

# Flexible selection of task-relevant features through across-area population gating

July 19, 2022

## Abstract

Brains can gracefully weed out irrelevant stimuli to guide behavior. This feat is believed to rely on a progressive selection of task-relevant stimuli across the cortical hierarchy, but the specific across-area interactions enabling stimulus selection are still unclear. Here, we propose that population gating, occurring within A1 but controlled by top-down inputs from mPFC, can support across-area stimulus selection.

Examining single-unit activity recorded while rats performed an auditory context-dependent task, we found that A1 encoded relevant and irrelevant stimuli along a common dimension of its neural space. Yet, the relevant stimulus encoding was enhanced along an extra dimension. In turn, mPFC encoded only the stimulus relevant to the ongoing context. To identify candidate mechanisms for stimulus selection within A1, we reverse-engineered low-rank RNNs trained on a similar task. Our analyses predicted that two context-modulated neural populations gated their preferred stimulus in opposite contexts, which we confirmed in further analyses of A1. We finally integrated our within-area observations in a two-region RNN and propose a novel mechanism for flexible across-area communication through fixed connectivity.

## Introduction

The informational value of different stimuli can change dramatically depending on the context, but animals can adapt with impressive flexibility to virtually any contingency change. A classical example of this feat is the so-called “cocktail party effect”, which refers to our ability to focus on a specific, currently relevant conversation while ignoring all the others. Understanding how stable neural circuits implement this kind of flexible, context-dependent behavior has proven challenging. While there is a growing consensus that it emerges from the interaction between different regions along the brain hierarchy (Brincat et al., 2018; Flesch et al., 2022; Panichello and Buschman, 2021; Siegel et al., 2015), the specific interactions are unclear.

One possibility is that regions early in the hierarchy merely represent the incoming stimuli and propagate their representations downstream, where context-dependent rules are applied to effectively guide behavior (Birman and Gardner, 2019; Li et al., 2009; Sasaki and Uka, 2009; Uka et al., 2012). In line with this view, pioneering work combining artificial neural networks and neurophysiological recordings from monkeys performing a canonical context-dependent task (Mante et al., 2013), shows that both relevant and irrelevant stimuli are encoded as late as the frontal cortex, suggesting that the selection of relevant stimuli indeed occurs late in the cortical hierarchy. Empirical evidence demonstrates however that primary sensory areas are modulated by behavioral context (Hajnal et al., 2021; Maunsell and Treue, 2006; Paneri and Gregoriou,

2017; Rodgers and DeWeese, 2014; Siegel et al., 2015), potentially through feedback interactions with downstream areas that could control the selection of the relevant stimulus already at an early stage of the cortical hierarchy (Fritz et al., 2010; Winkowski and Kanold, 2013). While an attractive possibility, the specific mechanisms through which different cortical areas cooperate to select the relevant stimuli earlier in the cortex are unclear.

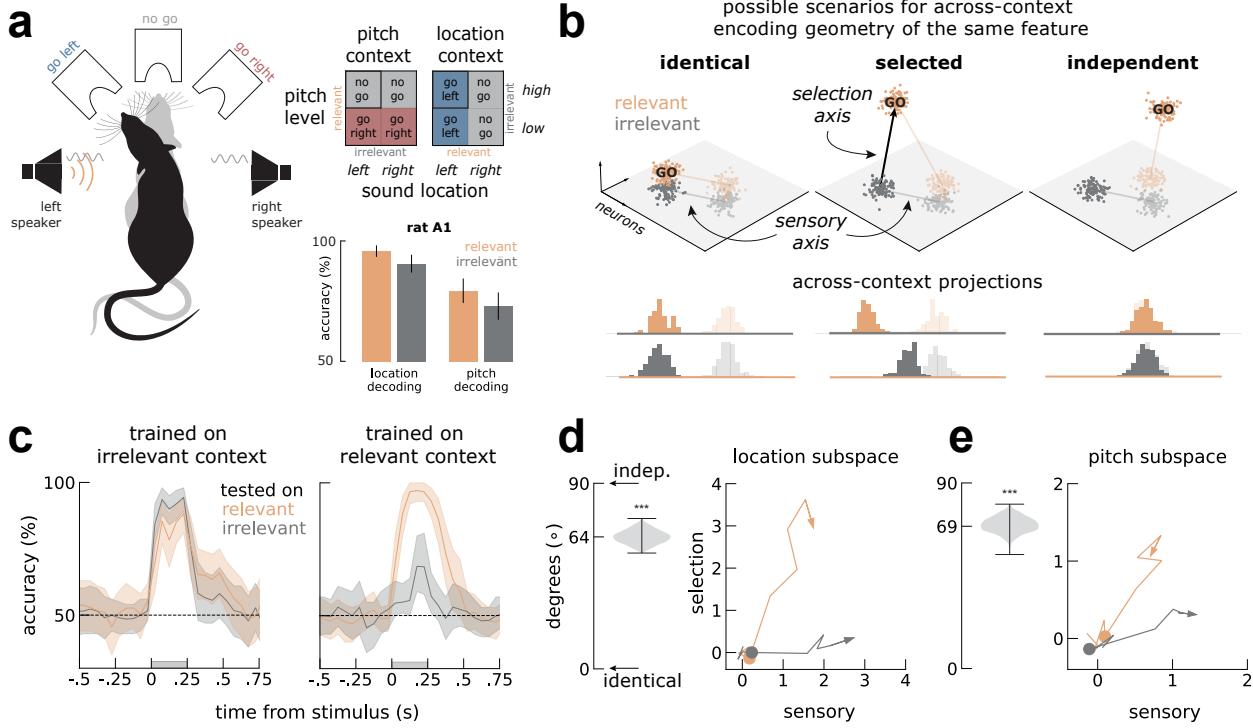
Here, we examine the population dynamics in the rat primary auditory cortex (A1) and the prelimbic region of medial prefrontal cortex (PFC), and propose a mechanism through which interactions between these two areas flexibly select relevant stimuli within A1 in a context-dependent task (Rodgers and DeWeese, 2014). We found that both relevant and irrelevant stimuli were encoded within a sensory subspace of A1, in line with other studies of humans and other animals performing context-dependent tasks (Flesch et al., 2022; Rodgers and DeWeese, 2014; Siegel et al., 2015). However, we found that the relevant stimuli were furthermore projected along an additional dimension, which we named ‘selection axis’. On the other hand, PFC encoded only the decision, fully determined by the selected stimuli. Both areas encoded context robustly throughout the trial (Rodgers and DeWeese, 2014). To investigate how this contextual information could drive stimulus selection in A1, we trained recurrent neural networks (RNN) on a similar task. Using the same analyses, we found that the geometry of the relevant and irrelevant stimuli representations resembled those of the rat’s A1. Reverse-engineering the mechanisms employed by these networks (Beiran et al., 2021; Dubreuil et al., 2022; Mastrogiovanni and Ostojic, 2018) predicted that context-modulated populations selectively gate the relevant stimuli in a context-dependent fashion, with different populations selecting specific stimuli in their preferred context. Further analyses of neural recordings revealed a similar population structure in A1, validating the model prediction and suggesting it could subserve the flexible communication of the selected stimulus with mPFC.

Integrating the neural computations observed within A1 and PFC in a two-area RNN model, we show how the two areas could solve the task by interacting through low-rank communication subspaces (Semedo et al., 2019, 2022). Despite fixed inter-area connectivity, the relevant stimulus was transmitted between A1 and PFC in a context-dependent manner. Our model suggests a novel mechanism through which areas interact flexibly along fixed subspaces: top-down contextual inputs from mPFC control communication by modulating the gain of separate populations in A1.

## Results

### Context-dependent stimulus representations in A1

To investigate how relevant stimuli are selected to guide flexible behavior, we analyzed neural activity previously collected (Rodgers and DeWeese, 2014) while rats performed a context-dependent, go/no-go auditory task (Methods). The animals were presented with an auditory stimulus (250 ms) consisting of a pitch warble from both speakers mixed with a broad-band noise lateralized to just one speaker (Fig. 1a, left). Contexts were alternated in blocks and indicated the relevant stimulus feature, i.e. pitch level (high or low) or noise location (left or right). The relevant feature (e.g. left/right and go-left/no-go, in the location task; Fig. 1a) indicated to the animal which port it had to lick to obtain a reward in each context (e.g. left/right and go-left/no-go, in the location task; Fig. 1a). Single-unit spike trains were collected either from the primary auditory cortex (A1) or medial prefrontal cortex (PFC) while the animals performed the task (Methods).



**Figure 1: Relevant and irrelevant stimuli are encoded in different subspaces in A1.**

**a)** Left, schematics of the auditory discrimination go/no-go task. Rats were presented with an auditory stimulus with two features (pitch and location). Two example trials (black vs gray rat) for the same stimuli (in orange, a noise burst on the left speaker and in gray a high pitch warble on both speakers) in different contexts (location, black; pitch, gray). Depending on the context, the animals had to attend to one of the stimulus features and respond accordingly: go left in location context or no go in the pitch context, for this stimulus pair. Right top, context-dependent go/no-go task rules specifying rewarded behavior for all stimulus pairs. Highlighted (black box) is the stimulus pair illustrated on the left. Bottom, both features are significantly decodable from A1, whether relevant (orange) or irrelevant (gray) (Rodgers and DeWeese, 2014).

**b)** Three possible scenarios for the encoding of the same feature depending on its relevance in each of the two contexts (orange and gray), as characterized by the geometric relationship of the coding axes across contexts. Different transparency levels refer to different conditions (e.g. left vs right location). Left: Identical encoding, where the coding axes are parallel in the two contexts; middle: enhanced encoding, where the go stimulus is enhanced by adding activity along a selection axis; right: independent encoding, corresponding to orthogonal coding axes. Bottom: to distinguish between scenarios, we project the trials in one context (colored histograms) onto the decoding axis (colored line) determined in the other context and inspect the resulting discrimination performance.

**c)** Across-context decoding (Across-context decoding in Methods) of location during pitch context and during location context. Left, irrelevant decoders work well both on irrelevant (gray) and relevant trials (orange). Right, relevant decoders work substantially better in relevant trials than in irrelevant trials. Shaded area marks the stimulus presentation period.

**d)** On the left, the angles between sensory and selection axis (before orthogonalization, Across-context decoding in Methods) estimated during location blocks. On the right, visualization of the activity elicited by relevant and irrelevant stimuli within the sensory-selection subspace. Colored circles mark the stimulus onset.

**e)** Same as d, but for the pitch context. Error-bars are bootstrapped 90% C.I., except in d,e) where they mark the extrema bootstrap.

