

Reviewed Preprint

v1 • December 12, 2025

Not revised

Reviewed Preprint

v2 • May 11, 2026

Revised by authors

✉ For correspondence:

giancarlo.lacamera@stonybrook.edualfredo.fontanini@stonybrook.edu

† These authors contributed equally to this work

Competing interests: No competing interests declared

Funding: See [page 38](#)

Reviewing editor: Thorsten Kahnt, National Institute on Drug Abuse Intramural Research Program, United States

© 2025, Lang et al. This article is distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use and redistribution provided that the original author and source are credited.

Linear and categorical coding units in the mouse gustatory cortex drive population dynamics and behavior in taste decision-making

Liam Lang^{1,2,4}, Camelia Yuejiao Zheng^{1,2,3,4,†}, Jennifer M Blackwell^{1,4,†}, Giancarlo La Camera^{1,2,4}✉, Alfredo Fontanini^{1,2,3,4}✉

¹Department of Neurobiology and Behavior, Stony Brook University, Stony Brook, United States • ²Graduate Program in Neuroscience, Stony Brook University, Stony Brook, United States • ³Medical Scientist Training Program, Stony Brook University, Stony Brook, United States • ⁴Center for Neural Circuit Dynamics, Stony Brook University, Stony Brook, United States

eLife Assessment

This **important** work advances our understanding of the single neuron coding types in the mouse gustatory cortex and the functional roles of these neurons for perceptual decision-making. The conclusions are based on **compelling** evidence from rigorous behavioral experiments, high-density electrophysiology, sophisticated data analysis, and neural network modeling with in silico perturbations of functionally-identified units. This work will be of broad interest to systems neuroscientists.

<https://doi.org/10.7554/eLife.109313.2.sa4>

Abstract

Cortical circuits produce time-varying patterns of population and single neuron activity that play a fundamental role in perceptual and behavioral processes. However, the functional contributions of individual neuron activity to population dynamics and behavior remain unclear. Here we addressed this issue focusing on the mouse gustatory cortex (GC) and using a taste mixture-based decision-making task, high-density electrophysiology, and computational modeling. GC population dynamics represented stimuli linearly during taste sampling and choices categorically before decisions. Single neurons were classified by their linear and categorical activity patterns, revealing subpopulations encoding sensory, perceptual, and decisional variables. To test their functional role, we built a recurrent neural network model of GC. Model perturbations showed linear and categorical neurons were essential for driving normal population dynamics and behavioral performance, whereas many units with other activity patterns could be silenced without consequence. These results have implications that extend beyond GC, and demonstrate the role of linear and categorical coding neurons in cortical dynamics and behavior during perceptual decision-making.

Introduction

Cortical circuits produce time-varying patterns of population and single neuron activity that play a fundamental role in perceptual and behavioral processes (Shuler and Bear, 2006 [↗](#); Guo et al., 2014b [↗](#); Buetfering et al., 2022 [↗](#)). Over the past decade, the gustatory cortex (GC) has emerged as a model for studying such cortical dynamics (Arieli et al., 2022 [↗](#); Livneh et al., 2020 [↗](#); Mahmood et al., 2025 [↗](#); Vincis et al., 2020 [↗](#); Vincis and Fontanini, 2016 [↗](#)). Population and single neuron activity in GC show characteristic time-varying modulations encoding multiple variables associated with a gustatory experience. Early studies on GC dynamics provide evidence for three

temporal epochs following intraoral delivery of tastants in rats. Right after taste delivery, GC neurons encode the somatosensory component of the stimulus hitting the tongue. After a few hundred milliseconds, neurons begin to encode the chemical identity of the tastant and, ultimately, its palatability (Katz et al., 2001). A similar sequence of activity has been described also in the GC of mice actively sampling tastants by licking a spout (Bouaichi and Vincis, 2020). These dynamics are not limited to the processing of taste; GC can also represent signals related to expectation (Mazzucato et al., 2019; Livneh and Andermann, 2021) and decision-making (Vincis et al., 2020; Lang et al., 2023).

Recent work relying on delayed response decision-making paradigms shows that populations as well as single neurons in GC can sequentially encode sensory, perceptual, and decisional variables. Head restrained mice were trained in a two-alternative choice (2AC) task (Guo et al., 2014a; Churchland and Ditterich, 2012) to sample gustatory stimuli, wait during a delay period, and lick either a left or a right spout based on specific taste-direction associations. In the context of this task, GC activity progresses from taste-coding during the sampling period to representing the licking direction predicted by each taste during the delay period (Vincis et al., 2020; Lang et al., 2023). Consistent results were observed in a similar task relying on taste mixtures to guide directional licking decisions. Two-photon calcium imaging has revealed that GC activity first encodes mixture components linearly and then, during the delay period, represents the binary decision to lick left or lick right (Kogan and Fontanini, 2024).

