).
1247 The importance of mixed selectivity in complex cognitive tasks. *Nature* *497*, 585-590.
- 1248 Rishel, C.A., Huang, G., and Freedman, D.J. (2013). Independent category and spatial
1249 encoding in parietal cortex. *Neuron* *77*, 969-979.
- 1250 Rust, N.C., and DiCarlo, J.J. (2010). Selectivity and tolerance ("invariance") both increase as
1251 visual information propagates from cortical area V4 to IT. *J Neurosci* *30*, 12978-12995.
- 1252 Woloszyn, L., and Sheinberg, D.L. (2009). Neural dynamics in inferior temporal cortex during a
1253 visual working memory task. *J Neurosci* *29*, 5494-5507.
- 1254 Zoccolan, D., Kouh, M., Poggio, T., and DiCarlo, J.J. (2007). Trade-off between object
1255 selectivity and tolerance in monkey inferotemporal cortex. *J Neurosci* *27*, 12292-12307.
- 1256

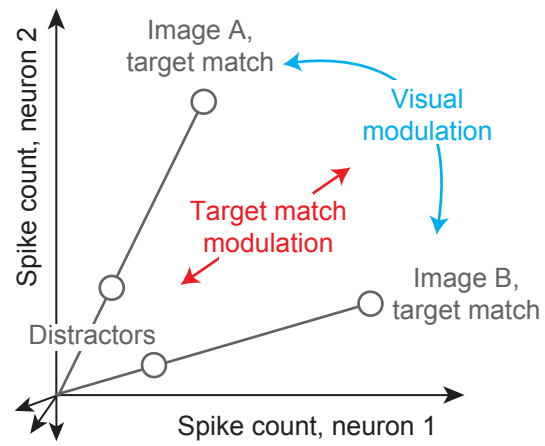


Figure 1