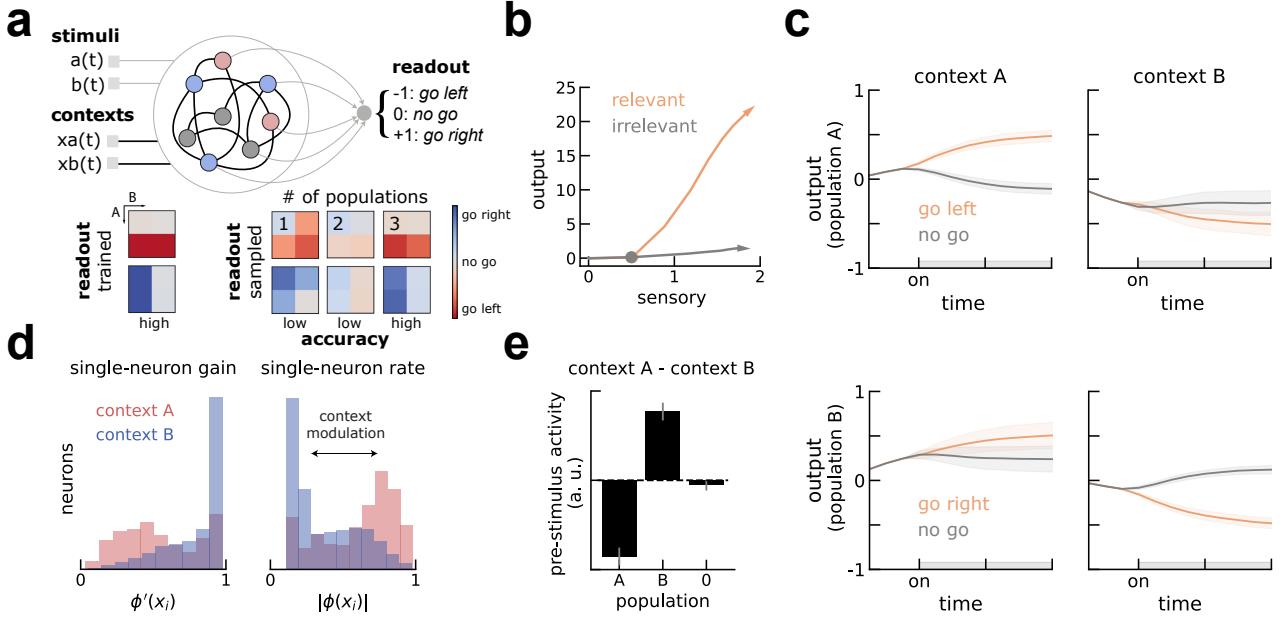
Previous decoding analyses of this dataset (Rodgers and DeWeese, 2014) showed that A1 represents context (Fig. S1a) and both stimulus features, regardless of their behavioral relevance (Fig. 1a). Here, we investigated if a given feature (pitch or localization) was encoded in the same format across contexts (Bernardi et al., 2020; Saez et al., 2015), i.e. depending on whether it was relevant (“relevant context”) or irrelevant (“irrelevant context”). In Fig. 1b, we illustrate three possible encoding scenarios in the neural activity state space, where each dimension represents a different neuron: (i) identical encoding, where the coding axes for the same feature are parallel between the relevant and irrelevant contexts; (ii) selection encoding, where the relevant go stimulus is enhanced by adding activity along a selection axis; and (iii) independent encoding, corresponding to orthogonal coding axes for the same feature across relevant and irrelevant contexts. If an auditory feature is encoded in similar formats across contexts (Fig. 1b, identical), projecting the activity collected during one context onto the decoding axis determined in the other context leads to similar separability between conditions (Fig. 1b, identical, bottom). On the other extreme, if the same feature is encoded in orthogonal formats in the two contexts, across-context projections are not separable (Fig. 1b, independent). In between these two extremes, for selection encoding (Fig. 1b, selected), the two conditions are equally separable along the decoding axes determined in the irrelevant context (Fig. 1b, sensory axis), but not along the decoder determined in the relevant context. Different codes can therefore be distinguished by their across-context decoding performance (Fig. 1, bottom).

To distinguish these possibilities, we trained stimuli-decoders on trials collected during the irrelevant context and tested their performance on trials during either context. We found that decoders trained on irrelevant trials performed well in both relevant and irrelevant contexts (Fig. 1c, left panel), evidence against an independent code and instead suggesting a sensory axis (Fig. 1b) that is shared across contexts. In contrast, the decoding accuracy of irrelevant trials was substantially reduced when tested with relevant decoders (Fig. 1c, right panel), discarding an identical code and suggesting a selection axis (Fig. 1b) along which a specific condition was enhanced in the relevant context. We quantified the angle between relevant and irrelevant decoding axes and found that they were correlated, but not parallel, as expected in a selection code (Fig. 1d,e, insets on the left). We therefore estimated the selection axis as the component of the relevant decoding axis that was orthogonal to the sensory axis (Fig. 1b). To visualize this particular encoding geometry, we then projected the trajectories of activity elicited by identical stimuli in the two contexts along the selection and sensory axes (Across-context decoding in Methods). We found that stimuli elicited activity mostly along the sensory axis when the stimuli were irrelevant (Fig. 1d,e, gray lines), but also along the selection axis when the same stimuli were presented in the relevant context (Fig. 1d,e, orange lines).

To elucidate how these context-dependent transformations could emerge in A1, we next trained a single-area RNN in a similar task.

## Single-area RNN predicts a non-random population structure

We implemented the context-dependent task using the NeuroGym toolbox (Molano-Mazon et al., 2022). Stimuli (A and B) were delivered transiently (gray bar, Fig. 2c), while the context was delivered throughout the whole trial (Methods). Aiming to replicate the stimulus selection seen in A1, we trained the network to select the relevant go stimulus along a readout vector that was fixed across contexts. To mimic our observations of mixed selectivity in A1 (Fig. S1b), the input and the readout weights on individual neurons were generated randomly (Hirokawa



**Figure 2: Trained RNN replicates A1 dynamics and predicts that population gating supports flexible selection of the relevant stimuli.** **a)** Top, schematics of the RNN. On each trial, the RNN receives 4 inputs (stimuli and contexts) and must output the correct choice (-1, 0 or +1, representing go left, no go or go right) onto a fixed readout axis. Depicted in black are the weights that are trained with backpropagation (i.e. contextual inputs and recurrent weights) and in gray those that remain fixed (i.e. stimuli input and readout). Bottom, average responses of trained (left) and resampled (right) networks separated by conditions and context (compare with schematics in Fig. 1a). Left, trained networks achieve perfect accuracy in both contexts. Right, clustering and resampling connectivity (Inferring populations in Methods) from a distribution based on an increasing number of populations shows at least 3 populations (population A, B and 0) are necessary to solve the task with comparable accuracy to trained networks (left). **b)** Similar to A1, the network represents both stimuli but enhances the relevant go stimulus along an additional axis. Colored circles mark the stimulus onset. **c)** Dynamics of activity of populations A and B projected on the output axis is reduced in opposite contexts, effectively gating the relevant go stimulus into the output axis (b). **d)** Single neurons in each population have different gains in the two contexts (here shown only for population B). This is reflected both in the slope of the transfer function (left,  $\phi'$ ) and in the single-neuron firing rate before the stimulus (right). **e)** The model predicts that the pre-stimulus firing rate can identify the 3 populations.

et al., 2019) and fixed during training (Trained A1 network in Methods).

To obtain easily interpretable RNNs, we constrained the recurrent connectivity matrix to be of low rank, allowing us to reverse-engineer the mechanisms employed by the trained networks (Beiran et al., 2021; Dubreuil et al., 2022). We found that a rank-one network was able to solve the task (Fig. 2a, bottom left), so that the connectivity matrix was defined by the outer product of two vectors, the output and input-selection vectors (Low-rank theory in Methods). After training, we froze the weights and collected the dynamics of all units during all types of trials. As we did with the biological units recorded from A1 (Fig. 1d,e), we projected the activity of the same stimuli separately when they were relevant (Fig. 2b, orange) or irrelevant (Fig. 2b, gray) onto the output and sensory axes (Trained A1 network in Methods). This confirmed that, as in A1, the network represented both stimuli along the same sensory axis, independently on whether they were relevant or irrelevant in the current context, but the relevant stimulus was enhanced along an additional axis.

We then used recently developed methods to reverse-engineer the mechanism through which the network learned to solve the task. Recent theoretical work has shown that context-dependent tasks such as the one considered here require neurons to be organized in different populations, each characterized by its joint statistics of connectivity parameters (Beiran et al., 2021; Dubreuil et al., 2022). A key empirical test of this finding is that networks generated by resampling connectivity parameters should solve the task with an accuracy similar to trained ones, as long as the statistics of connectivity within each population are preserved (Dubreuil et al., 2022). Performing this analysis (Inferring populations in Methods), we found that our trained networks relied on three populations, as resampling the connectivity vectors from the corresponding distribution led to high performance (Fig. 2a, bottom right). Due to differences in how these populations select distinct stimuli that we describe below, we labeled these 3 populations post hoc as A, B and 0.

We explored the contribution of each population to the overall dynamics leading to stimulus-selection along the output axis. With this aim, we examined separately the dynamics of the 3 populations individually by projecting their activity on the output axis in the two contexts. We observed that two out of three populations showed different stimulus-specific dynamics in the two contexts (Fig. 2c, Fig. S2). While population A selected the go stimulus A along the output axis during context A, it did not select this stimulus during context B (Fig. 2c, top), and vice versa for population and stimulus B (Fig. 2c, bottom). Note that neurons within each population were selective to both stimuli along the sensory axis, but collectively selected the go stimulus along an additional axis. We quantified this context-dependent dynamics by computing the context-dependent activation of each population (output gating, Methods) and compared it against randomly chosen populations. We found that populations A and B showed substantially more context-dependent activity along the readout axis than population 0, which reflected global dynamics (Fig. S2d). As found in a recent study (Dubreuil et al., 2022), this context-dependent modulation at the population level relied on selective gain modulation at the single-neuron level. This can be seen in Fig. 2d, in which we calculated the single-neuron gain as the slope of the transfer function (Methods) before the stimulus presentation. While neurons from population B were operating at higher levels of gain during context B, their gain was much lower during context A. In turn, population 0 did not show any gain modulation (Fig. S2d).

These analyses point to population gating through gain modulation as a candidate mechanism for solving the context-dependent task. To test whether population gating selects stim-

uli in the neural data, we sought a procedure to identify the two relevant populations from single-unit recordings. In the network model, the two populations are characterized by their connectivity and gain modulation, but this information is not directly accessible from extracellular recordings. However, since gain modulation arises from contextual inputs that shift the working point of individual neurons on their input-output function (see different context weight strengths to each population in Fig. S2c, (Dubreuil et al., 2022)), we found that gain modulation was reflected in the neuron's firing rate before stimulus onset (Fig. 2d,e). Specifically, we found that the two key populations in the model had decreased pre-stimulus firing rates in their preferred context. Therefore, the network model predicted that the single-neuron pre-stimulus firing rate would allow us to discriminate between neurons that perform the stimulus selection in the two contexts. We next tested this prediction in the A1 data.

