



# Superior colliculus neuronal ensemble activity signals optimal rather than subjective confidence

Brian Odegaard<sup>a,1,2</sup>, Piercesare Grimaldi<sup>b,c,1</sup>, Seong Hah Cho<sup>d,1</sup>, Megan A. K. Peters<sup>a,e</sup>, Hakwan Lau<sup>a,f,g</sup>, and Michele A. Basso<sup>b,c,f,h</sup>

<sup>a</sup>Department of Psychology, University of California, Los Angeles, CA 90095; <sup>b</sup>Department of Psychiatry and Biobehavioral Sciences, University of California, Los Angeles, CA 90095; <sup>c</sup>Semel Institute for Neuroscience and Human Behavior, University of California, Los Angeles, CA 90095; <sup>d</sup>Department of Integrative Physiology, University of California, Los Angeles, CA 90095; <sup>e</sup>Department of Bioengineering, University of California, Riverside, CA, 92521; <sup>f</sup>Brain Research Institute, University of California, Los Angeles, CA 90095; <sup>g</sup>Department of Psychology, University of Hong Kong, Hong Kong; and <sup>h</sup>Department of Neurobiology, University of California, Los Angeles, CA 90095

Edited by Michael E. Goldberg, Columbia University College of Physicians, New York, NY, and approved December 29, 2017 (received for review June 29, 2017)

Recent studies suggest that neurons in sensorimotor circuits involved in perceptual decision-making also play a role in decision confidence. In these studies, confidence is often considered to be an optimal readout of the probability that a decision is correct. However, the information leading to decision accuracy and the report of confidence often covaried, leaving open the possibility that there are actually two dissociable signal types in the brain: signals that correlate with decision accuracy (optimal confidence) and signals that correlate with subjects' behavioral reports of confidence (subjective confidence). We recorded neuronal activity from a sensorimotor decision area, the superior colliculus (SC) of monkeys, while they performed two different tasks. In our first task, decision accuracy and confidence covaried, as in previous studies. In our second task, we implemented a motion discrimination task with stimuli that were matched for decision accuracy but produced different levels of confidence, as reflected by behavioral reports. We used a multivariate decoder to predict monkeys' choices from neuronal population activity. As in previous studies on perceptual decision-making mechanisms, we found that neuronal decoding performance increased as decision accuracy increased. However, when decision accuracy was matched, performance of the decoder was similar between high and low subjective confidence conditions. These results show that the SC likely signals optimal decision confidence similar to previously reported cortical mechanisms, but is unlikely to play a critical role in subjective confidence. The results also motivate future investigations to determine where in the brain signals related to subjective confidence reside.

perceptual decision-making | multineuron recording | decoding | monkey | signal detection theory

When we view the world, our experience often includes an assessment of how confident we are in our perceptual decisions. For example, when driving on a foggy morning, there are moments when we can readily identify elements in our surroundings, and other moments when we are less sure about what lies ahead. Survival in any dynamic environment depends on being able to accurately assess how reliable our perceptions and decisions are in a given instance. Here, we ask: How is this subjective sense of confidence in our perceptual decisions represented in the brain?

Work in awake macaques reveals neuronal correlates of confidence in sensorimotor circuits involved in decision-making and action generation, such as the lateral intraparietal area (LIP) (1) and the supplementary eye fields (SEFs) (2). One pioneering study of the neurophysiological underpinnings of confidence employed an “opt-out” perceptual decision-making task (1). In this task, monkeys made decisions about the primary direction of motion in random dot displays and reported those decisions by making a saccade to one of two targets located in the visual field that corresponded to the dominant dot motion directions (right or left). On some trials, an opt-out option appeared orthogonal

to the other targets and was associated with a smaller, but guaranteed, reward; choosing the opt-out option indicates less confidence in the decision (3–6). In this task, neurons recorded from area LIP discharged with the highest rates when monkeys correctly chose targets associated with motion toward the response field (RF) and discharged with the lowest rates for correct, opposite RF choices (1). When monkeys chose to opt out, LIP neurons discharged at intermediate levels; these results were interpreted in support of the idea that LIP neurons encode a signal of decision confidence.

An issue arising from this LIP study and most other previous studies of confidence is that decision accuracy and confidence covary. That is, since subjects are usually more confident when they perform better on a given task, purported neuronal correlates of confidence may signal decision accuracy rather than subjective confidence per se. To make progress, two contributions may be needed: (i) a distinction between different notions of confidence, and (ii) a paradigm that dissociates subjective confidence and accuracy. Here, we address both needs.

First, according to one influential theoretical framework (7), confidence can be defined as the probability that a perceptual decision is correct, and this probability can be formalized by using signal detection theory (SDT) (8–11). For example, in thinking about the dot motion task within this framework, one

## Significance

Previously, the neuronal correlates of perceptual confidence have been identified in neural circuits responsible for deciding what an animal sees. However, behaviorally, confidence and perceptual decision accuracy are confounded; we are usually more confident about perceptual decisions when they are accurate. To tease them apart, we introduced a task with stimulus conditions that produced similar decision accuracy but different reports of subjective confidence. We decoded decision performance from neuronal signals in nonhuman primates in a subcortical region involved in decision-making, the superior colliculus (SC), and found that SC ensemble activity tracks decision accuracy, but not subjective confidence. These results challenge current ideas about how to measure subjective confidence in experiments and inspire ways to study its neuronal mechanisms.

Author contributions: B.O., P.G., S.H.C., H.L., and M.A.B. designed research; P.G., S.H.C., and M.A.B. performed research; B.O. analyzed data; and B.O., P.G., S.H.C., M.A.K.P., H.L., and M.A.B. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Published under the PNAS license.

<sup>1</sup>B.O., P.G., and S.H.C. contributed equally to this work.

<sup>2</sup>To whom correspondence should be addressed. Email: odegaard.brian@gmail.com.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1711628115/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1711628115/-DCSupplemental).

can assume there is a binary variable “ $s$ ” governing the true state of the world;  $s = -1$  or “+1” depending on whether the motion was primarily to the left (“L”) or right (“R”). On each trial, the subject makes a noisy measurement “ $m$ ” from the distribution  $P(m|s)$  —i.e., the distribution over the sensory measurement given the L or R stimulus (reflecting noise in the stimulus itself as well as in the sensory system). Confidence can then be defined as the distance between the measurement  $m$  and the decision boundary used to judge whether the dominant dot motion was to the left or right, which will monotonically reflect the probability or frequency that the decision made is correct (12). In other words, in this framework, confidence and accuracy are always correlated: As the distance from the measurement to the decision boundary increases, so will decision accuracy.

However, as has been noted (7), observers’ confidence judgments do not always correlate with task performance (13–17). Thus, a distinction must be made between a type of confidence that correlates with task accuracy, and a type of confidence that does not. We introduce the terms “optimal confidence” and “subjective confidence” to refer to these two types of confidence, respectively. Optimal confidence refers to the above SDT-based definition, in which an ideal observer’s confidence always monotonically tracks the accuracy of decisions. It is optimal in the sense that a subject making opt-out decisions based on optimal confidence will be able to maximize reward. Because optimal confidence always correlates with task accuracy, a neuronal correlate of optimal confidence can be found in population-level activity that effectively distinguishes between conditions yielding different levels of accuracy in a task.

On the other hand, subjective confidence is an actual (i.e., not necessarily ideal) observer’s internal belief about a perceptual decision, which is potentially prone to error. In SDT terms, this form of confidence depends not strictly on the distance from the sensory measurement to the decision boundary, but on the distance from the measurement to specific confidence criteria (see Fig. S1, which provides a formulation of the two types of confidence in terms of SDT), which may be arbitrarily placed based on some kind of heuristic. Subjective confidence can be indexed by behavioral reports which may or may not track decision accuracy perfectly. Because behavioral reports of confidence tend to track optimal confidence at least to some extent, isolating a true neuronal signal of subjective confidence can be difficult. However, this can be facilitated by paradigms that match decision accuracy across conditions, yet yield different behavioral reports about confidence. In this situation, a neuronal correlate of subjective confidence can be found in neuronal activity that tracks the behavioral reports of varying degrees of confidence, amid constant decision accuracy.

Recent work indicates that it is possible to dissociate the capacity to perform perceptual tasks from confidence reports by chemically inactivating the pulvinar (18) or orbitofrontal cortex (19), or psychophysically in humans (20–22). Therefore, we reasoned that we could develop visual stimuli that would lead to similar decision accuracy (and therefore, similar levels of optimal confidence), but yield different levels of confidence as measured by behavioral reports on individual trials (i.e., subjective confidence). Creation of these stimuli would allow us to investigate the neuronal mechanisms of confidence by determining whether activity in a given area signals optimal confidence, subjective confidence, or both (Fig. S1).

Therefore, in this study, monkeys performed two sets of experiments. The first was an opt-out task in which decision accuracy covaried with confidence similar to that performed previously for recordings in LIP (1), allowing us to search for neural correlates of optimal confidence. In the second experiment, building on innovative psychophysical work done in humans (20–23), we introduced a version of the dot-motion direction discrimination task in which we dissociated reports of

confidence from decision accuracy on individual trials. Using this task, we were able to successfully match decision accuracy [as defined by the SDT measure  $d'$  (8–10)], but produce different levels of confidence (defined as the probability of selecting the opt-out target when it was available), so that we could investigate the neuronal correlates of subjective confidence.