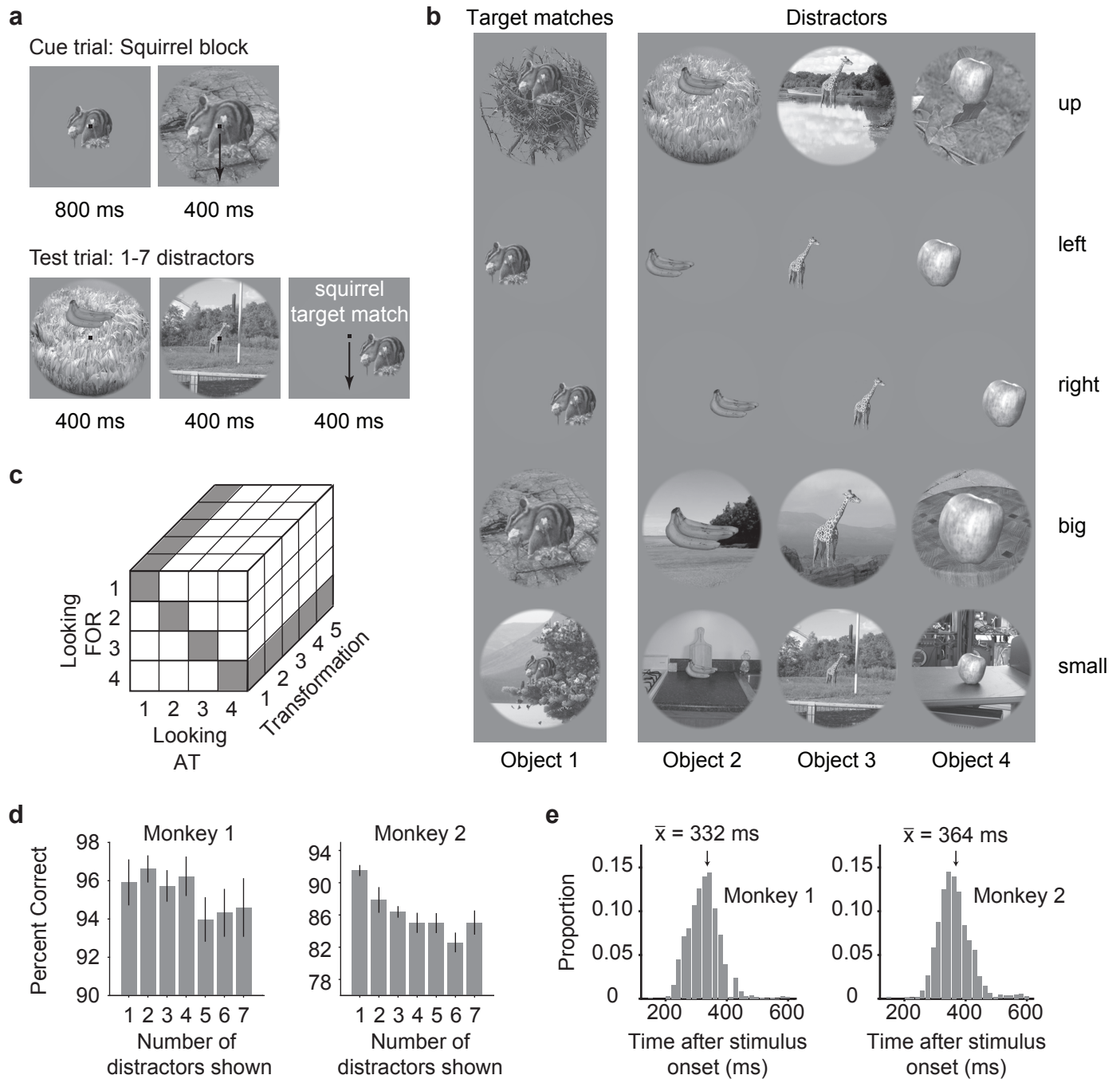


Figure 2

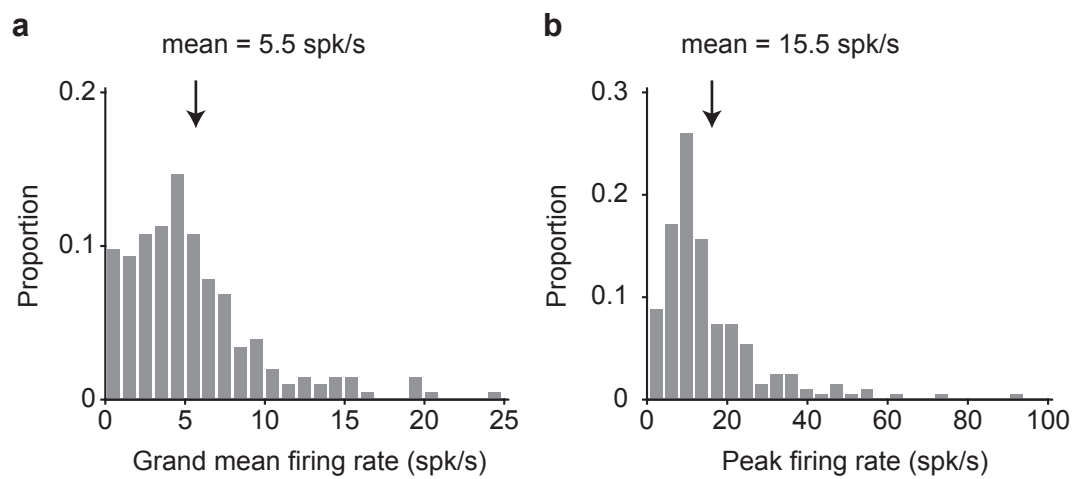


Figure 3

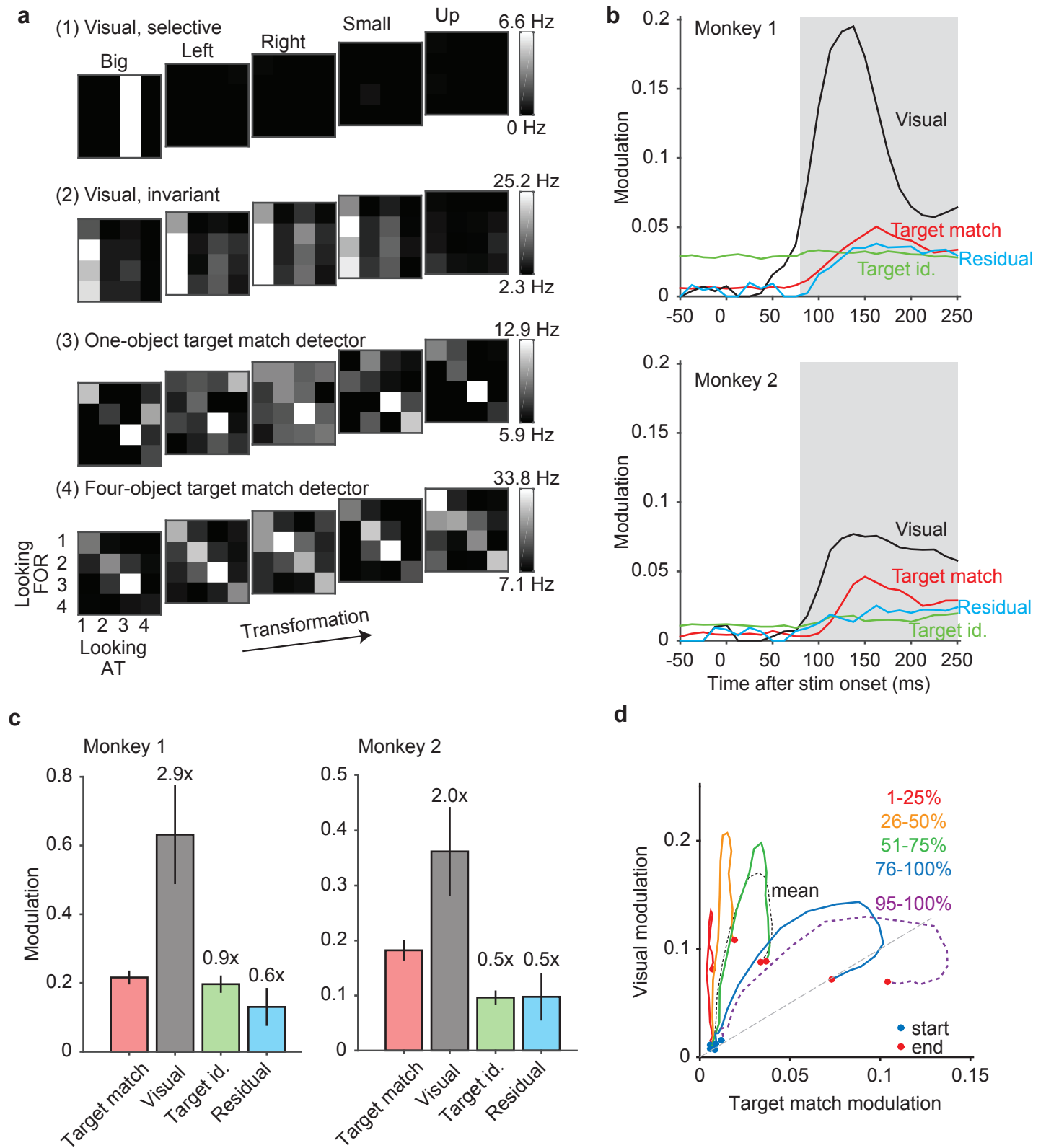


Figure 4

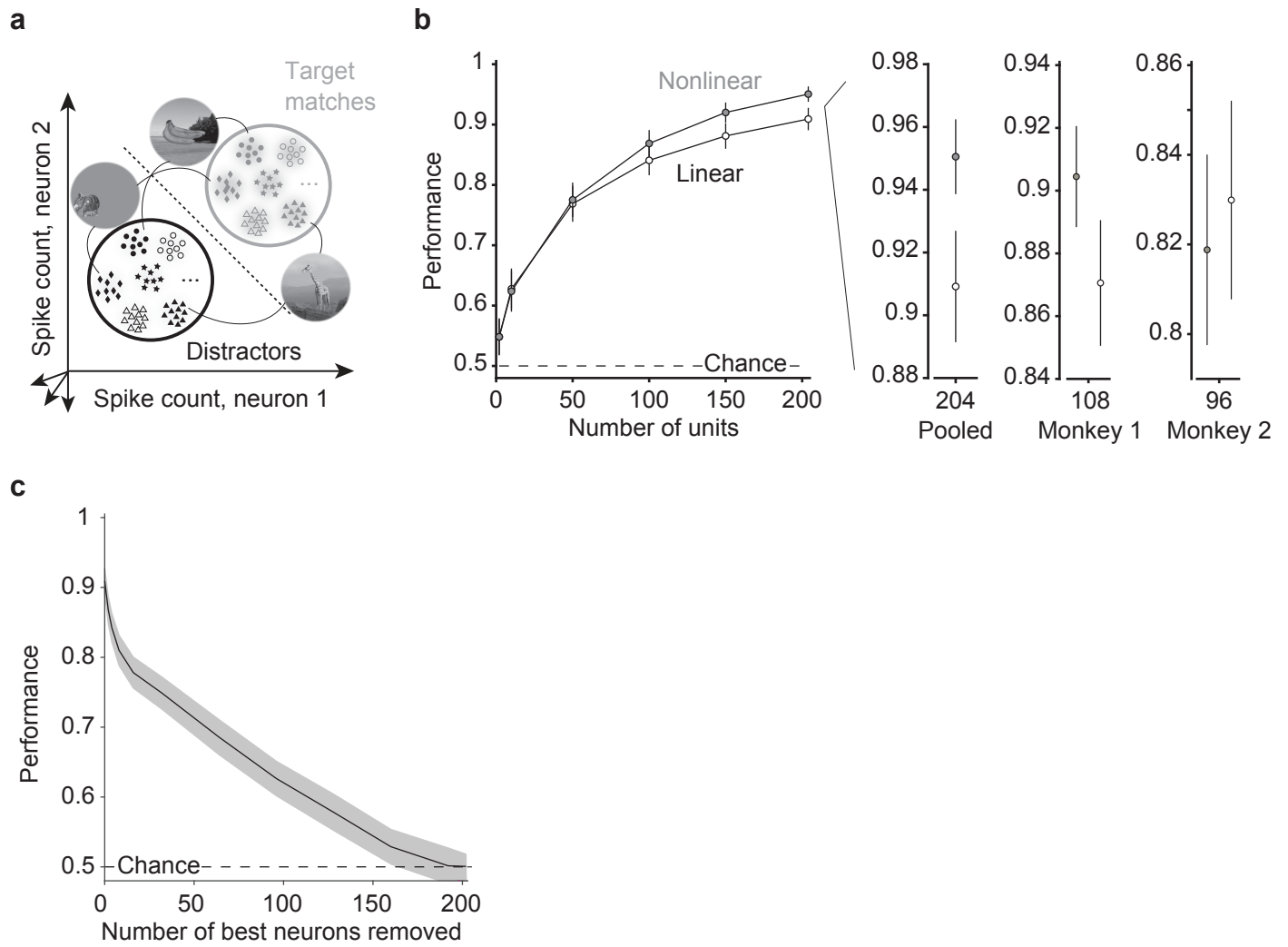