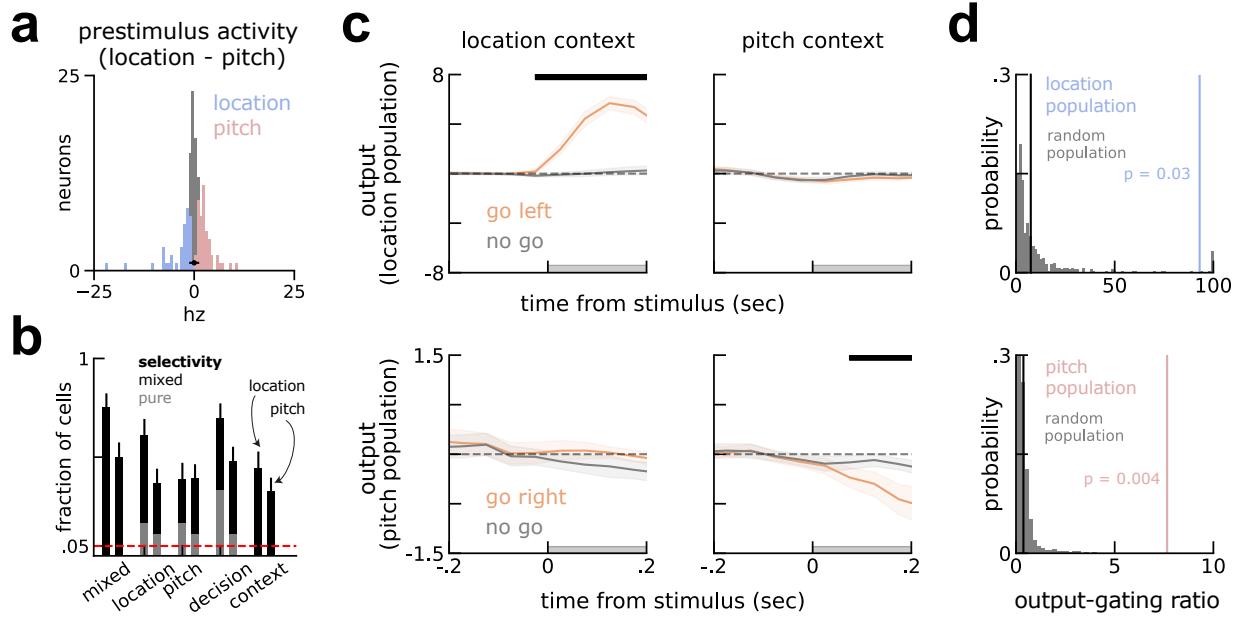
## The pre-stimulus activity of A1 neurons predicts their population structure

To test the prediction of different context-modulated neuronal populations selecting different stimuli, we grouped all the neurons recorded in A1 ( $n=130$ ) based on their context modulation during the pre-stimulus period (Fig. 3a, Mann-Whitney U test corrected for multiple comparisons). About a third of the neurons ( $n=48/130$ ) showed significantly lower spontaneous activity in the location context (henceforth location population) while another third ( $n=36/130$ ) showed decreased spontaneous activity in the pitch context (pitch population). The remaining ( $n=46/130$ ) neurons were not significantly modulated by context during the pre-stimulus (population 0). Despite being differentially modulated by context, neurons in either population had non-random mixed selectivity to the stimuli (Fig. 3b).

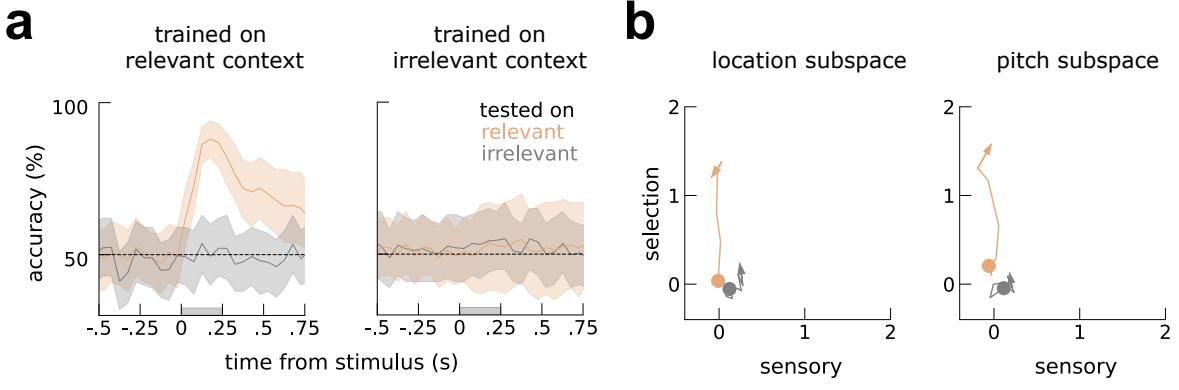
As in the model, we then inspected the projection of activity on the output axis for each of the two context-modulated populations. We estimated the output axis by decoding the two possible outputs (go left vs go right; Figure 1a and Methods) from each population and projected its activity along this axis, grouping trials by their context and output. As predicted by the model (Fig. 2), we found that the population of neurons with low spontaneous activity in a specific context gated the relevant go stimulus and ignored the irrelevant stimuli (Fig. 3c, top left, and bottom right). Conversely, in the opposite context, the output projection of the same population was essentially identical for all conditions (Fig. 3c, top right, and bottom left; see also Fig. S3). We found that the output gating of both populations, but not population 0 ( $p<0.25$ , population 0 vs shuffle;  $p<0.0025$ , 0 vs location;  $p<0.075$  population 0 vs pitch; Methods), was significantly higher than in randomly selected populations (Fig. 3d, Methods). In sum, we found that neurons grouped by their pre-stimulus context-modulation collectively select different stimuli, as was predicted by reverse-engineering the RNNs. Respectively, individual populations in A1 output the go stimuli in their preferred context but do not in the opposite context.

## mPFC encodes only the relevant stimulus along a selection axis

After characterizing a potential mechanism for the stimulus gating observed in A1, we investigated the existence of context-dependent neural dynamics within mPFC. Previous work has shown causal involvement of mPFC in action-selection during flexible behavior (Riaz et al., 2019; Rodgers and DeWeese, 2014), so we expected to see strong encoding of the stimulus relevant for the decision. As similarly done in A1, we tried decoding both relevant and ir-



**Figure 3: Pre-stimulus context-dependent activity reveals a subpopulation structure in A1, as predicted by the model.** **a)** Neurons are grouped in three populations, based on their pre-stimulus firing rate modulation to context. In black, the mean (circle) and [2.5, 97.5] percentiles (bar) of firing rate modulation to context of shuffled trials within neurons. **b)** Neurons in the pitch- and location-population (right and left slim bars, respectively) are mixed selective (Single-neuron selectivity in Methods), although a small fraction showed pure selectivity to some of the task variables (gray). **c)** Projection of stimulus responses of go left / right (orange) and no go (gray) onto the output axis for the location population (top) and pitch population (bottom). Different populations select the relevant go stimuli in different contexts. Left, projections of trials recorded during the location context. Right, equivalently for the pitch context. See Fig. S3 for the projections of individual stimuli here grouped as go left/right and no go. **d)** Permutation test shows output gating in c) is not visible in randomly picked populations (gray, Methods). Top, location-population has an output gating ratio (blue vertical bar) higher than chance ( $p=0.03$ ). Bottom, pitch-population has gain modulation higher than chance ( $p=0.004$ ). Population 0 (black vertical bar in both plots) did not show above chance gain modulation ( $p \leq 0.25$ ). All error-bars are bootstrapped standard errors of the mean.



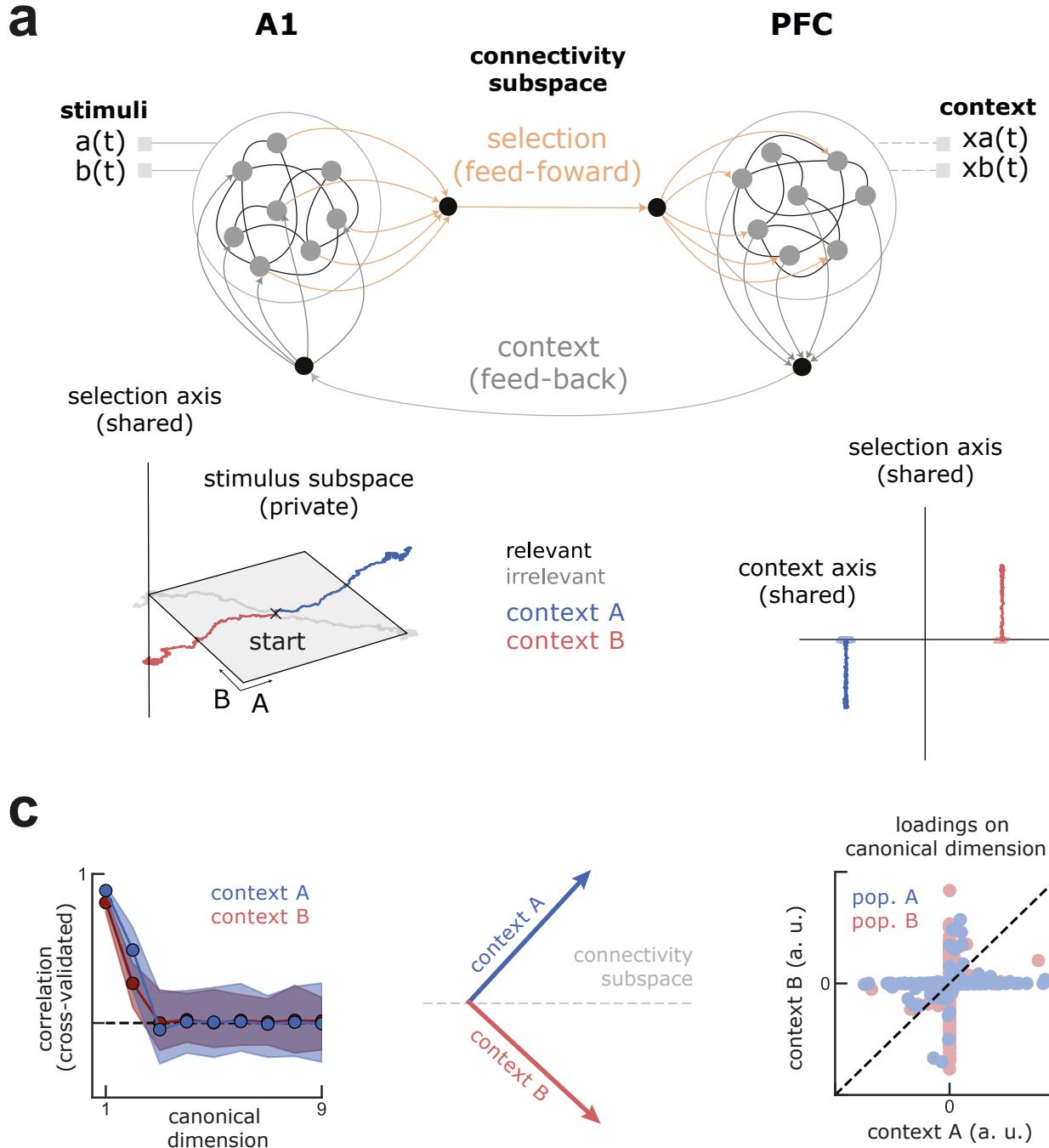
**Figure 4: PFC encodes only the selected stimuli along the selection axis.** **a)** Across-context decoding (Across-context decoding in Methods) of location during pitch context and during location context. Left, relevant decoders work well in relevant trials but not irrelevant trials. Right, irrelevant decoders fail both on irrelevant (gray) and relevant trials (orange). Shaded area marks the stimulus presentation period. Error-bars are bootstrap 90% C.I. See also Fig. 1c. **b)** visualization of the activity elicited by relevant and irrelevant stimuli within the sensory-selection subspace. Colored circles mark the stimulus onset. See also Fig. 1d,e.

relevant stimuli (Methods). As expected, mPFC indeed encoded the relevant stimulus (Fig. 4a, left orange), but in contrast to A1, it did not encode the irrelevant stimulus (Fig. 4a, gray). We also visualized mPFC context-dependent dynamics along a sensory and selection axis, estimated similarly to A1 (Fig 1). In contrast to A1, mPFC dynamics evolved exclusively along one axis encoding the decision (Figure 4b). Furthermore, separating neurons according to their pre stimulus activity did not reveal robust population gating, with all populations indistinguishable from randomly selected populations ( $p=0.1$ ,  $p=0.22$ ,  $p=0.083$  for population 0, 1 and 2, respectively; output-gating permutation test; Methods).