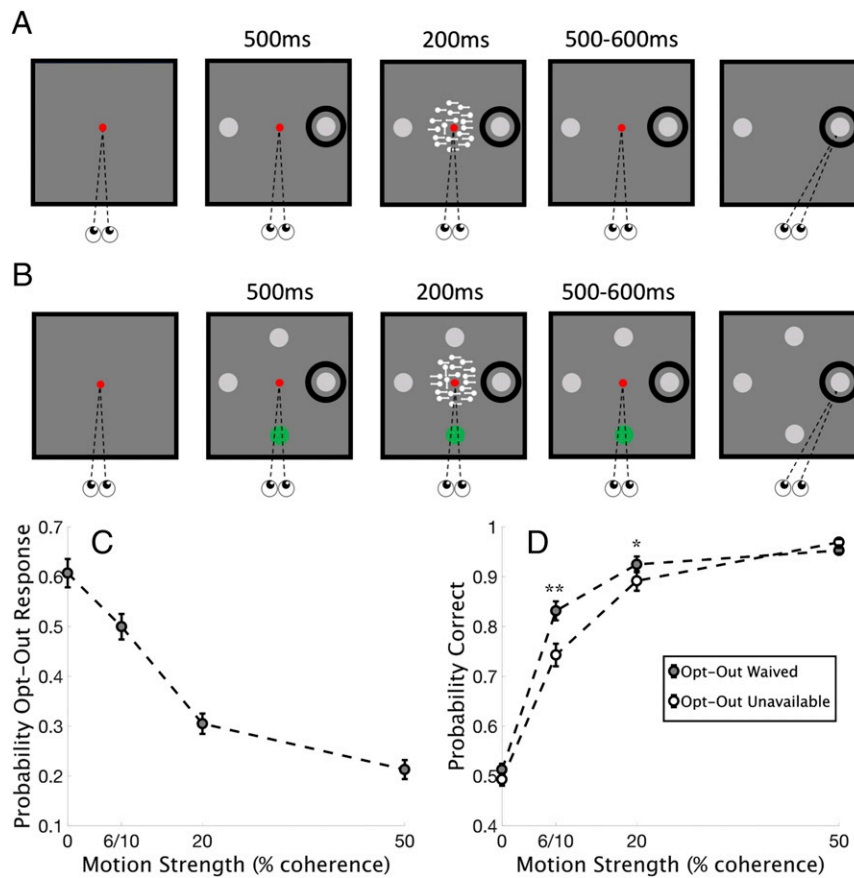
As monkeys performed these tasks, we recorded from multiple neurons simultaneously in the superior colliculus (colliculus), a subcortical structure that receives input from LIP and SEF and is involved in decision-making (24–30). We combined these behavioral paradigms and multineuron recordings with a machine learning approach (31) to decode population-level activity from hundreds of neurons recorded from the colliculus. We found that in the first task, a population decoder distinguished between high- and low-confidence trials in much the same way as LIP (1), providing strong evidence that the colliculus contributes to decision-making and optimal confidence in a manner similar to LIP. However, in our task in which visual stimuli were matched for sensitivity ( $d'$ ) but resulted in different reports of confidence, population-level activity in the colliculus failed to distinguish between conditions with different degrees of subjective confidence. Together, these findings support the hypothesis that the colliculus signals optimal confidence in dot-motion discrimination tasks, rather than subjective confidence. These results also reveal important considerations for the interpretation of existing data on decision-making confidence in other brain regions, too.

## Results

We used a multivariate decoding approach to assess population-level representations of perceptual decisions and confidence in the superior colliculus using random dot-motion discrimination tasks. We had two aims. Our first aim was to determine whether activity measured in the colliculus was similar to that observed previously in area LIP during performance of a confidence task (1). Our second aim was to arbitrate between two competing hypotheses: that neuronal activity in the colliculus primarily signals optimal confidence, as signals about confidence may correlate with decision accuracy, or, alternatively, that activity in the colliculus signals subjective confidence, as neuronal signals may differentiate between conditions where  $d'$  is matched, but confidence reports vary. We focus here on results obtained from a population decoding method.

We recorded neuronal activity in the colliculus using V-probe laminar electrodes containing 16 recording contacts (*Methods*). We measured both single-neuron and multineuron activity while monkeys performed a dot-motion discrimination task (Fig. 1 *A* and *B*). Each trial began when the animal established fixation on a central dot. Then, either two or four choice targets appeared for 500 ms. After this delay, the dot motion stimulus appeared at the center of the screen for 200 ms. When the motion stimulus disappeared, a delay period, selected randomly from between 500 and 600 ms, ensued. The fixation dot then disappeared, and monkeys indicated their motion direction decision by making a saccade to one of the choice targets, and they received a reward (sip of juice) for correct decisions. Importantly, on some trials there was an opt-out option. Choosing this target bypassed the motion discrimination question and led to a guaranteed, but smaller, reward compared with that received for correct decisions.

On trials when the opt-out option was available (Fig. 1*B*), we also included a fourth choice option which was opposite in location to the opt-out location to control for possible lateral interactions (see *Methods* for details). The fourth option never led to reward and was rarely chosen (~6.3% of all trials in stimulus-matched sessions). For each session, at least one of the choice targets appeared in the RF of at least one neuron recorded from the 16 contacts (black circle, Fig. 1*A*). The two trial types with (Fig. 1*B*) and without (Fig. 1*A*) the opt-out option available were randomly interleaved; because the properties of the random



**Fig. 1.** Stimulus-matched assessment of decision confidence in monkeys. (A and B) The behavioral task showing a trial in which the opt-out option was unavailable (A) and available (B). The trial types shown in A and B were randomly interleaved in each of the 19 stimulus-matched sessions. The red dot shows the fixation point, the gray dots show the possible choice targets, and the green dot shows the opt-out option. The black circle shows the RF. (C) The probability of choosing the opt-out option on trials when it was available (shown in B) is plotted as a function of motion coherence. Circles show means across sessions, and bars show SEM across sessions. Note that monkeys chose the opt-out option more often when motion coherence was low, indicating they were less confident about the motion direction decision. The number of trials making up this dataset is 14,642, as it includes all trials where the opt-out was offered. (D) The probability of correct choices is plotted against motion coherence for two monkeys using the same set of data as in C, but now plotting trials where an explicit decision about the motion direction was made (i.e., including trials with the opt-out unavailable, and excluding aborted trials, trials where the opt-out was selected, and trials where the lateral inhibition target was selected). The number of trials making up this dataset is 13,346. Circles show means across sessions, and bars show SEM. Gray filled circles show data when the opt-out option was available but waived (trials shown in B), and open circles show data when the opt-out option was unavailable (trials shown in A). Decision accuracy is higher for intermediate motion strengths when the opt-out target was available but waived, presumably reflecting higher confidence (t tests, Bonferroni corrected). \* $P < 0.01$ ; \*\* $P < 0.001$ .

dot-motion stimulus were identical between these trial types, we call these “stimulus-matched” sessions.

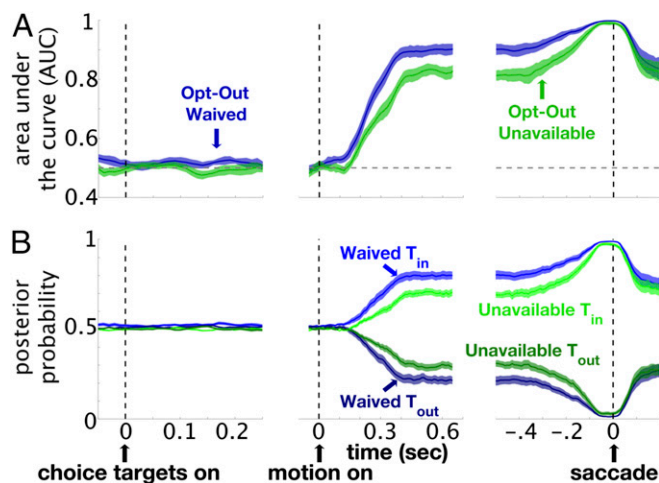
We reasoned that choices made with the opt-out unavailable occurred with a mix of high and low confidence, as monkeys were forced to choose one of the two targets. In trials with the opt-out option available, monkeys could report their level of confidence: Trials in which monkeys chose the opt-out target indicated low confidence, whereas trials in which monkeys waived the opt-out option and chose one of the targets corresponding to a direction of motion instead indicated high confidence (1, 3–5). Fig. 1 C and D show the behavior measured in trials with and without the opt-out option available. The probability of selecting the opt-out option, when available, decreased as a function of motion coherence, consistent with higher confidence on higher motion coherence trials (all t tests between conditions  $P < 0.05$ , Bonferroni corrected; Fig. 1C). Comparing trials in which the opt-out option was available and unavailable showed that at intermediate motion strengths, monkeys had a higher probability of being correct when the opt-out option was available and waived compared with when it was unavailable, indicating higher confidence (Fig. 1D).

