

# Task interference as a neuronal basis for the cost of cognitive flexibility

## Authors:

Cheng Xue\*<sup>1</sup>, Sol K. Markman\*<sup>1,2</sup>, Ruoyi Chen<sup>3</sup>, Lily E. Kramer<sup>1</sup>, Marlene R. Cohen<sup>1</sup>

<sup>1</sup>Department of Neurobiology, University of Chicago, IL, USA

<sup>2</sup>Department of Brain & Cognitive Sciences, Massachusetts Institute of Technology, MA, USA

<sup>3</sup>Department of Biological Sciences, Carnegie Mellon University, PA, USA

## Disclosures: None

Correspondence can be addressed to Cheng Xue ([cxue@uchicago.edu](mailto:cxue@uchicago.edu)) and Marlene R. Cohen (lead contact, [marlenecohen@uchicago.edu](mailto:marlenecohen@uchicago.edu))

## Abstract:

Humans and animals have an impressive ability to juggle multiple tasks in a constantly changing environment. This flexibility, however, leads to decreased performance under uncertain task conditions. Here, we combined monkey electrophysiology, human psychophysics, and artificial neural network modeling to investigate the neuronal mechanisms of this performance cost. We developed a behavioural paradigm to measure and influence participants' decision-making and perception in two distinct perceptual tasks. Our data revealed that both humans and monkeys, unlike an artificial neural network trained for the same tasks, make less accurate perceptual decisions when the task is uncertain. We generated a mechanistic hypothesis by comparing this neural network trained to produce correct choices with another network trained to replicate the participants' choices. We hypothesized, and confirmed with further behavioural, physiological, and causal experiments, that the cost of task flexibility comes from what we term task interference. Under uncertain conditions, interference between different tasks causes errors because it results in a stronger representation of irrelevant task features and entangled neuronal representations of different features. Our results suggest a tantalizing, general hypothesis: that cognitive capacity limitations, both in health and disease, stem from interference between neural representations of different stimuli, tasks, or memories.

## Introduction:

In the face of the inherent uncertainty and unexpected changes in natural environments, humans and animals have evolved remarkable cognitive flexibility, enabling them to seamlessly switch between tasks as the situation demands. However, this cognitive flexibility comes with a cost: when we are uncertain which task should be performed (typically during and shortly before or after task switching), we take longer to perform the task we choose and do so with lower accuracy (1-4).

1 Numerous non-exclusive hypotheses have been proposed based on behavioural evidence to  
2 explain the cost of task flexibility. These include limits on processing capacity (4), excessive  
3 attention to irrelevant stimuli or features (5-8) and proactive interference between competing  
4 tasks (1, 2, 9-12). These phenomenological explanations of behaviour do not directly address the  
5 neural mechanisms that underly task switching costs and related limits on other forms of  
6 cognition.

7 Here, we propose a novel neural mechanism that integrates and elucidates the roles and  
8 connections among various hypotheses addressing cognitive flexibility cost. We used a  
9 multidisciplinary approach including behavioural experiments in humans and rhesus monkeys,  
10 multi-neuron, multi-area recordings and causal manipulations in monkeys, and artificial neural  
11 network modeling to generate and test neural hypotheses. The core of our approach is our two-  
12 feature visual discrimination task, in which participants must discriminate one visual feature and  
13 ignore the other, inferring the implicitly changing task rule that determines which feature is  
14 relevant, based on their history of stimuli, choices and rewards.

15 Our findings suggest that the cost of task flexibility arises from stronger task interference under  
16 uncertain conditions. Such interference is a joint product of retaining task-irrelevant visual  
17 information and the entanglement of the neural representations of unrelated features. This  
18 phenomenon likely underlies the cost of cognitive flexibility across species, systems, and health  
19 states, offering promising avenues for advancing our understanding and addressing cognitive  
20 limitations in both basic science and translational research.

21

## 22 **Results:**

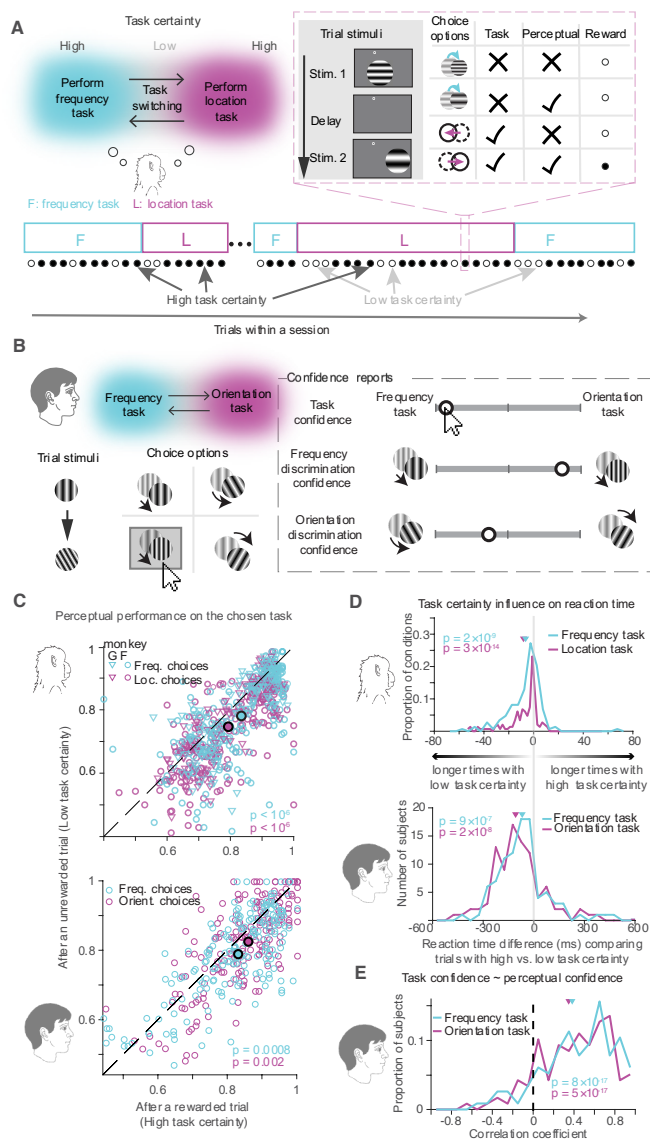
### 23 **Task switching can have behavioural costs.**

24 We designed a two-feature discrimination task that allows us to measure and manipulate  
25 perception and task certainty, pushing participants to flexible (task uncertain) and engaged (task  
26 certain) cognitive states(3) (Figure 1A). We trained 220 online human participants and two rhesus  
27 monkeys (*Macaca mulatta*; both male, 12 and 9 kg) to discriminate changes in two visual features  
28 of subsequently presented Gabor stimuli (spatial frequency and orientation for human  
29 participants; spatial frequency and location for monkey participants). Both features changed on  
30 each trial near perceptual threshold (average perceptual performance is 79% for monkeys, and  
31 82% for humans), but only one was relevant. Participants indicated both their chosen task and  
32 perceptual discrimination on every trial, and only correct discriminations of the relevant feature  
33 were rewarded (Figure 1A-B). The relevant feature switched without a cue with a constant low  
34 probability on each trial (2.5% for monkeys and 10% for humans).

35 The task structure meant that the participants' certainty about which task they would be  
36 rewarded for performing fluctuated throughout the session based on their history of stimuli,  
37 choices, and rewards. A rewarded trial suggested a high likelihood of the same task on the next  
38 trial (97.5% for monkeys, 90% for humans), thus leaving participants with higher task certainty.  
39 An unrewarded trial, however, led to task uncertainty because participants received no explicit  
40 feedback about the error type (i.e., whether the task had switched, they had made a perceptual

1 error, or both). For simplicity, we will refer to trials as having high or low task certainty based on  
 2 whether they were immediately preceded by a rewarded trial or an unrewarded trial. This  
 3 dichotomy proves to be an effective separation between high and low task certainty, supported  
 4 by task-related neuronal representations in the parietal cortex, and distinguished behavioural  
 5 patterns not solely attributable to recent rewards or trial position in a certain block (3). In the  
 6 following results, we will compare the trial condition under low task certainty (following an  
 7 unrewarded trial) vs. high task certainty (following a rewarded trial) to evaluate the cost for  
 8 flexibility.