## Multi-area RNN replicates A1 and PFC dynamics and predicts their interaction through fixed communication subspaces

In a final step, we integrated the neural dynamics observed in the rat's A1 and PFC into a single model, a two-area RNN that performed the task (Fig. 5a, Fig. S4). The goal of this model was to identify possible mechanisms by which interactions between brain areas could lead to the context-dependent population activity observed during the task. Our model relied on the recent hypothesis that different regions communicate through low-dimensional subspaces (Kohn et al., 2020; Semedo et al., 2019, 2022). Under this view, some information within an area is transmitted to some other area through a so-called communication subspace, while the rest remains in a private subspace. To directly implement this hypothesis, we engineered a network model by starting from two low-rank networks that represented A1 and PFC, and then connected them by adding low-rank interactions between areas (A1-PFC network in Methods).

Specifically, we set the connectivity geometry of A1 similarly to the trained network (Fig. 2), meaning that when biased by contextual inputs it selected the relevant stimuli along the A1 output axis. In turn, PFC was set up to store the current context in its persistent activity. Given our empirical observations (Fig. 1,3), our hypothesis was that the sensory stimuli information remains private within A1, while the selected stimulus information is communicated



**Figure 5: Figure 5. Engineered multi-area model replicates A1 and PFC dynamics and produces predictions for across-area interactions.** **a)** The interaction between A1 and PFC was set to occur through low-rank connectivity in opposite directions (feedforward in orange and feedback in gray). In contrast to the trained network, context is delivered transiently to PFC (dashed), stored in persistent activity and fed back to A1 (context, gray). In turn, stimuli are delivered to A1 and are not communicated to PFC, thus remaining in a “private subspace” of A1 (inset, bottom left). The relevant stimulus, which is selected within A1 by integrating the stimuli and context, is communicated to PFC along the selection axis (selection, orange). **b)** Estimation of the communication subspaces using canonical correlation analyses separately for each context (Communication subspace estimation in Methods). On the left, cross-validated correlation along different canonical dimensions is significant for two dimensions (error bars are 95% C.I.). A1 and PFC communicate through orthogonal communication subspaces in opposite contexts (red and blue), despite the model having fixed connectivity (gray). On the right, different populations participate in the communication subspace in different contexts. Population A has 0 loadings on<sup>10</sup> the communication subspace during context B (here shown the first canonical dimension); conversely for population B.

to PFC. To set up this communication channel, we implemented the connectivity from A1 to PFC as a rank-one matrix  $J_{A \rightarrow P} = I^{A \rightarrow P} \otimes n^{A \rightarrow P}$ , setting  $n_{A \rightarrow P}$  to be aligned with the output axis of A1 (feedforward in Fig. 4b, orange). On the other hand, context is stored in working memory within PFC, but communicated to A1 along another dimension, also implemented by a rank-one matrix  $J_{P \rightarrow A} = I^{P \rightarrow A} \otimes n^{P \rightarrow A}$  representing the connections from PFC to A1 (feedback in Fig. 4b, gray). As in the feedforward case, we set  $n^{P \rightarrow A}$  to be aligned with the output axis of PFC. Altogether, the connectivity between A1 and PFC was defined by two axes within each area, therefore its dimensionality was  $n=2$  (Fig. 4b,d).

After setting the two-area network to interact through low-rank subspaces, we tested its performance. This was done similarly to the trained network, but now stimuli were presented only to A1 and context to PFC. At the end of each trial, we read out the final response from PFC. In isolation, A1 represents all stimuli and PFC stores the current context in persistent activity (Fig. S4a), but they do not select the relevant stimulus in a context-dependent fashion (Fig. S4a,b). When we connect the two areas as described above, contextual information is propagated from PFC to A1, targeting specific neural populations to select the relevant stimulus within A1 (Fig. S4a, bottom). The selected stimulus is then propagated downstream to PFC, from which the final response was readout, effectively solving the context dependent task (Fig. S4b).

We then developed predictions about the interaction between A1 and PFC that could be tested in neurophysiological recordings. To this end, we attempted to estimate the communication subspace as done on experimental data. Following (Semedo et al., 2019, 2022), we computed the canonical correlation dimensions (cross-validated, Methods) between A1 and PFC on the residual activity, but separately for each context. Using this approach, we estimated the dimensionality of the communication subspace to be two-dimensional in both contexts (Fig. 4d, left). However, we found that the networks communicate through orthogonal subspaces in opposite contexts (Fig. 4d, middle; Fig. S4b), demonstrating that the communication between A1 and PFC alternates flexibly along different channels in opposite contexts, despite their fixed across-area connectivity. Moreover, the switch between different communication channels is controlled by top-down inputs from PFC to A1, which selects which population participates in the communication subspace (Fig. 4d, right) thereby determine what information is selected in A1 for propagation downstream.

All together, our model replicates the dynamics within A1 and PFC and proposes a novel mechanism for flexible communication across these areas within fixed connectivity.

## Discussion

Previous studies of neural activity during context-dependent behavior have found that both relevant and irrelevant stimuli are encoded in the cortex (Aoi et al., 2020; Mante et al., 2013; Rodgers and DeWeese, 2014; Siegel et al., 2015). Here, we used across-context decoding to characterize the specific encoding geometry of these stimuli in the rat’s auditory cortex (A1). We found evidence for the selection of the relevant stimuli along an axis (“selection axis”) orthogonal to the axis encoding the stimuli (“sensory axis”). This encoding geometry, which we term selection code, is related to previous work on the encoding of abstract variables (Bernardi et al., 2020). As it happens, this encoding geometry has several advantages relative to the alternatives illustrated in Fig. 1. First, it allows for sensory information invariance along the sensory axis, even across potentially very different contexts. In view of this, a pear will look like a pear, regardless of your current appetite. Then, despite encoding similar stimuli along a

common axis, it allows for their flexible selection depending on their current relevance, generalizing the previous finding in the ferret A1 of enhancement of go stimuli upon task engagement (Bagur et al., 2018).

By reverse-engineering RNNs trained with backpropagation to employ a selection code, we postulate specific mechanisms that could support this code in A1. We found a non-random population structure in the trained RNNs, with two populations selecting different go stimuli. In our model, context-dependent gating of the relevant stimulus was accomplished through gain modulation of specific populations. The model predicted that this population structure could be inferred from pre-stimulus firing rates in electrophysiological recordings. Indeed, we found evidence for such a structure in A1, but not in PFC where the irrelevant stimulus was not encoded. Note that in contrast to our model, which was perfectly symmetric, our decoding analyses of A1 revealed that pitch-related activity was weaker than location, but this does not change the interpretation of our results (see supplementary note).

Our final contribution is to incorporate the empirical findings within A1 and PFC in a multi-area network that postulate their interactions through low-rank communication subspaces. Previous work modeling communication subspaces have focused on noise correlations in spontaneous activity and feedforward interactions (Gozel and Doiron, 2022; Thivierge and Pilzak, 2022); see also (Perich et al., 2018) for a model of ‘output-null’ subspaces in the context of motor preparation. In contrast, we now propose a multi-region network with interactions in both feedforward and feedback directions, where PFC acts as a controller of A1, dynamically selecting the appropriate communication subspace for the ongoing context. Our model complements a large body of computational work focusing on multi-area interactions (see (Perich and Rajan, 2020) for a recent review). Our major contribution is to propose a single-neuron mechanism (i.e. gain modulation) for flexible selection of different subspaces through population neural dynamics (Javadzadeh and Hofer, 2022; Panichello and Buschman, 2021; Yoo and Hayden, 2020), despite fixed connectivity.

## Biological implementation of gain-modulation

In our model, gain modulation is accomplished by selectively targeting specific units with different contextual inputs (Fig. S2c, (Dubreuil et al., 2022)), pushing individual units to the non-linear regime of their input-output function. How gain modulation is accomplished in A1 remains to be fully elucidated (Carandini and Heeger, 2011), but possible ways in which neuronal populations in A1 could have reduced gain after increased activity include synaptic non-linearities, such as depression (Carandini and Heeger, 2011), or (loose) balance between inhibitory and excitatory neurons (see (Ahmadian and Miller, 2021) for a recent review). Future work is necessary to favor one of these, or others, specific mechanisms. In some studies probing context-dependent behavior, actions are decoupled from stimuli with careful task design (Mante et al., 2013; Siegel et al., 2015), while here we tackle it only indirectly (Fig S5; see also supplementary note); notwithstanding, the mechanisms proposed here for stimulus selection through gain-modulation would straightforwardly apply to decision selection within A1.

## Division of labor and implication to the communication subspaces hypothesis

In our model, there is a division of labor: A1 represents all stimuli, but when biased by the current context it selects the relevant stimulus; in turn, PFC reads out the relevant stimulus and uses it to effectively guide behavior and could in principle infer the current context through

trial and error. Crucially, these two areas communicate the key task variables, context and the selected stimulus, through rank-one subspaces in opposite directions. Here, we have assumed that mPFC interacts directly with A1, but the control of stimulus selection in A1 could be accomplished through a third area, such as the thalamus (Jaramillo et al., 2019) or the amygdala (Jercog et al., 2021). To our knowledge, this model is the first neural implementation of the communication subspace hypothesis (Semedo et al., 2019) that performs a cognitive task (but see (Gozel and Doiron, 2022; Thivierge and Pilzak, 2022)). While future theoretical work will be necessary to fully flesh out the implications of the communication subspace hypothesis (Semedo et al., 2019, 2022), our model already reveals several interesting insights. First, it shows that communication subspaces, empirically shown to play a role during passive viewing (Semedo et al., 2019, 2022), can be exploited during flexible behavior. In the model, PFC controls A1 with contextual inputs, biasing it to select and communicate the relevant stimulus. Using canonical correlations analyses (CCA) (Semedo et al., 2019, 2022), we show that this communication occurs along orthogonal subspaces that are explored flexibly in different contexts, but within the fixed subspace set by the network connectivity. This analysis is a specific prediction that can be tested in multi-area recordings from animals performing context-dependent behavior. Second, while the subspaces estimated with CCA are aligned with those implemented in the model, we used decoding analyses as baselines to show that this estimation is imperfect, mixing feedforward and feedback communication along the same dimensions (Fig. S4b).

