

## RESEARCH ARTICLE | *Sensory Processing*

# The integration of visual and target signals in V4 and IT during visual object search

 **Noam Roth** and  **Nicole C. Rust**

*Department of Psychology, University of Pennsylvania, Philadelphia, Pennsylvania*

Submitted 14 January 2019; accepted in final form 15 October 2019

**Roth N, Rust NC.** The integration of visual and target signals in V4 and IT during visual object search. *J Neurophysiol* 122: 2522–2540, 2019. First published October 16, 2019; doi:10.1152/jn.00024.2019.—Searching for a specific visual object requires our brain to compare the items in view with a remembered representation of the sought target to determine whether a target match is present. This comparison is thought to be implemented, in part, via the combination of top-down modulations reflecting target identity with feed-forward visual representations. However, it remains unclear whether top-down signals are integrated at a single locus within the ventral visual pathway (e.g., V4) or at multiple stages [e.g., both V4 and inferotemporal cortex (IT)]. To investigate, we recorded neural responses in V4 and IT as rhesus monkeys performed a task that required them to identify when a target object appeared across variation in position, size, and background context. We found nonvisual, task-specific signals in both V4 and IT. To evaluate whether V4 was the only locus for the integration of top-down signals, we evaluated several feed-forward accounts of processing from V4 to IT, including a model in which IT preferentially sampled from the best V4 units and a model that allowed for nonlinear IT computation. IT task-specific modulation was not accounted for by any of these feed-forward descriptions, suggesting that during object search, top-down signals are integrated directly within IT.

**NEW & NOTEWORTHY** To find specific objects, the brain must integrate top-down, target-specific signals with visual information about objects in view. However, the exact route of this integration in the ventral visual pathway is unclear. In the first study to systematically compare V4 and inferotemporal cortex (IT) during an invariant object search task, we demonstrate that top-down signals found in IT cannot be described as being inherited from V4 but rather must be integrated directly within IT itself.

invariant object recognition; object search; top-down signals; ventral visual pathway; visual attention

## INTRODUCTION

Finding a sought object requires our brains to perform at least two nontrivial computations. First, we must determine the identities of the objects in view, across variation in details such as their position, size, and background context. Second, we must compare this visual representation (of what we are looking at) with a remembered representation (of what we are looking for) to determine whether our target is in view. Con-

siderable evidence suggests that computations in the primate ventral visual pathway, including brain areas V1, V2, V4, and IT, support the process of invariant object recognition (reviewed by DiCarlo et al. 2012). Within V4 and IT, many neurons are also modulated by information about target identity as well as whether an image is a target match (Bichot et al. 2005; Chelazzi et al. 1998, 2001; Eskandar et al. 1992; Gibson and Maunsell 1997; Haenny et al. 1988; Kosai et al. 2014; Lueschow et al. 1994; Maunsell et al. 1991; Pagan et al. 2013; Roth and Rust 2018a). However, the route by which these signals arrive in V4 and IT remains unclear.

Here we present two proposals for how top-down signals reflecting the identity of a sought target and/or whether the object in view is a target match might arrive within V4 and IT during object search. In the first proposal (Fig. 1A), V4 serves as the sole locus of the combination of visual and top-down information and IT receives this information via feed-forward propagation from V4. In the second (Fig. 1B), top-down information is integrated directly in IT, either exclusively or in addition to its integration in V4. Our use of the term “integrated” is functional (as opposed to anatomical) and refers to the means by which top-down signals are reflected in the hierarchically arranged ventral visual pathway. We note that the anatomical locus of integration could be these brain areas or within other structures that they connect to, such as the pulvinar.

A number of studies report that task-relevant signals increase in a gradient-like fashion across the early visual hierarchy (i.e., V1, V2, and V4) during covert spatial attention and feature-based attention tasks (reviewed by Noudoost et al. 2010), consistent with the integration of top-down signals at multiple, early stages. However, the few studies that have compared top-down modulation at higher stages of the pathway, including V4 and IT, report it to be matched, both during visual target search (Chelazzi et al. 1998, 2001) as well as one covert spatial attention task (Moran and Desimone 1985). Additionally, while one might postulate that because receptive fields are smaller in V4, the brain would be more likely to integrate top-down signals in IT when a task requires spatial invariance (such as detecting the change in a feature despite its position), V4 feature-based attention effects have in fact been demonstrated to extend globally across the visual field (reviewed by Cohen and Maunsell 2011 and Maunsell and Treue 2006). Within V4, the responses of neurons have been demonstrated to be modulated by top-down attentional effects that extend well beyond the sizes of individual V4 receptive fields. These global effects are consistent with the idea that top-down

Address for reprint requests and other correspondence: N. Rust, Dept. of Psychology, Univ. of Pennsylvania, Goddard Laboratories, 428, Philadelphia, PA 19104 (e-mail: nrust@psych.upenn.edu).

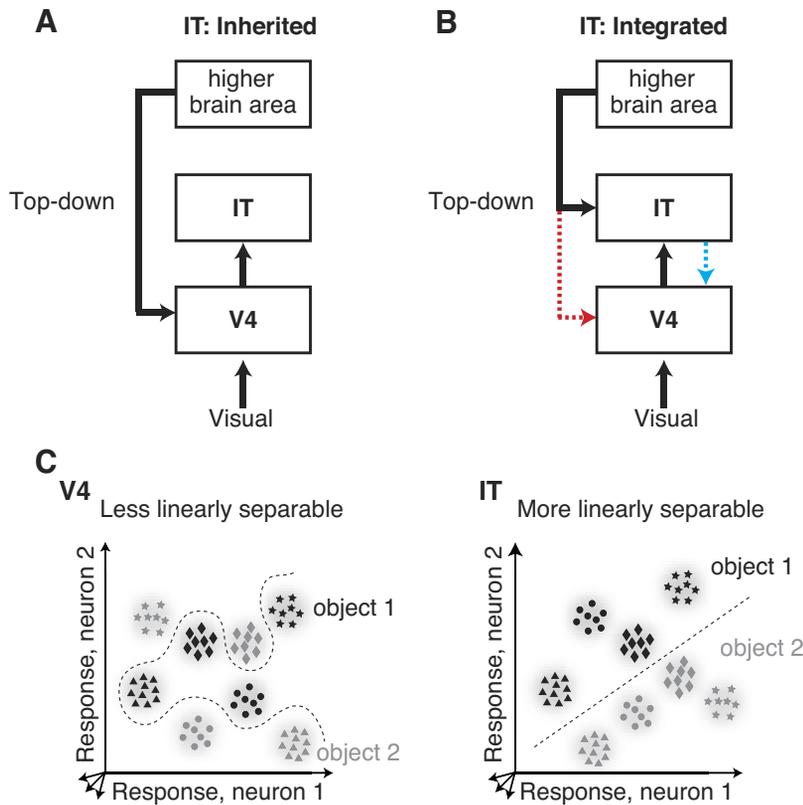


Fig. 1. Proposals for how top-down information might arrive within V4 and inferotemporal cortex (IT) during object search. *A*: the class of IT:inherited proposals predict that top-down information is integrated only in V4, and this information is then inherited by IT via feed-forward propagation. *B*: the class of IT:integrated proposals predict that top-down information is integrated directly in IT. This class includes proposals in which top-down information is integrated in both IT and V4 (red) as well as proposals in which top-down information is integrated exclusively in IT but is then fed back to V4 (cyan). *C*: cartoon depiction of the population representation of object identity in V4 and IT plotted as the response of 2 hypothetical neurons. In these plots, each point indicates the population response on a single trial, each cloud of points depicts the dispersion of responses to the same image across repeated trials, the shade of the points indicate the identity of the object contained in each image, and the shape of the points indicate the identity-preserving transformation at which each object appears. In V4, representations of object identity are not linearly separable, as indicated by the nonlinear boundary required to separate the 2 objects. In IT, representations of object identity are more linearly separable, as indicated by the linear boundary.

feature-based attention effects are integrated into V4 in a global manner, despite small sizes of V4 receptive fields. Together, the combined evidence that modulation magnitudes are matched in V4 and IT, coupled with the fact that top-down modulation can happen globally across the visual field in V4, suggests that global, top-down modulation may happen exclusively in V4 (Fig. 1A).

The means by which the brain integrates top-down information in the ventral visual pathway may very well depend on the specific task and its computational demands. The experiments described here were targeted at challenging the hypothesis that top-down integration happens exclusively in the ventral visual pathway within or before V4 (Fig. 1A) with a task that is seemingly better optimized for top-down integration in IT. Specifically, our experiments are designed to exploit differences in the visual representation of object identity between V4 and IT, where both brain areas have been demonstrated to reflect similar total amounts of information for object identification tasks, but in V4 this information is more implicitly formatted (i.e., requiring a nonlinear population readout; Fig. 1C, left) whereas in IT this information is more explicitly formatted (i.e., more accessible to a linear population readout; Fig. 1C, right; reviewed by DiCarlo et al. 2012). In our experiments, monkeys performed a delayed match to sample task in which they were cued to search for specific objects but those objects could appear at different positions, sizes and background contexts. In question is whether the more linearly formatted object representations in IT coincided with the preferential integration of top-down information in IT (Fig. 1B) or whether top-down information could be described as integrated in V4 (Fig. 1A), where top-down modulation has been demonstrated to occur globally across the visual field despite the small sizes of V4 receptive fields.

## MATERIALS AND METHODS

### Experimental Design

Experiments were performed on two adult male rhesus macaque monkeys (*Macaca mulatta*) with implanted head posts and recording chambers. All procedures were reviewed and approved by the University of Pennsylvania Institutional Animal Care and Use Committee.

### Visual Stimuli

Images were presented in a circular aperture with a radius of  $5^\circ$ , centered at fixation (Fig. 2A). Images contained objects presented at different positions, sizes, and background contexts. Sizes included the following: size-1 $\times$  ( $\sim 1.2^\circ$ ), size-1.5 $\times$  ( $\sim 1.8^\circ$ ), and size-2.25 $\times$  ( $\sim 2.7^\circ$ ). Positions included position-fixation ( $0^\circ$  horizontal and vertical), position-right ( $1.55^\circ$  to the right,  $0.67^\circ$  below fixation), position-left ( $1.55^\circ$  to the left,  $0.67^\circ$  below fixation), and position-up ( $0^\circ$  horizontal;  $1^\circ$  above fixation). Changes in size and position were combined to produce the following five transformations (Fig. 2C): “Up” (position-up, size-1.5 $\times$ ); “Left” (position-left; size-1.5 $\times$ ); “Right” (position-right; size-1.5 $\times$ ); “Big” (position-fixation; size-2.25 $\times$ ); and “Small” (position-fixation; size-1 $\times$ ). In addition, a different natural image background was chosen for each of the transformations Up, Big, and Small, whereas Left and Right were presented on a gray background (Fig. 2C). The complete image set included 4 objects, each presented at the 5 transformations described above, for a total of 20 images. The rationale behind selecting these particular transformations was to make the task of object identification (invariant to identity-preserving transformation) challenging for both V4 and IT, based on results from our previous work (Rust and DiCarlo 2010).

As described below, our experiment included two types of trials, cue trials and test trials. The set of 20 images described above were the only images presented on test trials. Cue trials also included a

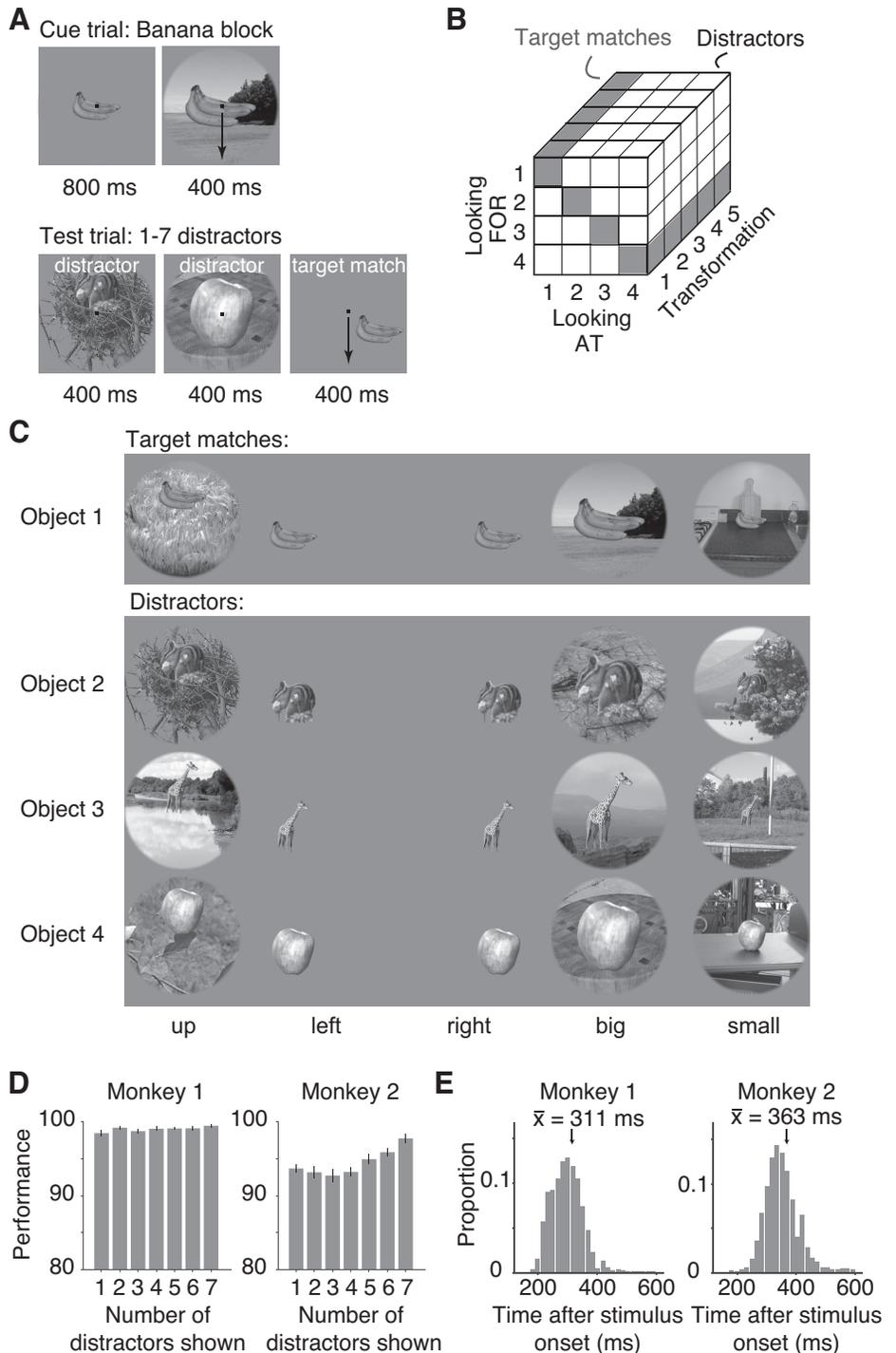


Fig. 2. The invariant delayed-match-to-sample task. *A*: monkeys initiated trials by fixating on a small dot. Each block (~3 min in duration) began with a cue trial that indicated the target object. On subsequent trials, a random number (1–7) of distractors were presented, and on most trials, this was followed by the target match. Monkeys were required to maintain fixation throughout the distractors and make a saccade to a response dot within a window 75–600 ms following the onset of the target match to receive a reward. In cases where the target match was presented for 400 ms and the monkey had still not broken fixation, a distractor stimulus was immediately presented. *B*: a schematic of the full experimental design, which included 80 conditions: “Looking AT” each of 4 objects, each presented at 5 identity-preserving transformations (for 20 images in total), viewed in the context of “Looking FOR” each object as a target. In this design, target matches (gray) fall along the diagonal of each Looking AT/Looking FOR transformation slice whereas distractors (white) fall off the diagonal. *C*: images used in the task: 4 objects were presented at each of 5 identity-preserving transformations (“up,” “left,” “right,” “big,” and “small”), for 20 images in total. In any given block, 5 of the images were presented as target matches and 15 were distractors. *D*: percent correct for each monkey, calculated based on both misses and false alarms (but disregarding fixation breaks), shown as a function of the number of distractors preceding the target match. Error bars indicate SE across experimental sessions. *E*: histograms of reaction times during correct trials (ms after stimulus onset), with means labeled.

version of each object (size-1.5 $\times$ ) presented at position-fixation on a gray background.