9



10

11 **FIG 1: Tasks and behavioural cost of flexibility.** (A) Monkeys discriminated either the spatial  
 12 frequency (cyan) or location (magenta) of two sequentially presented Gabor stimuli while  
 13 ignoring changes in the irrelevant feature. The relevant feature is not explicitly cued. There  
 14 was a fixed, low probability of a task switch after each trial, leading to task blocks of variable

1 length (cyan and magenta rectangles). Participants had to infer the relevant feature from  
2 their history of choices and feedback. Participants were rewarded only if they made the  
3 correct task choice and perceptual judgment (they indicated both on each trial by choosing  
4 one of four possible targets that represent the choices schematized in the first column of the  
5 table in A). Therefore, task certainty was generally high following rewarded trials and low  
6 following unrewarded trials. (B) Human participants performed the same task as the monkeys,  
7 alternating between orientation discrimination and frequency discrimination. On a small  
8 proportion of trials, some human participants were asked to report their confidence about  
9 their task choice and perceptual discrimination of both features on a continuous scale. (C)  
10 Behavioural cost of flexibility in perceptual performance for monkeys (top) and humans  
11 (bottom). Each open symbol represents perceptual performance for a set of trials with  
12 identical stimuli. Under low task certainty (following an unrewarded trial), both species show  
13 worse perceptual performance under low task certainty (following an unrewarded trial;  
14 ordinates) than under high task certainty (following a rewarded trial; abscissae;  $p=0.002$  for  
15 human orientation choices,  $p=0.0008$  for human frequency choices;  $p<10^{-6}$  for monkey  
16 location and frequency choices. Filled symbols denote mean performance. (D) Behavioural  
17 cost of flexibility in reaction time. The upper (monkey) and lower (human) histograms show  
18 distributions of the differences between reaction times under high and low task certainty for  
19 matched stimuli. The histograms are shifted to the left, indicating that participants reported  
20 their perceptual discriminations more slowly under low task certainty than under high task  
21 certainty.  $p<10^6$  for both species and both tasks. (E) Behavioural cost of flexibility in  
22 confidence. Cyan and magenta curves show the distribution of correlation coefficients  
23 between task confidence and perceptual confidence on the chosen task for human  
24 participants.  $p<10^6$  for both tasks.

25

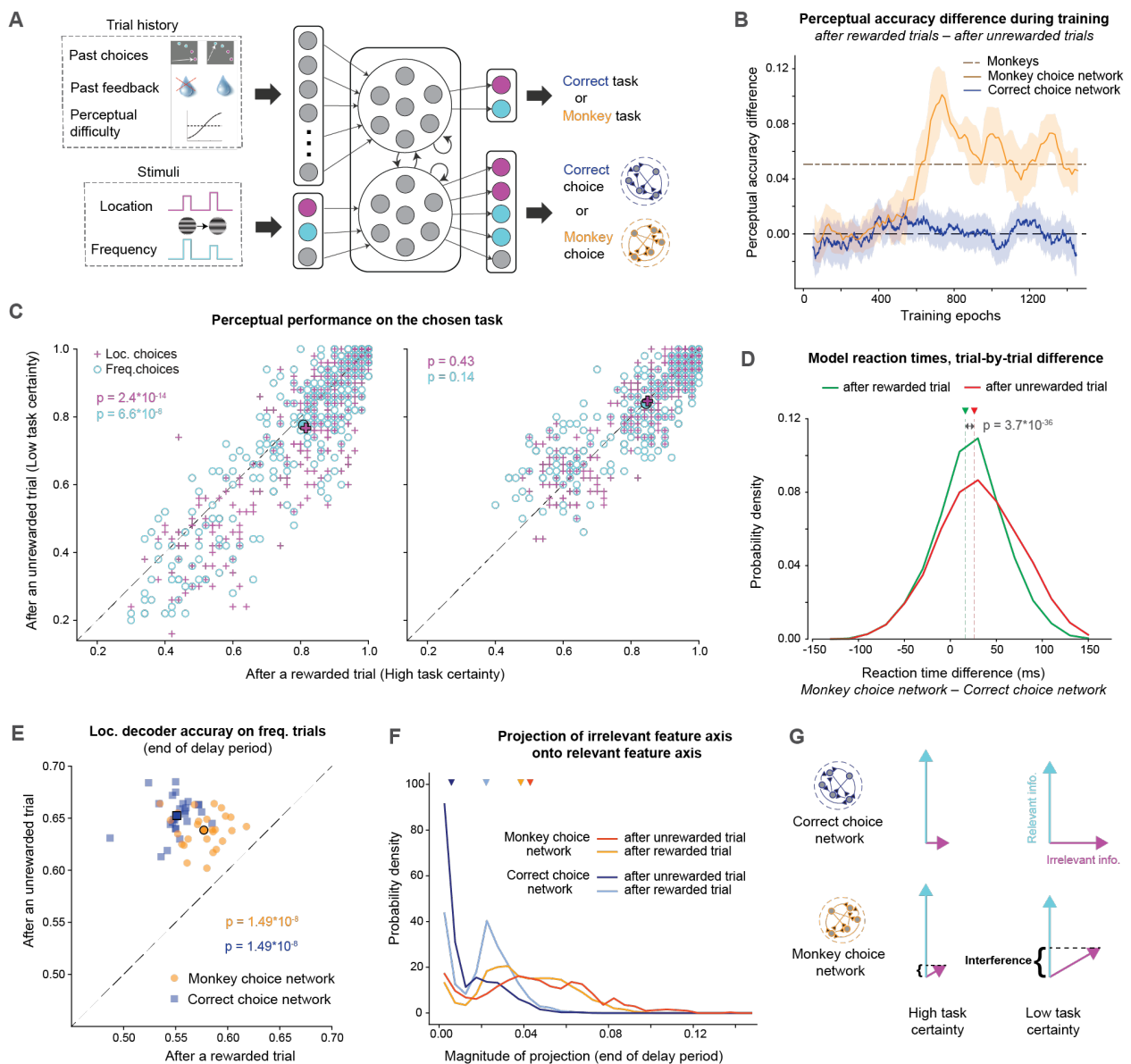
26 Both humans and monkeys demonstrate a cost of flexibility. Our participants perform the  
27 perceptual task better (Figure 1C), faster (Figure 1D), and with greater confidence (Figure 1E)  
28 following a rewarded trial, when task certainty is high.

29

### 30 **Recurrent neural networks trained on behavioural data can generate mechanistic hypotheses**

31 To generate hypotheses about the neural mechanisms underlying this cost of flexibility, we  
32 trained recurrent neural networks (RNNs) to perform the same two-feature discrimination task  
33 and then analyzed the internal dynamics of the models. We used two training strategies to create  
34 two distinct, comparable models. Both models received as inputs information about features of  
35 the stimuli on the current trial and the history of stimuli, choices, and rewards over the past nine  
36 trials (see Methods: RNN models and analyses). The first model, the “correct choice network,”  
37 was trained to produce the correct choice given those inputs (Figure 2A; see Methods). The  
38 correct choice network did not exhibit a cost of flexibility at any point during training (Figure 2B),  
39 which indicates that the cost of flexibility is not somehow related to the structure or statistics of  
40 our task; nor does it necessarily appear during the intermediate stage of task learning. The  
41 second model, the “monkey choice network,” was trained to predict the monkeys’ actual choices

1 (Figure 2A). Because of the way it was trained, unlike the correct choice network, the monkey  
 2 choice network exhibited the cost of flexibility characteristic of the behaviour of our monkeys  
 3 and human participants, which persisted from about halfway through training (Figure 2B and 2C).  
 4 To model reaction times, we applied a collapsing decision boundary to the four choice output  
 5 units such that the reaction time is given by the first timepoint in which any choice unit's activity  
 6 crosses the boundary. The boundary parameters were chosen such that the mean reaction times  
 7 of the monkey choice network after rewarded and unrewarded trials were comparable to those  
 8 of the monkeys. We compared these reaction times with those of the correct choice network  
 9 using the same decision boundary (Figure 2D). The correct choice network was faster than the  
 10 monkey choice network both after rewarded trials and unrewarded trials. Notably, this  
 11 difference was greater after unrewarded trials, indicating that task certainty has a greater effect  
 12 on the time course of choice output activity in the monkey choice network.