To determine if neuronal ensemble activity in the colliculus correlates with optimal confidence, we used multivariate classifiers to evaluate how population-level activity emerged over time as monkeys made decisions in the opt-out available and unavailable trials. Previous work showed that LIP discharge rates differed when a correct choice was reported by making a saccade toward the target in the RF (target-in, or “ $T_{in}$ ”) or away from the RF (target-out, or “ $T_{out}$ ”) (1). Here, we used a similar approach by evaluating the classifier’s ability to predict correct  $T_{in}$  and  $T_{out}$  choices with the opt-out choice available (but waived) and unavailable.

First, we assessed the area under the ROC curve (AUC) for the classifier (Methods) to evaluate the degree to which population-level activity in the colliculus may be informative for trial-by-trial predictions of particular behavioral responses (e.g.,  $T_{in}$  vs.  $T_{out}$  choices). This provided us with a measure of how effectively neuronal population activity discriminated between specific perceptual decisions. We focused on a comparison between the opt-out waived and the opt-out unavailable conditions, because, in both conditions, the motoric behavior is similar (i.e., the monkeys made saccades to choose one of the options to

reflect a perceptual decision rather than the opt-out option), and yet optimal confidence is expected to be higher in the opt-out waived condition; higher confidence and expected accuracy presumably cause the monkeys to waive the opt-out option when it is available.

Fig. 2 shows that neuronal activity in the colliculus signals correct  $T_{in}$  and  $T_{out}$  choices and that decoder performance (based on average performance on test sets using fivefold cross-validation) is higher when the opt-out option is available but waived. Here, we show the combined results across sessions for two monkeys, but we note that the decoding performance for both monkeys in this task was quite similar (Figs. S2 and S3). Fig. 2A shows that neuronal activity more accurately discriminates correct  $T_{in}$  choices vs. correct  $T_{out}$  choices when the opt-out choice was available but waived, compared with trials where opt-out was unavailable [ $t$  tests for all time windows  $>230$  ms after motion onset in Fig. 2A, Center,  $t(18) > 2.8$ ,  $P < 0.05$ ]. Following previous research (1), we interpreted this as evidence that the information contained in the neuronal population activity signals optimal confidence, as the population activity correlates with the differences in decision accuracy across these two conditions. To control for multiple comparisons throughout the entire motion onset period, we used the false discovery rate



**Fig. 2.** Decoding perceptual decisions made with different levels of confidence for the same motion stimuli (stimulus-matched). We trained and tested a decoding model using a 100-ms sliding window (step size = 10 ms) beginning 50 ms before the choice targets appeared through 200 ms after the choice report, to predict whether a given correct trial involved a choice in the RF ( $T_{in}$ ) or outside of the RF ( $T_{out}$ ). A total of 354 collicular neurons were used in this analysis, but the decoder was run independently by using fivefold cross-validation on data from each session (which included 9–26 simultaneously recorded neurons; *Methods*). *Left* is aligned to the onset of the choice targets, indicated by the dashed vertical line and upward arrow. *Center* is aligned to the onset of the motion stimulus, and *Right* is aligned to the onset of the saccade. Each data point represents classification performance of the midpoint of a given 100-ms time window (from 50 ms before to 50 ms after); smoothed data using a five-point moving average are represented. (A) Mean (thin solid lines) and SEM (shaded areas) classifier performance across sessions shown as the AUC plotted against time for opt-out waived and opt-out unavailable conditions. The ability of the classifier to predict a correct  $T_{in}$  or  $T_{out}$  choice was better on trials in which the opt-out option was available but waived (blue) and monkeys were more confident, compared with when the opt-out option was unavailable (green) and monkeys had a mix of higher and lower confidence in their decisions. (B) Similar to A, but plotting the average posterior probability over time. The y axis is the posterior probability of predicting that a given trial contains a correct  $T_{in}$  choice. This analysis is similar to the “decision variable” used in a previous study (31) and provides an estimate of the strength of the classifier’s predictions.

(FDR) method (32) to evaluate significance at each time point. With a FDR of 0.01, while four time windows between 100 and 230 ms were marginally significant, all time windows  $>230$  ms after motion onset were highly significant. We also performed further analyses to verify that decoding performance was not just driven by a few single neurons containing strong decision-related activity (as have been identified previously in the colliculus) and that population-level analyses added information over and above what single units can indicate (Fig. S4).

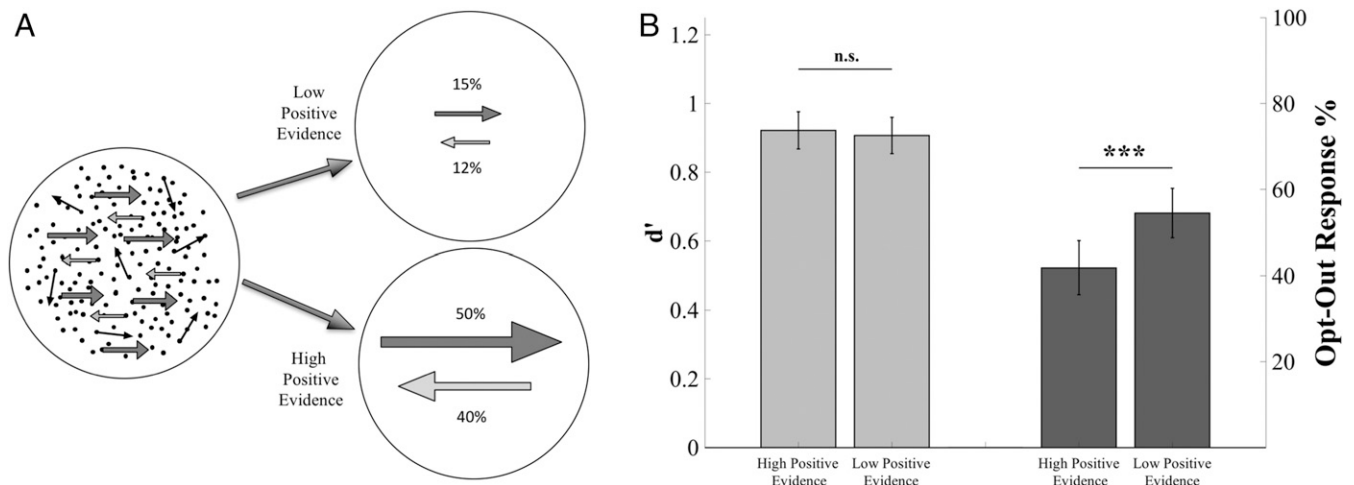
We also computed the posterior probabilities for trial-by-trial predictions made by the classifier. Roughly, following previous studies, these probabilities are interpreted as reflecting the strength of the internal decision evidence available to the animal (29–31). We sorted the classifier’s predictions by trial type over time (Fig. 2B). Similar to what was shown by using the AUC metric, differences between the posterior probabilities between  $T_{in}$  and  $T_{out}$  correct choices were significantly greater for the opt-out waived trials starting  $\sim 230$  ms after motion onset [all time windows  $>230$  ms after motion onset in Fig. 2B, Center,  $t(18) > 2.8$ ,  $P < 0.05$ ].

While the results from decoding  $T_{in}$  vs.  $T_{out}$  correct trials as a function of opt-out availability provide evidence to show that colliculus activity reflects perceptual decisions and their accuracy, does it reflect confidence behavior in any way for it to warrant the label optimal confidence? As shown in Fig. S5, the absolute magnitude of the strength of this signal decoded from neuronal populations in the colliculus was indeed correlated with actual opt-out rates. However, the relatively small magnitude of this effect suggests that the full story about confidence behavior may be more complicated.

The results described above provide evidence that the neuronal activity in the colliculus contains information about decision-making and decision confidence in much the same way as reported for area LIP (1). However, as noted, the task design used for both the colliculus and the LIP experiments leaves open the possible interpretation that the activity signals decision accuracy (and optimal confidence) rather than subjective confidence, since monkeys also performed better on the opt-out waived trials than on the opt-out unavailable trials. Therefore, we created a version of the dot-motion discrimination task in which decision accuracy was matched while confidence varied by manipulating the ratio of “positive evidence” (the amount of motion evidence toward the correct choice) to “negative evidence” (the amount of motion evidence toward the incorrect choice). Previous work showed that, while decision accuracy depends on the ratio of positive to negative evidence, subjective confidence depends on the overall magnitude of positive evidence (20–23). Thus, we presented monkeys with trials containing different ratios of positive and negative evidence to match decision accuracy (defined as perceptual sensitivity, or  $d'$ ; *Methods*) across two conditions (Fig. 3A) while attaining different levels of subjective confidence, as measured by their reports of confidence by choosing to opt out or not.