#### The Invariant Delayed-Match-to-Sample Task

Monkey behavioral training and testing utilized standard operant conditioning, head stabilization, and infrared video eye tracking. Custom software (<https://mworks.github.io/>) was used to present stimuli on an LCD monitor with an 85-Hz refresh rate.

The monkeys performed an invariant delayed-match-to-sample (IDMS) task (Fig. 2). As an overview, the task required the monkeys to make a saccade when a target object appeared within a sequence of

distractor images (Fig. 2A). Objects were presented at differing positions, sizes, and background contexts as described above and shown in Fig. 2C. Stimuli consisted of a fixed set of 20 images that included 4 objects, each presented at 5 different identity-preserving transformations (Fig. 2B). Each short block (~3 min) was run with a fixed target object before another target was pseudorandomly selected. Our design included two types of trials: cue trials and test trials (Fig. 2A). Only test trials were analyzed for this report.

A trial began when the monkey fixated on a red dot (0.15°) in the center of a gray screen, within a square window of  $\pm 1.5^\circ$ . Fixation was followed by a 250-ms delay before a stimulus appeared. Cue trials, which indicated the current target object, were presented at the

beginning of each short block or after three subsequent error trials. To minimize confusion, cue trials were designed to be distinct from test trials and began with the presentation of an image of each object that was distinct from the images used on test trials (a large version of the object presented at the center of gaze on a gray background; Fig. 2A). Test trials began with a distractor image, and neural responses to the first distractor were discarded to minimize nonstationarities such as stimulus onset effects. During the IDMS task, all images were presented at the center of gaze, in a circular aperture that blended into a gray background (Fig. 2C).

In each block, 5 images were presented as target matches and the other 15 as distractors. Distractor images were drawn randomly without replacement until each distractor was presented once on a correct trial, and the images were then re-randomized. Within each block, four repeated presentations of each of the 20 images were collected, and a new target object was then pseudorandomly selected. Each block lasted ~3 min. Following the presentation of all 4 objects as targets, the targets were rerandomized. At least 10 repeats of each condition were collected on correct trials. When more than 10 repeats were collected, the first 10 were used for analysis.

On most test trials, a target match followed the presentation of a random number of 1–6 distractors (probability of a target match at each position: 0.278, 0.151, 0.142, 0.128, 0.115, and 0.099; Fig. 2A). On a small fraction of trials (0.086), seven distractors were shown, and the monkeys were rewarded for fixating through all distractors. This translated into a function whereby if the monkey had not observed a target match by *position N–1* in the trial, the probability that the target match would appear at *position N* was for *positions 1–6*: 0.278, 0.209, 0.249, 0.299, 0.383, and 0.536. Each image was presented for 400 ms (or until the monkeys' eyes left the fixation window) and was immediately followed by the presentation of the next stimulus. Monkeys were rewarded for making a saccade to a response target within a window of 75–600 ms after the target match onset. In *monkey 1*, the response target was positioned 10° above fixation; in *monkey 2*, it was 10° below fixation. If the monkeys had not yet moved their eyes after 400 ms following target onset, a distractor stimulus was immediately presented. A trial was classified as a “false alarm” if the eyes left the fixation window via the top (*monkey 1*) or bottom (*monkey 2*) outside the allowable correct response period and traveled >0.5°. In contrast, all other instances in which the eyes left the fixation window during the presentation of distractors were characterized as fixation breaks. A trial was classified as a “miss” when the monkey continued fixating beyond 600 ms following the onset of the target match. Overall, monkeys performed this task with high accuracy. Disregarding fixation breaks (*monkey 1*: 11% of trials; *monkey 2*: 8% of trials), percent correct on the remaining trials was as follows: *monkey 1*: 98% correct, ~1% false alarms, and ~1% misses; *monkey 2*: 94% correct, 2% false alarms, and 4% misses. Behavioral performance was comparable for the sessions corresponding to recordings from the two areas (V4 percent correct overall = 96.5%; IT percent correct overall = 91.4%).

V4 receptive fields at and near the center of gaze are small: on average they have radii of 0.56° at the fovea, extending to radii of 1.4 at an eccentricity of 2.5° (Desimone and Schein 1987; Gattass et al. 1988). We thus took considerable care to ensure that the images were approximately placed in the same region of these receptive fields across repeated trials. In the second monkey, adequate fixational control could not be achieved through training. We thus applied a procedure in which we shifted each image at stimulus onset 25% toward the center of gaze (e.g., if the eyes were displaced 0.5° to the left, the image was repositioned 0.125° to the left and thus 0.375° from fixation). Image position then remained fixed until the onset of the next stimulus. The value of 25% was determined during training as an amount that enabled us to achieve better consistency of placement of the images on the receptive fields, while maintaining high behavioral performance.

### Neural Recording

The activity of neurons in V4 and IT was recorded via a single recording chamber for each brain area in each monkey. In both monkeys, chamber implantation and recording in IT preceded V4, and the IT recording chamber was implanted on the left hemisphere whereas the V4 recording chamber was implanted on the right hemisphere. While IT receptive fields span the vertical meridian, thus allowing us to access the visual representation of both sides with a single chamber, V4 receptive fields are confined to the contralateral hemifield. To simulate V4 coverage of the ipsilateral visual field, on roughly half of the V4 recording sessions, ( $n = 11/21$  sessions in *monkey 1*;  $n = 7/15$  sessions in *monkey 2*), we presented the images reflected across the vertical axis. We reflected the images, rather than, e.g., shifting them, to preserve the foveal-to-peripheral organization of the images. As an example, consider two units recorded in the right hemisphere with receptive fields centered at 1° to the left of fixation, one recorded on a standard session and the other recorded on a reflected session. The unit recorded on a reflected session would be presented with an object feature that typically would be shown at 1° to the right of fixation and was thus simulated as a unit in the left hemisphere. We then treated all V4 units recorded during these reflected sessions as if they were in the left hemisphere (and thus as receptive fields that were located in the right visual field).

Chamber placement for the IT chambers was guided by anatomical magnetic resonance images in both monkeys, and in one monkey, Brainsight neuronavigation (<https://www.rogue-research.com/>). Both V4 chambers were guided by Brainsight neuronavigation. The region of IT recorded was located on the ventral surface of the brain, over an area that spanned 4 mm lateral to the anterior middle temporal sulcus and 15–19 mm anterior to the ear canals, consistent with CIT/AIT. Both V4 chambers were centered 1 mm posterior to the ear canals and 29 mm lateral to the midline, positioned at a 30° angle. V4 recording sites were confirmed by a combination of receptive field location and position in the chamber, corresponding to results reported previously (Gattass et al. 1988). Specifically, we recorded from units within and around the inferior occipital sulcus, between the lunate sulcus and superior temporal sulcus. V4 units in lower visual field were confirmed as having receptive field centers that traversed from the vertical to horizontal meridian as recordings shifted from posterior to anterior. As expected, V4 units in the fovea and near the upper visual field were found lateral to those in the lower visual field and had receptive field centers that traversed from the horizontal meridian to the vertical meridian as recordings traversed medial to lateral and increased in depth.

Neural activity was recorded with 24-channel U-probes and V-probes (Plexon) with linearly arranged recording sites spaced with 100- $\mu$ m intervals. Continuous, wideband neural signals were amplified, digitized at 40 kHz, and stored using the OmniPlex Data Acquisition System (Plexon). Spike sorting was done manually offline (Plexon Offline Sorter). At least one candidate unit was identified on each recording channel, and two to three units were occasionally identified on the same channel. Spike sorting was performed blind to any experimental conditions to avoid bias. A multichannel recording session was included in the analysis if the animal performed the task until the completion of at least 10 correct trials per stimulus condition, there was no external noise source confounding the detection of spike waveforms, and the session included a threshold number of task-modulated units (>4 on 24 channels). The sample size for IT (number of units recorded) was chosen to approximately match our previous work (Pagan et al. 2013; Pagan and Rust 2014a). The sample size for V4 was selected to be threefold that number, to match the ratio between numbers of units estimated in V4 as compared with IT (DiCarlo et al. 2012).

For many of the analyses presented in this paper, we measured neural responses by counting spikes in a window that began 40 ms after stimulus onset in V4 and 80 ms after stimulus onset in IT. We

counted spikes in a 170-ms window in both areas, such that the spike counting windows were of equal length. Counting windows always preceded the monkeys' reaction times. On 7.7% of all correct target match presentations, the monkeys had reaction times faster than 250 ms, and those instances were excluded from analysis to ensure that spikes in both V4 and IT were only counted during periods of fixation.

In IT, we recorded neural responses across 20 experimental sessions (*monkey 1*: 10 sessions; *monkey 2*: 10 sessions). In V4, we recorded neural responses across 36 experimental sessions (*monkey 1*: 21 sessions; *monkey 2*: 15 sessions). When combining the units recorded across sessions into a larger "pseudopopulation" (a population of units combined across sessions and, when relevant, combined across monkeys), we began by screening for units that met three criteria. First, units needed to be modulated by our task, as quantified by a one-way ANOVA applied to our neural responses (80 conditions  $\times$  10 repeats,  $P < 0.01$ ). Second, units needed to pass a loose criterion on recording stability, as quantified by calculating the variance-to-mean ratio (Fano factor) for each unit, computed by fitting the relationship between the mean and variance of spike count across the 80 conditions (Fano factor  $< 5$ ). Finally, units needed to pass a loose criterion on unit recording isolation, quantified by calculating the signal-to-noise ratio (SNR) of the waveform as the difference between the maximum and minimum points of the average waveform, divided by twice the standard deviation across the differences between each waveform and the mean waveform (SNR  $> 2$ ). In IT, this yielded a pseudopopulation of 204 units (of 563 possible units), including 108 units from *monkey 1* and 96 units from *monkey 2*. In V4, this yielded a pseudopopulation of 650 units (of 970 possible units), including 382 units from *monkey 1* and 268 units from *monkey 2*.

#### V4 Receptive Field Mapping

To measure the location and extent of V4 receptive fields, bars were presented for 500 ms, 1 per trial, centered on a  $5 \times 5$  invisible grid. Bar orientation, length, and width as well as the grid center and extent were adjusted for each recording session based on preliminary hand mapping. On each trial, the monkey was required to maintain fixation on a small response dot ( $0.125^\circ$ ) to receive a reward. The responses to at least five repeats were collected at each position for each recording session. Only those units that produced clear visually evoked responses at a minimum of one position were considered for receptive field position analysis. The center of the receptive field was estimated by the maximum of the response across the  $5 \times 5$  grid of oriented bar stimuli and confirmed by visual inspection.

#### Quantifying Single-Unit Modulations

To quantify the degree to which individual V4 and IT units were modulated by task-relevant variables (Figs. 5, 6, and 7), such as changes in visual and target identity, we applied a bias-corrected, ANOVA-like procedure described in detail by Pagan and Rust (2014b) and summarized here. As an overview, this procedure is designed to parse each unit's total response variance into variance that can be attributed to each type of experimental parameter as well as variance that can be attributed to trial variability. Total variance is computed across the spike count responses for each unit across 16 conditions (4 images  $\times$  4 targets separately for each transformation) and 10 trials. Variances are then transformed into measures of spike count modulation (in the units of standard deviation around each unit's grand mean spike count) via a procedure that includes bias correction for overestimates in modulation due to noise.

To capture all types of modulation (such as modulation by changes in visual or target identity), the procedure begins by developing an orthonormal basis of 16 vectors. The number of basis vectors for each type of modulation is imposed by the experimental design. In particular, this basis  $\mathbf{b}$  includes vectors  $\mathbf{b}_i$  that reflect 1) the grand mean spike count across all conditions, 2) whether the object in view is a

target or a distractor ("target match"), 3) visual image identity ("visual"), 4) target object identity ("target identity"), and 5) nonlinear interactions between target and object identity not captured by target match modulation ("residual"). The initially designed set of vectors is then converted into an orthonormal basis via a Gram-Schmidt orthogonalization process.

The resulting basis spans the space of all possible responses for our task. Consequently, we can reexpress each trial-averaged vector of spike count responses to the 16 experimental conditions for each transformation,  $\mathbf{R}$ , as a weighted sum of these basis vectors. The weight corresponding to a basis vector for each unit reflects modulation of that unit's responses by that experimental parameter. To quantify the amounts of each type of modulation reflected by each unit, we began by computing the squared projection of each basis vector  $\mathbf{b}_i$  and  $\mathbf{R}$ . To correct for bias caused by overestimates in modulation due to noise, an analytical bias correction, described and verified in Pagan and Rust (2014b), was then subtracted from this value. The squared weight for each basis vector  $\mathbf{b}_i$  is thus calculated as:

$$w_i^2 = (\mathbf{R} \cdot \mathbf{b}_i^T)^2 - \frac{\sigma_i^2 \cdot (\mathbf{b}_i^T)^2}{m} \quad (1)$$

where  $\sigma_i^2$  indicates the trial variance, averaged across conditions ( $n = 16$ ), and  $m$  indicates the number of trials ( $m = 10$ ). If more than one dimension existed for a type of modulation, we summed values of the same type (Eq. 2). Next, we applied a normalization factor  $\{[1/(n - 1)]$  where  $n = 16\}$  to convert these summed values into variances. As a final step, we computed the square root of these quantities to convert them into modulation measures that reflected the number of spike count standard deviations around each unit's grand mean spike count. Modulation for each parameter type  $X$  was thus computed as:

$$\sigma_X = \sqrt{\frac{1}{n-1} \cdot \sum_{i=j}^k w_i^2} \quad (2)$$

for the weights  $w_j$  through  $w_k$  corresponding to basis vectors  $\mathbf{b}_j$  through  $\mathbf{b}_k$  for that parameter type, where the number of basis vectors corresponding to each parameter type were as follows: target match = 1; visual = 3; target identity = 3; residual = 8.