1 **FIG 2:** Generating mechanistic hypotheses using recurrent neural networks. (A) Model  
2 schematic for RNNs designed to either predict the correct choice (correct choice network) or  
3 predict the monkeys' choice on each trial (monkey choice network). (B) The difference in  
4 perceptual accuracy (after rewarded trials – after unrewarded trials) throughout training for  
5 both models. The two monkeys' difference in perceptual accuracy was on average about 5%,  
6 as indicated by the brown dashed line. The monkey choice network approaches this  
7 difference, while the correct choice network remains around zero (shaded regions show 95%  
8 confidence intervals from sliding window averaging). (C) The monkey choice network (left)  
9 displays a cost of flexibility after training, as shown by worse perceptual performance after  
10 an unrewarded, low task certainty trial (location choices:  $p=2.4\times 10^{-14}$ , frequency choices:  
11  $p=6.6\times 10^{-8}$ ). In contrast, the correct choice network (right) did not show a significant cost of  
12 flexibility (location choices:  $p=0.43$ , frequency choices:  $p=0.14$ ). Each point corresponds to a  
13 set of accuracies calculated from 50 trials with identical feature change amounts. (D) The  
14 monkey choice network displays a larger effect of task uncertainty on reaction time  
15 compared to the correct choice network. The plot depicts histograms of reaction times when  
16 the same collapsing decision boundary is applied to the outputs of both RNNs (see Methods).  
17 Overall, the monkey choice network has longer reaction times than the correct choice  
18 network given the same inputs (after unrewarded trials:  $p=1.7\times 10^{-271}$ , after rewarded trials:  
19  $p<10^{-300}$ ). The difference between models is greater after unrewarded than after rewarded  
20 trials ( $p=3.7\times 10^{-36}$ , Wilcoxon rank sum test). (E) At the end of the delay period between  
21 stimulus presentations (400 ms after first stimulus presentation), information about the  
22 irrelevant feature (e.g. location on trials when the model chose the spatial frequency task) is  
23 encoded better following non-rewarded than rewarded trials (points are above the diagonal  
24 for both models;  $p=1.49\times 10^{-8}$  for both). Each point represents a set of trials with an identical  
25 first stimulus and shows the performance of a corresponding cross-validated linear decoder.  
26 (F) The two features are represented more orthogonally in the correct choice network than  
27 in the monkey choice network at the end of the delay period (400 ms after first stimulus  
28 presentation). Histograms of the magnitude of the projection of the irrelevant feature axis  
29 onto the relevant feature axis for the monkey choice and correct choice networks after  
30 unrewarded and rewarded trials. The unit axis for each feature was computed using distance  
31 covariance analysis and scaled by the corresponding distance covariance (see Methods for  
32 details). All four distributions are significantly different (Wilcoxon rank sum tests: monkey  
33 choice network after non-reward > correct choice network after non-reward:  $p=6.1\times 10^{-106}$ ;  
34 monkey choice network after reward > correct choice network after reward:  $p=3.8\times 10^{-142}$ ;  
35 monkey choice network after reward > after non-reward:  $p=0.001$ , correct choice network  
36 after non-reward > after reward:  $p=1.8\times 10^{-35}$ ). Triangle markers represent the medians. (G)  
37 A schematic of how overrepresentation of the irrelevant feature when the task is uncertain  
38 (depicted in E) combined with non-orthogonal representations of task variables (depicted in  
39 F) leads to interference and thus a cost of flexibility. In an idealized correct choice network  
40 (top), the irrelevant feature is more strongly encoded (longer magenta axis) under low task  
41 certainty, however, it varies orthogonally to the relevant feature, thus avoiding interference.

1 In the monkey choice network (bottom), the feature axes are no longer orthogonal, leading  
2 to more interference under low task certainty.

3

#### 4 **Task uncertainty leads to task interference**

5 The observation that our monkey choice network, but not our correct choice network,  
6 demonstrates costs of flexibility gives us a platform for generating hypotheses about the  
7 mechanistic origins of this cost. We investigated the dynamics of the two networks and identified  
8 two features that could in principle underlie the cost of flexibility. First, consistent with the  
9 concepts of attentional and task-set inertia (1, 2, 5-7, 9, 10, 12), both models represented the  
10 irrelevant feature better under low task certainty than under high task certainty (orange data  
11 points, Figure 2E). When the task was more certain (following rewarded trials), the relevant  
12 feature was encoded throughout the entire delay period, while information about the irrelevant  
13 feature decayed nearly completely (Figure S2). This certainty-dependent ‘forgetting’ of irrelevant  
14 information occurred even in the correct choice network (blue data points, Figure 2E), which did  
15 not show a cost of flexibility, so it is on its own not a sufficient explanation of the cost of flexibility  
16 in biological organisms.

17 The second feature of the models is not predicted by common accounts of task switch costs. In  
18 the monkey choice network, the representations of the two visual features of the first stimulus  
19 become non-orthogonal before the onset of the second stimulus presentation (Figures 2F, S3).  
20 This tangling of the representations causes the monkey choice network to mix up information  
21 about the two features (e.g. judgments of spatial frequency are influenced by the location of the  
22 stimulus), which reduces accuracy. We can visualize the problem by identifying the axes that best  
23 represent information about each visual feature in a space in which the activity of each unit in  
24 the hidden layer is one dimension. These axes are the directions in population space along which  
25 there is most variation in responses to different values of that feature. If information about the  
26 two features is encoded independently, the axes best representing them will be orthogonal.  
27 However, in the monkey choice network, those axes rotate such that by the end of the delay  
28 period they are no longer orthogonal (Figures 2F, S3(B)). This non-orthogonality means that  
29 information about the two features is no longer independent, leading to ‘switch errors’ that  
30 worsen perceptual performance. In contrast, the correct choice network maintains nearly  
31 orthogonal feature axes, especially after unrewarded trials, such that their dot product remains  
32 closer to zero (Figures 2F, S3).

33 We hypothesize that the cost of flexibility is the joint product of these two mechanisms (Figure  
34 2G). The increased representation of irrelevant information when the task is uncertain is not a  
35 problem when information about the two features is encoded independently (as in the correct  
36 choice network). The independent representations make it straightforward to ignore irrelevant  
37 information, even when it is robustly represented. However, when information about the two  
38 features is not independent, as in the monkey choice network, the irrelevant information is  
39 confounded with information about the relevant feature, worsening perceptual performance.

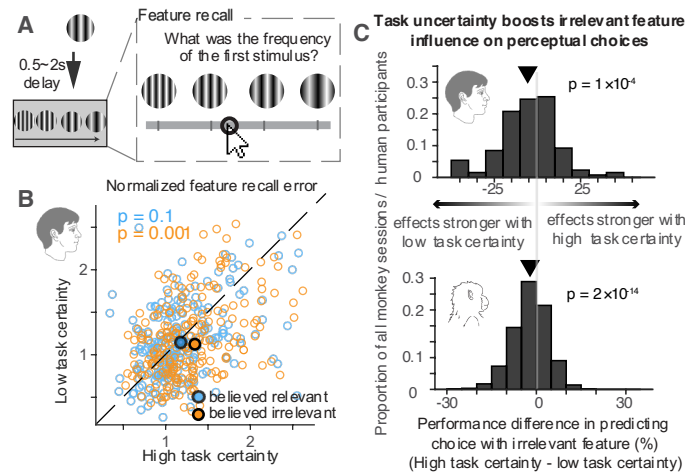
40 We tested predictions of this task interference hypothesis using new behavioural, neuronal, and  
41 causal experiments designed to identify the types of errors associated with task uncertainty.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17

## Behavioural evidence for task interference

We first looked for behavioural evidence supporting the first part of our task interference hypothesis, the information deemed irrelevant for the chosen task is more strongly retained under low task certainty. On a small subset of trials, we asked human participants to report one of the features (orientation or frequency) of the first stimulus after a delay (Figure 3A). Consistent with the idea that irrelevant information is maintained longer in short-term memory when the task is uncertain, participants had higher accuracy (smaller errors) when estimating the believed-irrelevant, but not the relevant, feature on low task certainty trials (Figure 3B).

Next, we looked for behavioural evidence supporting the hypothesis that the two features are entangled when the task is uncertain. The monkey choice network predicted that under uncertain task conditions, participants' choices about one feature (e.g., spatial frequency) are influenced by the other feature (e.g., location). Consistent with this hypothesis, the choices of humans (Figure 3C upper panel) and monkeys (Figure 3C lower panel) were more strongly related to the irrelevant stimulus feature following a non-reward (when the task was relatively uncertain) than a reward (when the task was more certain).



18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

**FIG 3:** Behavioural evidence for the task interference hypothesis. (A) In a subset of the trials, we asked human participants to reproduce one of the features (spatial frequency or orientation) of the first stimulus according to their memory following a delay. The plot depicts the experimental design, including the slider used to reproduce the stimulus. (B) When participants reported that a feature was behaviourally irrelevant, their recall error was better (smaller) for that feature if the preceding trial had been incorrect (low task certainty, y-axis) than correct (high task certainty, x-axis). The orange data points lie significantly below the identity line ( $p=0.001$ ), indicating that believed-irrelevant information was reported with smaller errors under low task certainty. In contrast, there was no significant relationship between task certainty and participants' ability to remember the believed-relevant feature ( $p>0.05$ ), although the deviations from the identity line are not significantly different between believed relevant and believed irrelevant features (Wilcoxon rank sum test,  $p>0.05$ ).