Fig. 3B shows that this manipulation yielded statistically similar levels of decision sensitivity as measured by  $d'$  (sign test,  $z = 0.83$ ,  $P = 0.40$ ), but different degrees of confidence, as indicated by the percentage of trials in which monkeys chose to opt out (sign test,  $z = -4.59$ ,  $P < 10^{-5}$ ). In this new behavioral task, trials with and without the opt-out option were randomly interleaved, allowing us to compute  $d'$  from trials without the opt-out (demonstrating that performance is adequately matched with these stimuli), while evaluating possible differences in subjective confidence from the proportion of trials in which the opt-out was selected when it was available. Data from individual sessions is shown in Fig. S6.

With this new task, we evaluated whether the same  $T_{in}$  vs.  $T_{out}$  decision-related activity differed between the two “sensitivity-matched” condition types (high positive evidence vs. low positive



**Fig. 3.** A task for dissociating sensitivity and confidence (sensitivity-matched). (A) By manipulating the ratio of positive evidence (dot motion toward the correct decision; dark gray rightward arrows) to negative evidence (motion incompatible with the correct decision; light gray leftward arrows), it is possible to match sensitivity across two conditions, as measured by  $d'$ , but achieve different levels of confidence, as indexed by the proportion of trials the monkeys chose to opt out. Shown here is a representative example of two conditions that could achieve this result; please note that random dot motion is also included in these conditions, and the exact ratio of positive to negative evidence varied slightly in each session, but the overall number of dots remained constant (*Methods*). The sequence of events for this paradigm was identical to the stimulus-matched paradigm described in Fig. 1, but we refer to this task as sensitivity-matched. (B)  $d'$  and the percentage of opt-out choices (when the opt-out was available and was selected) are plotted for high and low positive evidence conditions. Across 23 behavioral sessions from two monkeys, the results show statistically indistinguishable sensitivity (light gray bars) between high and low positive evidence conditions (10,197 trials in the  $d'$  analysis), but different percentages of opt-out choices (16,471 trials in the opt-out response analysis; dark gray bars). \*\*\* $P < 10^{-5}$ ; n.s., not significant. Bars show averages across sessions, and error bars are SEM.

evidence); that is, whether it passes the test to be considered a neural correlate of subjective confidence. Fig. 4 shows the decoding results for the sensitivity-matched task. The neurons recorded in each of the sensitivity-matched sessions used in this decoding analysis were different from the neurons used in decoding the stimulus-matched sessions. While trials with and without the opt-out were randomly interleaved in the sensitivity-matched sessions, we focused our initial decoding analyses solely on the opt-out unavailable trials, excluding opt-out waived trials (Fig. 4). This was to ensure that, should the decoder identify a difference between the two conditions, this difference would not be driven by a sheer difference in internal perceptual response, as the difference in decision criteria for opt-out behavior between the high and low positive evidence conditions means that opt-out waived trials will be more frequent in the high positive evidence condition, and as such, the average internal response strength will not be matched between the conditions (Fig. S1).

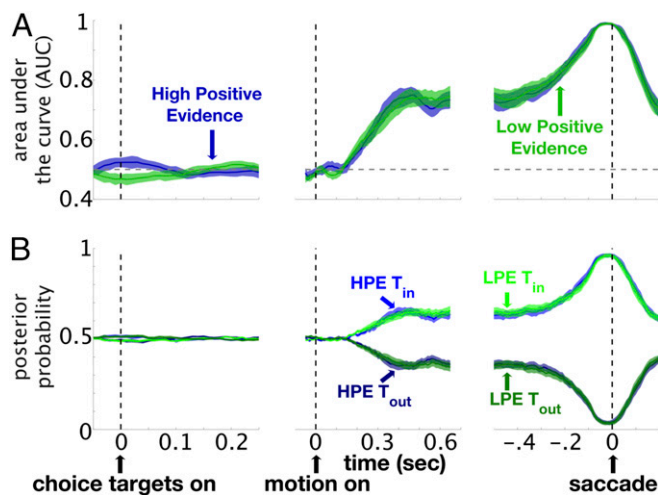
Following motion onset, the decoder performance was statistically indistinguishable for both the high positive evidence (higher confidence) and low positive evidence (lower confidence) conditions for nearly all time points [59 of 65,  $t$  tests,  $t(22) < 2.1$ ,  $P > 0.05$ ]. Importantly, by using a FDR of 0.01 to correct for multiple comparisons, none of the time windows reached significance. We observed a similar pattern when comparing the posterior probabilities for the high- and low-confidence trials from the sensitivity-matched task (Fig. 4B). The temporal evolution of the strength of the predictions produced by the classifier was statistically indistinguishable for almost all time points [59 of 65,  $t$  tests,  $t(22) < 2.1$ ,  $P > 0.05$ ], and by using the FDR method to account for false positives, no significant differences were found for any of the time windows following motion onset.

To further assess whether there were distinct signals for subjective confidence in the activity of colliculus neurons, we performed a cross-generalization analysis. Although the magnitude of decodability between  $T_{in}$  and  $T_{out}$  choices was similar between the high and low positive evidence conditions, it was possible that the difference in confidence was reflected by different neurons

contributing to this same level of decodability. In that case, some neurons would have still meaningfully reflected the different levels of subjective confidence. If that were true, the performance of a classifier that was trained on trials from the high positive evidence condition and tested on trials from the low positive evidence condition (or vice versa) should have been reduced compared with the performance of classifiers trained and tested within the same condition. That is, if we observed that information was substantially lost through the cross-generalization process, it would have provided evidence for distinct neuronal signals for high and low confidence. Fig. 5 shows the performance of a classifier trained on trials from the high positive evidence condition and tested on trials from the low positive evidence condition as measured by the AUC and posterior probability. This classifier showed similar performance to classifiers trained and tested on trials from a single condition [one-way ANOVA, 64 of 65 time windows following motion onset,  $F(66) < 1$ ,  $P > 0.05$ ]; the ability to decode was roughly equivalent across the two conditions, and a comparison between training on low positive evidence and testing on high positive evidence yielded statistically indistinguishable results.

Taken together, with differences in population neuronal activity in the colliculus during a stimulus-matched confidence task (Fig. 2) but similarity across conditions in a sensitivity-matched confidence task (Figs. 4 and 5), these results provide evidence that colliculus activity likely reflects optimal confidence and not subjective confidence per se.

Although our decoding results suggest that, at the population level, collicular activity discriminates different perceptual decisions at similar levels for the high- and low-confidence conditions when sensitivity is matched, there may still be individual neurons that reflect subjective confidence in other ways. To investigate this possibility, we computed a normalized “discriminability index” (*Methods*) to determine how effectively individual neurons could discriminate between  $T_{in}$  and  $T_{out}$  choices as a function of confidence in the sensitivity-matched task, which included conditions that varied in terms of positive evidence level (high vs. low), as well as conditions that varied in terms of opt-out availability,



**Fig. 4.** Decoding perceptual decisions made with different levels of confidence and the same level of sensitivity. We trained and tested a decoding model using a 100-ms sliding window (step size = 10 ms) beginning 50 ms before the choice targets appeared through 200 ms after the choice report, to predict whether a given correct trial included a saccade toward the choice target in the RF ( $T_{in}$ ) or outside of the RF ( $T_{out}$ ). The data are from the sensitivity-matched task shown in Fig. 3 and contain 6,910 trials from 421 neurons from two monkeys (23 total sessions). The decoder was run separately on neurons from each recording session. Each data point represents the classification performance of the midpoint of a given 100-ms time window (from 50 ms before to 50 ms after); smoothed data using a five-point moving average are represented. (A) The mean classifier performance as AUC plotted against time in seconds. (B) The mean posterior probability for  $T_{in}$  and  $T_{out}$  choices plotted against time in seconds; the y axis reflects the posterior probability that a given trial contains a correct  $T_{in}$  choice. The blue lines and shaded areas show the mean and SEM from the high positive evidence (HPE) condition (high confidence), and the green lines and shaded areas show the mean and SEM from the low positive evidence (LPE) condition (low confidence).

like in the original task (available vs. unavailable). For a neuron to signal subjective confidence, it should show greater discriminability of  $T_{in}$  vs.  $T_{out}$  choices, not only on trials in which the opt-out choice was available but waived (compared with when it was unavailable), but also on trials from the high positive evidence condition; put simply, neurons that care about optimal confidence as defined by opt-out availability should also care about subjective confidence based on evidence ratios if activity signals subjective confidence at all. Including both of these condition types in the sensitivity-matched task allowed us to assess this.