When estimating modulation for individual units (Fig. 5), the bias-corrected squared values were rectified for each unit before taking the square root. When estimating modulation population means (Fig. 6, B–E, and Fig. 7), the bias-corrected squared values were averaged across units before taking the square root. When estimating modulation population means within the broader 170-ms bins for V4 and IT, respectively (bar graphs shown in Fig. 6, B–E, and Fig. 7, A and B, right), the modulations shown are roughly equal to the sum (or integral) of the modulations in the 50-ms sliding bins (Fig. 6, B–E, and Fig. 7, A and B, left) summed across the 40- to 210-ms or 80- to 250-ms bins in V4 and IT, respectively. However, because we computed the modulations separately in the different bin sizes before bias correcting, averaging across units, and taking the square root, this relationship is not exactly equivalent to an integral. Because these measures were not normally distributed, standard error about the mean was computed via a bootstrap procedure. On each iteration of the bootstrap (across 1,000 iterations), we randomly sampled values from the modulation values for each unit in the population, with replacement. Standard error was computed as the standard deviation across the means of these resampled populations.

#### Population Performance: Target Match Information

To determine the ability of the V4 and IT populations to classify target matches versus distractors (Fig. 8), we applied two types of decoders: a Fisher linear discriminant (FLD; a linear classifier) and a maximum likelihood decoder (a decoder that can classify based on

linear as well as nonlinearly formatted target match information). Both decoders were cross validated with the same resampling procedure. On each iteration of the resampling, we randomly shuffled the trials for each condition and for each unit and (for numbers of units less than the full population size) randomly selected units (with the exception of Fig. 8D, cyan, where we selected the “best” units, as described below). On each iteration, eight (of the 10 total) trials from each condition were used for training the decoder, one trial from each condition was used to determine a value for regularization of the FLD linear classifier, and one trial from each condition was used for a cross-validated measurement of performance.

To circumvent issues related to the format of visual information, classifier analyses were performed per transformation (for the 4 of the 5 transformations used, see Fig. 5; Big, Up, Left, and Small). The data for each transformation consisted of 16 conditions (4 visual objects viewed under 4 different target contexts). To ensure that decoder performance relied only on target match information and was not biased due to other factors, such as differences in the numbers of each class, each classification was computed for four target matches versus four (of 12 possible) distractors. Each set of four distractors was selected to span all possible combinations of mismatched object and target identities (e.g., *objects 1, 2, 3, and 4* paired with *targets 4, 3, 2, and 1*), of which there are nine possible sets. Performance was computed on each resampling iteration by averaging the binary performance outcomes across the nine possible sets of target matches and distractors, each of which contained eight cross-validated test trials, and across the four transformations used. For both types of classifiers, means  $\pm$  SE of performance was computed as the means  $\pm$  SD of performance across 1,000 resampling iterations. Standard error thus reflected the variability due to the specific trials assigned to training and testing and, for populations smaller than the full size, the specific units chosen.

To compute linear classifier performance (Fig. 8D), we used a two-way FLD as described below.

The general form of a linear decoding axis is:

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b, \quad (3)$$

where  $\mathbf{w}$  is an  $N$ -dimensional vector containing the linear weights applied to each of  $N$  units and  $b$  is a scalar value. We fit these parameters using an FLD, where the vector of linear weights was calculated as:

$$\mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2) \quad (4)$$

and  $b$  was calculated as:

$$b = \mathbf{w} \cdot \frac{1}{2}(\mu_1 + \mu_2) = \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 - \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 \quad (5)$$

Here  $\mu_1$  and  $\mu_2$  are the means of two classes (target matches and distractors) and the mean covariance matrix is calculated as:

$$\Sigma = \frac{\Sigma_1 + \Sigma_2}{2} \quad (6)$$

where  $\Sigma_1$  and  $\Sigma_2$  are the regularized covariance matrices of the two classes. These covariance matrices were computed using a regularized estimate equal to a linear combination of the sample covariance and the identity matrix  $I$  (Pagan and Rust 2014a):

$$\Sigma_i = \gamma \Sigma_i + (1 - \gamma) \cdot I \quad (7)$$

We determined  $\gamma$  by exploring a range of values from 0.01 to 0.99, and we selected the value that maximized average performance across all iterations, measured with the cross-validation regularization trials set aside for this purpose (see below). We then computed performance for that value of  $\gamma$  with separately measured test trials, to ensure a fully cross-validated measure. Because this calculation of the FLD parameters incorporates the off-diagonal terms of the covariance

matrix, FLD weights are optimized for both the information conveyed by individual units as well as their pairwise interactions.

To compute neural population performance, we began by computing the dot product of the test data  $\mathbf{x}$  and the linear weights  $\mathbf{w}$ , adjusted by  $b$  (Eq. 3). Each test trial was then assigned to one class, and proportion correct was then computed as the fraction of test trials that were correctly assigned, according to their true labels. To compute linear classifier performance for the best V4 units (Fig. 8D, cyan), we ranked units by their  $d'$  based on the training data and subselected top-ranked units to measure cross-validated performance. Unit  $d'$  was computed as:

$$d' = \frac{|\mu_{\text{Match}} - \mu_{\text{Distractor}}|}{\sigma_{\text{pooled}}}, \quad (8)$$

where  $\mu_{\text{Match}}$  and  $\mu_{\text{Distractor}}$  correspond to the mean response across the set of target match and distractors,  $\sigma_{\text{pooled}} = \sqrt{\frac{\sigma_{\text{Match}}^2 + \sigma_{\text{Distractor}}^2}{2}}$ , and  $\sigma_{\text{Match}}$  and  $\sigma_{\text{Distractor}}$  correspond to the standard deviation of responses across the set of target matches and distractors, respectively.

As a measure of total target match information (Fig. 8F; combined linear and nonlinear), we implemented a maximum likelihood decoder (Pagan et al. 2013, 2016). We began by using the set of training trials to compute the average response  $r_{uc}$  of each unit  $u$  to each of the two conditions  $c$  (target matches vs. distractors). We then computed the likelihood that a test response  $k$  was generated from a particular condition as a Poisson-distributed variable:

$$lik_{u,c}(k) = \frac{(r_{uc})^k \cdot e^{-r_{uc}}}{k!} \quad (9)$$

The likelihood that a population response vector was generated in response to each condition was then computed as the product of the likelihoods of the individual units. We assigned the population response to the category with the maximum likelihood, and we computed performance as the fraction of trials in which the classification was correct based on the true labels of the test data.

### Statistical Analysis

Because our measures were not normally distributed, we computed  $P$  values via resampling procedures. When comparing the magnitudes of single unit modulation values between V4 and IT (Fig. 5, Fig. 6, B–E, and Fig. 7), a bootstrap procedure was applied in which values were randomly sampled from the values for each unit, with replacement, across many iterations. We calculated  $P$  values as the fraction of resampling iterations on which the difference between the means of the resampled data was flipped in sign relative to the actual difference between the means of the full data set (for example, if the mean of visual modulation in V4 was larger than the mean of visual modulation in IT, the fraction of iterations in which the mean of visual modulation in IT was larger than the mean of visual modulation in V4).

When comparing population decoding measures (Fig. 8, D and F), 1,000 iterations of cross-validated population performance were computed, and  $P$  values were calculated as the fraction of classifier iterations on which the difference was flipped in sign relative to the actual difference between the means across classifier iterations (for example, if the mean of decoding *measure 1* was larger than the mean of decoding *measure 2*, the fraction of iterations in which the mean of *measure 2* was larger than the mean of *measure 1*). When evaluating whether a population decoding measure was different from chance (Fig. 8, D and F),  $P$  values were calculated as the fraction of classifier iterations on which performance was greater than chance performance (50%).

## RESULTS

Understanding the path by which top-down information arrives in the ventral visual pathway during visual object search is challenging, due to several factors. First, receptive fields at different stages of the pathway have markedly different sizes and it is unclear how to best compare them. Second, during tasks that approximate real-world object search, the responses of individual V4 and IT neurons operate in a low spike count regime in which the amounts and types of signals within individual units are difficult to measure (Roth and Rust 2018a, 2018b). This noisy, low-spike-count regime is a combined consequence of spike count windows that are short, as implied by reaction times that are fast (~250 ms) and individual stimuli that drive only a subset of neurons robustly (Roth and Rust 2018b). These challenges have traditionally been addressed using approaches that seek to reduce the noise in a manner that implicitly makes unrealistic assumptions about the brain. For example, tailoring stimuli to fit within the small sizes of V4 receptive fields and/or aligning stimuli with the peak of each neuron's preferences disregards the contributions of the receptive field surround and/or neurons that are activated along their curve flank. Similarly, counting spikes in long windows that exceed natural reaction times assumes that neural responses are stationary, whereas their responses are not (e.g., Pagan and Rust 2014b). Here we apply a complementary approach that involves studying V4 and IT in a manner analogous to the way that the brain addresses the challenge of noisy individual neuron responses during real-world object search: by combining the noisy responses of individual neurons across a neural population (i.e., via weighted population decoding schemes).

As an overview of our experiments, we trained two monkeys to perform an invariant delayed-match-to-sample (IDMS) task that required them to report when target objects appeared. On any given trial of the IDMS task, monkeys were instructed about the object that they were searching for, but they did not know the specific position, size, and background context with which it would appear. The task was designed around previously established differences in the nature of visual object representations between V4 and IT, where visual information in V4 is more implicitly formatted (i.e., less linearly separable; Fig. 1C, *left*) whereas in IT this information is more explicitly formatted (i.e., more accessible to a linear population readout; Fig. 1C, *right*; reviewed by DiCarlo et al. 2012). At the level of individual neurons, the response property supporting IT population linear separability is thought to be object "tolerance" or the preservation of rank-order tuning for objects across identity-preserving transformations (Ito et al. 1995; Li et al. 2009). Object tolerance is most intuitively described for changes in object position and spatially large IT receptive fields that have the same tuning for object identity across all regions of their receptive field. In contrast, V4 receptive fields are smaller and are stimulated by different regions of the image when an object moves to a new position, and this can lead to a change in the rank-order tuning for objects at different positions. Consequently, individual V4 neurons are less object tolerant, thereby leading to a more nonlinear population representation of object identity.

In the context of the IDMS task where monkeys were searching for an object that could appear at different transfor-

mations, the visual information required to solve the task (i.e., identity of the object in view) is expected to exist in both V4 and IT, but this information is expected to be formatted in a more linear manner in IT (Fig. 1C). In question is whether this difference in information format between V4 and IT coincides with the preferential integration of top-down information in IT (Fig. 1B). Equivalently, in question is whether during the IDMS task the brain somehow manages to integrate invariant information about object identity in V4 where receptive fields are small, which it may in fact be able to do, given the spatially global top-down modulations documented in V4 feature-based attention experiments (Cohen and Maunsell 2011; Maunsell and Treue 2006;), or whether it preferentially targets IT, where receptive fields are larger.

The IT data presented here were also included in two earlier publications (Roth and Rust 2018a, 2018b). There we reported that during IDMS neural signals in IT reflected behavioral confusions on the trials in which the monkeys made errors and IT target match signals were configured in a manner that minimized their interference with IT visual representations. The focus of the current report is a determination of how these signals arrive in IT via a systematic comparison between IT and its input brain area, V4; the V4 data have not been published previously.

### *The Invariant Delayed-Match-to-Sample Task*

Monkeys performed an invariant delayed-match-to-sample (IDMS) task in short blocks of trials (~3 min on average) with a fixed target object. Each block began with a cue trial that indicated the target for that block (Fig. 2A, "Cue trial"). The remainder of the block was comprised primarily of test trials (Fig. 2A, "Test trial"). Test trials began with the presentation of a distractor, and on most trials, this was followed by zero to five additional distractors (for a total of 1–6 distractor images) and then an image containing the target match. The monkeys' task required them to maintain fixation during the presentation of distractors and make a saccade in response to the appearance of a target match to receive a juice reward. To minimize the possibility that monkeys would predict the target match, on a small fraction of the trials the target match did not appear and the monkeys were rewarded for maintaining fixation through seven distractors. Unlike other classic DMS tasks (Chelazzi et al. 1993; Eskandar et al. 1992; Lueschow et al. 1994; Miller and Desimone 1994; Pagan et al. 2013), our experimental design does not incorporate a sample at the beginning of each test trial, to better mimic real-world object search, where target matches are not repeats of the same image presented shortly before. Rather, our cue trials contain a sample stimulus that is not an exact match to the targets, and cue trials are only presented at the beginning of each block. One benefit of this task design is that it better isolates target match modulation from other types of stimulus repetition effects, including repetition suppression (Miller and Desimone 1994). This distinction is important for this study, as our goal was to systematically compare the magnitudes of top-down modulation in V4 and IT and repetition suppression is likely to be a largely feed-forward process.

Our experimental stimuli consisted of a fixed set of 20 images: 4 objects presented at each of 5 transformations (Fig. 2C). These specific images were selected to make the task of

classifying object identity challenging for the IT population, and these specific transformations were selected based on findings from our previous work (Rust and DiCarlo 2010). In a given target block (e.g., a “banana block”), a subset of 5 of the images were target matches and the remaining 15 were distractors (Fig. 2C). The full experimental design amounted to 20 images (4 objects presented at 5 identity-preserving transformations), all viewed in the context of each of the 4 objects as a target, resulting in 80 experimental conditions (Fig. 2B). In this design, target matches fall along the diagonal of each “looking at”/“looking for” matrix slice (where a matrix slice corresponds to the conditions at one fixed transformation; Fig. 2B, gray). For each of the 80 conditions, we collected at least 10 repeats of correct trials. Behavioral performance was high overall (Fig. 2D). The monkeys’ mean reaction times (computed as the time their eyes left the fixation window relative to the target match stimulus onset) were 311 and 363 ms for *monkeys 1* and 2, respectively (Fig. 2E).

To systematically compare the responses of V4 and IT during this task, we applied a population-based approach in which we fixed the images and their placement in the visual field across all the units that we studied, and we sampled from units whose receptive fields overlapped the stimuli. Specifically, we presented images at the center of gaze, with a diameter of  $5^\circ$ . Neurons in IT typically have receptive fields that extend beyond  $5^\circ$  and extend into all four quadrants (Op De Beeck and Vogels 2000). In contrast, V4 receptive fields are smaller, retinotopically organized, and confined to the contralateral hemifield (Desimone and Schein 1987; Gattass et al. 1988). To compare these two brain areas, we applied extensions of approaches developed in our earlier work in which we compared the responses of a set of randomly sampled IT units with a population of V4 units whose receptive fields tiled the image (Rust and DiCarlo 2010). This required sampling V4 units with receptive fields in both upper and lower visual fields (Fig. 3), which we achieved through recording at different positions within and around the inferior occipital sulcus. This also required measuring units with receptive fields on both sides of the vertical meridian, which we approximated by isolating our recordings to one hemisphere but reflecting the images along the vertical axis in approximately half the sessions.