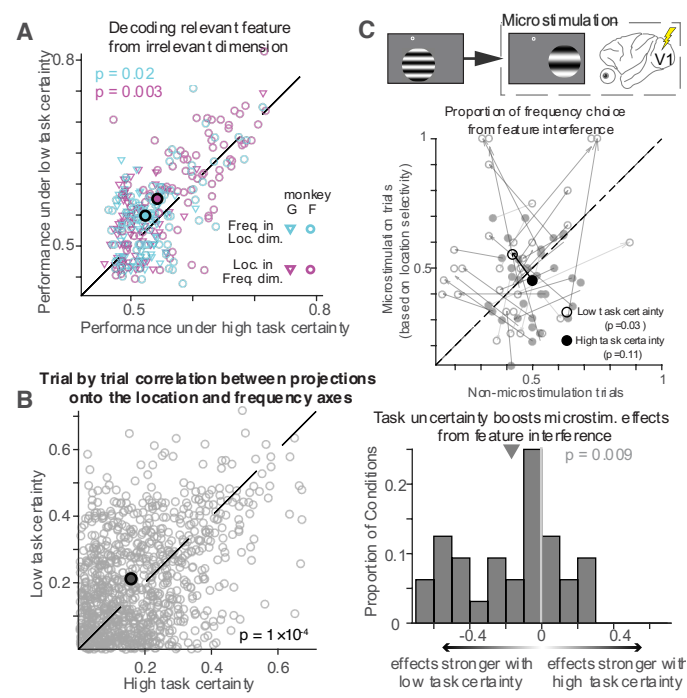
1 (C) Perceptual choices are significantly more predictable from the irrelevant feature under  
 2 lower task certainty. For each monkey experiment session and human participant, we predict  
 3 the perceptual choices on one feature based on the changes in the other feature for trial-  
 4 number matched conditions of low and high task certainty. Histograms (top: humans, bottom:  
 5 monkeys) show that across monkey experimental sessions / human participants, the  
 6 distributions of differences in prediction performances under high and low task certainty  
 7 were significantly below zero ( $p=1\times 10^{-4}$  for humans,  $p=2\times 10^{-14}$  for monkeys), indicating that  
 8 irrelevant feature is more predictive of perceptual choices under task uncertainty.

9

10 **Neuronal evidence for task interference**

11 To test the predictions of the feature interference hypothesis on neural representations, we  
 12 recorded groups of neurons in monkey primary visual cortex (V1), which is known to encode the  
 13 visual features (frequency and location) that the monkeys discriminated. We analyzed the  
 14 neurons in a similar way as the units in our model, beginning by using linear regression to identify  
 15 the axes in V1 population activity space that best encode each visual feature. We found evidence  
 16 for task interference in both signal and noise. To look for entangled representations in the signal,  
 17 we tried to decode the frequency value from location encoding axis, and vice versa. Consistent  
 18 with our hypothesis, under low task certainty, the representations of spatial frequency and  
 19 location were more entangled, as indicated by better decoder performance under low task  
 20 certainty (Figure 4A). To investigate noise-based evidence for task interference, we investigated  
 21 the V1 population responses to the same stimulus across trials. The trial-to-trial fluctuations  
 22 between projections onto the spatial location and frequency axes were more correlated under  
 23 low task certainty (Figure 4B).

24



1 **FIG 4:** Neuronal and causal evidence for task interference from neural population recordings  
2 and electrical microstimulation in monkey V1. (A) Stronger feature entanglement under low  
3 task certainty. For each monkey recording session, we identified the dimensions in V1  
4 population space along which we could best linearly decode each visual feature. Using linear  
5 decoders, we assessed our ability to decode the relevant feature in the encoding dimension  
6 of the irrelevant feature under high (abscissa) and low task certainty (ordinate). Across  
7 sessions, for both features, the points are distributed significantly above the identity line  
8 ( $p=0.003$  for location classifiers in frequency dimension,  $p=0.02$  for frequency classifiers in  
9 location dimension), indicating increased feature entanglement under low task certainty. (B)  
10 Projections of V1 population responses onto the spatial frequency and location axes are more  
11 correlated under low task certainty, indicating that the two axes are less orthogonal. The  
12 abscissa and ordinate of each data point show the absolute value of the correlation  
13 coefficient relating projections onto the two axes for trials with identical stimuli under high  
14 and low task certainty. We use the absolute value because there is no a priori prediction  
15 about the sign of the correlation, only that nonzero correlation indicates non-orthogonality.  
16 The mean of the distribution is above the identity line, indicating that decoding of the two  
17 features is more correlated under low task certainty ( $p=1\times 10^{-4}$ ). (C) Causal tests of the task  
18 interference hypothesis using V1 microstimulation. On a subset of randomly interleaved trials,  
19 we microstimulated V1 neurons during the display of the second stimulus. Each point in the  
20 scatterplot in the upper panel represents a set of trials with identical stimuli, displaying the  
21 proportion of a given spatial frequency choices on trials with microstimulation (ordinate) or  
22 without microstimulation (abscissa). If the neuronal population representations of the two  
23 features were independent, the points would align along the identity line, while deviation  
24 indicates feature interference. The observation that microstimulation has a bigger impact on  
25 behaviour under low task certainty is evidence for task interference (Compare  
26 microstimulation effects in high certainty: filled grey symbols, not significantly away from  
27 identity line,  $p=0.11$ ; vs. low task certainty: open grey symbols, above identity line,  $p=0.03$ ).  
28 Black symbols indicate the means of each task certainty condition. The bottom histogram  
29 plots the differences between microstimulation effects under low and high task certainty (i.e.  
30 how much the arrows in the scatter moved the data points away from the diagonal). The  
31 distribution across conditions reveals a significantly negative distribution ( $p=0.009$ ) that  
32 confirms stronger feature interference in V1 under low task certainty.

33

#### 34 **Causal evidence for task interference**

35 Finally, we used a causal manipulation to test the prediction of the task interference hypothesis  
36 that the representations of the two features are more entangled under low task certainty. If the  
37 representations of the two features are orthogonal, then a causal manipulation that biases  
38 judgments of one feature should not impact judgments of the other. If, on the other hand, the  
39 representations of the two features are entangled, then biasing judgments of one feature should  
40 necessarily impact judgments of the other in the way predicted by the entanglement.

41 We therefore used electrical microstimulation in V1 to bias judgments of one feature (location)  
42 and measured the impact on judgments of the other feature (spatial frequency) under different

1 task certainty conditions. Thanks to the retinotopic organization of V1, we were able to induce a  
2 bias in the perceived location change by microstimulating V1 during the presentation of the  
3 second stimulus. The bias was predicted by the tuning of the stimulated site: when the receptive  
4 field recorded on the microstimulated channel was located to the right of the first stimulus, the  
5 monkey reported rightward location shifts more frequently, and vice versa (Supplementary  
6 Figure S4(A)). On trials when the monkey performed the spatial frequency task, microstimulation  
7 (which affects essentially only location judgments; Fig. S4(B)) should not have influenced choices  
8 if the representations were orthogonal. Consistent with the task interference hypothesis,  
9 microstimulation had a bigger impact on spatial frequency judgments on trials with low than high  
10 task certainty (Fig. 4C open symbols significantly deviate more from the diagonal line than filled  
11 symbols), suggesting that the representations of location and spatial frequency were less  
12 orthogonal under low certainty.

13

## 14 **Discussion:**

15 Humans and animals must flexibly switch between tasks to survive in an uncertain world. This  
16 flexibility has long been known to come at a cost. The cost of flexibility has been explained by a  
17 diverse set of ideas, which are mostly modeled and tested using behaviour (13-20). We aimed to  
18 unify and mechanistically explain the extensive body of behavioural literature that has explored  
19 this intriguing phenomenon.

20 We used a multidisciplinary approach to identify neural population reasons for the cost  
21 associated with switching between tasks. By grounding our investigation in the neural  
22 representations of tasks and their associated variables, we uncovered a task interference  
23 neuronal mechanism that unifies previous findings and has implications for our understanding of  
24 many cognitive processes.

### 25 *Integrative approach reveals neuronal mechanism for complex behaviour*

26 In our research, we presented a novel research paradigm combining artificial neural networks,  
27 monkey electrophysiology, and human psychophysics, to understand the neuronal mechanism  
28 for advanced cognition and complex behaviour. By constraining artificial neural networks using  
29 actual behavioural observations, we generate data-driven hypotheses about the neuronal  
30 dynamics that underlie the intricacies in real behaviour. Large-scale electrophysiology in  
31 behaving monkeys allows for these hypotheses to be empirically tested at a single-neuron level,  
32 providing a direct examination of the neural correlates of hypothesized cognitive processes.  
33 Human psychophysics complements this, extending the relevance of our findings beyond a single  
34 species and ensuring applicability to broader cognitive phenomena.

35 Combined, our approach holds potential to extract meaningful insights from complex data, a  
36 critical challenge in an era when datasets become increasingly high dimensional (39).

### 37 *Task interference may impact neural processes throughout the brain*

38 Remarkably, we observed neural signatures of task interference even in V1, a primary sensory  
39 area not typically associated with strong modulations by cognitive processes (21, 22). This  
40 suggests that the phenomenon of interfering neural representations is not confined to specific

1 brain regions but rather extends throughout the brain. Future research should investigate  
2 whether task interference poses an even stronger limit in higher-order visual areas and prefrontal  
3 cortex, which are more intimately linked with more advanced cognitive processing.