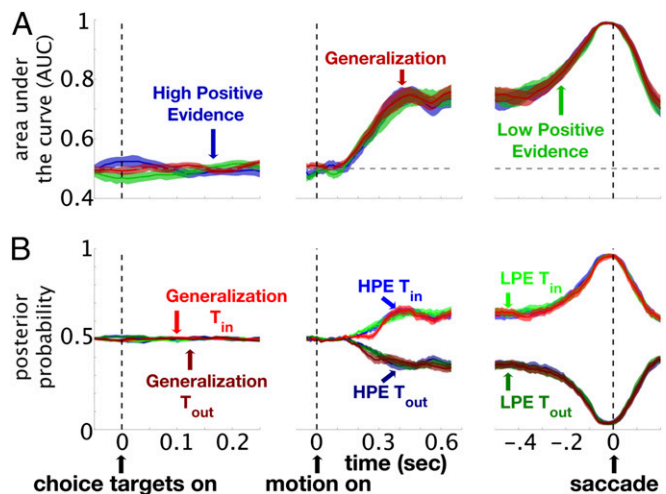
The discriminability index ranges from  $-1$  to  $+1$ . In the opt-out available and unavailable conditions (the stimulus-matched conditions), neurons that maximally discriminate  $T_{in}$  and  $T_{out}$  when the opt-out is available but waived have a value of  $+1$ ; neurons that maximally discriminate  $T_{in}$  and  $T_{out}$  when the opt-out is unavailable have a value of  $-1$ . In the sensitivity-matched conditions, neurons that maximally discriminate  $T_{in}$  and  $T_{out}$  in the high positive evidence condition have a value of  $+1$ , whereas neurons that maximally discriminate  $T_{in}$  and  $T_{out}$  in the low positive evidence condition have a value of  $-1$ . Neurons with values near  $1$  on both discriminability indices are neurons that signal confidence.

In line with the decoding results obtained for the original stimulus-matched task, which only included two conditions that varied in terms of opt-out availability, we found a significant number of neurons with higher discriminability indices when monkeys waived the opt-out choice vs. when the opt-out choice was unavailable (sign rank test,  $z = 13.87$ ,  $P < 10^{-42}$ ). For the sensitivity-matched conditions, however, discriminability indices regarding capacities as a function of evidence level were distributed symmetrically around  $0$  (sign rank test,  $z = -0.78$ ,  $P =$

$0.44$ ), indicating a lack of a significant number of neurons that discriminated choices more effectively in the high positive evidence trials. Fig. S7A shows histograms of the discriminability indices for neurons recorded in sensitivity-matched sessions.

Because there were some collicular neurons that gave the appearance of signaling confidence based on their discriminability index (Fig. S7A), we also analyzed neurons with discriminability indices  $>0$  on both measures, to determine if this ability to discriminate  $T_{in}$  and  $T_{out}$  choices was stable across different trials. We divided each session's dataset into "odd" and "even" trials and computed two discriminability indices; one for odd trials and one for even trials for each neuron. If a neuron signals subjective confidence, it should do so for all trials and should do so consistently. If, however, the signal is inconsistent across trials, it is unlikely to provide a signal that is usable by the brain; the odd-even index assesses this. Fig. S8A shows the possible confidence neurons—those falling in the upper right quadrant when the discriminability index for the stimulus- and sensitivity-matched conditions were computed from odd trials only. By computing these same two discriminability indices for the same neurons from even trials (Fig. S8B and C), it became clear that while neurons were stable in their increased capacity to discriminate  $T_{in}$  and  $T_{out}$  as a function of opt-out availability (sign rank test,  $z = 9.82$ ,  $P < 10^{-23}$ ), they were not stable in their capacity to discriminate these trial types in sensitivity-matched conditions (sign rank test,  $z = 0.29$ ,  $P = 0.76$ ). Thus, consistent with the population-level analysis, collicular activity appears better explained by optimal confidence than subjective confidence.

We also assessed whether subjective confidence might be encoded in the colliculus in other ways, beyond what is reflected by the ability to distinguish between correct  $T_{in}$  and  $T_{out}$  choices. To that end, we analyzed whether we could effectively decode the probability of opting out in the high and low positive evidence conditions in our sensitivity-matched task. Specifically, in decoding trials when the monkeys chose one of the rewarded targets (both correct and incorrect trials) compared with trials when the



**Fig. 5.** Generalization analysis reveals little evidence for subjective confidence signals in the superior colliculus. (A) The mean classifier performance as AUC plotted against time in seconds. Data are the same as in Fig. 4, with the addition of results from a linear classifier trained on trials from the high positive evidence condition (HPE; high confidence) and tested on trials from the low positive evidence condition (LPE; low confidence), shown in red. Lines show averages, and shaded areas show SEM. (B) The mean posterior probability for  $T_{in}$  and  $T_{out}$  choices plotted against time in seconds.  $T_{in}$  indicates correct trials in which the monkeys made a saccade toward the choice targets in the RF, and  $T_{out}$  indicates correct trials in which the monkeys made saccades toward choice targets outside of the RF.

monkeys made the decision to opt out, we could evaluate whether this comparison could effectively discriminate between the high and low positive evidence conditions, which produced different levels of subjective confidence. Since the purpose here is to actually predict opt-out rate, in this analysis, we only included trials where the opt-out option was available. Importantly, to make sure the neuronal activity reflected the decision to opt out rather than just the motoric signal for saccade toward the opt-out option, the neurons included here never had the opt-out option presented in their RFs. As shown in Fig. S9, this analysis showed that our decoder performed similarly in both the high and low positive evidence conditions, as both the AUC [64 of 65 time windows following motion onset,  $t(22) < 2.1$ ,  $P > 0.05$ ] and posterior probability [all time windows following motion onset,  $t(22) < 2.1$ ,  $P > 0.05$ ] metrics were quite similar across time.

Finally, to further investigate other ways that subjective confidence may be coded in the colliculus, we also trained the decoder to directly classify any differences between the high and low positive evidence conditions (Fig. S10). Despite the various differences in sheer physical stimulus properties and levels of reward expectation, we found that colliculus activity did not strongly distinguish between them.

## Discussion

We combined psychophysics with multineuron recordings and population-decoding methods to determine whether activity in the superior colliculus of monkeys signals decision confidence. Using a task similar to that used in conjunction with recordings in area LIP (1), we identified population-level activity in the colliculus that distinguished between different choices and different levels of confidence in much the same way as LIP. That is, when decision accuracy and decision confidence covary, the colliculus signals confidence in a manner similar to LIP. This is consistent with an interpretation that the colliculus, like LIP, signals more than just eye movements and plays an important role in perceptual decision-making (24, 25, 28, 33–35). However, when comparing collicular activity using a task that dissociated optimal from subjective confidence, we found that both population and single neuron activity was indistinguishable between high- vs. low-confidence conditions. Further analyses also failed to find strong evidence in favor of the claim that the colliculus signals subjective confidence per se. Thus, we conclude that the role of the colliculus in decision confidence likely primarily concerns optimal confidence.

These findings raise interesting questions regarding previous interpretations of studies. By using an opt-out task (1), neuronal correlates of confidence have been found in LIP (1) and in the pulvinar (18). Even though our study also used the opt-out design, in a previous investigation (1), monkeys were informed about the opt-out option only after the motion stimulus appeared and presumably after they made their decision. In our paradigm, the choice options appeared before the onset of the motion stimulus to avoid visual contamination of the neuronal activity during the stimulus period. Despite this difference, the ability of collicular neurons to distinguish  $T_{in}$  and  $T_{out}$  choices with different levels of confidence was surprisingly similar to the activity patterns seen in LIP. Similar findings were also obtained in the SEF by using a wagering task in which monkeys reported their confidence by making “bets” after each perceptual decision (2). To the extent that the colliculus may signal primarily optimal confidence rather than subjective confidence, this open question may apply to those other regions, too. Further research is needed to answer this question.

Despite similarities to previous studies, the neuronal ensemble activity in the colliculus did not pass our sensitivity-matched tests for subjective confidence. However, there could be other neuronal signatures that differ between our two confidence conditions (such as those involving temporal patterns) that our analyses

were unable to identify. However, to the extent that confidence is reflected by firing rate differences between  $T_{in}$  and  $T_{out}$ , as has been assessed by previous studies (1), such activity patterns across the population of neurons assessed seem highly similar between the high and low positive evidence conditions, as the decoders generalized remarkably well between them (Fig. 5). To exercise further caution, we also conducted analysis of individual neurons (Figs. S7 and S8). We found that to the extent that some neurons might have shown any difference in discriminability between these sensitivity-matched conditions, such differences were unlikely to be stable properties of the neurons.

Based on the human literature (36–38) as well as animal studies (19, 39), one intriguing possibility is that subjective confidence may reside in prefrontal cortex, even under sensitivity-matched conditions. Although one previous study (2) recorded from the lateral prefrontal cortex as well as the frontal eye fields and did not find neurons reflecting optimal confidence in these areas as defined above, it remains to be tested whether such neuronal signatures for subjective confidence may emerge when confidence is dissociated from sensitivity, or when an opt-out task rather than a wagering task is adopted. In humans, under sensitivity-matched conditions, hemodynamic activity differs between conditions involving different levels of reported confidence (36–38). Applying magnetic stimulation or chemical inactivation to the prefrontal cortex alters confidence reports while sensitivity remains unchanged (19, 37, 39). In another study in monkeys, muscimol injection to the pulvinar impaired confidence reports, as assessed by an opt-out task, while leaving decision accuracy unchanged (18). Such effects may involve the interactions between the known projections from the dorsal central pulvinar to the prefrontal cortex (40–42). The work in prefrontal cortex and pulvinar, like our work reported here, also argues strongly for a distinction between optimal confidence based on perceptual decisions and subjective confidence that is dissociable from perceptual decisions. We propose that combining our behavioral task with multineuron recordings in the prefrontal cortex and pulvinar may uncover representations of subjective confidence independent of optimal confidence.