## Towards a multi-area perspective on context-dependent behavior

With the advent of large-scale recordings, it is becoming clear that animal behavior implicates multiple areas. In a rare tour-de-force, a recent study recorded simultaneously from six areas along the primate visual pathway, while subjects were engaged in a visual context-dependent task (Brincat et al., 2018; Siegel et al., 2015). This study shows clearly, perhaps unsurprisingly, that visual sensory information is quickly and more strongly encoded in the visual cortex (V4) than in associative areas, indicative of the feedforward flow of sensory information. On the other hand, the current context and the monkey’s decision are encoded earlier and more prominently in higher-order areas, such as PFC, consistent with the feedback flow of contextual information in our model. Both types of information are eventually encoded in all of the recorded areas, suggesting inter-area communication in feedforward and feedback directions. Interestingly, both relevant and irrelevant stimuli were decoded with comparable accuracies across the brain hierarchy, generalizing a previous finding in the monkey frontal eye field (Aoi et al., 2020; Mante et al., 2013). In contrast, we did not find encoding of irrelevant stimuli in mPFC, consistent instead with early selection of the relevant stimuli. This discrepancy might be due first and foremost to differences in animal species, but also to task differences. However, another recent study recording simultaneously from several areas across the monkey brain (V4, FEF, Parietal, and PFC), shows that visual areas (V4) encode strongly both relevant and irrelevant stimuli, but areas downstream such as FEF or PFC give clear preference to the relevant stimulus and are more predictive of the upcoming action – in line with the view taken here. Similarly, a recent MEG study of humans performing a context-dependent task, shows that decoding of irrelevant features from the dorsal premotor cortex, to which the prelimbic part of the rat mPFC is arguably reminiscent (Uylings and van Eden, 1990), is substantially lower than the decoding of relevant features. Similar findings have also been reported in human fMRI (Flesch et al., 2022). Our model proposes that different areas, with different computational roles, interact through low-rank subspaces (Semedo et al., 2019, 2022). Together with recent work on multi-area interactions (Perich and Rajan, 2020; Semedo et al., 2020) our work motivates an exciting new look on previous and future multi-region recordings (Kohn et al., 2020).

# Methods

## Animal training and electrophysiology

All procedures were approved by the Animal Care and Use Committee at the University of California, Berkeley. We are reanalyzing a previously collected dataset, so we are describing the experimental procedures here only briefly. For a complete description, we refer the reader to the original publication (Rodgers and DeWeese, 2014).

**Task.** Six rats were trained to respond to either of two simultaneously presented sounds, in a context-dependent fashion. The rats initiated each trial by holding their nose in the center port of a three-port behavior box. Each stimulus was 250 ms in duration and consisted of two different features — location and pitch. More specifically, the stimulus consisted of a noise burst played either from the left or right speaker (location feature), and a high or low pitched frequency-modulated tone (pitch feature), played from both speakers simultaneously. On each contextual block, one of the features was the relevant feature and its value determined the correct response, while the other feature was deemed irrelevant. During location blocks, the reward could be collected on the left port (go left) when the stimulus was presented on the left speaker; no reward could be collected when the stimulus was presented on the right (no go). Conversely, during pitch block, the reward could be collected on the right port (go right) when the pitch was of low frequency; no reward could be collected for high frequency (no go). Correct responses were rewarded with water, while mistakes were penalized with a 2–6 s timeout. Before contextual block changes, the rats performed 20 “cue trials”, in which the rat heard only relevant sounds without the irrelevant feature. Incorrect responses and ‘cue trials’ were excluded from all the analyses.

**Single-unit recordings.** After training, tetrodes were implanted into the rats’ brains, targeting either A1 or the prelimbic region of mPFC and single-unit spike trains were collected while the animals performed the task. Here, we only analyzed units with a sort quality defined as ‘great’ or ‘good’ in the CSV file (<https://github.com/exrodgers/Rodgers2014>) recorded during at least 10 correct trials of each kind. All analyses were performed on raw spike counts computed within windows of 50ms. See the original publication for more details (Rodgers and DeWeese, 2014).

## Single-cell analyses

**Single-neuron selectivity.** To estimate single-neuron selectivity we regressed single neuron spike counts during the stimulus presentation against a linear combination of all task variables of interest (Mante et al., 2013), namely location, pitch, decision, and context. We considered a task variable to be significantly encoded by a neuron if its regression weights were significantly different than 0, as accessed with the statsmodels python package. Neurons with only one significant weight were considered to have pure selectivity and otherwise mixed selective (Fig. 3).

**Identification of subpopulations by pre-stimulus modulation.** To test the RNN prediction laid out in Fig. 2, we averaged each neuron’s firing rate before stimulus onset separately in each context. For each neuron, we then tested for their different firing rates in the two contexts (Mann-Whitney U test), corrected for multiple comparisons (Benjamini/Hochberg). Out of  $n=130/131$  neurons in A1/PFC, some neurons had significantly lower firing rates during the location context ( $n=48/58$  in A1/PFC), others during the pitch context ( $n=36/36$ ) and were thus labeled as location and pitch population, respectively. Neurons that did not show

significant context modulation ( $n=46/37$ ) were labeled as “population 0”.

## Population analyses

**Pseudo-population decoding.** All decoding analyses were performed on ‘pseudo-trials’, pooling across all animals (Meyers et al., 2008). Specifically, we build pseudo-simultaneous populations by resampling with repetition 50 pseudo-trials from each condition and neuron. We repeated this process 500 times, leading to 500 folds across which we computed decoding variability. All decoding performances were cross-validated by splitting the training and testing dataset in two halves (50% trials for testing). Importantly, the dataset splitting was performed independently for each fold. We decoded the variable of interest – context, location or pitch – using the scikit-learn package `sklearn.linear_model.LogisticRegression`. To estimate the output axis in Fig. 3 and communication subspaces in the multi-area network (Fig. 4), we also used linear discriminant analysis (Bagur et al., 2018).

**Across-context decoding.** To investigate the stimuli encoding geometry within and across contexts, we performed across-context decoding (Bernardi et al., 2020; Saez et al., 2015). In this case, we also used pseudo-populations, but training and testing was done with datasets collected during different contexts. For instance, we trained logistic regression decoders to discriminate the location of the stimuli (left vs right) during pitch (location) blocks and then tested these decoders either on pitch blocks or on location blocks (Fig. 1c). When training and testing within the same context, we set aside 50% of trials for cross-validation. This was not necessary when this process was done across contexts, but we also subsampled 50% of the trials within each fold to avoid unfair comparisons. We repeated this process for all time points (Fig. S5) and found that selection and sensory axes were stable during stimulus presentation. For all decoding analyses we therefore used the average weights during this period.

For the visualization of activity along a decoding axis, we removed the non-linearity of logistic regression. Specifically, we collected the weights trained with logistic regression and projected the activity elicited by go and no-go stimuli on these weights. We then plotted the distance between these two conditions (without applying the logistic non-linearity). Importantly, before projecting on these weights, we orthogonalized the sensory and selection axis using QR decomposition (Mante et al., 2013).

**Output-gating.** To quantify the degree of output gating seen in A1 (Fig. 3), we calculated the following ratio during stimulus:

$$\frac{|GO_{ctx=pop} - NoGO_{ctx=pop}|}{|GO_{ctx\neq pop} - NoGO_{ctx\neq pop}|}$$

With GO (NoGO) corresponding to the average activity elicited along the decision axis for the go (no-go) stimuli and  $pop \in \{location, pitch\}$ . This value was high when a specific population was strongly modulated by context, i.e. with large activity values along the output axis for its corresponding context and low activity values in the opposite context. We also computed the same ratio for population 0 (considering either contexts as the relevant context) and for randomly labeled neurons. In the latter case, we permuted trial labels for each neuron and relabeled them based on the recomputed pre-stimulus activity with permuted trials. We then used the distribution of output gating calculated on permuted trials to evaluate a permutation

test (Fig. 3).

**Communication subspace estimation.** For the multi-area network simulations (see below), we estimated the communication subspaces using Canonical Correlation Analyses (CCA), which is a common approach for aligning neural representations (Gallego et al., 2018, 2020; Susillo et al., 2015) and more recently to study multi-area interactions (Semedo et al., 2019, 2022). Here, as done previously for studying multi-area interactions (Semedo et al., 2019, 2022), we focused on noise correlations. Specifically, we started by running 1000 trials of the go/no-go context-dependent task. We then focused on the activity during the stimulus presentation, where stimulus selection and context information was flowing feedforward and feedback, respectively. We then removed the mean activity of each neuron and stimulus conditions (Semedo et al., 2019, 2022) and computed the canonical dimensions in the following way. First, to avoid overfitting we reduced the dimensionality of the neural activity collected from both areas using PCA (scikit-learn python package, (Pedregosa et al., 2011)), keeping only the 10 dimensions with the most variability. We then used CCA (scikit-learn python package, (Pedregosa et al., 2011)) to find the canonical dimensions, along which the activity from the two areas were maximally correlated. We did this on one half of the trials and then computed the Pearson correlation with the other half of the trials and repeated this process 250 times (folds). When keeping 10 dimensions, we found that the communication subspace was 2D, as expected. However, we noticed that the number of correlated dimensions was sensitive to the number of principal components that we kept in the preprocessing step (Fig. S4c). We estimated the canonical dimensions either using data from all trials or separating by context.

**Angle between subspaces.** To quantify the alignment between the estimated communication subspaces, we computed the subspace overlap (Bondanelli et al., 2021). Specifically, we computed the arccosine of the largest singular value of  $B^T \hat{B}$ , where  $B$  and  $\hat{B}$  are the basis defined by the across-area connectivity vectors and the estimated subspace, respectively.

## Recurrent Neural Networks

**Go/no-go context-dependent decision-making task.** We implemented an abstraction of (Rodgers and DeWeese, 2014) task using the NeuroGym toolbox (Molano-Mazon et al., 2022). Briefly, the stimulus was 4-dimensional, reflecting the pitch and location feature in the rat’s experiment, in this case called A and B, and the two contexts, context A and context B. During stimulus presentation, we added gaussian noise with  $\sigma = 1$  on top of the stimuli mean levels. The stimuli features had two levels (-1,1) as well as the contextual inputs (0,1). Before the stimulus presentation, which lasted 10 timesteps, there was a pre-stimulus period of 4 timesteps. Contextual inputs were delivered during both periods, in contrast to the features that were delivered exclusively during stimulus presentation. As was the case for the rats, the network had to select the relevant stimuli level and ignore the irrelevant stimuli, depending on the context level (A=-1, B=+1 in context A, and A=+1, B=0 in context B).

**Continuous-time RNN.** The dynamics of each unit  $i$  were determined by the sum of recurrent weights  $J_{ij}$  and feedforward inputs weights  $I_i^l$ :

$$\tau \dot{x}_i(t) = -x_i(t) + \frac{1}{N} \sum_j J_{ij} \phi(x_j(t)) + \sum_l^{N_{input}} u^l(t) I_i^l + \eta_i(t)$$

With  $\phi = \tanh$ . The time constant  $\tau = 100ms$  was the same for all neurons  $i$ . For simulation and training, the equation was solved using Euler’s method with a time step  $\Delta t = 20ms$ .

The independent white noise term  $\eta_i$  was simulated by drawing at each time step from a gaussian with mean 0 and standard deviation 0.05. To calculate the gain of each neuron  $\phi'(x_i)$ , we passed each neuron activity through  $\phi'(x_i) = 1 - \tanh(x_i)^2$ .