Because V4 receptive fields in the region of the field that we recorded are small, one issue of concern is the replicability of retinal image placement across trials. We quantified the stability of monkeys’ eye positions across repeated trials as the spatial deviation in retinal image placement across trials, measured relative to the mean position across trials, and we determined the proportion of eye positions that were within windows corresponding to V4 receptive field sizes at the range of eccentricities we recorded (Gattass et al. 1988). In *monkey 1*, 85% of eye positions fell within windows corresponding to the average RF sizes at the fovea (average foveal receptive field size =  $0.56^\circ$ ), and 97% of eye positions were within windows corresponding to RF sizes at an eccentricity of  $2.5^\circ$  (average receptive field size at  $2.5^\circ = 1.4^\circ$ ). To achieve similar precision in *monkey 2*, fixational control was improved by employing a procedure in which eye position was determined just before stimulus onset and the image was shifted at stimulus onset such that it was positioned closer to the center of gaze (see MATERIALS AND METHODS). The image then remained in the

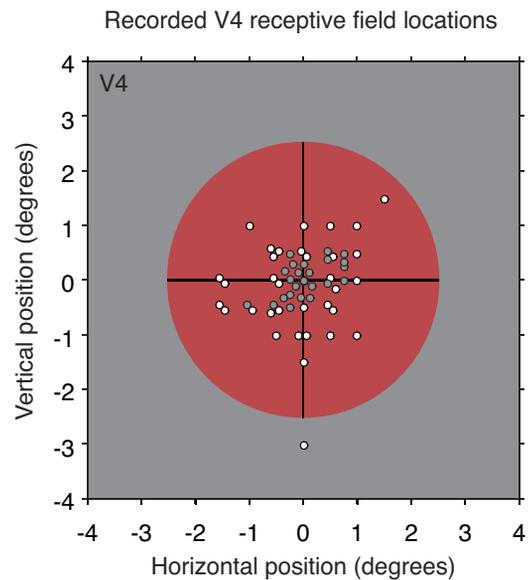


Fig. 3. V4 receptive field locations. Images were displayed at the center of gaze and were  $5^\circ$  in diameter. Red circle indicates the location and size of the images. We targeted V4 units with receptive fields that tiled the images. After approximate receptive field localization with hand mapping, receptive field locations were determined with oriented bar stimuli presented in a  $5 \times 5$  grid of different positions (see MATERIALS AND METHODS). Shown are the receptive field centers of a subset of recorded V4 units; 1 dot is shown for each unique receptive field location recorded. On approximately half of the sessions, images were reflected across the vertical axis, and for these sessions, the receptive field centers are plotted in the ipsilateral visual field. *Monkey 1*: white; *monkey 2*: gray.

same, fixed position for the duration of the image viewing period. The resulting retinal stability was comparable to *monkey 1*: on average, 95 and 99% of presentations occurred within windows with a radius of  $0.56$  and  $1.4^\circ$ , respectively. These approaches were also effective in producing similar distributions of neural trial-by-trial variability between the two monkeys and between the two brain areas, as measured by the means  $\pm$  SD of the variance-to-mean ratio (Fano factor) across units (*monkey 1*: V4 =  $1.61 \pm 0.65$ ; IT =  $1.62 \pm 0.68$ ; *monkey 2*: V4 =  $1.49 \pm 0.54$ ; IT =  $1.28 \pm 0.36$ ).

As two monkeys performed this task, we recorded neural activity from small populations using 24-channel probes that were acutely lowered into V4 or IT before each session. Our goal was to record the data required to allow us to evaluate the hypotheses presented in Fig. 1, A and B, in the context of estimates of the amount of convergence between V4 and IT, which range from onefold based on measures of neural signals (Rust and DiCarlo 2010) to threefold based on anatomical estimates (DiCarlo et al. 2012). Consequently, we aimed to collect threefold more units from V4. Following a screen for units based on their stability, isolation, and task modulation (see MATERIALS AND METHODS), our data included 650 V4 units and 204 IT units (*monkey 1*: 382 units in V4 and 108 in IT; *monkey 2*: 268 units in V4 and 96 in IT). The data reported here were extracted from trials with correct responses. For most of our analyses, we counted spikes in equal length windows in V4 and IT but adjusted for the difference in latency between the two brain areas (170 ms, V4: 40–210 ms; IT: 80–250 ms following stimulus onset). These windows always preceded the monkeys’ reaction times and thus corresponded to periods of fixation. Consistent with a previous report (Rust and DiCarlo

2012), firing rates in V4 and IT were moderate, and counting spikes in the short windows implied by the monkeys' fast reaction times only produced a few spikes per trial on average (mean firing rate was ~1 spike/trial; V4: 1.34 spikes/trial; IT: 0.92 spikes/trial). In the face of these low spike counts, the statistical power in our analyses arises from combining the responses across many neurons using weighted decoding schemes.

*Different Types of Signals Are Reflected in V4 and IT Units*

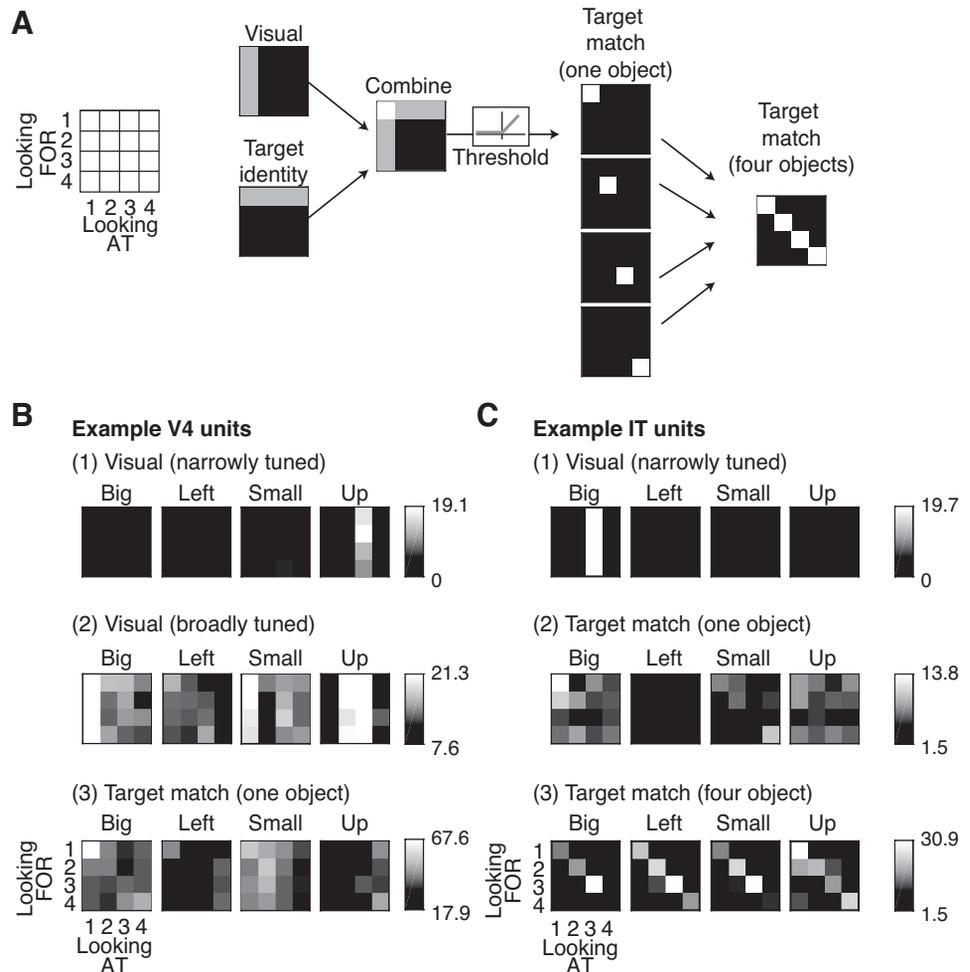
To interpret the different types of signals that might be reflected in V4 and IT during object search, it is useful to conceptualize how signals that differentiate between target matches versus distractors ("target match signals"), which reflect the solution to the task, might be computed. When considered in terms of a single 4 × 4 "looking at" versus "looking for" matrix (Fig. 4A, left), target match signals are reflected as diagonal structure [i.e., you are looking at the same image you are looking for, Fig. 4A, right, "Target match (four objects)"]. In the most straightforward description of target match computation, congruent "visual" information (what you are "looking at"; vertical structure) and "target identity" information (what you are "looking for"; horizontal structure) combine in a nonlinear fashion to compute target match detectors that are selective for one object presented as a target match ["Target match (one object)"]. Finally, these are pooled across

the four different objects to create "Four object target match detectors" that respond whenever a target is in view (Fig. 4A, right). We found examples of these types of idealized units in V4 and/or IT (Fig. 4, B and C). In both areas, we found "purely visual" units that responded selectively to images but were not modulated by other factors, such as target identity or whether an image was presented as a target match (Fig. 4, B and C, "Visual"). In contrast, one notable difference between V4 and IT was the existence of a handful of IT units (~10/204) that reflected the remarkable property of responding to nearly every image presented as a target match (every object at every transformation) but not when those same images were presented as distractors [Fig. 4C, "Target match (four objects)"]. We did not find any such units in V4. However, in both V4 and IT, we found units that responded preferentially to individual objects presented as target matches as compared with distractors [Fig. 4, B and C, "Target match (one object)"]. We note that while these illustrative examples were chosen because they reflect intuitive forms of pure selectivity, many (if not most) units tended to reflect less intuitive mixtures of visual and task-relevant modulation.

*Visual Representations in V4 and IT During the IDMS Task*

When making systematic comparisons between V4 and IT, there are important factors to consider. For example, should the information contained in the V4 and IT populations be com-

Fig. 4. Conceptualizing invariant delayed-match-to-sample (IDMS) task computation. A: an idealized depiction of how target match signals, which reflect the solution to the IDMS task, might be computed. For simplicity, the computation is described for one 4 × 4 slice of the experimental design matrix, which corresponds to viewing each of 4 objects ("Looking AT") in the context of each of 4 objects as a target ("Looking FOR") at one transformation. In the first stage of this idealization of target match computation, a unit reflecting visual information and a unit reflecting persistent target identity information (i.e., working memory) are combined, and the result is passed through a threshold. The resulting unit reflects target match information for one object. Next, 4 of these units (each with a different object preference) are linearly combined to produce a unit that signals whether a target match is present, regardless of the identity of the object. B and C: example single-unit responses. Response matrixes are plotted as four 4 × 4 slices of the experimental design matrix. Each slice corresponds to viewing each of 4 objects (Looking AT) in the context of each of 4 objects (Looking FOR) at 1 of the 4 transformations used. Shown are the response matrixes corresponding to 3 example units from V4 and inferotemporal cortex (IT). Response matrixes were plotted as the average firing rates across trials and rescaled from the minimum (black) to maximum (white) response across all experimental conditions. Scale bars indicate each unit's minimum and maximum average firing rate response in spikes/s.



pared with equal numbers of units? Similarly, what are appropriate benchmarks for determining whether the samples recorded from each brain area are representative? As an example, imagine a scenario in which the information about whether an image is a target match or a distractor is reflected in both V4 and IT to the same degree, but the V4 neurons recorded in an experiment all have small, overlapping receptive fields confined to the same, small region of the visual field. In contrast, IT neurons, by virtue of their large receptive fields, would have access to much more of the visual field. From this data we might erroneously find that the magnitude of total target match information is larger in IT than V4 by way of nonrepresentative sampling.

As a benchmark for assessing whether the data we recorded from each brain area were representative, we compared the amount of visual modulation present in each brain area, at each transformation, with the following rationale. First, all the visual information contained in IT is thought to arrive there after first traveling through V4 (Felleman and Van Essen 1991), and consequently, samples of V4 and IT are comparable only if visual information is equal or higher in the V4 sample. Second, by comparing visual information at each transformation separately, we circumvent issues related to the differences in the format of visual information between the two brain areas described in Fig. 1C.

To compare the amounts of visual information in our recorded V4 and IT populations, we computed a single-unit measure of visual modulation described by (Pagan and Rust 2014b) and included in a number of our earlier reports (Pagan et al. 2013; Pagan and Rust 2014a; Roth and Rust 2018a). The advantage of this measure is that it disentangles modulations due to changes in visual identity from other factors, such as top-down target modulation, whereas many other traditional measures (e.g., single-neuron  $d'$  or receiver operating characteristic) do not. More specifically, while measures like single-unit  $d'$  are in fact proportional to modulation by the variable of interest (e.g., visual modulation), they are also inversely proportional to other types of modulation (e.g., target modulation), when they exist, because these other types of modulation act as a form of noise or nuisance variability that lowers  $d'$  (Pagan and Rust 2014b). Consequently, for a given value of single-neuron  $d'$ , it remains ambiguous whether that value arose from a particular amount of visual modulation in the absence of other modulation types or by larger visual modulation in the presence of other modulation types, and this disambiguation has important consequences for comparing the magnitudes of visual modulation in our recorded V4 and IT data. To resolve this ambiguity, our measure applies an extension of an ANOVA and, like the ANOVA, parses a unit's total response variance into that which can be attributed to modulation by each type of stimulus variable (e.g., changes in the visual image, changes in target identity, etc.), as well as their nonlinear interactions. We also apply a bias-correction procedure to correct for the overestimation of these variances due to limited samples, developed and tested extensively in simulation (Pagan and Rust 2014b). Finally, because variance non-intuitively scales as the square of response changes (i.e., when firing rates double, variances quadruple), we calculate modulation in units of standard deviation, computed as the square root of the bias-corrected variance quantities. In sum, this measure of visual modulation quantifies the modulation in a

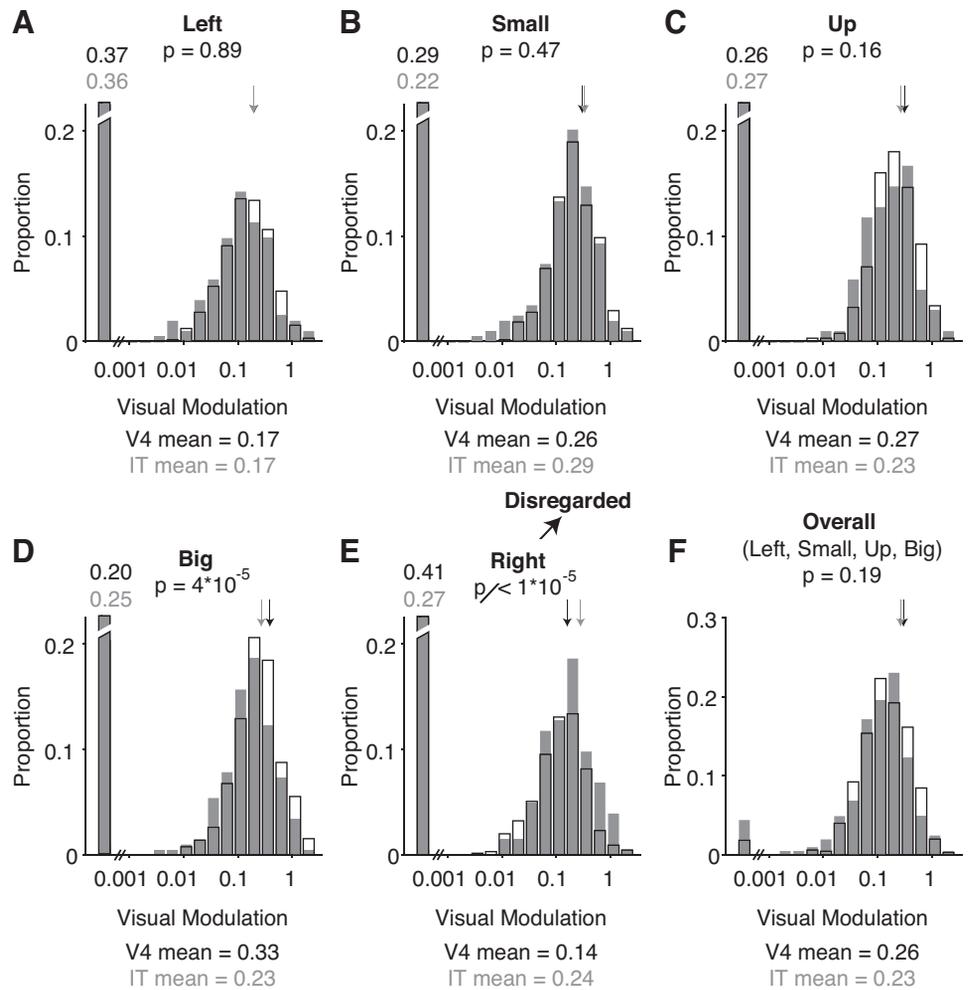
unit's spike count that can be attributed to changes in the identity of the object in view, and it is computed separately for each of the five transformations.