#### 4 *Dimensionality as a fundamental limit on cognition*

5 Limited capacity is a repeatedly reported constraint on multiple cognitive processes including  
6 working memory (23), attention (24), decision-making (25), and learning (26). Our findings raise  
7 the intriguing possibility that the dimensionality of a neural population representation is a  
8 fundamental limit on cognition because it limits the behavioural repertoire that the population  
9 can enable. The dimensionality limits the number of quantities (e.g. sensory features, number of  
10 objects, task variables, timescales, etc.) that can be represented independently. In theory, the  
11 dimensionality of a neural population representation could approach the number of neurons in  
12 that population, but the true dimensionality is typically much lower (27, 28).

13 Observations in diverse experimental systems are broadly consistent with the possibility that the  
14 dimensionality of neural population representations limits cognitive capacity. For instance,  
15 associating multiple simple features (such as the colour and surface texture of a banana) to aid  
16 in perceiving a complex feature (such as the ripeness of the banana), might combine, or reduce  
17 the dimensionality of, the representations of those associated features (29). However, reversing  
18 such learned associations would in this scenario require the brain to increase the dimensionality  
19 again, which could account for the challenges in reversal learning in operant conditioning in mice  
20 (30), rule learning in children with bipolar disorder (31), and language learning in neurotypical  
21 people (32).

22 Similarly, limitations in working memory capacity (23, 33) might be thought of as a limit on the  
23 number of items that can be represented independently in a population of neurons. People can  
24 overcome those limits by ‘chunking’ the memory of different items into larger, familiar units (e.g.,  
25 memorizing the digits of a credit card or phone number in groups of three or four (33, 34)). This  
26 strategy may effectively reduce the neural dimensions needed to represent and maintain the  
27 whole sequence. Interestingly, artificial neural networks exhibit human-like working memory  
28 capacity limitations only after being pre-trained on large-scale natural image tasks (35), which  
29 may indicate that the mechanisms adopted to process high-dimensional inputs compromise  
30 performance on reduced laboratory tasks, perhaps due to an increase in interference between  
31 stimuli. The possibility that limitations on diverse cognitive processes share a neuronal  
32 mechanism offers a path toward understanding and alleviating those limitations.

#### 33 *Translational implications of task interference*

34 Deficits in cognitive flexibility are an especially stark consequence of the many brain disorders,  
35 including dementia, substance abuse, and developmental disorders, that are associated with  
36 diminished cognitive capacity (36, 37). If task interference indeed underlies many cognitive  
37 limitations, our results offer hope for improving and repairing cognitive capacity. For example,  
38 we recently demonstrated that methylphenidate (trade name Ritalin, a stimulant used to treat  
39 disorders of attention; (38)) effectively increases the dimensionality of neural population  
40 representations in the visual cortex and that this increase is correlated with improvement in a  
41 visual task (38). This is a proof of principle supporting the idea that repurposing existing drugs to

1 address other cognitive abilities that might share a neural mechanism with their primary target  
2 could be an effective way to improve cognition in health and disease. The suggestion in our V1  
3 results that task interference affects brain-wide mechanisms means that systemic drugs and/or  
4 therapies that affect widespread neuromodulatory systems might be especially effective.

5 Together, our results hold promising implications for the future of basic science and translational  
6 efforts to mechanistically understand, improve, and repair cognition. Our task interference  
7 hypothesis provides a concrete set of testable predictions with implications for theoretical and  
8 experimental work spanning many species and systems.

9

## 10 **Methods:**

### 11 **Behavioural task for non-human primates**

12 We designed a two-interval, two-feature discrimination task with stochastic task switching that  
13 allows us to measure and manipulate task certainty and perceptual decision-making (Xue et al,  
14 2022). Briefly, two rhesus monkeys (both male, 12 and 9 kg) fixated a central spot while two  
15 Gabor stimuli ( $\sim 2^\circ$  visual angle in diameter) were displayed in series (200 ms each), separated by  
16 a variable delay (300-500 ms). The second Gabor stimulus differed from the first in its spatial  
17 location (shifted left or right) and spatial frequency (higher or lower), with independently  
18 randomized change amounts in the two features. Small titrations of the change amounts in the  
19 two features were made before each experiment to keep the perceptual performance of the  
20 animals in both tasks at roughly 75%. After a 150 ms second delay, the fixation dot disappeared,  
21 and the animals looked at one of four peripheral targets to indicate both the inferred relevant  
22 feature and the direction of change in that feature (spatial frequency increase or decrease and  
23 spatial location left-shift or right-shift). The relevant feature changed stochastically on 2.5% of  
24 trials, and the monkeys were rewarded only if they correctly discriminated the relevant feature.  
25 The visual display contained no information about the behavioural relevance of features, so the  
26 animals needed to infer the relevant feature based on their stimulus, choice, and reward history.  
27 The median number of trials in each daily session was 1086. All animal procedures were approved  
28 by the Institutional Animal Care and Use Committees of the University of Pittsburgh and Carnegie  
29 Mellon University, where the animal experiments were conducted.

### 30 **Behavioural task for the human participants**

31 We adapted our task for online human psychophysics (built with Gorilla experiment building tools  
32 (40) and customized scripts). We recruited a total of 220 adult human participants through  
33 Prolific (prolific.co), who performed slightly different combinations of behavioural experiments.

34 In the main experiments, we used a version of the two-interval, two-feature discrimination task  
35 with stochastic uncued task switching that we adapted to better suit human participants and  
36 online experiments (Figure 1B). The timing was slightly different (500 ms initial fixation period,  
37 300 ms stimulus presentations, and 500 ms delay period). The two features to be discriminated  
38 were orientation and spatial frequency. As in the monkey experiments, the changes in the two  
39 features were independently varied. Similar to the monkeys, the change amounts were set at  
40 each participant's perceptual threshold, at which the participants performed perceptual

1 discriminations in both tasks at roughly 75% correct, determined by a psychophysical staircase  
2 procedure before the main task session. Participants indicated their choice by clicking on one of  
3 four buttons representing each feature and its change direction (indicated by a combination of  
4 text and graphical icon). As in the monkey experiments, the human participants received  
5 feedback following each choice (a green tick indicated a correct report of the direction of change  
6 in the task-relevant feature, and a red cross indicated an incorrect perceptual judgment, task  
7 choice, or both). As in the monkey experiments, the relevant feature switched stochastically, but  
8 the frequency was different (with a probability of 0.1 on each trial for humans, 0.025 for  
9 monkeys). The overall length of a typical session varied between 100 to 150 trials and around 15-  
10 20 minutes.

### 11 *Task-belief and feature perception confidence report*

12 Some task sessions required participants to report their confidence in their perceptual judgment  
13 on a random 20% of trials. After reporting their decision on the main task but before receiving  
14 feedback, participants were presented with three confidence sliders. Participants were  
15 instructed to indicate their task-belief on the first slider, ranging from a strong belief that spatial  
16 frequency was relevant to a strong belief that orientation was relevant. An intermediate value  
17 indicated low task certainty. The second and third sliders were used to report the perceived  
18 direction of change for each feature, where the absolute value reflected confidence in the  
19 perceptual judgment.

### 20 *Variable delay and feature recall question*

21 In some sessions, we explored memory for believed-relevant and believed-irrelevant stimulus  
22 features by interrupting a randomly selected 20% of trials after either 500 ms (short delay) or  
23 2000 ms (long delay) to ask participants to recall a feature of the first Gabor. Participants  
24 indicated the feature memory using a slider, facilitated by four reference Gabor patches that  
25 spanned the possible range. We quantified estimation error as the difference between the actual  
26 stimulus feature and the recalled feature.

## 27 **RNN models and analyses**

### 28 *Training and architecture*

29 We trained RNN models using custom code based on PsychRNN (Ehrlich et al., 2021) and  
30 TensorFlow (Abadi et al., 2016). The models were governed by the following dynamical equations:

$$\begin{aligned} 31 \quad \tau dx &= (-x + W_{rec} r + b_{rec} + W_{in} u) dt + \sigma_{rec} \sqrt{2\tau} d\xi \\ 32 \quad r &= f(x) = \max(0, x) \\ 33 \quad z &= W_{out} r + b_{out} \end{aligned}$$

34 where  $x$  is the recurrent state, which when passed through the rectified linear unit (ReLU)  
35 activation function gives the firing rate  $r$ ,  $u$  is the input vector, and  $z$  is the output vector.  $W_{in}$ ,  
36  $W_{rec}$ , and  $W_{out}$  are the input, recurrent, and output weight matrices, and  $b_{rec}$  and  $b_{out}$  are constant  
37 biases added to the recurrent and output units.  $dt$  is the discrete time step,  $\tau$  is the neuronal time  
38 constant, and  $\sigma_{rec}$  scales recurrent noise, which is a Gaussian noise process ( $d\xi$ ).