**Trained A1 network.** For the A1 network in isolation, the connectivity matrix  $J$  was constrained to be low-rank during training. We found empirically during training that a rank-one network could solve the task, meaning that  $J_{ij} = m_i n_j$ . The network received  $N_{input} = 4$ ,  $(u^A, u^B, u^{xA}, u^{xB})$ , corresponding to stimulus  $A, B$  and context  $A, B$ . The input vectors  $I^i$  defined the sensory axes. We trained the networks using backpropagation through time to minimize the following mean squared error loss function during the last timestep of each trial t:

$$L = \sum_t (z_t - X^T w)^2$$

Where  $z_t$  is the correct response on trial  $t$ ,  $X$  the network activity during the last timestep and  $w$  the readout vector. Only the contextual inputs ( $I^{xA}, I^{xB}$ ) and recurrent weights ( $m, n$ ) were optimized during 64000 trials (in batches of 160 trials each). Optimization was carried out using Adam (Kingma and Ba, 2014) in pytorch (Paszke et al., 2017) with the decay rates of the first and second moments of 0.9 and 0.999, and learning rate of 0.001.

**Low-rank theory.** We found empirically that rank-one connectivity (i.e.  $J = m \otimes n$ ) was enough to solve the task of interest. It can therefore be shown that the network activity is constrained to be at most  $(1 + N_{inputs})$ :

$$x(t) = \kappa(t)m + \sum_i^{N_{input}} v(t)I_i$$

Where  $v(t)$  is the low-pass filter version of the input  $u(t)$  (Dubreuil et al., 2022). In this setting,  $n$  can be seen as the input selection vector and  $m$  as the output vector of a single latent variable  $\kappa$ . Previous theoretical work has shown that computations performed by low-rank networks, including those trained through back-propagation, are fully determined by the rank of their connectivity matrix and the geometric relationship of its connectivity vectors —  $m, n, I^A, I^B, I^{xA}, I^{xB}$  for the case of the trained A1 network. This relationship is characterized by the overlaps between the different connectivity vectors (e.g.  $\sigma_{mn} = m^T n$ ), which can be subdivided into an arbitrary number of subpopulations (Dubreuil et al., 2022). While the rank determines the number of latent variables  $\kappa$  that can be manipulated, the number of populations  $\mathcal{P}$  constrains the possible computations on the latent variables. For a given rank-one network, with  $\mathcal{P}$  populations, it can be shown that in the limit of  $N \rightarrow \infty$  the dynamics of the latent variables  $\kappa$  is described by:

$$\dot{\kappa} = -\kappa + \sum_p^{\mathcal{P}} \left[ \tilde{\sigma}_{nm}^{(p)} \kappa + \sum_i^{N_{input}} \tilde{\sigma}_{nI_i}^{(p)} v_i \right]$$

With  $\tilde{\sigma}_{mn}^{(p)} = \sigma_{mn}^{(p)} \langle \phi' \rangle^{(p)}$ , which can be seen as the functional connectivity – i.e. a function of the effective connectivity  $\sigma_{mn}^{(p)}$  and the population average gain  $\langle \phi' \rangle$ . Different populations can have different functional connectivity, depending on their average gain, which is itself a recurrent function of  $x$  and the active inputs (Beiran et al., 2021; Dubreuil et al., 2022).

**Inferring populations.** To infer the minimal number of populations necessary to solve the task, we followed a previously proposed approach (Dubreuil et al., 2022). Briefly, we used the method BayesianGaussianMixture from the scikit-learn python package (Pedregosa et al.,

2011) to cluster neurons in an increasing number of independent populations. After clustering, we calculated the empirical means and covariance matrices of each cluster (i.e. population) independently. We then sampled new connectivity vectors from multivariate gaussian distributions defined by these mean and covariance matrices and concatenated across populations. Finally, we evaluated the performance of networks with the sampled connectivity.

**One solution for context-dependent, go/no-go tasks.** We found out that rank-one connectivity with 3 populations solves our task. Neurons in all populations are selective to all external stimuli, but they differ in which stimulus is integrated into the latent variable (or output axis). While the first 2 populations select 1 stimulus (i.e.  $\sigma_{nI_A}^{(1)} > 0, \sigma_{nI_B}^{(2)} > 0$ ) and should ignore the other (i.e.  $\sigma_{nI_B}^{(1)} = 0, \sigma_{nI_A}^{(2)} = 0$ ), the third population must have negative feedback ( $\sigma_{nm}^{(3)} < 0$ ), which we found out to be essential to implement the no-go condition ( $\kappa = 0$ ) (Beiran et al., 2021). See Fig. S2e,f for the dynamics of this model.

**A1-PFC network.** In contrast to the A1 RNN, in which we trained the connectivity vectors, we directly engineered the A1-PFC RNN. Specifically, we adapted the low-rank framework to describe across-area dynamics. To model multi-area interactions, we represent the connectivity matrix  $J$  in terms of a block structure:

Recurrent connectivity  $A$  and  $P$  populates the diagonal and feedforward  $J_{A \rightarrow P}$  and feedback  $J_{A \rightarrow P}$  the off diagonals. Our key assumption is that each block has a rank-one structure, and is thus defined by the outer product of two connectivity vectors (e.g.  $A = m_A \otimes n_A$ ). Under this constraint, we can separate the recurrent, feedforward, and feedback inputs in a compact form for the dynamics within A1 and PFC:

$$\begin{aligned}\dot{x}_{A_i} &= -x_i + \frac{m_i^A}{N_A} \sum_{j \in A} n_j^A \phi(x_j) + \frac{I_i^{P \rightarrow A}}{N_P} \sum_{k \in P} n_k^{P \rightarrow A} \phi(x_k) + \sum_l^{N_{input}^A} u_l^A(t) I_i^l. \\ \dot{x}_{P_i} &= -x_i + \frac{m_i^P}{N_P} \sum_{j \in P} n_j^P \phi(x_j) + \frac{I_i^{A \rightarrow P}}{N_A} \sum_{k \in A} n_k^{A \rightarrow P} \phi(x_k) + \sum_l^{N_{input}^P} u_l^P(t) I_i^l\end{aligned}$$

For the purpose of this study, we assume that context information was delivered only to PFC. Moreover, to ensure a rank-one feedback communication subspace from PFC to A1, we assume A1 receives an additional constant input  $I^k$ , the details of which are described in the next section. In turn, stimuli  $(u^A, u^B)$  are delivered exclusively to A1. Thus,  $N_{input}^A = 2$  and  $N_{input}^P = 1$ . Under these assumptions, in the limit of  $N \rightarrow \infty$  and assuming again  $\mathcal{P}$  populations, the dynamics of the high-dimensional A1-PFC network can be reduced to the dynamics of the following latent variables in A1:

$$\begin{aligned}\dot{\kappa}_A &= -\kappa_A + \sum_p^{\mathcal{P}} \left[ \tilde{\sigma}_{n^A m^A}^{(p)} \kappa_A + \sum_i^{N_{input}} \tilde{\sigma}_{n^A I^i}^{(p)} v^i + \tilde{\sigma}_{n^A I^P \rightarrow A}^{(p)} v_{P \rightarrow A}^{(p)} \right] + \tilde{\sigma}_n^{A I^k} \\ \dot{v}_{P \rightarrow A} &= -v_{P \rightarrow A} + \sum_p^{\mathcal{P}} \tilde{\sigma}_{m^P n^P \rightarrow A}^{(p)} \kappa_P\end{aligned}$$

And in PFC:

$$\begin{aligned}\dot{\kappa}_P &= -\kappa_P + \sum_p^{\mathcal{P}} \left[ \tilde{\sigma}_{n^P m^P}^{(p)} \kappa_P + \tilde{\sigma}_{n^P I^x}^{(p)} v_x + \tilde{\sigma}_{n^P I^A \rightarrow P}^{(p)} v_{A \rightarrow P} \right] \\ \dot{v}_{A \rightarrow P} &= -v_{A \rightarrow P} + \sum_p^{\mathcal{P}} \tilde{\sigma}_{m^A n^A \rightarrow P}^{(p)} \kappa_A\end{aligned}$$

In addition to the internal latent variables  $\kappa_A, \kappa_P$ , there are now two extra latent variables corresponding to the communication subspace  $v_{A \rightarrow P}, v_{P \rightarrow A}$ . The key new elements in this formulation are the overlaps (e.g.  $\sigma_{n^{A \rightarrow P} m}$ ) between the output vectors within an area, such as  $m$  in  $A1$ , and the corresponding input selection vectors populating the off diagonals of  $J$ , such as  $n^{A \rightarrow P}$ . Reminiscent of the case for within area dynamics (Beiran et al.; Dubreuil et al., 2022; Mastrogiuseppe and Ostojic, 2018), non-negative overlap leads to the communication of the corresponding latent variables. Generally speaking, there is one of these overlaps for each within-area variable. In the simplified case addressed here, we set all the overlaps to zero, except for those related to the within-area latent variables. For simplicity, we set  $n^{A \rightarrow P} = m$  and  $n^{P \rightarrow A} = m^P$ , ensuring both across-area overlaps were non-zero. The overlaps within  $A1$  were set similar to the trained RNN, and when modulated by context, this network integrated the relevant stimuli along the output axis  $m$  with the same mechanism as the trained network. In turn, the geometry of PFC was set so it integrated its inputs (i.e. context) into one of the two fixed points. Specifically, we set  $\sigma_{mn} > 0$  and  $\sigma_{nI} > 0$ .

To encode context through a one-dimensional communication subspace in the feedback direction, we introduced a bias term. This bias term, which was fixed and present in all trials and timepoints was defined as  $I_k = \frac{I_{ctxA} + I_{ctxB}}{2}$ . This way the context value (-1 or 1) was projected along one dimension, which we conveniently defined as  $I_x = \frac{I_{ctxA} - I_{ctxB}}{2}$ , making the net input  $I_{ctxA}$  when context = 1 and  $I_{ctxB}$  when context = -1.

## Data code availability

Data is currently available at <https://crcns.org/data-sets/pfc/pfc-1> and the code necessary to replicate all figures will be available at <https://jmourabarbosa.github.io/publications/> upon publication.

## Author contributions

J.B., R.P, S.O. and Y.B analyzed the data. J.B. and S.O. conceived and performed the modeling research. C.R. conceived and performed the experiments. J.B., S.O. and Y.B. wrote the manuscript. All authors critically revised and edited the manuscript.

## Competing interests

The authors declare no competing interests.

## Acknowledgements

We thank Heike Stein, Manuel Molano and Ramon Nogueira for feedback on the manuscript. This work was supported by FrontCog grant ANR-17-EURE-0017, ANR-JCJC-DynaMiC and the NIH Brain Initiative project U01-NS122123 (SO, JB). JB was supported by the Fyssen Foundation.