For three of the five transformations (Left, Small, and Up), mean visual modulation was statistically indistinguishable between V4 and IT (Fig. 5, A–C). For one transformation (Big; Fig. 5D) mean visual modulation was larger in V4, but we retained this transformation for subsequent analyses because its incorporation reflected a worst-case scenario against the sampling problem of concern (i.e., one in which V4 has been inadequately sampled). In contrast, for the final transformation (Fig. 5E, Right), the V4 population had significantly lower performance than IT ( $P < 1e10^{-5}$ ), and investigation of the recorded receptive field locations (Fig. 3) revealed that this was likely due to incomplete sampling at that location. As such, we disregarded this transformation from further analyses. Subsequent analyses are focused on the four of five transformations in which visual modulation, averaged across transformations, was not statistically distinguishable in V4 as compared with IT, either in the pooled data or in either monkey (Fig. 5F; *monkey 1*: V4 mean = 0.30, IT mean = 0.27,  $P = 0.07$ ; *monkey 2*: V4 mean = 0.18, IT mean = 0.19,  $P = 0.93$ ). The fact that visual modulation is matched between V4 and IT across these four transformations suggests that the two populations can and should be compared with approximately matched numbers of units, consistent with previous reports (Rust and DiCarlo 2010). We have confirmed via simulation that this approach is capable of detecting scenarios in which the visual information reflected across the V4 population converges to produce a considerably smaller sized IT population; as expected, these cases produce average visual modulation per unit that is larger in IT, as a consequence of concentrating the visual information reflected in V4 in a smaller number of IT units (not shown).

#### Quantifying Top-Down Modulation in V4 and IT

To quantify the magnitudes of these different types of task relevant signals in V4 and IT, we extended the procedure presented in Fig. 5 to not only quantify “visual” modulation (i.e., modulation that can be attributed to changes in the identity of the visual image) but also other types of nonoverlapping modulations that could be attributed to “target identity” modulation: changes in the identity of a sought target; “target match” modulation: changes in whether an image was a target match or a distractor; and residual” modulation: nonlinear interactions between visual and target identity that are not target match modulation including, for example, incongruent nonlinear combinations of visual and target identity (e.g., a unit visually responsive to the combination of *object 1* and *target 2* or equivalently, selectivity for a particular distractor condition). When considered in terms of a single  $4 \times 4$  “looking at” versus “looking for” slice of the experimental design matrix (Fig. 2B), these modulations produce vertical, horizontal, diagonal, and off-diagonal structure, respectively (Fig. 6A). More specifically, this procedure decomposes the total variance into a linear sum of these components, and then transforms variance into units of spike count modulation (as described above). Consequently, all of the response variance of a unit is accounted for with our technique. In these analyses, modulation magnitudes were computed per transformation (i.e., per slice of

Fig. 5. Comparison of visual modulation in V4 and inferotemporal cortex (IT). Shown are distributions of visual modulation magnitudes across units, parsed by transformation, for V4 (open bars,  $n = 650$  units) and IT (gray,  $n = 204$  units) and plotted on a log axis. Following a bias correction to remove the impact of trial variability, visual modulation was computed in units of standard deviation around each unit's grand mean spike count. The first bin includes units with negligible visual modulation (modulation  $< 0.001$ ), and the broken axis indicates that these bars should extend to the proportions labeled just above. Means of each distribution, including units with negligible visual modulation, are indicated by arrows, and values are indicated at the bottom. The  $P$  values at the top were computed via a bootstrap significance test evaluating the probability that differences in the means between V4 and IT can be attributed to chance. A–E: distributions parsed by transformation. Visual modulation corresponding to the transformation “Right” was higher in IT as compared with V4, due to incomplete sampling of receptive fields at this location (Fig. 3) and was thus disregarded from further analyses. F: distributions of visual modulation, averaged for each unit across the transformations “Left,” “Small,” “Up,” and “Big.”



the experimental design matrix) and then averaged to compute mean modulations per transformation. Notably, this analysis defines target match modulation as a differential response to the same images presented as target matches versus distractors, or equivalently, diagonal structure in the transformation slices presented in Fig. 6A. Consequently, units similar to both the “Target match (one object)” unit as well as the “Target match (four objects)” unit (Fig. 4, B and C) reflect target match modulation, as both units have a diagonal component to their responses. What differentiates these two types of units is that the “Target match (one object)” unit also reflects selectivity for image and target identity, which is reflected in this analysis as a mixture of target match, visual, and target identity modulation. Below we both quantify each type of modulation separately, as well as combine target identity, target match and residual modulation into one combined measure of “top-down” modulation.

We applied this analysis to spike count windows positioned at sliding locations relative to stimulus onset, as well as the same counting windows described above (170 ms; V4: 40–210 ms; IT: 80–250 ms; Fig. 6, B–E). As expected, visual modulation did not exist before stimulus onset, and visual signals arrived in V4 ~40 ms earlier than in IT in both animals (Fig. 6B). In contrast, modulations reflecting information about whether an image was a target match or a distractor (‘target match’ modulation) were considerably smaller in V4 as com-

pared with IT in both animals (Fig. 6C; *monkey 1*:  $P < 0.001$ ; *monkey 2*:  $P < 0.001$ ). In *monkey 1*, V4 target match modulations increased throughout the viewing period and reached levels that were similar to those found in IT, but this rise occurred with a delay in V4 relative to IT. This was not replicated in *monkey 2*, where target match modulations in V4 were small throughout the viewing period. When the V4 spike count window was extended to a longer window such that its offset was matched to IT (40–250 ms), target match information remained smaller in V4 than IT (V4: 0.12 vs. IT: 0.20,  $P < 0.001$ ).

Modulations reflecting information about the identity of the target (“target identity” modulation) were present in both V4 and IT before stimulus onset (Fig. 6D), consistent with persistent working memory signals in both brain areas. These persistent signals were stronger in IT as compared with V4 in *monkey 1* ( $P < 0.001$ ) but comparable in size between V4 and IT in *monkey 2* ( $P = 0.23$ ). Lastly, we found that in both V4 and IT, residual modulation was small relative to the other types of modulations (Fig. 6E). Residual modulation was larger in IT than V4 in both monkeys (*monkey 1*:  $P = 0.048$ ; *monkey 2*:  $P = 0.006$ ). To summarize these results, we found that in both monkeys, visual modulation was matched between V4 and IT whereas target match signals were weaker in V4. We also found persistent target identity signals that were reflected

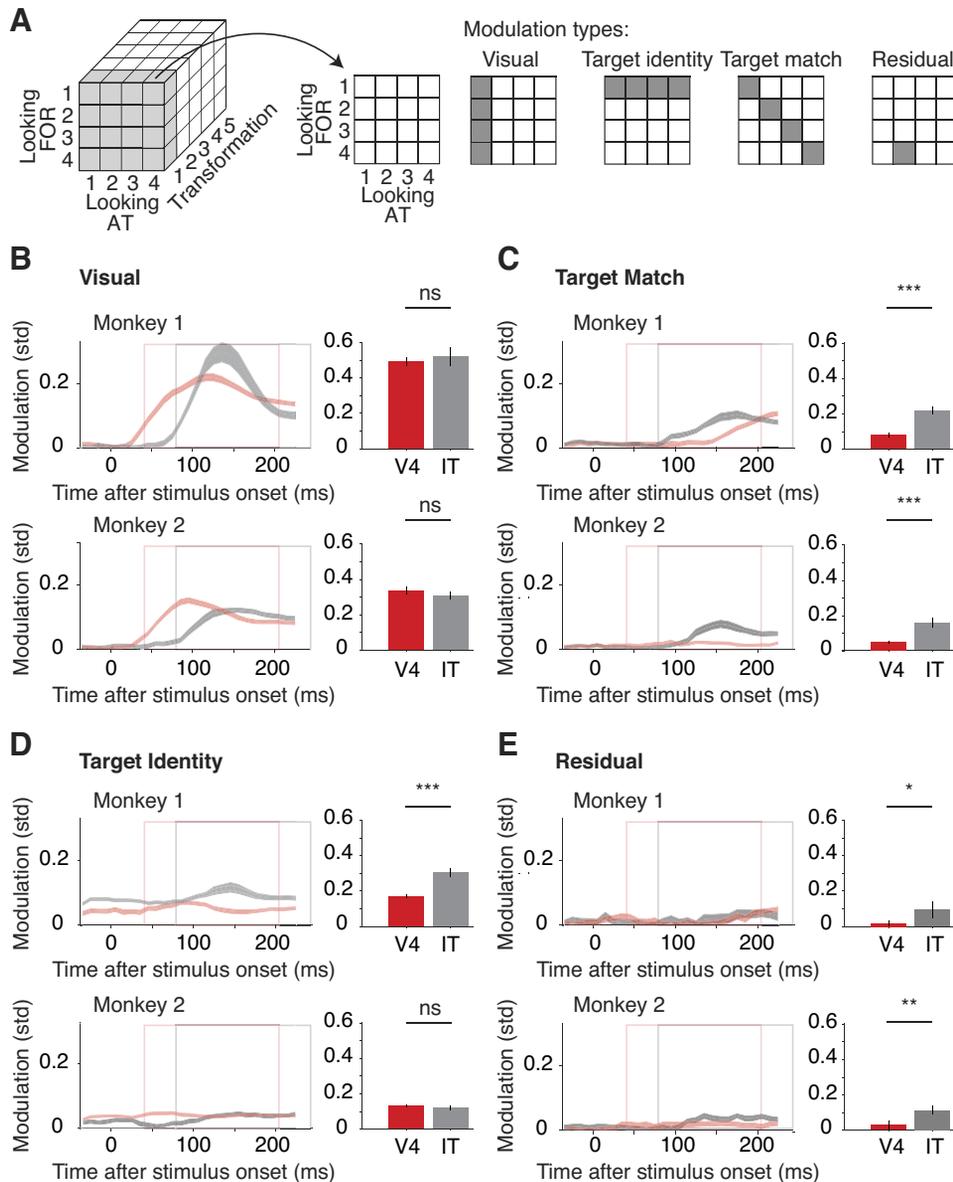


Fig. 6. Time course of different types of single unit modulations in V4 and inferotemporal cortex (IT). **A**: to illustrate the different types of task-relevant signals that could be present in V4 and IT, shown is a slice through the invariant delayed-match-to-sample (IDMS) experimental design (Fig. 2B), corresponding to 1 transformation. Shown are visual modulations, which differentiate between different objects in view (vertical structure); target identity modulations, which differentiate between different target objects (horizontal structure); target match modulations, which differentiate between whether objects appear as a target match vs. a distractor (diagonal structure); and residual modulations, which differentiate between any other types of conditions (e.g., a response to a particular distractor condition such as looking for *object 4* when looking at *object 2*). **B–E**: modulations were computed for each type of experimental parameter in units of the SDs around each unit's grand mean spike count (see RESULTS). In **B–E**, average modulation magnitudes across units in V4 (red;  $n = 650$ ) and IT (gray;  $n = 204$ ) shown on the left as a function of time (ms after stimulus onset). Modulation magnitudes, computed in spike count bins 50-ms wide and shifted by 10 ms, are plotted corresponding to the midpoint of each bin, and consequently, end 25 ms before the termination of the last bin. The bar plots show average signal magnitudes quantified within broader spike counting windows indicated by the rectangles on the left (V4: 40–210 ms, red rectangle; IT: 80–250 ms, gray rectangle). Modulation in broader bins can loosely be envisioned as the integral of modulation in narrower bins; however, due to bias correction, this is not exactly the case (see MATERIALS AND METHODS). \* $P < 0.05$ ; \*\* $P < 0.01$ ; \*\*\* $P < 0.001$ ; ns $P > 0.05$ . Error bars reflect the SE of modulation across units, computed via a bootstrap procedure.

in both areas before and throughout the stimulus-evoked period.

As a complementary analysis, we also quantified the total amount of top-down modulation (combined target match, target identity, and residual modulation) and compared it to the evolution of the visual modulation. In both brain areas, top-down modulation was considerable throughout the analysis window (Fig. 7, *A* and *B*, left). During the latency-corrected stimulus-evoked period, top-down modulations were 34 and 77% the size of the visual modulations in V4 and IT, respectively (Fig. 7, *A* and *B*, right). These results demonstrate that considerable nonvisual, task-relevant modulations exist in both brain areas, and they also suggest that these are smaller in V4 as compared with IT. Total top-down information also remained smaller in V4 than IT even when the V4 spike count window was extended to a longer window such that its offset was matched to IT (40–250 ms; V4: 0.22 vs. IT: 0.32,  $P = 0.002$ ). Additionally, because the results presented thus far were computed per transformation and then averaged, we

examined the degree to which they held for each transformation individually (Fig. 7C). Total top-down modulations were significantly smaller in V4 than IT in all cases (Big:  $P = 0.004$ ; Left:  $P < 0.001$ ; Small:  $P < 0.001$ ; Up:  $P < 0.001$ ), verifying that our results are not specific to a subset of the transformations. To address the concern that lower top-down modulation in V4 might be a consequence of smaller V4 receptive fields and consequently averaging across transformations at which V4 units are not visually responsive, we performed a control analysis in which we recomputed top-down modulation after screening for units with visual modulation  $> 0$  at each transformation independently (Fig. 7D). Top-down modulation remained significantly smaller in V4 as compared with IT at all transformations (Big:  $P = 0.018$ ; Left:  $P = 0.007$ ; Small:  $P = 0.024$ ; Up:  $P < 0.001$ ).