Figure 5

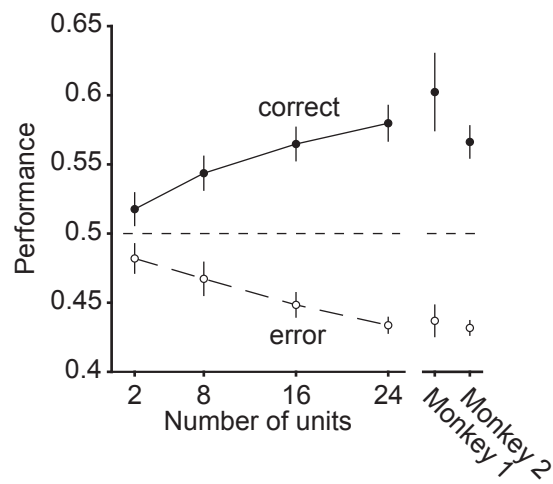


Figure 6

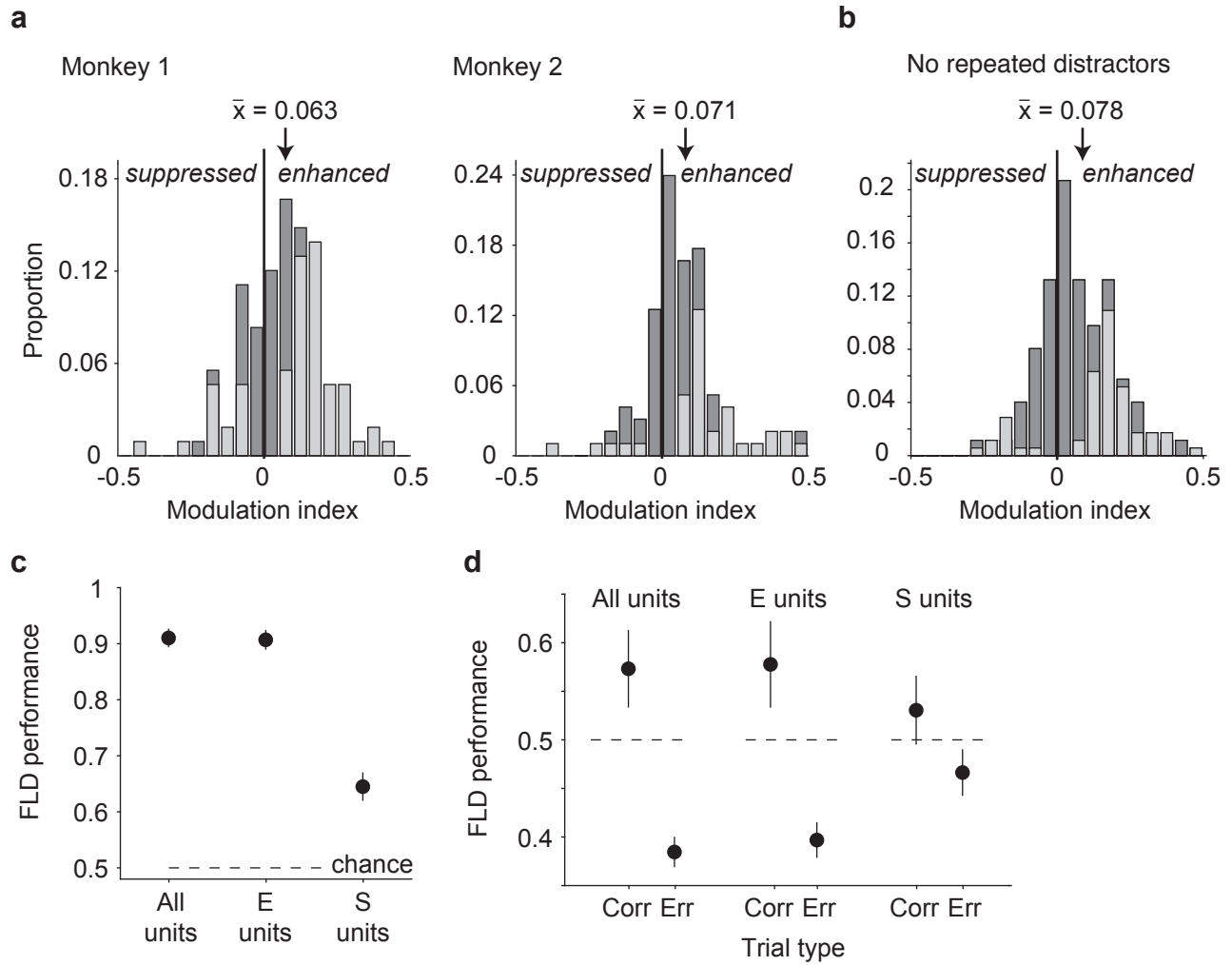


Figure 7

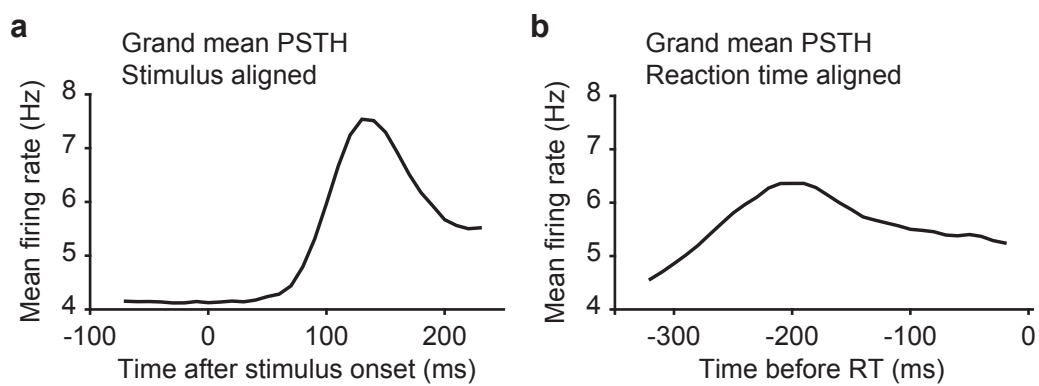