## Supplementary Notes

### Discarding motor execution confounds

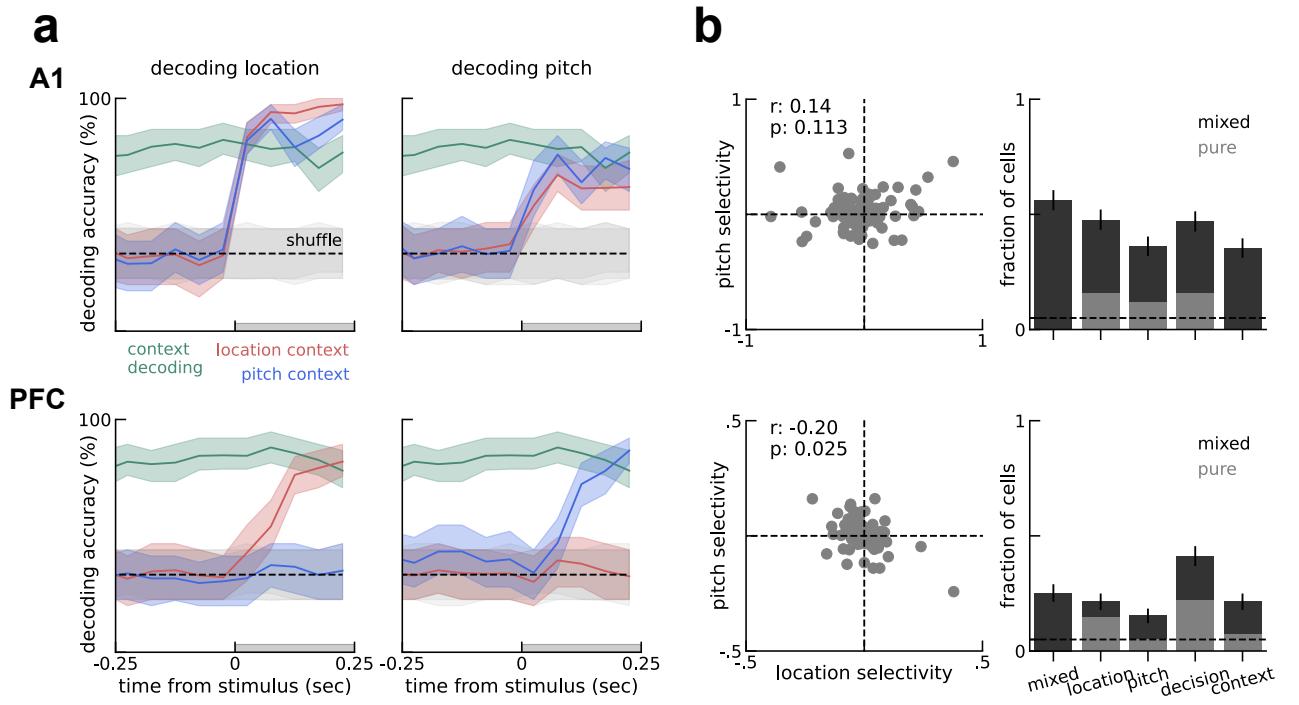
At the level of the task design, action (go left / go right) and the relevant go stimuli (left / low pitch) were not dissociable, which could in principle confound the abstract representation of the relevant stimuli with activity related to motor execution. To investigate this potential

confound, we cross-correlated the decoding weights estimated during different time points across the trial (Fig. S5). These analyses revealed that relevant and irrelevant decoding dimensions were constrained to stimulus presentation and that a third dimension was explored after the stimulus offset. This axis was orthogonal to the aforementioned stimuli axis (Fig. S5), likely encoding the rat's ongoing motor execution. To avoid confounding motor execution with the encoding of relevant stimuli, this period was ignored from all analyses and only the sensory and selection axes — orthogonal to the motor axis — were interpreted. In contrast, PFC did not encode the irrelevant stimulus and the decoding of relevant stimulus was well aligned with decision decoding (Fig. S5), suggesting that the decoding of relevant stimulus from PFC was in fact capturing motor variability. Importantly, in contrast to A1, selecting PFC neurons based on their pre-stimulus modulation to context did not reveal a population structure with output-gating (Fig. 3), indicating further that grouping A1 neurons by their pre-stimulus firing rate did not select motor related neurons. In a recent experiment (Yin et al., 2020) that decoupled sounds from motor output (i.e. go vs no-go), Yin and colleagues found that sound encoding emerged earlier in the ferret A1 (25 ms) than in PFC (50-100 ms). On the other hand, motor output related information emerged earlier in PFC (50 ms) and feedback information appeared in A1 around 600ms (their Fig. 4), one order of magnitude later than what we here interpret as stimulus encoding (1-2 bins,  $\pm$ 25-50 ms; Fig. 1,3) but in line with the aforementioned motor-related variability in A1 that we discarded from our analyses (Fig. S5). Ultimately, only task designs decoupling decision and stimulus will be able to determine if decision is already computed in A1 before reaching PFC; notwithstanding, the mechanisms proposed here for stimulus selection would straightforwardly apply to decision selection within A1.

### **Slight asymmetry between contexts / location more strongly encoded than pitch**

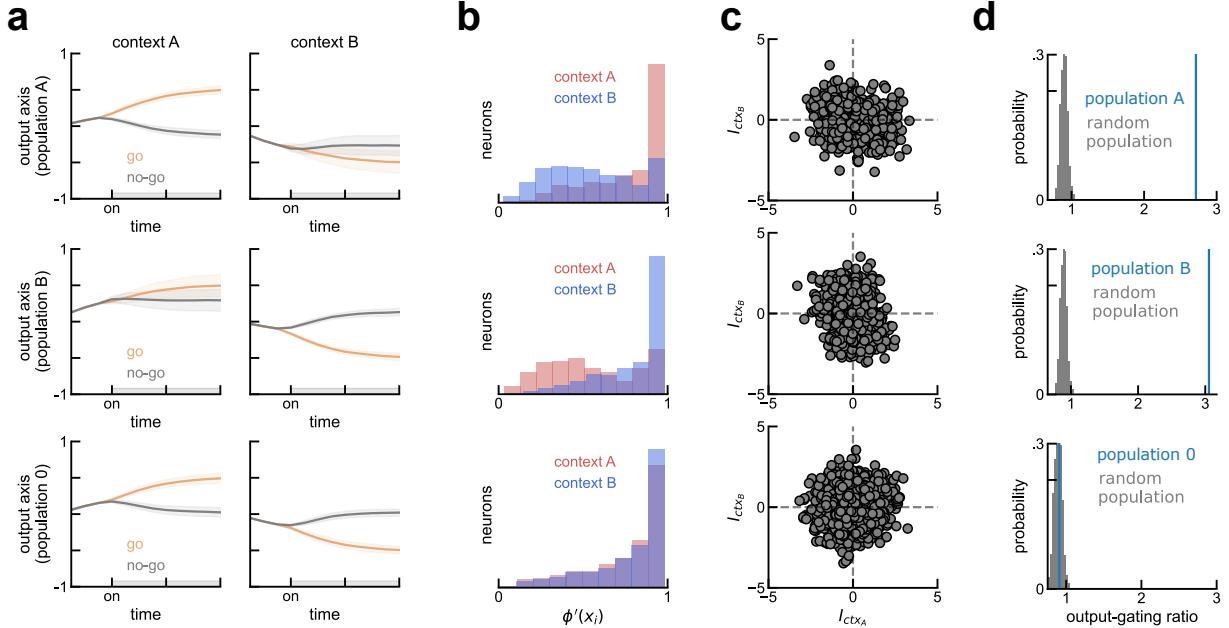
We found stronger encoding of location than of pitch in A1 (Fig. 1 & 3), but we are reluctant to interpret this as a general finding. Instead, we speculate that this slight asymmetry between the two contexts was due to all recordings being performed in the left brain hemisphere. Crucially, the noise bursts indicating "nogo" came from the right, contralateral to the recordings and expected to cause stronger responses than stimuli presented on the ipsilateral hemifield. In addition, the animals showed lower performance in pitch blocks (Rodgers and DeWeese, 2014). Together, these asymmetries in the task may explain the differences in effect size between the two contexts. Future experiments with perfectly symmetric tasks and/or bilateral recordings are necessary to validate this possibility, but they will not change the qualitative interpretation of our results.

## **Supplementary Figures**



**Figure 6: A1 and PFC encode different task variables.** **a)** Logistic regression decoding (Methods) of location (left), pitch (right) in each context (red and blue) and decoding of context overlaid on both plots for comparison. Top, A1 encodes both stimulus' features in either contexts. Bottom, PFC encodes only the relevant stimulus' features for the ongoing context. **b)** Left, feature-selectivity is mixed in both areas. Right, fraction of cells with significant task-variables regressors (Methods). Error-bars are bootstrapped SEM.

## Network model (trained)



## Low-dimensional Model

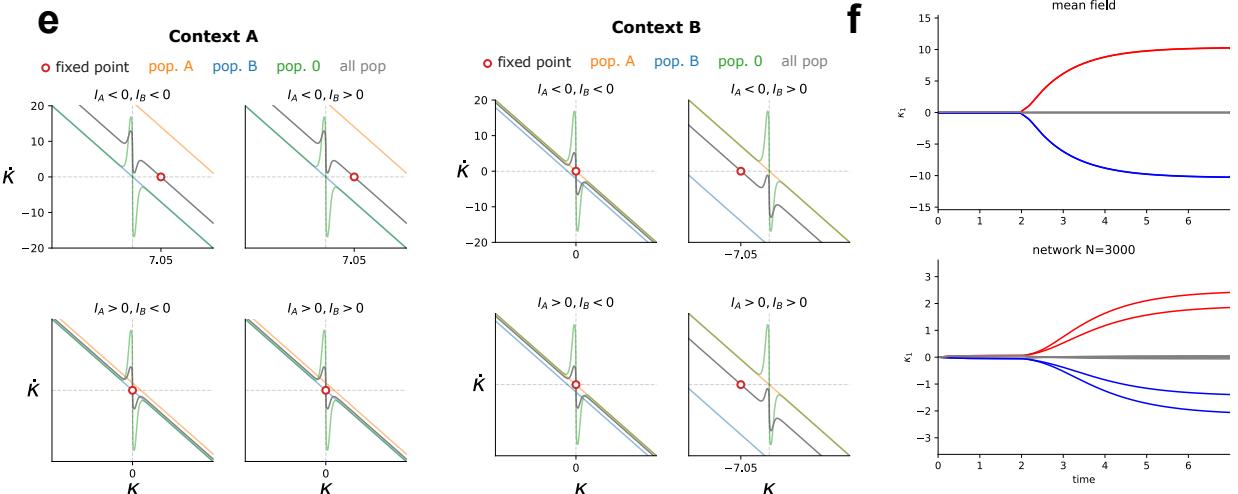


Figure 7: Network and low-dimensional model of a context-dependent go/no-go task. **a,b)** Same as Fig. 2c and d, but showing data from all three populations. **c)** Context weights optimized through training have different variances for each population, the mechanism supporting gain modulation (Dubreuil et al., 2022). Population A and B have a larger range of weights for context B and A, respectively. Population 0 has the same range of weights for both contextual inputs. **d)** Populations A and B showed substantially more context-dependent activity along the readout axis than population 0 as measured with output-gating ratio (Methods). **e)** Dynamics of kappa separated for each population (color lines) and collectively for all populations (gray) for both contexts and all stimulus combinations. Here it can be seen that population 0 (green) contributes equally for all contexts and stimulus conditions. Namely, it pushes the dynamics of kappa towards 0, essential to have a fixed corresponding to the no go conditions (two bottom conditions for context A, on the left; and two left conditions for context B, on the right). In contrast, population A (orange) and B (blue) are inactive during context B and A, respectively, and do not contribute to the dynamics during those conditions. On the other hand, the same populations are active in opposite contexts and integrate the relevant stimulus into kappa dynamics (e.g. orange lines in context A when input A  $\neq 0$ ). **f)** Top, dynamics of the low-dimensional model for all trials. Bottom, dynamics a network with weights sampled from the distribution defined in f (Methods). Quantitative differences due to