39 Our monkey choice network (Figure 2A) was trained to predict the choices of two macaques  
40 based on the same history of stimuli, choices, and rewards experienced by our monkey

1 participants. The RNN had two modules, which we term the ‘task module’ and the ‘perception  
2 module.’ The two modules were fully connected in the recurrent layer but received different  
3 inputs and mapped onto different outputs. The task module received the history of monkeys’  
4 choices, stimulus feature change amounts (signaling perceptual difficulty), and rewards for nine  
5 preceding trials. These inputs were constant throughout the trial, based on the assumption that  
6 task rule dynamics are slow compared to single-trial stimulus response dynamics. The task  
7 module was trained to output the task that the monkey chose on the current trial as a one-hot  
8 vector. The perception module received the spatial frequency (SF) and spatial location (SL) values  
9 of the first stimulus and, after a variable delay (300 to 500 ms, uniform distribution), of the  
10 second stimulus. Units in the perception module displayed mixed selectivity for the two features.  
11 The perception module was trained to output one of the four choices (SF increase, SF decrease,  
12 SL decrease, SL increase) as a one-hot vector.

13 We compare this monkey choice network to a correct choice network (Figure 2A), which was  
14 trained to output the correct choice (rather than the monkeys’ choices) given the perceptual  
15 inputs on the current trial and its own trial history (feedback based on the model’s choices rather  
16 than the monkeys’). Stimulus values were drawn from the same data used to train the monkey  
17 choice network. After the models were trained, we used the monkeys’ trial history inputs for  
18 both models so that we could directly compare trials (all model figures besides 2B). The correct  
19 choice network trace in figure 2B was calculated directly from its training batches (200 trials per  
20 epoch), and sliding window averaged over 100 epochs. The monkey choice network trace was  
21 calculated post-hoc from weights saved during training (every 15 epochs), and sliding window  
22 averaged over 6 datapoints (90 epochs). We chose to analyze the monkey choice network at the  
23 point during training with the highest monkey choice accuracy (see Figure S1), which was at  
24 training epoch 1245 out of 1500. All RNNs were trained using the Adam optimizer, with L1 and  
25 L2 regularization of weights and L2 regularization of firing rates. Note that to keep the model  
26 architecture consistent with the monkey choice network, the correct choice network here has  
27 key differences in the inputs from the model we used in our previous study (3).

### 28 *Reaction time calculation*

29 To calculate model reaction times, we applied a collapsing decision boundary of the form

30  $f(t) = a \cdot \left(1 - \frac{b \cdot t}{t+c}\right)$  to the choice output units of the RNNs, where  $t$  is the time after choice  
31 period onset. For figure 2D, the parameter values were  $a=1.5$  ms,  $b=1.8$ ,  $c=28$  ms.

### 32 *Distance covariance analysis for assessing orthogonality of feature axes*

33 To quantify the (non-)orthogonality of feature axes in the models, we used distance covariance  
34 analysis (DCA), a dimensionality reduction method that identifies linear projections that capture  
35 interactions between variables by maximizing the correlational statistic distance covariance (41).  
36 After reducing the dimensionality of the activity of all units in the hidden layer to the first 10  
37 principal components (capturing >95% of the variance), we calculated the distance covariance  
38 statistic and corresponding axis separately for each unique set of trial history inputs, feature  
39 (location and spatial frequency), and timepoint during the delay period, from 500 simulated trials  
40 with different feature values. The projection of the irrelevant feature axis onto the relevant

1 feature axis was calculated by taking the dot product of the unit vectors and scaling it by the  
2 smaller distance covariance (that of the believed-irrelevant feature) (Figures 2F and S3).

### 3 **Electrophysiology for the monkey experiments**

4 Different aspects of a subset of the electrophysiological data have been reported previously (Xue  
5 et al, 2022). The visual stimuli were displayed on a linearized CRT monitor (1,024 × 768 pixels,  
6 120-Hz refresh rate) placed 57 cm from the animal. We monitored eye position using an infrared  
7 eye tracker (Eyelink 1000, SR Research) and used custom software (written in MATLAB using the  
8 Psychophysics Toolbox, (25) to present stimuli and monitor behaviour. We recorded eye position  
9 and pupil diameter (1,000 samples per s), neuronal responses (30,000 samples per s), and the  
10 signal from a photodiode to align neuronal responses to stimulus presentation times (30,000  
11 samples per s) using hardware from Ripple. We recorded neuronal activity from chronically  
12 implanted Utah arrays (Blackrock Microsystems; 48 electrodes in V1) during daily experimental  
13 sessions for several months in each animal (90 sessions from monkey F and 68 sessions from  
14 monkey G). We set the threshold for each channel at three times the standard deviation and used  
15 threshold crossings as the multiunit activity on that unit. We positioned the stimuli to maximize  
16 the overlap between potential stimulus locations and the joint receptive fields of V1 units. The  
17 receptive fields were measured using separate mapping sessions, during which the monkeys  
18 fixated their gaze on a dot, while Gabor stimuli were subsequently flashed at random locations  
19 for 100ms each across a portion of the screen.

20 We included experimental sessions for analysis if the total number of completed trials was at  
21 least 480 (which was ten times the number of channels on the multielectrode array). During  
22 complete trials the monkeys successfully maintained fixation until they indicated their choice).  
23 We analyzed the activity of V1 units during stimulus display periods, shifted by 34 ms neuronal  
24 response latency (42). V1 units are included for analyses if the stimulus response is 25% larger  
25 than baseline activity, measured 100 ms before stimulus onset, and 2) larger than 5 sp/s. These  
26 procedures resulted in 89 sessions from Monkey F and 68 sessions from Monkey G; an average  
27 46 V1 units, and a median 1086 completed trials per session. Other aspects of these data have  
28 been reported in a previous study (3).

29 We did V1 microstimulation experiments during 6 sessions in one monkey (monkey F, Figure 4).  
30 We were unable to conduct these experiments in the other monkey because the experiments  
31 were interrupted by the COVID shutdown of the lab, after which there were no usable channels  
32 on the array. During microstimulation sessions, we electrically stimulated one channel on the V1  
33 array with 20, 30  $\mu$ A, 200 Hz biphasic pulses between 50 and 150 ms after the onset of the second  
34 Gabor stimulus on a randomly selected 50% of trials.

35 Microstimulation perturbed the monkey's reported location shift consistently based on the  
36 relative locations of the V1 receptive field and the first Gabor location (Supplementary Figure  
37 S4(A), but did not effectively manipulate the monkeys reported frequency change based on the  
38 preferred frequency relative to the first Gabor frequency (Supplementary Figure S4(B). This  
39 discrepancy may be attributed to the fact that microstimulation likely activates neurons through  
40 direct axonal activation within a volume of tens of microns in diameter (43). While these  
41 activated neurons may exhibit different frequency selectivity than the recorded neurons on the  
42 same channel, they likely share similar receptive fields due to V1's retinotopic structure. We



1 assessed the level of interference as the degree to which the receptive field location relative to  
2 the first Gabor location influenced the monkey's frequency choices (which should be zero if the  
3 representations of frequency and location are orthogonal). We then compared the level of  
4 interference induced by microstimulation following rewarded trials and following unrewarded  
5 trials.

## 6 **Statistical tests**

7 All p-values reported in this study are from Wilcoxon signed rank tests unless otherwise specified.

8

9

## 10 **Acknowledgments**

11 We are grateful to John Maunsell, Douglas Ruff, Ramanujan Srinath, and Pouya Bashivan for  
12 comments on this manuscript. This work was supported by the Simons Foundation (Simons  
13 Collaboration on the Global Brain award 542961SPI to MRC), the National Eye Institute of the  
14 National Institutes of Health (awards R01EY022930, R01EY034723, and RF1NS121913 to MRC).

15

## 16 **Author contributions**

17 C.X. and M.R.C. designed research; S.K.M. and C.X. trained recurrent neural nets and analyzed  
18 the data; C.X. and L.E.K. performed monkey training and neuronal recording; C.X. analyzed the  
19 neuronal data and performed the microstimulation experiments; R.C. and C.X. performed the  
20 online human psychophysics experiments and data analyses; M.R.C. supervised the findings of  
21 this work; C.X., S.K.M and M.R.C. wrote the paper.