Figure 8

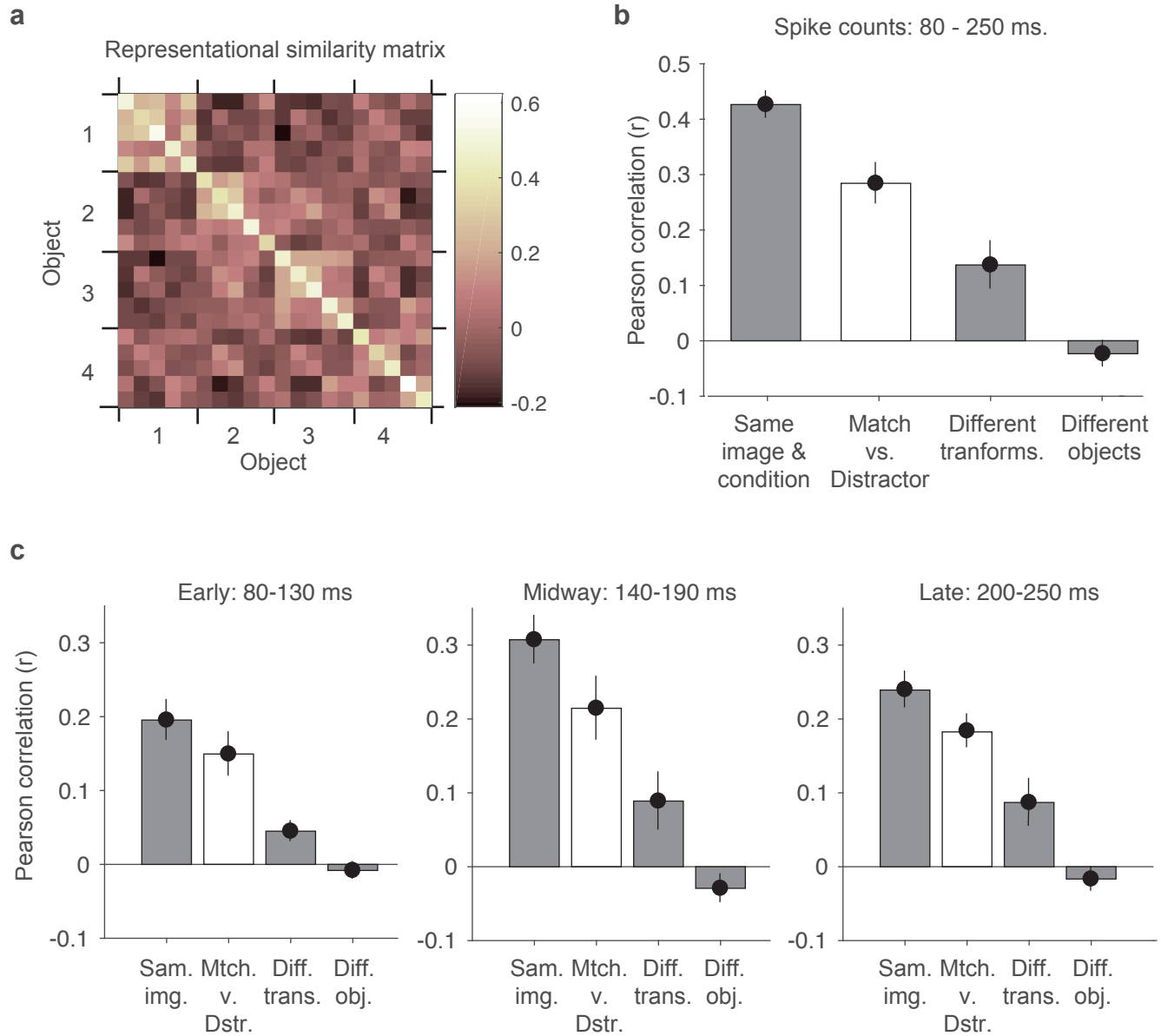


Figure 9

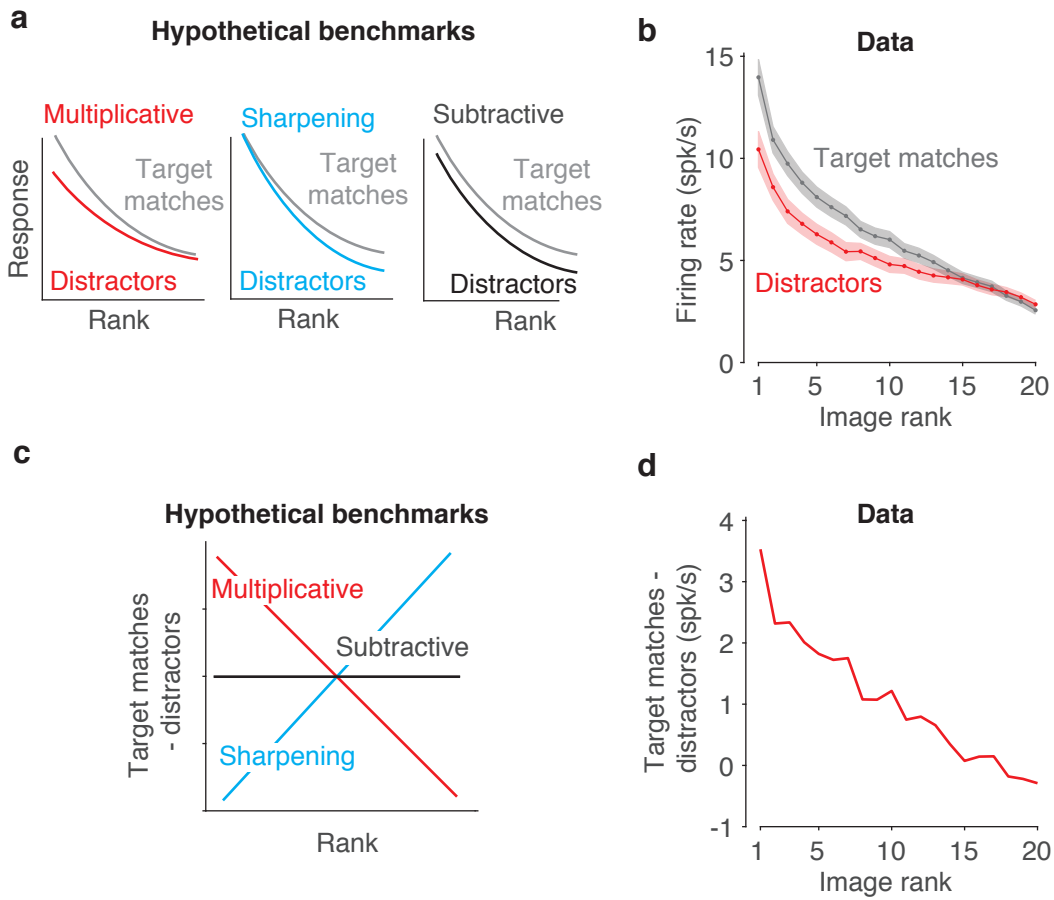


Figure 10

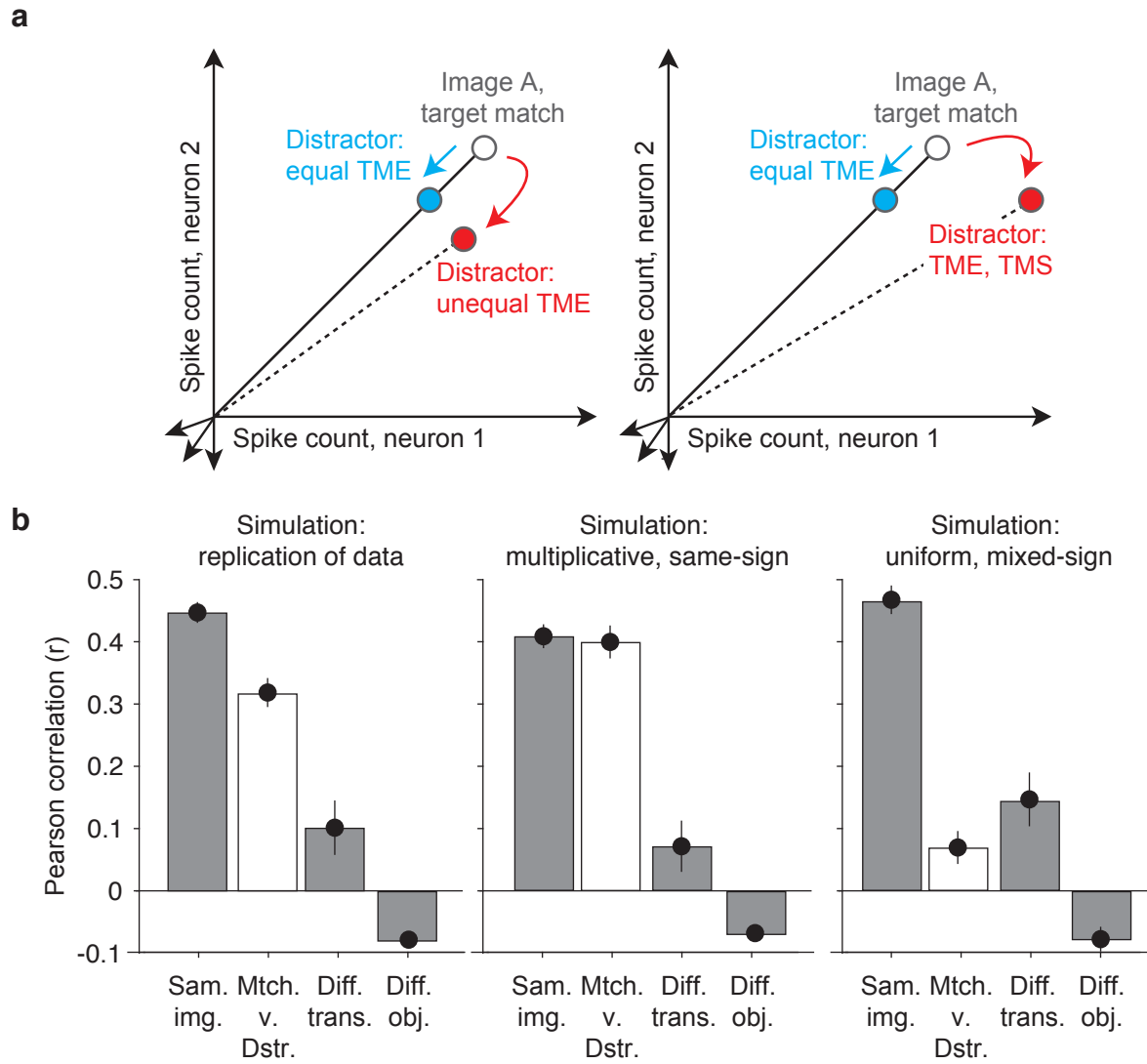


Figure 11