In summary, the fact that average visual modulation is matched in V4 and IT implies that the top-down modulation reflected in the two brain areas can and should be compared with equal numbers of units (or equivalently, by comparing

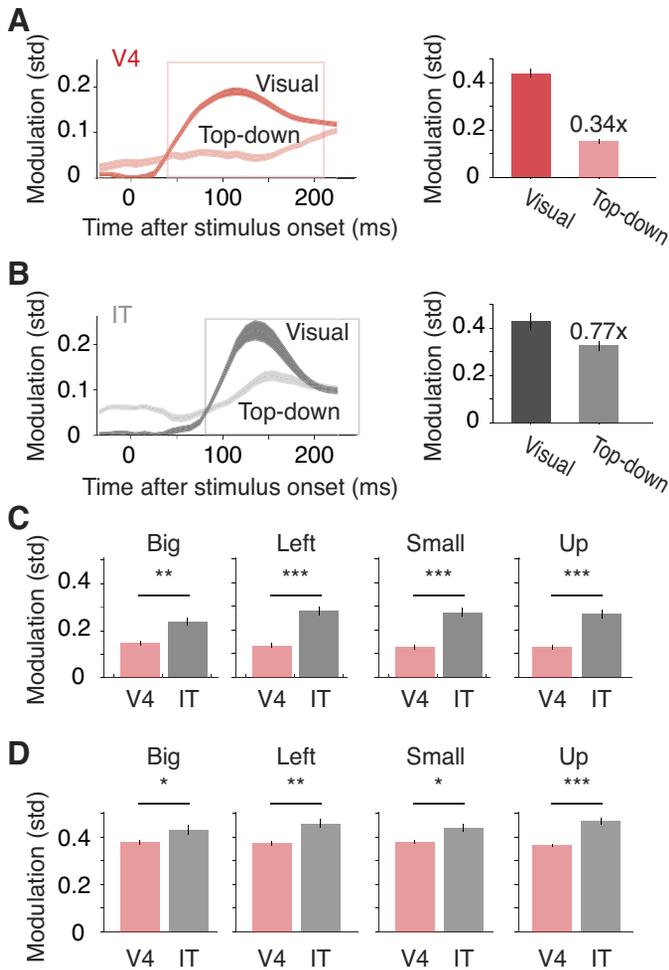


Fig. 7. Top-down modulations in V4 and inferotemporal cortex (IT). *A* and *B*: top-down modulations (V4: light red; IT: light gray) were computed as the sum of target identity, target match, and residual modulations and are shown alongside visual modulations (V4: dark red; IT: dark gray). Mean modulation magnitudes are computed in the same manner and shown with the same conventions as Fig. 6. Labels in the bar plots above the top-down modulation magnitudes indicate the proportional size of top-down relative to visual modulations in each brain area. *C*: mean top-down modulations in V4 and IT, computed as in *A* and *B*, but for each transformation separately. *D*: top-down modulations recomputed after screening for units that were visually responsive (visual modulation  $>0$ ). In *A–C*, error bars reflect the SE of modulation across units, computed via a bootstrap procedure. \* $P < 0.05$ ; \*\* $P < 0.01$ ; \*\*\* $P < 0.001$ .

average modulations per unit). When compared in this way, top-down modulation is consistently smaller in V4 than IT, even when it is pooled across all the different ways in which it can be reflected (Fig. 7, *A* and *B*) as well as when units are first screened for visual responsiveness (Fig. 7*D*).

#### Population-Based Comparisons of Top-Down Modulation in V4 and IT During the IDMS Task

The IDMS task involved the monkeys viewing images and determining whether they were target matches (requiring a saccade) or distractors (requiring fixation) and thus can be reconceptualized as a two-way classification of the same images presented as target matches versus distractors. When applied to neural data, the information available for this task could be linearly formatted (Fig. 8, *A* and *C*) or it could be

nonlinear (Fig. 8, *B* and *E*). Linear target match information is reflected as modulation along the diagonal of each unit's response matrix (Fig. 8*A*) whereas nonlinear target match modulation can be reflected as visual and top-down modulation reflected in the same or different units (Fig. 8*B*). Under the assumption that visual modulations are larger than top-down modulations (Fig. 7, *A* and *B*), top-down modulations serve as the bottleneck for target match task performance, and task performance can, in turn, be used as a proxy for comparing the magnitudes of top-down modulation in V4 and IT. To compare the amounts of linearly and nonlinearly formatted information available in V4 and IT during the IDMS task, we compared the performance of a handful of weighted linear and nonlinear decoders applied to each population. By weighting each neuron (e.g., in the case of linear decoding, proportional to the amount of task-relevant information that it carries), this process ensures that a unit's responses were not considered in situations where that unit was not informative. Additionally, we applied the decoders to the neural responses to each transformation separately and then averaged decoder performance across transformations to account for differences in the format of visual information between V4 and IT (Fig. 1*C*). Decoding analyses were always performed with equal numbers of target matches and distractors (see MATERIALS AND METHODS) to avoid bias.

*Linear target match information, uniform sampling.* To quantify the amount of linearly separable target match information in V4 and IT, we computed the cross-validated performance of a FLD, which weights each unit proportional to its  $d'$  for this task, adjusted for any correlations that exist between units (see MATERIALS AND METHODS). To determine whether weighted, uniform sampling of V4 could account for target match information in IT, we randomly selected IT units up to the total numbers of units that we recorded (Fig. 8*D*, gray) and compared this to a random selection of V4 units for matched sized populations (and thus always a subset of the V4 data Fig. 8*D*, red). Cross-validated population performance was higher than chance in V4 ( $P < 0.001$ ) but was significantly higher in IT as compared with V4 in both monkeys (Fig. 8*D*, gray vs. red;  $P < 0.001$ ). We also compared population performance when all V4 units were included for each monkey ( $\sim 3\times$  units in V4 as compared with IT; Fig. 8*D*, red open circles) and found that performance in V4 remained considerably lower than IT ( $P < 0.001$ ). These results suggest that IT target match information is not directly inherited from V4 under a model that proposes uniform sampling of V4 by IT and up to threefold convergence.

*Linear target match information, best unit sampling.* Evidence from other studies suggests that the brain can learn to preferentially readout the subset of neurons that carry the most task-relevant information with extensive training (Law and Gold 2009), and the monkeys involved in these experiments were trained extensively. Could a version of the feed-forward proposal in which IT preferentially samples the “best” V4 neurons account for our data? To allow us to address this question, we sampled threefold more units in V4 as compared with IT, consistent with anatomical estimates of the ratios of neurons between the two brain areas (DiCarlo et al. 2012). To assess whether a “best unit” sampling description of V4 by IT could account for our data, we recomputed performance for V4 and IT populations that were matched in size but when only the

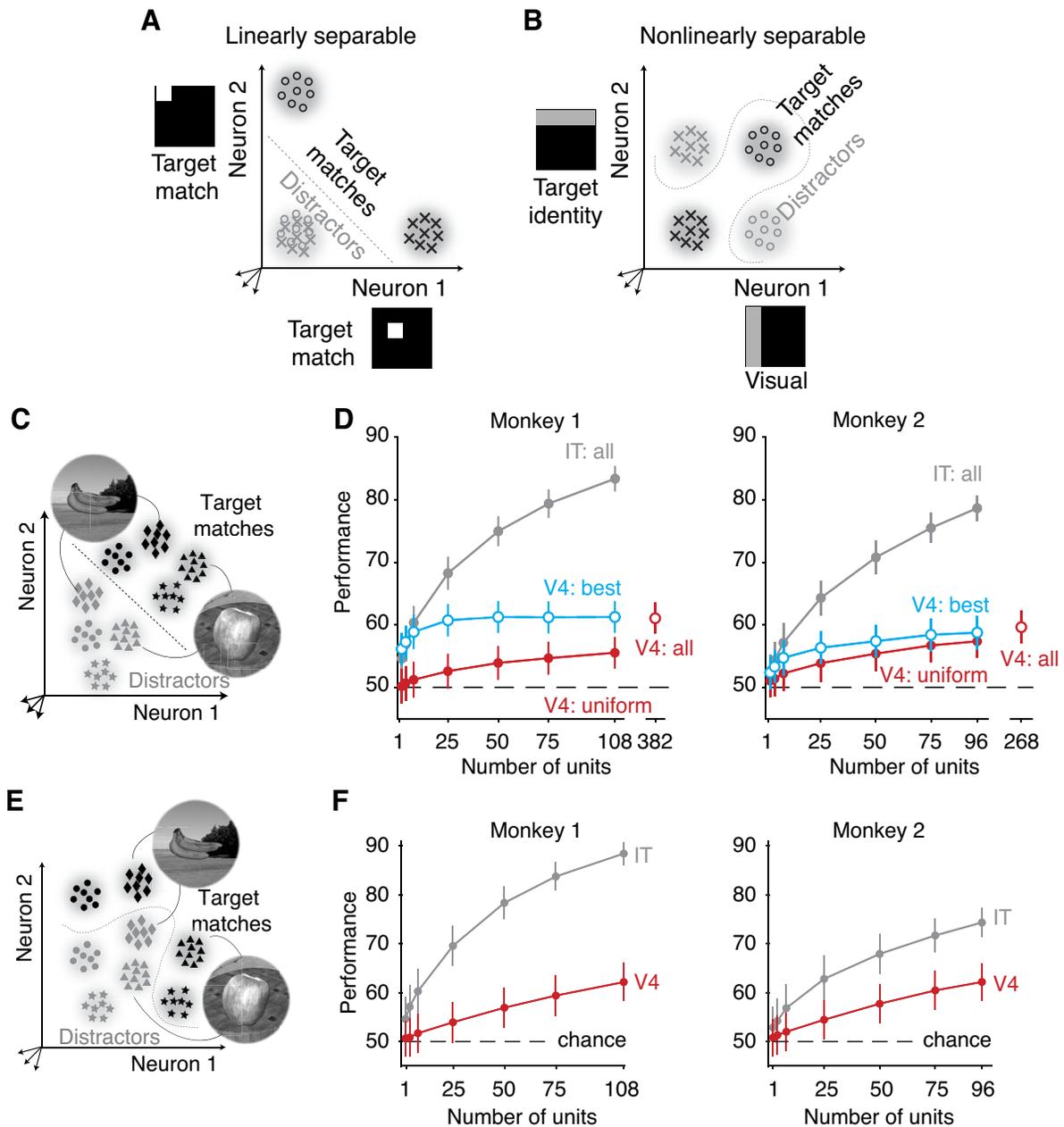


Fig. 8. Population-based comparisons of top-down information in V4 and inferotemporal cortex (IT). *A* and *B*: the invariant delayed-match-to-sample (IDMS) task is presented as a 2-way classification of the same images presented as target matches vs. as distractors with a linear decision boundary. Each point depicts a hypothetical population response for a population of 2 neurons on a single trial, and clusters of points depict the dispersion of responses across repeated trials for the same condition. Included are the hypothetical responses to the same images presented as target matches (black) and as distractors (gray). *A*: an illustration of the linearly separable target match information reflected in a population comprised of 2 hypothetical “Target match (one object)” units, selective for different objects. *B*: an illustration of the nonlinearly separable target match information reflected in a population comprised of 1 unit with visual identity modulation and the other with target identity modulation, both selective for the same object. *C*: an illustration of linear separability for 4 images presented as target matches vs. as distractors. *D*: performance of a linear classifier trained to classify whether an object was a target match or a distractor, invariant of object identity (at each transformation and averaged across transformations). Red: uniform subsampling of V4 units; cyan: sampling the best *N* V4 units based on the training data; gray: uniform sampling of IT units. Total numbers of IT units: *monkey 1*: *n* = 108 units, *monkey 2*: *n* = 96 units. Total numbers of V4 units: *monkey 1*: *n* = 382, *monkey 2*: *n* = 268. Error bars (SE) reflect the variability that can be attributed to the specific subset of trials chosen for training and testing and, for subsets of units smaller than the full population, the specific subset of units chosen. Dashed line indicates chance performance. *E*: an illustration of nonlinear separability for four images. *F*: performance of a nonlinear, maximum likelihood classifier trained to classify whether an object was a target match or a distractor, invariant of object identity. Performance was assessed at each identity-preserving transformation and then averaged. Error bars (SE) reflect the variability that can be attributed to the specific subset of trials chosen for training and testing, and, for subsets of units smaller than the full population, the specific subset of units chosen. Dashed line indicates chance performance.

$N$  top-ranked V4 units were included for different sized  $N$ . Notably, this analysis differs from the analysis presented above in which all V4 units were included insofar as the computation combined with the assignment of weights based on limited samples is a process that contains some noise, and in scenarios where only a subset of units carry a signal, performance is expected to increase until all signal-carrying units are included, but then can fall slightly with the inclusion of units that do not reflect any signal. In this analysis, units were ranked based on the training data before computing cross-validated performance. We found that V4 performance was slightly higher for the best units as compared with randomly selected units (Fig. 8D, cyan vs. red); however, performance for the best V4 units remained lower than IT performance in both monkeys (Fig. 8D, cyan vs. gray,  $P < 0.001$ ).

These results suggest that during IDMS IT target match modulation cannot be accounted for via feed-forward propagation of this modulation from V4, even if IT were to sample from the best V4 subset.

*Nonlinear target match information.* The results above suggest that IT did not inherit its top-down information from V4 in a manner which can be described by purely linear computation (either by sampling uniformly or preferentially sampling from the best V4 units). However, it could be the case that some of the information differentiating target matches and distractors was present at the level of V4 but in a nonlinear format. One example of nonlinearly formatted target match information is the case in which visual modulation and target identity modulation are reflected in different units (Fig. 8B). To quantify the “total” target match information in V4 and IT, regardless of its format, we measured cross-validated performance for a maximum likelihood (as opposed to linear) classifier (see MATERIALS AND METHODS). We confirmed via simulation that the nonlinear decoding analysis we apply is capable of decoding nonlinearly formatted information (e.g., Fig. 8B) to determine whether a target match is present (not shown). When applied to the data, cross-validated population performance was higher than chance in V4 (Fig. 8F, red; in both monkeys, compared at  $n = 108$  in *monkey 1* and  $n = 96$  in *monkey 2*,  $P < 0.001$ ) but was also higher in IT as compared with V4 (Fig. 8F, gray; in both monkeys, compared at  $n = 108$  in *monkey 1*,  $P < 0.001$  and  $n = 96$  in *monkey 2*,  $P = 0.009$ ).