These dynamics are not epiphenomenal, as perturbations of neural activity at specific times in the trial interfere with ingestive behaviors as well as taste-guided decision-making. For instance, silencing GC during the palatability epoch delays the onset of aversive reactions to taste (Mukherjee et al., 2019). In a 2AC task, silencing GC during the delay period affects task performance (Vincis et al., 2020). The importance of these dynamics has therefore spurred extensive investigations of their underlying mechanisms. Spiking network models have unveiled architectural features that are sufficient for producing population dynamics matching those seen in GC during taste-processing, expectation, and decision-making (Miller and Katz, 2010; Mazzucato et al., 2015, 2016, 2019; Lang et al., 2023). Furthermore, these studies have provided important information on how perturbing activity during different temporal epochs can impact behavior (Lang et al., 2023). Yet, as important as this work has been in advancing our understanding of cortical dynamics, many questions remain unanswered. The relationship among individual neuron activity patterns, population-level dynamics, and behavioral outcomes has yet to be fully elucidated. Most critically, we still do not understand how neurons encoding specific task features contribute to population dynamics and influence overall performance.

Previous computational studies were not designed to investigate the role of specific single neuron firing patterns as both network connectivity and the contribution of functional groups of neurons were established *a priori* and tuned to obtain the desired dynamics and behavior. In this study we overcome the limitations of previous approaches by combining high-density behavioral electrophysiology (Jun et al., 2017) with recurrent neural network (RNN) modeling (Cohen et al., 2020; Valente et al., 2022). Specifically, we recorded ensembles of GC neurons in mice performing a taste mixture directional task analogous to the task used by Kogan and Fontanini (2024). Population analyses revealed a progression from linear encoding of taste concentrations to categorical prospective encoding of directional licking decisions. Guided by population dynamics, we identified single units that either encoded taste concentration linearly, or task events categorically. A subpopulation of neurons linearly tracked the concentration of the components in the mixture, while others categorically encoded the prevailing taste quality or the imminent licking direction. Linear coding of the stimulus was more predominant early on, following taste sampling, while categorical coding of licking direction was more pronounced later in the trial, before lateral licking. Categorical coding of taste quality was present, albeit in small percentage, throughout the period from sampling to lateral licking. To test the functional significance of single neuron firing patterns, we trained an RNN on both the recorded neural activity and the behavioral performance of the mice for each session. A fraction of the network's units were trained to match the firing activity of recorded single units, while the rest of the units

were free from constraints during training. All the single units recorded in an individual session were fed to the network to ensure that the RNN would not be biased toward linear or categorical patterns or any experimenter-selected feature. After training, the RNNs exhibited linear and categorical patterns in both their constrained and unconstrained units. We hypothesized that these units with specific activity patterns, though relatively small compared to the overall population, were critical for population activity and behavior. We tested this hypothesis by re-running the networks after systematically removing these units that had emerged during training. The simulations revealed that linear coding neurons as well as categorical neurons representing both perceptual and decisional variables were indeed necessary for driving realistic population dynamics and behavioral performance, while a large fraction of the units, exhibiting different patterns of activity, could be silenced without consequence for performance.

Altogether, the results presented here demonstrate the functional significance of specific single neuron firing patterns observed in GC during a taste mixture 2AC task and establish an approach that successfully relates population-level dynamics, individual neuron activity patterns, and behavioral outcomes. Furthermore, since these dynamics are not unique to GC, but have been observed in a variety of cortical circuits (Goltstein et al., 2021 [↗](#); Niessing and Friedrich, 2010 [↗](#); Reinert et al., 2021 [↗](#)), our work provides insights that generalize far beyond taste.

Results

Behavioral task and electrophysiological recordings

Thirteen mice were tested on a binary sucrose/NaCl mixture-based perceptual decision-making task (Figure 1A [↗](#)). The task design is based on a two-alternative choice (2AC) paradigm (Guo et al., 2014a [↗](#); Churchland and Ditterich, 2012 [↗](#); Kogan and Fontanini, 2024 [↗](#)). Briefly, mice were trained to sample either sucrose (100 mM) or NaCl (100 mM) from a central spout (for ~0.9 s), wait for a delay period (~3.0 s), then lick one of two lateral spouts according to the task policy (e.g., sucrose → lick left; NaCl → lick right). Taste-side pairings were counterbalanced across subjects. Correct lateral licks were rewarded with a drop of water from the lateral spout. The delay period—important for temporally separating the sensory and cognitive processing during this task—was implemented such that the average time between the first lick to the central spout (time point “T”) and the first lick to a lateral spout (time point “D”) was ~3.9 s. Upon reaching criterion for discriminating sucrose from NaCl (i.e., 85% correct performance for three days in a row), mice were tested with the following mixtures (expressed in %Sucrose/%NaCl): 0/100, 25/75, 35/65, 45/55, 55/45, 65/35, 75/25, or 100/0. If the mixture was predominantly NaCl, mice had to lick on the side that was trained to be associated with NaCl; if it was predominantly sucrose, mice had to lick on the sucrose side.