Finally, one issue with the stimulus-matched task (and sensitivity-matched task) is that the condition with mixed high- and low-confidence trials contains only two visual stimuli, whereas the condition with the opt-out choice available contains four stimuli. It is well-documented that neuronal activity in the colliculus is modulated by choice target uncertainty; specifically, as the number of possible targets increases, activity in the colliculus decreases (22). Our results cannot be explained by this because as Fig. 2 Left shows, in the opt-out unavailable condition, there were fewer possible stimuli, yet that activity was indistinguishable from that seen in the opt-out waived condition in which there were more stimuli on the screen. This is opposite to what would be expected for an interpretation based on lateral inhibition or uncertainty.

In summary, our findings highlight the important roles played by the superior colliculus in decision-making, beyond its well-known role in eye movements (43), and, perhaps more importantly, they raise critical questions about the interpretation of previous findings and open up exciting possibilities for future studies of subjective confidence.

## Methods

**Surgical Procedures.** Two male rhesus monkeys (9–13 kg) were prepared for electrophysiological recordings and measurements of eye movements. Anesthesia was induced with an intramuscular injection of ketamine (5.0 mg/kg) and midazolam (0.2 mg/kg), and atropine (0.04 mg/kg) was provided to limit salivation. Monkeys were then intubated and maintained at a general anesthetic plane with isoflurane. One hour before the procedure, animals received buprenorphine (0.01 mg/kg) and the antibiotic Excede (20 mg/kg; 7 d slow release) and then meloxicam (0.3 mg/kg) at the conclusion of the

procedure, and meloxicam (0.2 mg/kg) and buprenorphine (0.01 mg/kg) for 3 d postsurgically as analgesia. Monkeys were implanted with MRI-compatible headposts, and one (monkey H) was implanted with eye loops (44, 45) to measure eye position. In the other monkey (monkey P), eye position was measured with an iView camera (Sensomotoric Instruments). Both monkeys received MRI-compatible recording chambers placed over the superior colliculus (anterior–posterior + 3, medial–lateral 0) and angled posteriorly at 38°. Precise placement of the post and chambers was performed by using MRI-guided surgical software (BrainSight; Rogue Research). All surgical procedures were performed under general anesthesia by using aseptic procedures. All experimental protocols were approved by the University of California, Los Angeles Chancellor's Animal Research Committee and complied with and generally exceeded standards set by the Public Health Service policy on the humane care and use of laboratory animals.

**Eye Movement Recording Procedures.** We used a QNX-based real-time experimental data-acquisition system and Windows-based visual stimulus generation system ("Rex" and "Vex"), developed and distributed by the Laboratory of Sensorimotor Research, National Eye Institute (Bethesda) (46) to create the behavioral paradigm, display the visual stimulus, and acquire two channels of eye position data. Voltage signals proportional to horizontal and vertical components of eye position were filtered (eight pole Bessel –3 dB, 180 Hz), digitized at 16-bit resolution, and sampled at 1 kHz (PCI-6036E; National Instruments). The camera-acquired eye position signals were filtered digitally by using a built-in bilateral filter. We used an automated procedure to define saccadic eye movements using eye velocity (20° per s) and acceleration criteria (5,000° per s<sup>2</sup>), respectively. The adequacy of the algorithm was verified and adjusted as necessary on a trial-by-trial basis by the experimenter.

**Electrophysiological Procedures.** We recorded multineuron activity from the intermediate layers of the superior colliculus using a platinum/iridium V-Probe coated with polyimide (Plexon) with an impedance of 275 ( $\pm$ 50) k $\Omega$ . The electrode was aimed at the colliculus perpendicular to its surface by using guide tubes positioned with a grid system (47) and advanced by using an electronic microdrive system controlled by a graphical user interface (Nan Instruments). Action potential waveforms were bandpass-filtered (250 Hz–5 kHz; four pole Butterworth), and amplified, by using the BlackRock NSP hardware system controlled by the Cerebus software suite (BlackRock Microsystems). The voltage data were sampled and digitalized at 30 kHz with 16-bit resolution and saved to disk for offline sorting. For isolating neurons online, we used time and amplitude windowing criteria (Cerebus; BlackRock Inc.). Waveforms satisfying these criteria generated transistor–transistor logic pulses indicating the time of occurrence of an action potential and were sampled and digitized at 1 kHz with 16-bit resolution and saved to disk.

Action potential waveforms were sorted offline by using the Plexon Offline Sorter (Plexon, Inc.) and classified into single neurons ( $n = 115$ ) and multineuron ( $n = 660$ ) activity. At the start of each recording session, we aimed to identify a recording site with at least one buildup neuron, in light of their established role in higher-level phenomena such as attention, selection, and decision-making (reviewed in ref. 48). We classified buildup neurons as those neurons having a significantly higher discharge rate during the stimulus period (200–600 ms after motion onset) compared with baseline (200–0 ms before the stimulus appears). While the recording procedure first focused on identifying buildup neurons before continuing with the experiment, all neurons that were recorded in a session (both buildup and non-buildup) were used in the decoding analysis for a given session.

RFs of collicular neurons were mapped online to provide an estimate of the center of the RF to place at least one choice target. We determined the general characteristics of the neuronal activity and an estimate of the center of the preferred RF by requiring monkeys to make saccades to different locations in the visual field. We made a qualitative assessment online about the preferred location on the basis of maximal discharge determined audibly. We confirmed the center of the RF by plotting the discharge as a heat map across visual space. Only neurons with RF eccentricities between 7° and 20° were studied to ensure no overlap of the RF with the centrally placed moving dot stimulus.

The neurons we recorded from were different in each recording session; the neurons from the 19 stimulus-matched sessions which were used were different from the neurons from the 23 sensitivity-matched sessions.

**Behavioral Task.** We used the same behavioral task in both the stimulus- and sensitivity-matched paradigms. Each trial in both paradigms began when monkeys acquired a centrally located spot and remained fixated for 500 ms. Then, the choice targets appeared. One choice target appeared in the center

of the RF of at least one of the recorded neurons ( $T_{in}$ ), and the other choice target appeared in the opposite hemifield ( $T_{out}$ ). These positions were randomized on each trial. For both the stimulus- and sensitivity-matched paradigms, half of the trials had only two choice targets (i.e., "opt-out unavailable") and half had an opt-out choice target available. These trial types were randomized in each session. All targets, including the opt-out, were isoluminant. The location of the opt-out choice was orthogonal to the two motion choice targets (90°), and on these trials, we also presented a fourth dot, irrelevant to the task, 180° opposite to the opt-out target location. This was included to control for possible lateral interactions (24, 49). That is, to ensure that any differences between the opt-out waived and the opt-out unavailable trials were not driven by introducing an additional response target in an orthogonal location, we introduced a fourth dot to make the stimulus symmetrical, so that each possible target in the opt-out available condition was surrounded by a isoluminant targets at the same distance and relative locations.

After the choice targets appeared and monkeys maintained fixation on the central spot for ~500 ms, the dot motion stimulus appeared centrally for 200 ms. Monkeys maintained fixation for another 500- to 600-ms interval (the exact time was randomly selected between those two times from a uniform distribution) and then were cued to report their decision by removal of the fixation point. If the correct choice occurred, monkeys received a juice reward (0.2 mL). If the incorrect choice occurred, monkeys received no reward and a timeout of 2,000 ms. On trials in which monkeys selected the opt-out choice, they received a smaller but guaranteed reward (80% of the correct choice reward amount).

**Stimuli.** For both tasks, the motion stimulus appeared on a CRT display operating at 60 Hz. The motion speed was 5° per s, and the same dots were maintained on the screen for the duration of the stimulus (200 ms). Some dots moved coherently in a single direction (coherence percentages described below), while the other dots moved with randomly selected trajectories. The radius of the motion stimulus was 3°, and the size of dots in the display were 0.05°. The dot density in both tasks was 50 dots per degree squared. Each dot moved in the same direction for the duration of a given trial. For all motion stimuli, the total number of dots appearing on the display was kept constant to maintain isoluminance.

For the stimulus-matched paradigm, four motion coherence levels were tested for each monkey. For monkey P, we tested performance with 20%, 10%, 6%, and 0% coherence. For monkey H, we tested performance with 50%, 10%, 6%, and 0% coherence. Different coherence levels were used to yield approximately equivalent performance levels across the two monkeys. Dots moving in random directions were also included, and the total number of dots in all displays was the same.

For the sensitivity-matched paradigm, the dot coherence ratios characterized by positive evidence (PE; motion favoring the correct choice) and negative evidence (NE; motion favoring the incorrect choice) were customized for each monkey in each session to yield similar  $d'$  values across two conditions on trials where the opt-out was unavailable, but different amounts of selecting the opt-out across those two conditions on trials when it was available. For the eight  $d'$ -matched sessions for monkey P, one  $d'$ -matched session included a 50%PE/30%NE coherence ratio for high positive evidence and 20%PE/17%NE coherence ratio for low positive evidence; two  $d'$ -matched sessions included 50%PE/30%NE coherence ratio for high positive evidence and 35%PE/21%NE coherence ratio for low positive evidence; four  $d'$ -matched sessions included a 50%PE/30%NE coherence ratio for high positive evidence and 20%PE/12%NE coherence ratio for low positive evidence; and one  $d'$ -matched session included a 50%PE/34%NE coherence ratio for high positive evidence and 20%PE/9%NE coherence ratio for low positive evidence.