22

## 23 **References:**

- 24 1. Monsell S. Task switching. *Trends Cogn Sci.* 2003;7(3):134-40.  
25 2. Kiesel A, Steinhauser M, Wendt M, Falkenstein M, Jost K, Philipp AM, et al. Control and  
26 interference in task switching--a review. *Psychol Bull.* 2010;136(5):849-74.  
27 3. Xue C, Kramer LE, Cohen MR. Dynamic task-belief is an integral part of decision-making. *Neuron.*  
28 2022.  
29 4. Koch I, Poljac E, Müller H, Kiesel A. Cognitive structure, flexibility, and plasticity in human  
30 multitasking-An integrative review of dual-task and task-switching research. *Psychol Bull.*  
31 2018;144(6):557-83.  
32 5. Longman CS, Lavric A, Munteanu C, Monsell S. Attentional inertia and delayed orienting of  
33 spatial attention in task-switching. *J Exp Psychol Hum Percept Perform.* 2014;40(4):1580-602.  
34 6. Mayr U, Kuhns D, Rieter M. Eye movements reveal dynamics of task control. *J Exp Psychol Gen.*  
35 2013;142(2):489-509.  
36 7. Longman CS, Lavric A, Monsell S. More attention to attention? An eye-tracking investigation of  
37 selection of perceptual attributes during a task switch. *J Exp Psychol Learn Mem Cogn.* 2013;39(4):1142-  
38 51.  
39 8. Pashler H. Dual-task interference in simple tasks: data and theory. *Psychol Bull.*  
40 1994;116(2):220-44.

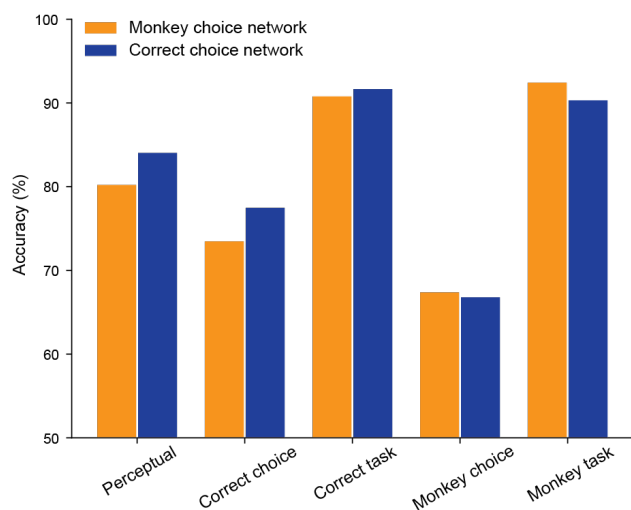
- 1 9. Allport A, Wylie G. Task-switching: Positive and negative priming of task-set. 1999.
- 2 10. Yeung N, Nystrom LE, Aronson JA, Cohen JD. Between-task competition and cognitive control in  
3 task switching. *J Neurosci*. 2006;26(5):1429-38.
- 4 11. Waszak F, Hommel B, Allport A. Task-switching and long-term priming: role of episodic stimulus-  
5 task bindings in task-shift costs. *Cogn Psychol*. 2003;46(4):361-413.
- 6 12. Allport A, Styles EA, Hsieh S. Shifting Intentional Set: Exploring the Dynamic Control of Tasks.  
7 1994.
- 8 13. Ardid S, Wang XJ. A tweaking principle for executive control: neuronal circuit mechanism for  
9 rule-based task switching and conflict resolution. *J Neurosci*. 2013;33(50):19504-17.
- 10 14. Driscoll L, Shenoy K, Sussillo D. Flexible multitask computation in recurrent networks utilizes  
11 shared dynamical motifs. *bioRxiv*. 2022:2022.08.15.503870.
- 12 15. Yang GR, Joglekar MR, Song HF, Newsome WT, Wang XJ. Task representations in neural  
13 networks trained to perform many cognitive tasks. *Nat Neurosci*. 2019;22(2):297-306.
- 14 16. Mante V, Sussillo D, Shenoy KV, Newsome WT. Context-dependent computation by recurrent  
15 dynamics in prefrontal cortex. *Nature*. 2013;503(7474):78-84.
- 16 17. Flesch T, Juechems K, Dumbalska T, Saxe A, Summerfield C. Orthogonal representations for  
17 robust context-dependent task performance in brains and neural networks. *Neuron*. 2022;110(24):4212-  
18 9.
- 19 18. Riveland R, Pouget A. Generalization in Sensorimotor Networks Configured with Natural  
20 Language Instructions. *bioRxiv*. 2023:2022.02.22.481293.
- 21 19. Deco G, Rolls ET. Attention and working memory: a dynamical model of neuronal activity in the  
22 prefrontal cortex. *Eur J Neurosci*. 2003;18(8):2374-90.
- 23 20. Duan CA, Pagan M, Piet AT, Kopec CD, Akrami A, Riordan AJ, et al. Collicular circuits for flexible  
24 sensorimotor routing. *Nat Neurosci*. 2021;24(8):1110-20.
- 25 21. Luck SJ, Chelazzi L, Hillyard SA, Desimone R. Neural mechanisms of spatial selective attention in  
26 areas V1, V2, and V4 of macaque visual cortex. *J Neurophysiol*. 1997;77(1):24-42.
- 27 22. McAdams CJ, Reid RC. Attention modulates the responses of simple cells in monkey primary  
28 visual cortex. *J Neurosci*. 2005;25(47):11023-33.
- 29 23. Cowan N. Working memory capacity: Classic edition: Psychology press; 2016.
- 30 24. Pylyshyn ZW, Storm RW. Tracking multiple independent targets: evidence for a parallel tracking  
31 mechanism. *Spa Vis*. 1988;3(3):179-97.
- 32 25. Kang YH, Löffler A, Jeurissen D, Zylberberg A, Wolpert DM, Shadlen MN. Multiple decisions  
33 about one object involve parallel sensory acquisition but time-multiplexed evidence incorporation.  
34 *bioRxiv*. 2020:2020.10.15.341008.
- 35 26. Van Merriënboer JJ, Sweller J. Cognitive load theory and complex learning: Recent  
36 developments and future directions. *Educational psychology review*. 2005;17:147-77.
- 37 27. Saxena S, Cunningham JP. Towards the neural population doctrine. *Curr Opin Neurobiol*.  
38 2019;55:103-11.
- 39 28. Chung S, Abbott LF. Neural population geometry: An approach for understanding biological and  
40 artificial neural networks. *Curr Opin Neurobiol*. 2021;70:137-44.
- 41 29. Libby A, Buschman TJ. Rotational dynamics reduce interference between sensory and memory  
42 representations. *Nat Neurosci*. 2021;24(5):715-26.
- 43 30. Meister M. Learning, fast and slow. *Curr Opin Neurobiol*. 2022;75:102555.
- 44 31. Wegbreit E, Cushman GK, Weissman AB, Bojanek E, Kim KL, Leibenluft E, et al. Reversal-learning  
45 deficits in childhood-onset bipolar disorder across the transition from childhood to young adulthood. *J*  
46 *Affect Disord*. 2016;203:46-54.
- 47 32. Idemaru K, Holt LL. Word recognition reflects dimension-based statistical learning. *J Exp Psychol*  
48 *Hum Percept Perform*. 2011;37(6):1939-56.

- 1 33. MILLER GA. The magical number seven plus or minus two: some limits on our capacity for  
2 processing information. *Psychol Rev.* 1956;63(2):81-97.
- 3 34. Thalmann M, Souza AS, Oberauer K. How does chunking help working memory? *J Exp Psychol*  
4 *Learn Mem Cogn.* 2019;45(1):37-55.
- 5 35. Xie Y, Duan Y, Cheng A, Jiang P, Cueva CJ, Yang GR. Natural constraints explain working memory  
6 capacity limitations in sensory-cognitive models. *bioRxiv.* 2023:2023.03.30.534982.
- 7 36. Ferreira CS, Maraver MJ, Hanslmayr S, Bajo T. Theta oscillations show impaired interference  
8 detection in older adults during selective memory retrieval. *Sci Rep.* 2019;9(1):9977.
- 9 37. Dajani DR, Uddin LQ. Demystifying cognitive flexibility: Implications for clinical and  
10 developmental neuroscience. *Trends Neurosci.* 2015;38(9):571-8.
- 11 38. Ni AM, Bowes BS, Ruff DA, Cohen MR. Methylphenidate as a causal test of translational and  
12 basic neural coding hypotheses. *Proc Natl Acad Sci U S A.* 2022;119(17):e2120529119.
- 13 39. Fetsch CR. The importance of task design and behavioural control for understanding the neural  
14 basis of cognitive functions. *Curr Opin Neurobiol.* 2016;37:16-22.
- 15 40. Anwyl-Irvine AL, Massonnié J, Flitton A, Kirkham N, Evershed JK. Gorilla in our midst: An online  
16 behavioural experiment builder. *Behav Res Methods.* 2020;52(1):388-407.
- 17 41. Cowley B, Semedo J, Zandvakili A, Smith M, Kohn A, Yu B. Distance Covariance Analysis. In: Aarti  
18 S, Jerry Z, editors. *Proceedings of the 20th International Conference on Artificial Intelligence and*  
19 *Statistics; Proceedings of Machine Learning Research: PMLR; 2017. p. 242--51.*
- 20 42. Schmolesky MT, Wang Y, Hanes DP, Thompson KG, Leutgeb S, Schall JD, et al. Signal timing  
21 across the macaque visual system. *J Neurophysiol.* 1998;79(6):3272-8.
- 22 43. Histed MH, Bonin V, Reid RC. Direct activation of sparse, distributed populations of cortical  
23 neurons by electrical microstimulation. *Neuron.* 2009;63(4):508-22.