Mice were tested on this task up to 3 times each (average ~1.8 sessions per animal), resulting in a total of 23 sessions. Mice performed an average of 137 trials per test session (range: 65 to 195) with an average overall accuracy of $77.2 \pm 4.3\%$ (range: 69.7% to 86.2%). As expected, performance was near chance level for difficult discriminations (average $56.9 \pm 8.9\%$ over sessions for 45/55 and 55/45 trials) and well above it for easy discriminations (average $93.1 \pm 3.5\%$ over sessions for 0/100 and 100/0 trials). The psychometric curve for session-averaged performance is plotted in Figure 1B [↗](#). There was a slight asymmetry in the curve (i.e., $P(\text{Sucrose choice} \mid \text{Stimulus} = 50/50) > 0.5$), which could be due to mismatches in perceived intensities of mixture components or to lateral motor biases, but given its small magnitude, we did not investigate it further. Just before each testing session, high-density Neuropixels probes were acutely inserted in the GC so that single unit electrophysiological recordings could be obtained during performance of the test (Figure 1C-E [↗](#)). Electrode positioning was reconstructed histologically upon the termination of the experiment. A Python-based GUI for Histological E-data Registration in Brain Space was used to register the slice images onto the Allen CCF mouse atlas (Fuglstad et al., 2023 [↗](#)). Figure S1 [↗](#) shows the reconstructed positioning of the probes. All channels were sorted, but only those mapped within GC were used for analysis. Simultaneously recorded GC ensembles had an average size of ~27 neurons (range: 7 to 68), for a total of 626 neurons across all sessions.

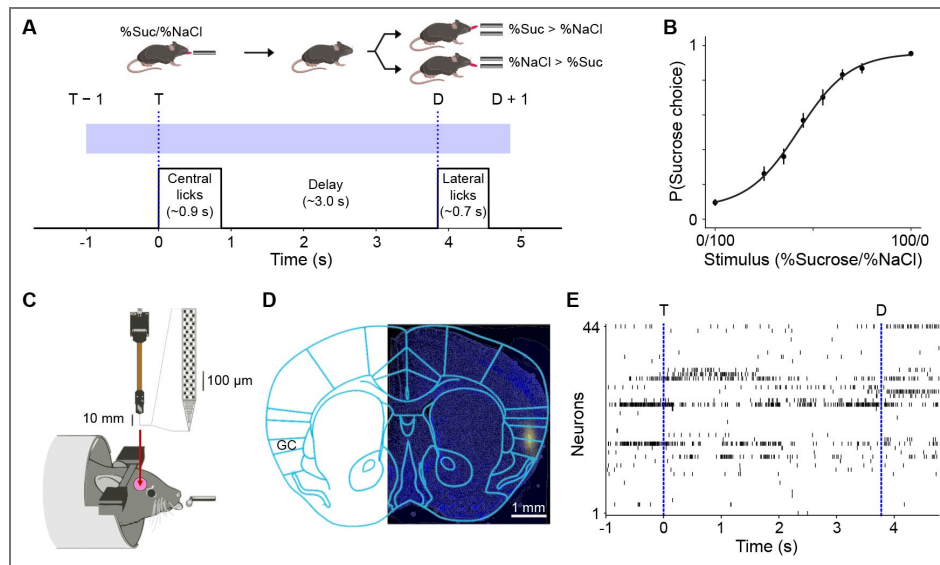


Figure 1. Experimental paradigm.

A: Schematic of the behavioral task. The blue rectangle indicates the temporal window over which all analyses were conducted, from 1 s before the first central lick, T, to 1 s after the first lateral lick, D. **B:** Psychometric curve averaged across sessions and subjects. Circles and error bars represent mean \pm s.e.m. ($N = 23$ sessions across 13 subjects) for the probability of a choice in the sucrose-associated direction for each stimulus value; the continuous curve is a sigmoidal curve fit to the means. **C:** Schematic of acute probe insertion. **D:** Example histological section indicating accurate probe placement in GC (blue: Hoechst; yellow: DiI applied to probe). **E:** Example spike raster plot for neurons simultaneously recorded within a single session from GC of behaving mice.