For the 15  $d'$ -matched sessions for monkey H, one  $d'$ -matched session included a 50%PE/30%NE coherence ratio for high positive evidence and 35%PE/21%NE coherence ratio for low positive evidence; two  $d'$ -matched sessions included 50%PE/33%NE coherence ratio for high positive evidence and 20%PE/5%NE coherence ratio for low positive evidence; one  $d'$ -matched session included a 50%PE/37%NE coherence ratio for high positive evidence and 20%PE/7%NE coherence ratio for low positive evidence; one  $d'$ -matched session included a 50%PE/37%NE coherence ratio for high positive evidence and 23%PE/7%NE coherence ratio for low positive evidence; nine  $d'$ -matched sessions included a 50%PE/37%NE coherence ratio for high positive evidence and 25%PE/7%NE coherence ratio for low positive evidence; and one  $d'$ -matched session included 35%PE/30%NE coherence ratio for high positive evidence and 20%PE/12%NE coherence ratio for low positive evidence.



**Behavioral Data Analysis.** We used SDT to quantify the decision sensitivity of the monkeys in our behavioral task. In this task, monkeys were presented with a dot motion stimulus and had to make a discrimination judgment as to whether the primary motion direction was to the right or left.  $d'$  is a measure of an observer's capacity to perform a sensory task: A  $d'$  score of 0 indicates a complete inability to discriminate left and right motion directions in this task, while  $d'$  scores  $>0$  quantify an observer's sensitivity to make this type of discrimination. As noted by Wickens (9),  $d'$  in discrimination tasks can be computed by adding the Z-transformed correct-response probabilities for both stimulus types (p. 116). Thus,  $d'$  was calculated as:

$$d' = Z(p_A) + Z(p_B), \quad [1]$$

where in this task,  $p_A$  refers to the probability of a correct judgment for trials where the primary motion direction was toward the left, and  $p_B$  refers to the probability of a correct judgment where the primary motion direction was to the right. This equation yields the exact same  $d'$  values as the standard  $d'$  equation for detection tasks [ $Z(\text{Hit Rate}) - Z(\text{False Alarm Rate})$ ] but provides a more accurate characterization for discrimination judgments, as "false alarms" are not possible in this type of task, since a primary motion direction is present on every trial.

The sensitivity-matched sessions included two different trial types: On some trials, the opt-out was unavailable, and monkeys' only choice was between the two response options. These trials allowed us to determine that our two evidence conditions were matched. On other trials, the opt-out was available but could be waived, and these trials allowed us to infer different levels of confidence across these two conditions. We only computed  $d'$  from trials where the opt-out choice was unavailable and focused the decoding analyses on these trials alone. This was done to ensure that, should our subsequent decoding analyses identify a difference across conditions, this difference would not be driven solely by differences in the perceptual criterion used for each condition. The two trial types were randomly interleaved in sensitivity-matched sessions, and data from these different trial types is shown in Fig. 3 and Fig. S7. We also note that in our sensitivity-matched sessions, we only analyzed days in which the  $d'$  scores between our high and low positive evidence trials were within 0.7 of one another (see Fig. S6 for individual session results).

**Decoding Analysis.** To investigate how population activity in the superior colliculus may be related to optimal and subjective confidence, we applied a decoding model to analyze time-varying neuronal activity and performed our decoding analyses separately on the data from each recording session. In each session, between 9 and 26 neurons were recorded from our V-Probe recording device, and all units used in decoding for a given session were recorded simultaneously.

We first quantified neuronal discharge rates across all electrodes with a sliding window analysis, computing the sum of action potentials occurring within 100-ms time windows (step size = 10 ms). Next, we applied a logistic regression model using the fitlinear function in MATLAB (Mathworks, 2016). The general idea behind this linear classification function is that on any given trial, the overall classification score  $f(x)$  can be predicted from the neuronal activity at a given time point using the following equation:

$$f(x) = \beta x + b. \quad [2]$$

In this equation,  $x$  is the vector of the summed spike counts for each neuron in a given time window,  $\beta$  is a vector representing the linear coefficient estimates for each neuron, and  $b$  is the scalar bias, reflecting the intercept estimate. However, since our decoding analyses focused on categorical outcomes instead of continuous measures, we applied the "logistic" learner from fitlinear, which implements the "logit" score transformation function to the raw classification scores to yield the probability of a given class (e.g.,  $X$ ), via the following equation:

$$p(X) = \frac{1}{(1 + e^{-\beta x + b})}, \quad [3]$$

with the following loss function for classification, where  $y \in \{\pm 1\}$ :

$$L(y, f(x)) = \log(1 + \exp(-yf(x))). \quad [4]$$

This implementation uses the following ridge regularization penalty to avoid overfitting in our procedure, with a lambda value of  $(1/\text{number of neurons})$  in a given session:

$$\frac{\lambda}{2} \sum_{j=1}^p \beta_j^2. \quad [5]$$

We also implemented a uniform prior in the fitlinear function over the two classes, which specified that the two classes being predicted were equally likely on each trial. Finally, we estimated the posterior probabilities for class predictions on each trial using the predict function from the Statistics and Machine Learning Toolbox in MATLAB.

As has been noted previously, logistic regression classifiers find the best hyperplane that separates the population response patterns associated with the two classes that are being predicted (29). Therefore, the essential idea behind the aforementioned analysis is that, for every trial, the decoder will produce not only a final prediction of, for example, whether the monkey chose the target located in the RF or the target out of the RF, but also a measure of the strength of the prediction (via the posterior probability metric), which corresponds to the prediction's distance from the hyperplane. We further explain the utility of these metrics in Results.

The model was implemented by using fivefold cross-validation at each time point, with 80% of the data as the "training set" for fitting the  $\beta$  coefficients and 20% of the data as the "test set." In all figures and results, we report the average performance across all five test sets. Two metrics enabled us to assess performance of the model: First, we used AUC as our method to assess decoder accuracy. Second, we sorted each model's predictions by trial type and evaluated the posterior probability of particular class predictions over time. This allowed us to assess the strength of the classifiers' prediction for each trial type across time, within a range of 0–1. Thus, the results we report are based on average AUC and posterior probabilities across the five test sets at each time point.

Three time periods were of particular interest for our decoding procedure. First, the time period around the onset of the targets, to determine whether the prestimulus activity held any predictive power for the monkeys' upcoming decisions. Second, the time period following onset of the motion stimulus, since this is the time when monkeys are forming their decisions, and, as such, the activity could signal both decision sensitivity and/or subjective confidence. Finally, the time period around the saccade is also informative, as this time window reflects the ceiling for classification performance based on the recorded neuronal activity.

We note that recent work has demonstrated the utility of decoding approaches compared with single-neuron analyses (31, 50), and, indeed, our own analysis revealed a stronger capacity to classify correct perceptual decisions by using population-level analyses compared with single neurons (Fig. S4). While we do think that single-neuron analysis of our data can also be informative, we think a machine-learning approach is particularly advantageous, as decision confidence may be encoded by complex patterns of neuronal activity distributed across many neurons within a brain region, as has been shown in other recent work (38).

**Discriminability Index.** To assess each neuron's discriminative capacity for  $T_{in}$  and  $T_{out}$  choices, we computed a "discriminability index." This metric produces a normalized value between  $-1$  and  $1$  specifying both the strength and direction of a neuron's predictive power for a given two-class discrimination problem. For example, in our initial analysis (Fig. 2), we classified whether a given correct choice would be toward the RF ( $T_{in}$ ) or away from the RF ( $T_{out}$ ). We hypothesized that the ability to discriminate would change as a function of opt-out availability. Thus, we computed the discriminability index for each neuron for the  $T_{in}$  vs.  $T_{out}$  classification procedures in the following manner:

$$\text{Stimulus-Matched Discriminability Index} = \frac{\frac{\text{Opt-Out Waived}}{|T_{in} - T_{out}|} - \frac{\text{Opt-Out Unavailable}}{|T_{in} - T_{out}|}}{\frac{\text{Opt-Out Waived}}{|T_{in} - T_{out}|} + \frac{\text{Opt-Out Unavailable}}{|T_{in} - T_{out}|}}. \quad [6]$$

Negative values mean that the neuronal activity is more discriminable for  $T_{in}$  compared with  $T_{out}$  when the opt-out is unavailable compared with when it is waived; positive values indicate that the neuronal activity is more discriminable between  $T_{in}$  compared with  $T_{out}$  when the opt-out is waived compared with when it is unavailable. With the sensitivity-matched data, we computed the same discriminability index for trials from the high positive evidence condition and from the low positive evidence condition using the following equation:

## Sensitivity-Matched Discriminability Index

$$= \frac{\begin{array}{c} \text{High Positive Evidence} \\ |T_{in} - T_{out}| \end{array} - \begin{array}{c} \text{Low Positive Evidence} \\ |T_{in} - T_{out}| \end{array}}{\begin{array}{c} |T_{in} - T_{out}| \\ \text{High Positive Evidence} \end{array} + \begin{array}{c} |T_{in} - T_{out}| \\ \text{Low Positive Evidence} \end{array}} \quad [7]$$

Positive values indicate the neuronal activity is more discriminable for  $T_{in}$  compared with  $T_{out}$  for trials in the high positive evidence condition compared with trials from the low positive evidence condition, and negative values indicate that the neuronal activity is more discriminable for  $T_{in}$  compared with  $T_{out}$  for trials in the low positive evidence condition compared with trials from the high positive evidence condition. We declared

- Kiani R, Shadlen MN (2009) Representation of confidence associated with a decision by neurons in the parietal cortex. *Science* 324:759–764.
- Middlebrooks PG, Sommer MA (2012) Neuronal correlates of metacognition in primate frontal cortex. *Neuron* 75:517–530.
- Kornell N, Son LK, Terrace HS (2007) Transfer of metacognitive skills and hint seeking in monkeys. *Psychol Sci* 18:64–71.
- Smith JD, et al. (1995) The uncertain response in the bottlenosed dolphin (*Tursiops truncatus*). *J Exp Psychol Gen* 124:391–408.
- Shields WE, Smith JD, Washburn DA (1997) Uncertain responses by humans and rhesus monkeys (*Macaca mulatta*) in a psychophysical same-different task. *J Exp Psychol Gen* 126:147–164.
- Grimaldi P, Lau H, Basso MA (2015) There are things that we know that we know, and there are things that we do not know that we do not know: Confidence in decision-making. *Neurosci Biobehav Rev* 55:88–97.
- Pouget A, Drugowitsch J, Kepecs A (2016) Confidence and certainty: Distinct probabilistic quantities for different goals. *Nat Neurosci* 19:366–374.
- Green DM, Swets JA (1966) *Signal Detection Theory and Psychophysics* (John Wiley and Sons, New York).
- Wickens TD (2002) *Elementary Signal Detection Theory* (Oxford Univ Press, Oxford).
- Macmillan NA, Creelman CD (2005) *Detection Theory: A User's Guide* (Lawrence Erlbaum Associates Publishers, Mahwah, NJ), 2nd Ed.
- Kepecs A, Mainen ZF (2012) A computational framework for the study of confidence in humans and animals. *Philos Trans R Soc Lond B Biol Sci* 367:1322–1337.
- Kepecs A, Uchida N, Zariwala HA, Mainen ZF (2008) Neural correlates, computation and behavioural impact of decision confidence. *Nature* 455:227–231.
- Camerer CF, Lovo D (1999) Overconfidence and excess entry: An experimental approach. *Am Econ Rev* 89:306–318.
- Baranski JV, Petrusic WM (1994) The calibration and resolution of confidence in perceptual judgments. *Percept Psychophys* 55:412–428.
- Adams JK (1957) A confidence scale defined in terms of expected percentages. *Am J Psychol* 70:432–436.
- Moore DA, Healy PJ (2008) The trouble with overconfidence. *Psychol Rev* 115: 502–517.
- Harvey N (1997) Confidence in judgment. *Trends Cogn Sci* 1:78–82.
- Komura Y, Nikkuni A, Hirashima N, Uetake T, Miyamoto A (2013) Responses of pulvinar neurons reflect a subject's confidence in visual categorization. *Nat Neurosci* 16:749–755.
- Lak A, et al. (2014) Orbitofrontal cortex is required for optimal waiting based on decision confidence. *Neuron* 84:190–201.
- Koizumi A, Maniscalco B, Lau H (2015) Does perceptual confidence facilitate cognitive control? *Atten Percept Psychophys* 77:1295–1306.
- Zylberberg A, Fetsch CR, Shadlen MN (2016) The influence of evidence volatility on choice, reaction time and confidence in a perceptual decision. *Elife* 5:e17688.
- Zylberberg A, Bartfeld P, Sigman M (2012) The construction of confidence in a perceptual decision. *Front Integr Neurosci* 6:79.
- Samaha J, Barrett JJ, Sheldon AD, LaRocque JJ, Postle BR (2016) Dissociating perceptual confidence from discrimination accuracy reveals no influence of metacognitive awareness on working memory. *Front Psychol* 7:851.
- Basso MA, Wurtz RH (1998) Modulation of neuronal activity in superior colliculus by changes in target probability. *J Neurosci* 18:7519–7534.
- Horwitz GD, Newsome WT (1999) Separate signals for target selection and movement specification in the superior colliculus. *Science* 284:1158–1161.
- Krauzlis R, Dill N (2002) Neural correlates of target choice for pursuit and saccades in the primate superior colliculus. *Neuron* 35:355–363.
- McPeck RM, Keller EL (2004) Deficits in saccade target selection after inactivation of superior colliculus. *Nat Neurosci* 7:757–763.
- Kim B, Basso MA (2008) Saccade target selection in the superior colliculus: A signal detection theory approach. *J Neurosci* 28:2991–3007.
- Fries W (1984) Cortical projections to the superior colliculus in the macaque monkey: A retrograde study using horseradish peroxidase. *J Comp Neurol* 230:55–76.
- Paré M, Wurtz RH (2001) Progression in neuronal processing for saccadic eye movements from parietal cortex area lip to superior colliculus. *J Neurophysiol* 85: 2545–2562.
- Kiani R, Cueva CJ, Reppas JB, Newsome WT (2014) Dynamics of neural population responses in prefrontal cortex indicate changes of mind on single trials. *Curr Biol* 24: 1542–1547.
- Benjamini Y, Krieger AM, Yekutieli D (2006) Adaptive linear step-up procedures that control the false discovery rate. *Biometrika* 93:491–507.
- Shadlen MN, Newsome WT (2001) Neural basis of a perceptual decision in the parietal cortex (area LIP) of the rhesus monkey. *J Neurophysiol* 86:1916–1936.
- Gold JI, Shadlen MN (2007) The neural basis of decision making. *Annu Rev Neurosci* 30:535–574.
- Horwitz GD, Newsome WT (2001) Target selection for saccadic eye movements: Prelude activity in the superior colliculus during a direction-discrimination task. *J Neurophysiol* 86:2543–2558.
- Lau H, Passingham RE (2006) Relative blindsight in normal observers and the neural correlate of visual consciousness. *Proc Natl Acad Sci USA* 103:18763–18768.
- Rounis E, Maniscalco B, Rothwell JC, Passingham RE, Lau H (2010) Theta-burst transcranial magnetic stimulation to the prefrontal cortex impairs metacognitive visual awareness. *Cogn Neurosci* 1:165–175.
- Cortese A, Amano K, Koizumi A, Kawato M, Lau H (2016) Multivoxel neurofeedback selectively modulates confidence without changing perceptual performance. *Nat Commun* 7:13669.
- Miyamoto K, et al. (2017) Causal neural network of metamemory for retrospection in primates. *Science* 355:188–193.
- Romanski LM, Giguere M, Bates JF, Goldman-Rakic PS (1997) Topographic organization of medial pulvinar connections with the prefrontal cortex in the rhesus monkey. *J Comp Neurol* 379:313–332.
- Shipp S (2003) The functional logic of cortico-pulvinar connections. *Philos Trans R Soc Lond B Biol Sci* 358:1605–1624.
- Pessoa L, Adolphs R (2010) Emotion processing and the amygdala: From a 'low road' to 'many roads' of evaluating biological significance. *Nat Rev Neurosci* 11:773–783.
- Sparks DL, Hartwich-Young R (1989) The deep layers of the superior colliculus. *Rev Oculomot Res* 3:213–255.
- Judge SJ, Richmond BJ, Chu FC (1980) Implantation of magnetic search coils for measurement of eye position: An improved method. *Vision Res* 20:535–538.
- Fuchs AF, Robinson DA (1966) A method for measuring horizontal and vertical eye movement chronically in the monkey. *J Appl Physiol* 21:1068–1070.
- Hays AV, Jr, Richmond BJ, Optican LM (1982) Unix-based multiple-process system, for real-time data acquisition and control. Available at <https://www.osti.gov/scitech/biblio/5213621>. Accessed March 24, 2017.
- Crist CF, Yamasaki DS, Komatsu H, Wurtz RH (1988) A grid system and a microsyringe for single cell recording. *J Neurosci Methods* 26:117–122.
- Basso MA, May PJ (2017) Circuits for action and cognition: A view from the superior colliculus. *Annu Rev Vis Sci* 3:197–226.
- Rizzolatti G, Camarda R, Grupp LA, Pisa M (1973) Inhibition of visual responses of single units in the cat superior colliculus by the introduction of a second visual stimulus. *Brain Res* 61:390–394.
- Leavitt ML, Pieper F, Sachs AJ, Martinez-Trujillo JC (2017) Correlated variability modifies working memory fidelity in primate prefrontal neuronal ensembles. *Proc Natl Acad Sci USA* 114:E2494–E2503.