*Population-based comparisons, summarized.* Taken together, these results rule out all variants of the IT:inherited proposal presented in Fig. 1A, including descriptions in which IT preferentially samples from the best V4 units (Fig. 8D, cyan vs gray), as well as descriptions that allow for nonlinear computation on the information arriving in IT from V4 (Fig. 8F). Because these results thus suggest that IT target match information is not exclusively inherited via feed-forward projections arriving from V4, we can conclude that, as suggested by the IT:integrated proposal, at least some component of top-down information is integrated directly in IT during the IDMS task (Fig. 1B).

#### *The Relationship Between the Visually Evoked Response and Target Match Modulation in Individual Units*

The results presented in Fig. 8 establish that during the IDMS task target match modulation is stronger in IT. How does the relationship between the visually evoked response and

the magnitude of target match modulation manifest in the responses of individual units in each brain area?

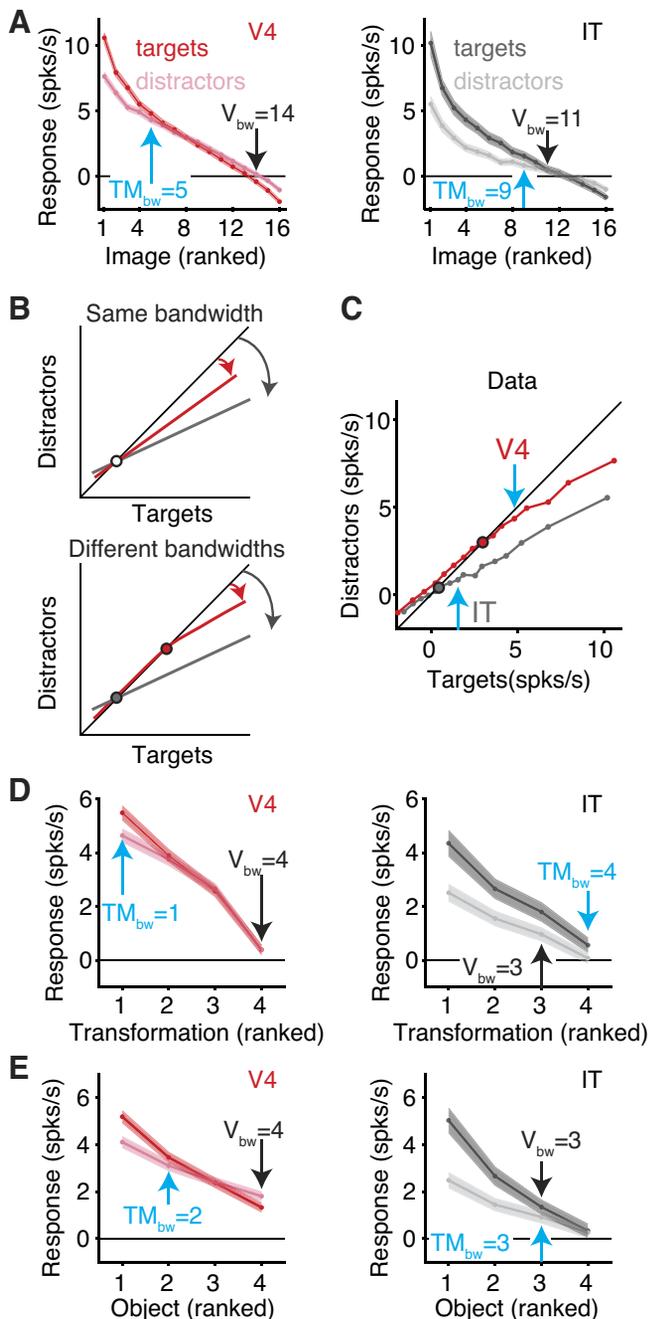
Conceptually, the increase in target match modulation in IT as compared with V4 could be reflected in a similar manner in the two brain areas, such as multiplicative rescaling, but with a larger multiplicative factor in IT. In this scenario, both brain areas would exhibit tuning for target match modulation that is matched in bandwidth to the visually evoked tuning for images. Alternatively, increases in target match modulation in IT over V4 could result from similar target match modulation magnitudes for each image at which it appears but with broader tuning across images in IT. We note that these two possibilities are not mutually exclusive.

To explore these issues, we computed the magnitude of target match modulation as a function of the effectiveness of each image at driving individual units to fire. Specifically, we ranked the responses of each unit to the 16 images (4 objects  $\times$  4 transformations; averaged across target matches and distractors) and then examined the responses to target matches and distractors separately, after averaging across all units (Fig. 9A). Shown are results computed after subtracting the baseline firing rate from each unit's responses, where negative values correspond to responses below baseline. In both V4 and IT, target match modulation for the highest ranked images was reflected as target match enhancement, and both brain areas exhibited weak target match suppression for the least effective images (Fig. 9A). The net impact of target match modulation (averaged across all ranks) was target match enhancement in both brain areas, with stronger net modulation in IT (the difference in the average response to target matches and distractors, averaged across all ranks, divided by the average response to target matches: V4 = +0.07 spks/s; IT = +0.46 spks/s; Fig. 9A).

Next, we computed tuning bandwidths separately for the visually evoked response and for target match enhancement. To determine the bandwidth of the visually evoked response, we determined the number of image ranks with responses significantly above the baseline response (for data averaged across target matches and distractors; Wilcoxon rank sum test; criterion  $P = 0.05$ ). We found the visually evoked bandwidths to be slightly broader in V4 (bandwidth = 14 images; Fig. 9A, left, black arrow) as compared with IT (bandwidth = 11 images; Fig. 9A, right, black arrow). To determine the bandwidths of target match enhancement, we computed the number of images for which the responses to target matches were significantly higher than the responses to distractors (Wilcoxon rank sum test; criterion  $P = 0.05$ ). In V4, target match enhancement was more narrowly tuned than the visually evoked response (target match enhancement bandwidth = 5 images; visually evoked response bandwidth = 14 images; Fig. 9A, left, cyan vs. black arrows). In IT, these bandwidths were more similar (target match enhancement bandwidth = 9 images; visually evoked response bandwidth = 11 images; Fig. 9A, right, cyan vs. black arrows).

As a complementary analysis to determine the degree to which target match modulation took on the same functional form in V4 and IT relative to the “multiplicative” benchmark, we plotted the responses to target matches against the responses to distractors, similar to a quantile-quantile or Q-Q plot. In these plots, rescaling the responses by different multiplicative factors but with a fixed bandwidth amounts to

changing the slope of the lines around a fixed pivot point (Fig. 9B, *top*) whereas changing the bandwidths results in different points of intersection with the unity line (Fig. 9B, *bottom*). Our data were not well described as a change in slope around a fixed pivot point but rather by more narrowly tuned target match enhancement in V4 (Fig. 9C). Notably, this narrower bandwidth cannot be explained by a more narrowly tuned visually evoked response in V4, because the visually evoked response was slightly broader in V4 as compared with IT (as described above, Fig. 9A). These results suggest that increased target match modulation in IT over V4 is not a result of the same multiplicative process but with different multiplicative factors in the two brain areas, rather, in IT, it is both larger for the image ranks at which it appears and it is more broadly tuned across images.



The 16 images included in this analysis differed from one another across changes in both object identity as well as object transformation. To determine whether broader target match enhancement in IT as compared with V4 depended differentially on these two factors, we performed the same analysis by ranking each unit by transformation (after averaging across object identity; Fig. 9D) or ranking by object identity (after averaging across transformations; Fig. 9E). Across transformations, the bandwidth of target match enhancement was considerably narrower in V4 (1 transformation; Fig. 9D, *left*, cyan arrow) as compared with IT (4 transformations; Fig. 9D, *right*, cyan arrow). Once again, this difference could not be accounted for by a difference in the bandwidth of the visually evoked response, as this was broader in V4 (4 transformations; Fig. 9D, *left*, black arrow) than IT (3 transformations; Fig. 9D, *right*, black arrow). Across object identity, the bandwidth of target match enhancement was slightly narrower in V4 (2 objects; Fig. 9E, *left*, cyan arrow) as compared with IT (3 objects; Fig. 9E, *right*, cyan arrow) whereas the visually evoked response remained slightly broader in V4 (4 objects; Fig. 9E, *left*, black arrow) as compared with IT (3 objects; Fig. 9E, *right*, black arrow).

Together, these results suggest that the increase in target match modulation in IT as compared with V4 can be attributed to both an increase in the magnitude of target match modulation for the images that are most effective at driving individual units to fire as well as an increase in bandwidth over the number of images that have this modulation. These IT bandwidth increases, in turn, appear to follow from broader tuning of target match enhancement across both identity-preserving transformations as well as object identity, with a larger influence of the former as compared with the latter.

Fig. 9. Comparison of the visually evoked response and target match modulation in individual units. **A**: the spike count responses to target matches [V4: *left*, red; inferotemporal cortex (IT): *right*, dark gray] and distractors (V4: *left*, pink; IT: *right*, light gray) to the 16 different images that made it through the screen presented in Fig. 5, ranked by image preference. After subtraction of the baseline responses (computed 170 ms before the onset of the first image in each trial), the responses of each unit to each image were ranked by their average response across target matches and distractors. Shown are means  $\pm$  SE across all units, after sorting each unit by image rank (V4:  $n = 650$ ; IT:  $n = 204$ ). The bandwidth of the visually evoked response ( $V_{bw}$ ; black arrows) was computed as the number of images for which the response (averaged across targets and distractors) was significantly larger than baseline (Wilcoxon rank sum test; criterion  $P = 0.05$ ). Target match bandwidth ( $TM_{bw}$ ; cyan arrows) was computed as the number of images for which the target match response was significantly higher than the distractor response (Wilcoxon rank sum test; criterion  $P = 0.05$ ). **B**: schematics showing the hypothetical responses to distractors plotted against hypothetical responses to target matches, similar to a quantile-quantile (Q-Q) plot. *Top*: shown are two hypothetical populations whose responses are both multiplicatively rescaled but with different factors. This manifests in these plots as changes in the slope of these lines around a fixed pivot point (white dot). *Bottom*: shown are 2 hypothetical populations whose responses differ by both their magnitudes of target match modulation as well as the bandwidth of this tuning; changes in bandwidth manifest as different points of intersection with the unity line (red and gray large dots). **C**: V4 and IT spike count responses, plotted as described for **B**. Points of intersection with the unity line are indicated with large red and gray dots and for comparison, the cyan arrows from **A** are included. **D**: responses to ranked images (averaged across object identity), plotted with the same conventions as **A**. **E**: responses to ranked objects (averaged across transformation identity), plotted with the same conventions as **A**.

## DISCUSSION

Finding sought objects requires the brain to compare visual information about the objects in view with information about the currently sought target to compute a signal that reports when a target match has been found. During object search, information about the identity of a sought target and/or whether it is a target match is thought to be fed back to mid to higher stages of the ventral visual pathway, including V4 and IT, but the specific path this information takes is unclear. In this study, we sought to differentiate between scenarios in which top-down information is integrated directly in IT (Fig. 1B) versus those in which it is integrated in V4 and arrives in IT via feed-forward propagation (Fig. 1A). We evaluated a number of feed-forward descriptions between V4 and IT and found none of them could account for the amount of nonvisual, task-relevant information present in IT. These included a model in which IT uniformly samples target match signals from V4 (Fig. 8D, red), a model in which IT preferentially samples target match signals from the best V4 units (Fig. 8D, cyan), and a model that allowed for IT nonlinear processing of inputs arriving from V4 (Fig. 8F). Together, these results suggest that during IDMS, top-down, task-specific signals in IT are not exclusively inherited from V4 but rather are integrated within IT, at least in part.

Our study employed a population-based approach in which we recorded neural responses in an unbiased manner while the monkeys performed the task and we counted spikes in the short windows implied by the monkeys' fast reaction times. Our approaches are a scientific advance over other approaches that artificially increase the signal-to-noise ratio in a data set in a manner that implicitly makes unrealistic assumptions about the brain. For example, tailoring stimuli to align with the peak of each neuron's preferences disregards the contributions of neurons that are activated along their curve flank, and counting spikes in long windows that exceed natural reaction times assumes that neural responses are stationary, but they are not (e.g., Pagan and Rust 2014a). We emphasize that while these types of "increase the SNR approaches" continue to play a fundamental role in the study of visual processing (and we continue to be champions of them in those contexts), the questions under investigation here, which relate to the path of information flow in the brain during the IDMS task, are better studied in a more assumption free manner, using the types of population-based approaches we apply here. In the context of the IDMS task, reaction times are naturally fast, implying spike count windows that are short (e.g., after accounting for latency, 170 ms). In such cases, the brain is expected to operate in a low spike count regime and modulation magnitudes are expected to be small. Similarly, when stimuli are not artificially optimized for each neuron's preferences, firing rate histograms are expected to have long tails. The way that the brain deals with these types of issues is, of course, by combining the neural responses across many (noisy) neurons via combination rules adjusted for the task at hand (e.g., a weighted linear population readout). The statistical power in our data arises, similarly, from population-based analyses.

In our study, we were careful to consider whether the larger IT top-down information we observed (relative to V4) was not being confounded with visual information (i.e., stimulus tuning). First, our analyses go to great lengths to demonstrate that

larger IT top-down information does not follow from differences in the total amount of visual information between V4 and IT. Specifically, we demonstrate that all the visual information available in the IT data is also present in randomly sampled populations of a brain area known to provide the input to IT, V4, with approximately matched numbers of units (Fig. 5). This establishes that we have recorded from comparable V4 and IT populations and that the visual information in IT can be described as being inherited from V4. If IT also inherited its top-down information from V4, we would expect that when visual information is matched between V4 and IT, top-down information would be matched as well. However, it was not: we found that total top-down information was considerably larger in IT than V4 even when visual information was matched, implying that while visual information can be described as feed-forward, top-down task-relevant information cannot and thus must be integrated directly within IT itself. It might be of further concern that the larger IT top-down information we observed (relative to V4) erroneously follows from differences in the format of visual information between V4 and IT (e.g., the fact that V4 neurons have smaller receptive fields). We emphasize that our experiment exploited rather than controlled for differences in the format of visual information between V4 and IT, by testing where top-down signals are integrated in a task that required spatial invariance. This experimental design was motivated by findings that top-down effects appear globally across the visual field in V4 despite the small sizes of V4 receptive fields, and we thus wondered whether spatially global, top-down integration in V4 could account for top-down signals in IT within a feed-forward framework.