Population dynamics during mixture-based decision-making

To assess the involvement of GC in representing different events within a taste mixture 2AC task, we analyzed task-related, population firing activity. We used a “warped” time scale to allow for consistent analysis of dynamics with respect to multiple behavioral events of interest that have trial-to-trial variability. Each neuron’s peristimulus time histograms (PSTHs) were constructed by aligning spike trains to the first central lick T, warping each trial’s duration between first central and lateral lick to the overall average (~3.9 s), calculating firing rates in ~50 ms bins, averaging over correct trials of each mixture separately, and smoothing with Gaussian kernels. Visual inspection of the population PSTH averaged across all recorded neurons and for all stimuli (Figure 2A) revealed increases in firing rates aligned with the central and lateral licks. Consistent with this, we found 57.5% (370/626) of neurons responded to taste and 43.9% (275/626) displayed preparatory activity before lateral licks (see **Methods: Responsivity and selectivity analyses**; Table S1 provides these counts across sessions). To identify when populations of GC neurons discriminate between trial types (correct predominantly-sucrose vs correct predominantly-NaCl trials), we computed auROCs and measured differences in firing rate distributions between the two types of trials (Figure 2B). Neurons could maximally discriminate between correct sucrose and NaCl trial types at all time points (white dots), with noticeable concentrations of differential activity around time points T and D (white traces overlaid on the heatmaps are neuron-averaged auROCs). When separating neurons based on whether their peak differential activity was in favor of predominantly-sucrose or predominantly-NaCl mixtures, we found no significant difference between the number of neurons that “preferred” one to the other (306 vs 320; 2-tailed binomial test, $p = 0.548$). Though there were slight qualitative differences between mean auROC traces, they had similar peaks during sampling (0.062 for NaCl-preferring, 0.052 for sucrose-preferring between T and D) and post-decision (0.077 for NaCl-preferring, 0.075 for sucrose-preferring after D). Moreover, the distribution of peak au-ROCs was not different between NaCl- and sucrose-preferring neurons (rank-sum test, $p = 0.789$).

To assess how such task-related firing encoded task variables (e.g., stimuli and decisions), we trained classifiers to decode stimulus and choice variables from activity vectors for each ensemble of simultaneously recorded neurons (23 ensembles with a median size of 27 neurons). For each recording session, all trials were labeled according to the delivered stimulus (8 possible labels depending on the mixture) and whether the animal chose the left or right lateral spout. At each time point, a decoder then predicted a label for each trial as the label whose trial-averaged activity vector was nearest (in terms of Euclidean distance) to the activity vector of the trial in question. Session-averaged decoding for both stimulus and choice was significantly above chance throughout the period between taste delivery and reward; however, the dynamics differed significantly (2-way within-subjects ANOVA with factors time [sampling/delay] and decoded variable [stimulus/choice]; interaction $p < 0.001$). Stimulus decoding peaked at 31.0% ~425 ms after stimulus delivery (Figure 2C), while choice decoding also rose with stimulus delivery but ramped before the lateral licks and peaked at 81.9% ~75 ms after the decision time (Figure 2D).

To further assess response dynamics, population activity trajectories across different trial types were analyzed with both unsupervised and supervised dimensionality reduction techniques. For the unsupervised method, we calculated pairwise Euclidean distances in 626-dimensional neural space between all stimulus types and time points (Figure 3A). We then used a t-distributed stochastic neighbor embedding (t-SNE) to non-linearly map the neural data into a 2-dimensional space while attempting to preserve the true distance structure (Figure 3B). Notably, activity trajectories diverged at the time of taste delivery according to stimulus type, with more distant stimuli (e.g., 0/100 vs 100/0) corresponding to more distant neural activities and more similar stimuli (e.g., 45/55 vs 55/45) corresponding to closer neural activities. The trajectories then binarized according to stimulus choice, with a large distance between trials predicting different licking directions regardless of exact stimuli. To separate the components associated with the stimulus from those associated with the choice, we employed a demixed principal component analysis (dPCA; Kobak et al., 2016) and found the dimensions of maximal data variance with respect to stimulus- and choice-specific variance. The projection of the population activity onto the