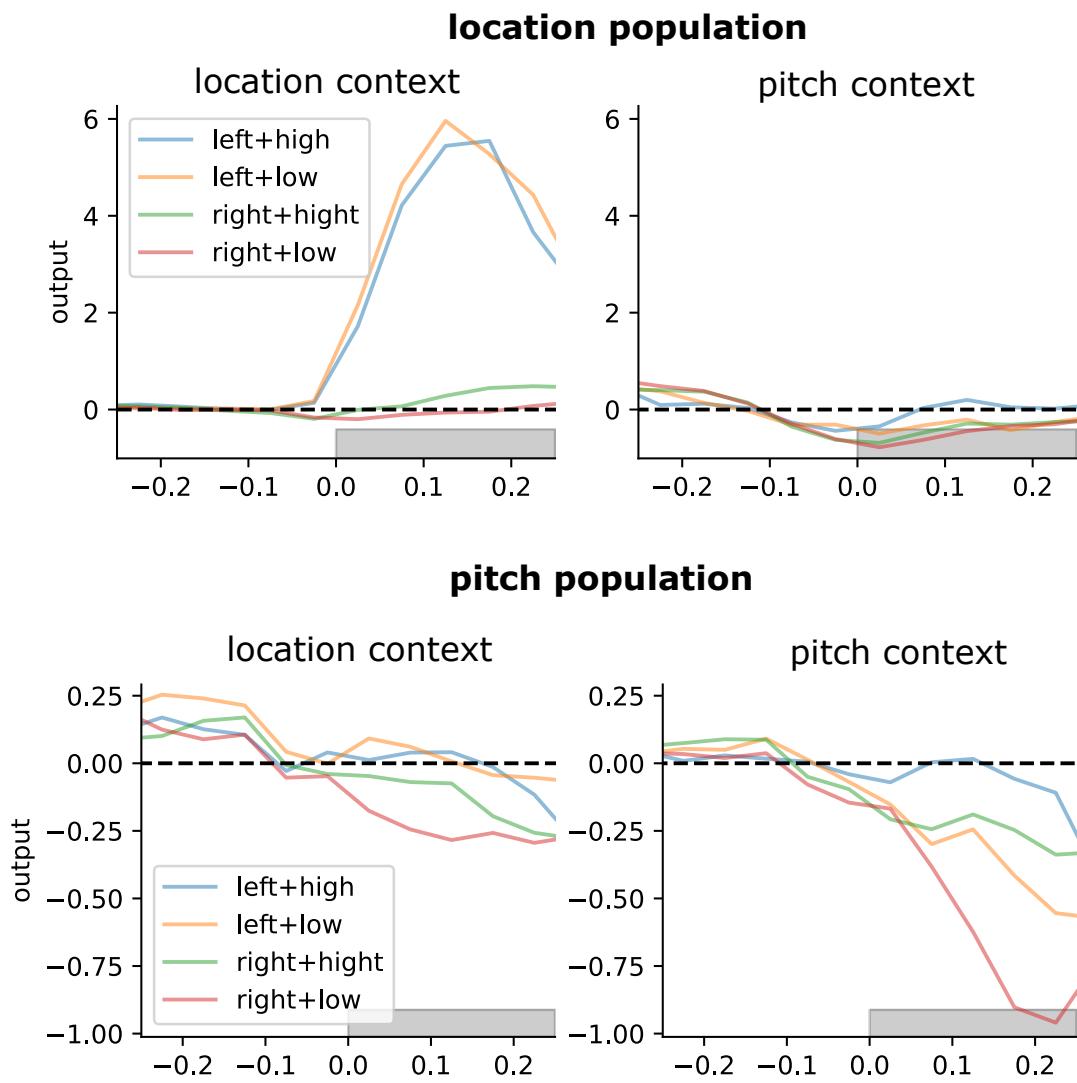


Figure 8: Same as Fig. 3, but separating by stimuli. Go stimuli for the location context (left) in orange and blue. Go stimuli for the pitch context (low) in red and orange.

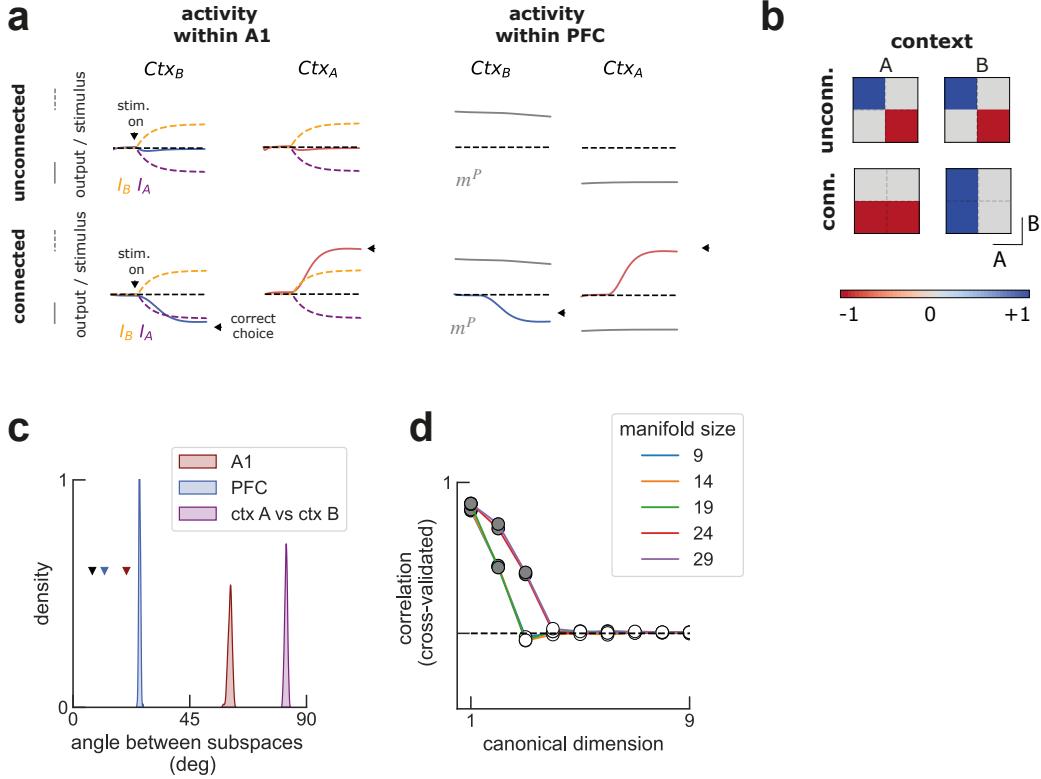


Figure 9: **a)** A1 and PFC in isolation did not integrate the relevant stimulus in a context-dependent fashion (unconnected), but they did when set up to cooperate through communication subspaces (connected). We found that connecting these two areas drove A1 to integrate into the recurrent dynamics the relevant, but ignore the irrelevant stimulus. Moreover, meaningful choices could be read out from PFC instead of A1 (black triangles). In the figure we illustrate two trials. Specifically, the projection of the network activity on different connectivity vectors:  $I_A, I_B$  in purple and yellow, respectively;  $m_A, m_p$  in red/blue and gray, respectively; and on the input-selector vector from A1 to PFC (Methods), red/blue. **b)** When unconnected, the two areas do not show context-dependent behavior, but they do so when set to interact through low-rank connectivity. **c)** principal angle between different subspaces (Methods). Communication subspaces inferred during opposite contexts are almost orthogonal (purple). In red and blue, the angle between the subspaces estimated with canonical correlation analysis and those defined by the network connectivity. For comparison, colored triangles mark the angle between connectivity subspaces and those determined by decoding context and decision from each area; black triangle marks the mean angle between the same subspaces estimated from different folds (Methods), the minimum empirical distance possible between subspaces. **d)** As described in the methods, we only used the first PCs of the neural to estimate the communication subspace using CCA. This is a typical preprocessing step (e.g. (Gallego et al., 2018)). We found empirically that the estimated subspace is sensible to the number of PCs we kept (manifold size). For the purpose of Fig. 4, we only kept the first 9 PCs.

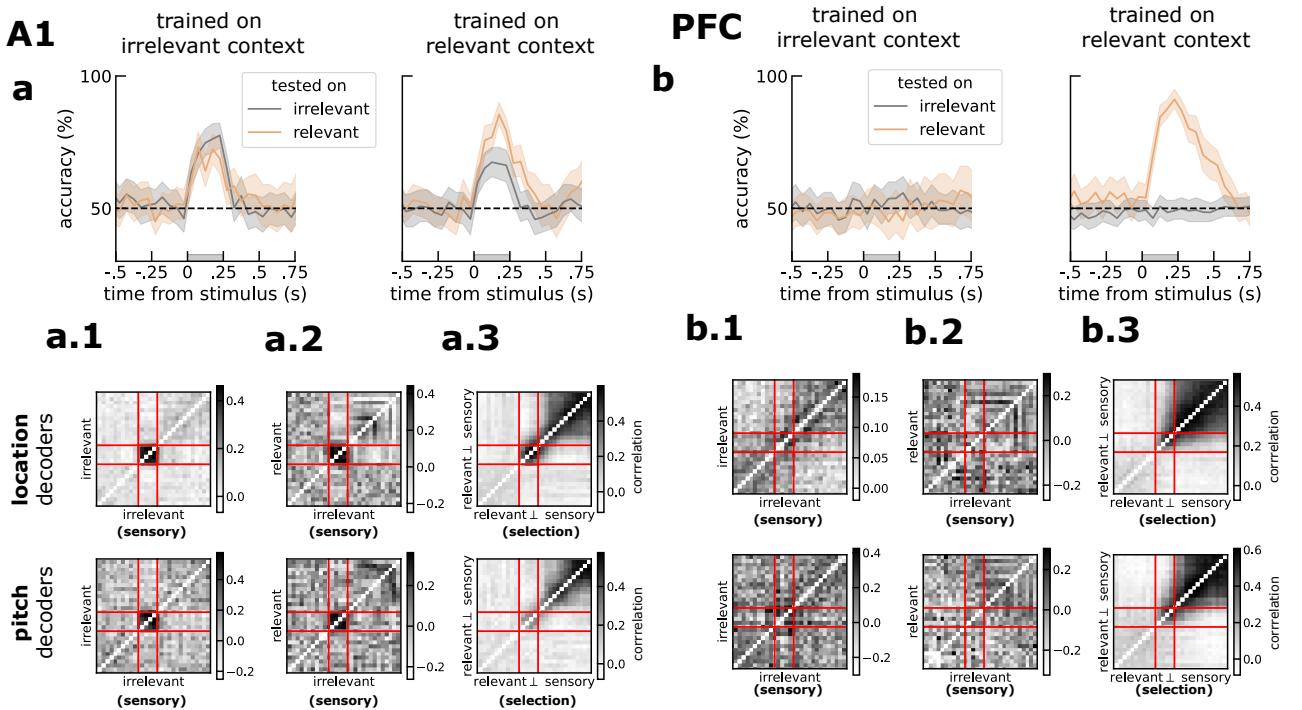


Figure 10: Top, same as Fig. 1 but for pitch instead of location. Error-bars are bootstrap SEM. Bottom, cross-correlation of decoding weights for different axes. In a.3, it can be seen that the relevant stimuli code has two components, one during the stimulus and another one during responses. On the other hand, for PFC (b.3), the code seems to unfold along similar codes during stimulus and response.