We found nonvisual, task-specific signals to be sizeable in V4 (~35% of the size of visual modulation), consistent with many other reports (Bichot et al. 2005; Chelazzi et al. 2001; Cohen and Maunsell 2009; Haenny et al. 1988; Hayden and Gallant 2005; Kosai et al. 2014; Luck et al. 1997; McAdams and Maunsell 1999, 2000; Mirabella et al. 2007; Moran and Desimone 1985; Motter 1994a, 1994b; Ogawa and Komatsu 2004). At the same time, we also found nonvisual, task-specific modulations to be even larger in IT (~75% the size of visual modulation). In a previous study, during a visual target search task in which monkeys made a saccade to a target match following the presentation of a sample image, nonvisual, task-specific signals were reported to be more similar in V4 and IT (63 and 70% of the visually evoked response in V4 and IT, respectively; Chelazzi et al. 1998, 2001). One notable difference between our study and this earlier work is that our study compared V4 and IT during a version of the delayed-match-to-sample task in which sought target objects could appear at different positions, sizes and background contexts. The fact that top-down, task-specific signals were considerably larger in IT versus V4 in our task may follow from the fact that IT contains a more explicit, linear representation of object identity across these transformations than V4 (reviewed by DiCarlo et al. 2012). Consequently, top-down modulation may be targeted directly to IT in situations that require an invariant object representation whereas the brain might target the pathway differently when tasks have different computational requirements. For example, because V4 receptive fields are smaller and retinotopically organized, V4 might serve as the primary locus for the integration of top-down signals for tasks that

require spatial specificity, such as covert spatial attention tasks, and in these tasks little top-down integration might occur in IT (Moran and Desimone 1985). Only one earlier study has reported on the responses of IT neurons in the context of a DMS task in which objects could appear at different identity-preserving transformations (Lueschow et al. 1994), but this study did not measure signals in V4.

Our results, which demonstrate larger nonvisual, task-relevant modulations in IT as compared with V4, are consistent with more general interpretations that the magnitudes of top-down modulation exist in a gradient-like fashion hierarchically along the ventral visual pathway (reviewed by Noudoost et al. 2010). As described above, such gradients are consistent both with the integration of top-down modulation at multiple stages of the pathway (Fig. 1B, red) as well as integration at a single locus, followed by feedback within the pathway itself (Fig. 1B, cyan). One study (Buffalo et al. 2010) provided evidence supporting the latter description in V1, V2, and V4 in the form of noting that not only the magnitude of modulation was greater in higher visual areas, but it also arrived earlier, consistent with a feed-back description. In our data, this issue was ambiguous: we found that in one monkey the arrival of target match modulation was delayed in V4 as compared with IT (Fig. 6C, *monkey 1*) whereas in the other monkey, target match modulation was small in V4 throughout the viewing period (Fig. 6C, *monkey 2*).

In an earlier series of reports, we compared the responses of IT and its projection area, perirhinal cortex, during a more classic version of the delayed-match-to-sample task (that did not incorporate variation in the objects' transformations; Pagan et al. 2013, 2016; Pagan and Rust 2014a). We found that the responses of perirhinal cortex were well-described by a model in which top-down, task-relevant signals were integrated within or before IT consistent with a feed-forward process between IT and perirhinal cortex. The results presented here extend this understanding to suggest that the locus of top-down integration during DMS search tasks is unlikely to exclusively be V4 and that some amount of top-down integration is likely to happen directly within IT itself.

Computing a target match signal requires the combination of the visual representation of the currently viewed scene with a remembered representation of the sought target (e.g., Fig. 4A). In an analysis of the same IT data presented here, we found that the IT population misclassified trials on which the monkeys made errors, supporting notions that the IT target match signal is in fact related to the neural signals used to make target match behavioral judgments (Roth and Rust 2018a). The additional target match information present in IT that is not also present in V4 could reflect the implementation of this comparison in IT itself, or alternatively, the comparison might be implemented in a higher order brain area and fed back to IT cortex. The timing of the arrival of this signal in IT (which peaks at ~150 ms; Fig. 6C) relative to the monkeys' median reaction times (~335 ms; Fig. 2E), does not rule out the former scenario, but with our data we cannot definitively distinguish between these alternatives. Additionally, in this study monkeys were trained extensively on the images used in these experiments and future experiments will be required to address the degree to which these results hold under more everyday conditions in which monkeys are viewing images and objects for the first time.

In this study, the IDMS task was implemented in short blocks with a fixed target such that a cue image did not need to appear at the beginning of each trial. This design was inspired by studies reporting that when a cue image is presented shortly before a target match (e.g., in the same trial), active target match modulation intermingles with passive, repetition suppression (Miller and Desimone 1994), which we have also found to be true (Pagan et al. 2013). For example, whereas approximately balanced numbers of IT units reflect net target match enhancement versus target match suppression in a classic DMS design (Pagan et al. 2013), IT units nearly universally reflect target match enhancement during the IDMS task (Roth and Rust 2018a). This distinction is important for this study, as our goal was to systematically compare the magnitudes of top-down modulation in V4 and IT whereas repetition suppression is likely to be a largely feed-forward process. At the same time, the block design we employed here also imposes some drawbacks, including less frequent changes of the target identity and thus less frequent changes of the issuance of new top-down information. Future experiments will be required to assess how these results compare in the context of more dynamic and more realistic object search conditions (where presumably one ceases to look for the same target upon finding it).

#### ACKNOWLEDGMENTS

We thank Margot P. Wohl and Krystal Henderson for technical contributions.

#### GRANTS

This work was supported by National Eye Institute Grant R01-EY-020851, Simons Foundation (through an award from the Simons Collaboration on the Global Brain), and McKnight Endowment for Neuroscience.

#### DISCLOSURES

No conflicts of interest, financial or otherwise, are declared by the authors.

#### AUTHOR CONTRIBUTIONS

N.R. and N.C.R. conceived and designed research; N.R. performed experiments; N.R. and N.C.R. analyzed data; N.R. and N.C.R. interpreted results of experiments; N.R. and N.C.R. prepared figures; N.R. and N.C.R. drafted manuscript; N.R. and N.C.R. edited and revised manuscript; N.R. and N.C.R. approved final version of manuscript.

#### REFERENCES

- Bichot NP, Rossi AF, Desimone R. Parallel and serial neural mechanisms for visual search in macaque area V4. *Science* 308: 529–534, 2005. doi:10.1126/science.1109676.
- Buffalo EA, Fries P, Landman R, Liang H, Desimone R. A backward progression of attentional effects in the ventral stream. *Proc Natl Acad Sci USA* 107: 361–365, 2010. doi:10.1073/pnas.0907658106.
- Chelazzi L, Duncan J, Miller EK, Desimone R. Responses of neurons in inferior temporal cortex during memory-guided visual search. *J Neurophysiol* 80: 2918–2940, 1998. doi:10.1152/jn.1998.80.6.2918.
- Chelazzi L, Miller EK, Duncan J, Desimone R. A neural basis for visual search in inferior temporal cortex. *Nature* 363: 345–347, 1993. doi:10.1038/363345a0.
- Chelazzi L, Miller EK, Duncan J, Desimone R. Responses of neurons in macaque area V4 during memory-guided visual search. *Cereb Cortex* 11: 761–772, 2001. doi:10.1093/cercor/11.8.761.
- Cohen MR, Maunsell JH. Attention improves performance primarily by reducing interneuronal correlations. *Nat Neurosci* 12: 1594–1600, 2009. doi:10.1038/nn.2439.

- Cohen MR, Maunsell JH.** Using neuronal populations to study the mechanisms underlying spatial and feature attention. *Neuron* 70: 1192–1204, 2011. doi:10.1016/j.neuron.2011.04.029.
- Desimone R, Schein SJ.** Visual properties of neurons in area V4 of the macaque: sensitivity to stimulus form. *J Neurophysiol* 57: 835–868, 1987. doi:10.1152/jn.1987.57.3.835.
- DiCarlo JJ, Zoccolan D, Rust NC.** How does the brain solve visual object recognition? *Neuron* 73: 415–434, 2012. doi:10.1016/j.neuron.2012.01.010.
- Eskandar EN, Richmond BJ, Optican LM.** Role of inferior temporal neurons in visual memory. I. Temporal encoding of information about visual images, recalled images, and behavioral context. *J Neurophysiol* 68: 1277–1295, 1992. doi:10.1152/jn.1992.68.4.1277.
- Felleman DJ, Van Essen DC.** Distributed hierarchical processing in the primate cerebral cortex. *Cereb Cortex* 1: 1–47, 1991. doi:10.1093/cercor/1.1.1.
- Gattass R, Sousa AP, Gross CG.** Visuotopic organization and extent of V3 and V4 of the macaque. *J Neurosci* 8: 1831–1845, 1988. doi:10.1523/JNEUROSCI.08-06-01831.1988.
- Gibson JR, Maunsell JH.** Sensory modality specificity of neural activity related to memory in visual cortex. *J Neurophysiol* 78: 1263–1275, 1997. doi:10.1152/jn.1997.78.3.1263.
- Haenny PE, Maunsell JH, Schiller PH.** State dependent activity in monkey visual cortex. II. Retinal and extraretinal factors in V4. *Exp Brain Res* 69: 245–259, 1988. doi:10.1007/BF00247570.
- Hayden BY, Gallant JL.** Time course of attention reveals different mechanisms for spatial and feature-based attention in area V4. *Neuron* 47: 637–643, 2005. doi:10.1016/j.neuron.2005.07.020.
- Ito M, Tamura H, Fujita I, Tanaka K.** Size and position invariance of neuronal responses in monkey inferotemporal cortex. *J Neurophysiol* 73: 218–226, 1995. doi:10.1152/jn.1995.73.1.218.
- Kosai Y, El-Shamayleh Y, Fyall AM, Pasupathy A.** The role of visual area V4 in the discrimination of partially occluded shapes. *J Neurosci* 34: 8570–8584, 2014. doi:10.1523/JNEUROSCI.1375-14.2014.
- Law CT, Gold JI.** Reinforcement learning can account for associative and perceptual learning on a visual-decision task. *Nat Neurosci* 12: 655–663, 2009. doi:10.1038/nn.2304.
- Li N, Cox DD, Zoccolan D, DiCarlo JJ.** What response properties do individual neurons need to underlie position and clutter “invariant” object recognition? *J Neurophysiol* 102: 360–376, 2009. doi:10.1152/jn.90745.2008.
- Luck SJ, Chelazzi L, Hillyard SA, Desimone R.** Neural mechanisms of spatial selective attention in areas V1, V2, and V4 of macaque visual cortex. *J Neurophysiol* 77: 24–42, 1997. doi:10.1152/jn.1997.77.1.24.
- Lueschow A, Miller EK, Desimone R.** Inferior temporal mechanisms for invariant object recognition. *Cereb Cortex* 4: 523–531, 1994. doi:10.1093/cercor/4.5.523.
- Maunsell JH, Sclar G, Nealey TA, DePriest DD.** Extraretinal representations in area V4 in the macaque monkey. *Vis Neurosci* 7: 561–573, 1991. doi:10.1017/S095252380001035X.
- Maunsell JH, Treue S.** Feature-based attention in visual cortex. *Trends Neurosci* 29: 317–322, 2006. doi:10.1016/j.tins.2006.04.001.
- McAdams CJ, Maunsell JH.** Effects of attention on orientation-tuning functions of single neurons in macaque cortical area V4. *J Neurosci* 19: 431–441, 1999. doi:10.1523/JNEUROSCI.19-01-00431.1999.
- McAdams CJ, Maunsell JH.** Attention to both space and feature modulates neuronal responses in macaque area V4. *J Neurophysiol* 83: 1751–1755, 2000. doi:10.1152/jn.2000.83.3.1751.
- Miller EK, Desimone R.** Parallel neuronal mechanisms for short-term memory. *Science* 263: 520–522, 1994. doi:10.1126/science.8290960.
- Mirabella G, Bertini G, Samengo I, Kilavik BE, Frilli D, Della Libera C, Chelazzi L.** Neurons in area V4 of the macaque translate attended visual features into behaviorally relevant categories. *Neuron* 54: 303–318, 2007. doi:10.1016/j.neuron.2007.04.007.
- Moran J, Desimone R.** Selective attention gates visual processing in the extrastriate cortex. *Science* 229: 782–784, 1985. doi:10.1126/science.4023713.
- Motter BC.** Neural correlates of attentive selection for color or luminance in extrastriate area V4. *J Neurosci* 14: 2178–2189, 1994a. doi:10.1523/JNEUROSCI.14-04-02178.1994.
- Motter BC.** Neural correlates of feature selective memory and pop-out in extrastriate area V4. *J Neurosci* 14: 2190–2199, 1994b. doi:10.1523/JNEUROSCI.14-04-02190.1994.
- Noudoost B, Chang MH, Steinmetz NA, Moore T.** Top-down control of visual attention. *Curr Opin Neurobiol* 20: 183–190, 2010. doi:10.1016/j.conb.2010.02.003.
- Ogawa T, Komatsu H.** Target selection in area V4 during a multidimensional visual search task. *J Neurosci* 24: 6371–6382, 2004. doi:10.1523/JNEUROSCI.0569-04.2004.
- Op De Beeck H, Vogels R.** Spatial sensitivity of macaque inferior temporal neurons. *J Comp Neurol* 426: 505–518, 2000. doi:10.1002/1096-9861(20001030)426:4<505:AID-CNE1>3.0.CO;2-M.
- Pagan M, Rust NC.** Dynamic target match signals in perirhinal cortex can be explained by instantaneous computations that act on dynamic input from inferotemporal cortex. *J Neurosci* 34: 11067–11084, 2014a. doi:10.1523/JNEUROSCI.4040-13.2014.
- Pagan M, Rust NC.** Quantifying the signals contained in heterogeneous neural responses and determining their relationships with task performance. *J Neurophysiol* 112: 1584–1598, 2014b. doi:10.1152/jn.00260.2014.
- Pagan M, Simoncelli EP, Rust NC.** Neural quadratic discriminant analysis: nonlinear decoding with V1-like computation. *Neural Comput* 28: 2291–2319, 2016. doi:10.1162/NECO\_a\_00890.
- Pagan M, Urban LS, Wohl MP, Rust NC.** Signals in inferotemporal and perirhinal cortex suggest an untangling of visual target information. *Nat Neurosci* 16: 1132–1139, 2013. doi:10.1038/nn.3433.
- Roth N, Rust NC.** Inferotemporal cortex multiplexes behaviorally relevant target match signals and visual representations in a manner that minimizes their interference. *PLoS One* 13: e0200528, 2018a. doi:10.1371/journal.pone.0200528.
- Roth N, Rust NC.** Rethinking assumptions about how trial and nuisance variability impact neural task performance in a fast processing regime. *J Neurophysiol* 121: 115–130, 2018b. doi:10.1152/jn.00503.2018.
- Rust NC, Dicarlo JJ.** Selectivity and tolerance (“invariance”) both increase as visual information propagates from cortical area V4 to IT. *J Neurosci* 30: 12978–12995, 2010. doi:10.1523/JNEUROSCI.0179-10.2010.
- Rust NC, DiCarlo JJ.** Balanced increases in selectivity and tolerance produce constant sparseness along the ventral visual stream. *J Neurosci* 32: 10170–10182, 2012. doi:10.1523/JNEUROSCI.6125-11.2012.