24

## 25 Supplementary Figures

26

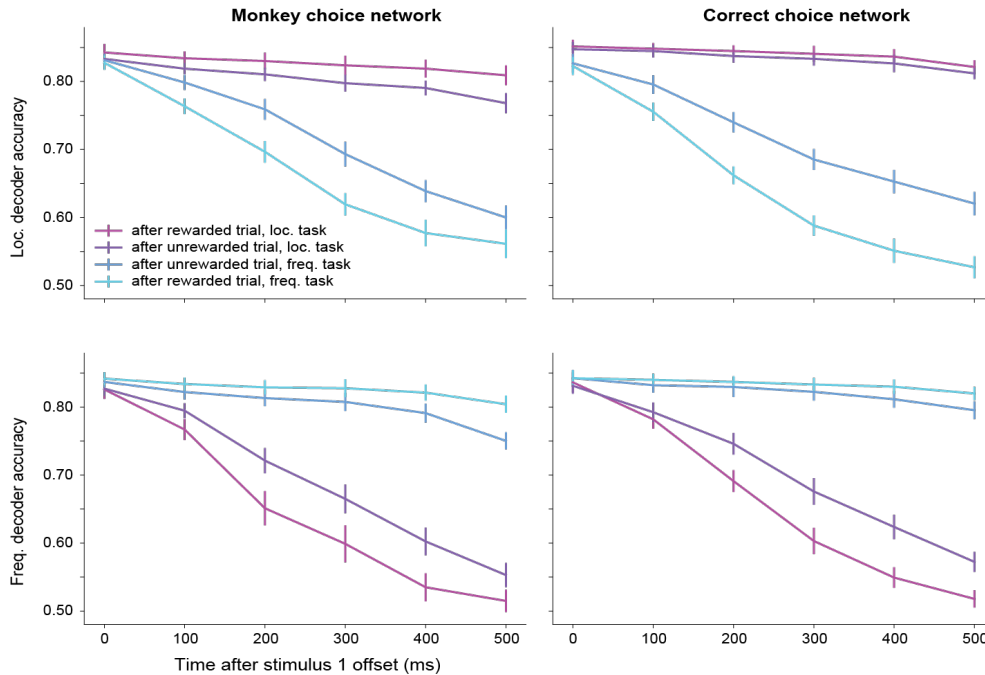


27

28 **FIG S1:** Accuracies of the monkey choice network (orange) and correct choice network (blue)  
29 after training. Perceptual: percentage of correct perceptual judgments, regardless of the  
30 chosen task (chance=50%); Correct choice: percentage of correct choices (chance=25%);  
31 Correct task: percentage of correct task choices, regardless of the perceptual judgment  
32 (chance=50%); Monkey choice: percentage of choices that match the monkeys' choices

1 (chance=25%). Monkey task: percentage of task choices that match the monkeys' chosen task,  
2 regardless of the perceptual judgment (chance=50%). All differences between models are  
3 statistically significant (two-sided z-tests,  $p < 0.0005$ ).

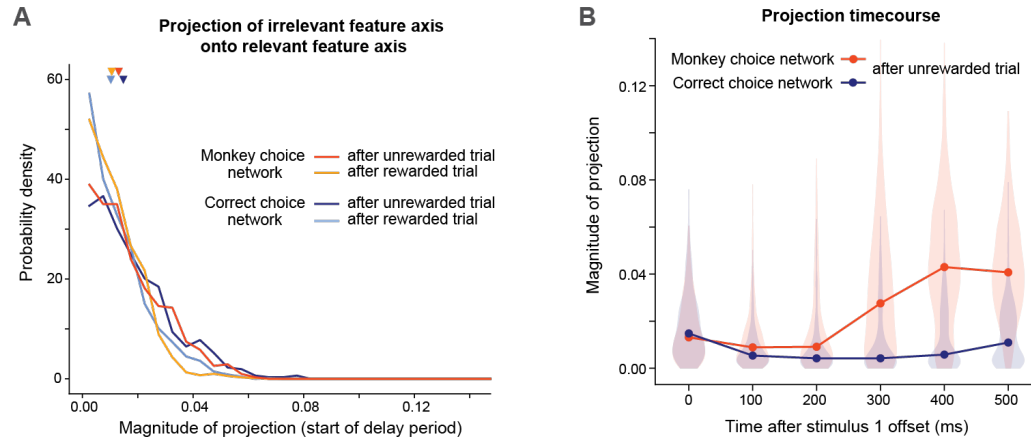
4



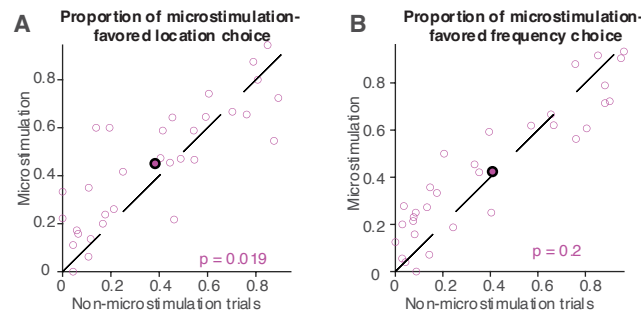
5

6 **FIG S2:** Feature (location, top row, and frequency, bottom row) decoder accuracies during a 500  
7 ms interstimulus delay period for the monkey choice network (left column) and correct choice  
8 network (right column) in four trial conditions. The conditions are as follows: 1) after a  
9 rewarded trial, location task chosen (high location task certainty); 2) after an unrewarded trial,  
10 location task chosen (low location task certainty); 3) after an unrewarded trial, frequency task  
11 chosen (low frequency task certainty); 4) after a rewarded trial, frequency task chosen (high  
12 frequency task certainty). In both models, the more a feature is believed to be irrelevant, the  
13 more quickly information about that feature decays following stimulus offset. Each point  
14 represents the mean performance of cross-validated linear decoders (same as Figure 2E), with  
15 error bars showing  $\pm 1$  standard deviation. A separate decoder was trained for each stimulus  
16 condition (identical first stimulus) and timepoint.

17



1  
 2 **FIG S3:** Feature axes become less orthogonal following stimulus offset in the monkey choice  
 3 network but not the correct choice network. Orthogonality is quantified by calculating the  
 4 projection of the irrelevant feature axis onto the relevant feature axis. (A) At the beginning of  
 5 the interstimulus delay period (0 ms after first stimulus offset), the monkey choice network  
 6 represents the two features (location and frequency) as orthogonally as the correct choice  
 7 network both after unrewarded trials ( $p=0.16$ ) and after rewarded trials ( $p=0.67$ , Wilcoxon rank  
 8 sum tests). The median projections (triangle markers) are all below 0.02. This contrasts with  
 9 Figure 2F, which shows the same histograms at the end of the delay period (400 ms after first  
 10 stimulus offset). (B) A violin plot showing the time course of projections during a 500 ms  
 11 interstimulus delay period. Points represent the medians of the projection distributions; shaded  
 12 regions show kernel density estimates. Feature axes exhibit a partial ‘collapse’ in the monkey  
 13 choice network, but not the correct choice network, such that the projection of the irrelevant  
 14 feature axis onto the relevant feature axis increases significantly.  
 15



16  
 17  
 18 **FIG S4:** Electrical microstimulation induced location choice bias predicted by the receptive  
 19 field (RF) of nearby recorded neuron (A) but did not significantly induce frequency choice bias  
 20 predicted by the frequency selectivity of nearby recorded neuron (B). We selectively  
 21 microstimulated V1 neurons during the second stimulus to observe the impact of  
 22 microstimulation on behavioural choices. Each point in the scatterplot represents a stimulus  
 23 condition. (A) In trials where the RF of the neuron found at the microstimulated site is to the  
 24 left side of the first stimulus, then we define location left-shift choice as the “microstimulation  
 25 favored location choice”; and vice versa for the location right-shift choice with RF to the right  
 26 side of the first stimulus. The scatter plot displays the behavioural effects of microstimulation

1 on location choices by plotting the proportion of the microstimulation favored location  
2 choices with microstimulation (ordinate of the points) against that without microstimulation  
3 (abscissa of the points). The dots lie significantly above the diagonal line, showing that  
4 microstimulation biased the monkey's location choices predicted by the RF location of nearby  
5 recorded cells. (B) Similar figure convention to (A), except it shows the absence of a  
6 systematic behavioural frequency choice bias predicted by the frequency selectivity of nearby  
7 recorded cells  
8