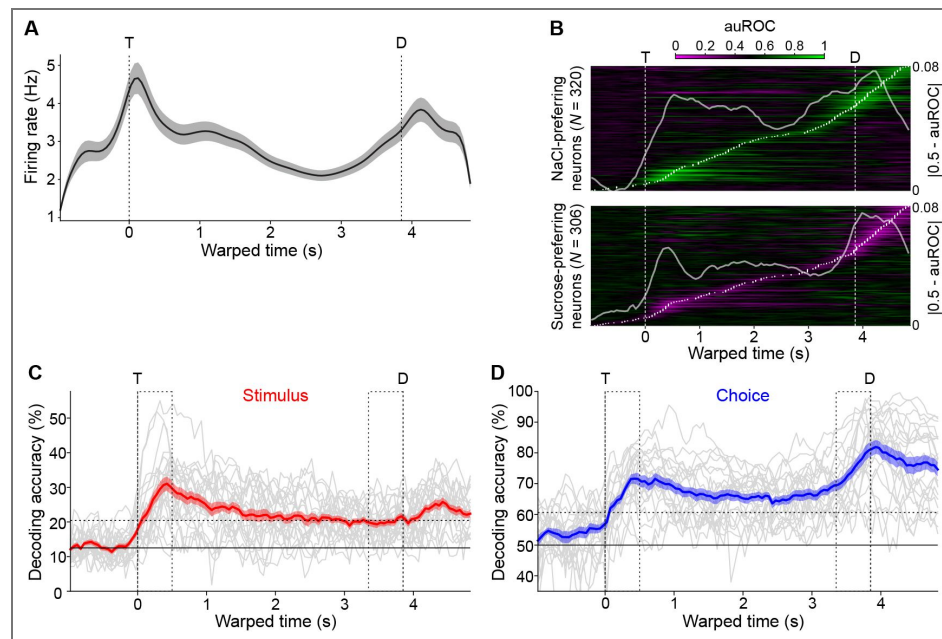


Figure 2. Population activity and information encoding during taste mixture-based decision-making.

A: Population PSTH ($N = 626$). Vertical dashed lines indicate the first central (T) and lateral licks (D), respectively. The trace represents mean firing rate; shading represents s.e.m. **B:** Population heatmaps for single unit differential activity between correct predominantly-sucrose and correct predominantly-NaCl trials. White dots indicate each unit's time of peak differential activity. Traces are ordered by peak time and separated by whether peak differential activity is in favor of NaCl ("NaCl-preferring") or sucrose ("Sucrose-preferring"). White trace is the mean auROC across neurons. **C-D:** Decoding of task-relevant variables. For each session (grey trace), accuracy is plotted over time with colored shaded traces representing the mean \pm s.e.m. over sessions. Trial labels to be decoded were mixtures (**C**) and choice (**D**). Horizontal solid line represents theoretical chance level. Horizontal dashed line represents theoretical significant decoding threshold ($\alpha = 0.01$). Dashed rectangles mark the "sampling" and "delay" analysis windows.

principal stimulus component (Figure 3C) shows a monotonic, linear separation of activity according to mixture, irrespective of the ultimate choice (note how each average error trial trajectory, the dotted lines, overlaps with its average correct trial trajectory, the solid lines). In contrast, the projection of population activity onto the principal choice component (Figure 3D) shows a binary separation of activity according to upcoming choice irrespective of mixture (note how average error trial trajectories for predominantly-sucrose mixtures overlap with average correct trial trajectories for predominantly-NaCl mixtures, and vice versa). The time courses of the projected activities suggest that a binarization of activity according to selected choice emerges just before lateral licking as the graded stimulus-based activity slowly collapses.

In summary, population analyses show a linear representation of taste mixture information toward the beginning of the trial and a categorical representation of decisions toward the end.

Single unit responses during mixture-based decision-making

Upon observing that population activity can represent task variables either in a linear or binary fashion, we investigated whether such representations might also be found at the level of single units. A response profile analysis was performed. To extract a response profile (essentially a tuning curve) a neuron's firing rate in a specified time window was averaged over correct trials of a particular stimulus and plotted as a function of the mixture. Two temporally separated windows of interest, [T, T + 500 ms] and [D – 500 ms, D], were chosen and, for each, the single unit response profiles for all neurons were constructed. Profiles were assigned a label, “linear” or “step”, after a least-squares regression statistically determined the shape that fit the profile best (if neither fit was significant, it was assigned the “other” label; see **Methods: Response profiles**). The shape templates were chosen based on previous published work (Kogan and Fontanini, 2024; Maier and Katz, 2013), with linear fits representing response profiles that track the concentration of one of the components in the mixture and step fits representing response profiles that change abruptly at the 50/50 mixture. Figure 4A shows examples of linear (left) and step (middle, right) single unit response profiles (bottom), as well as these neurons' corresponding PSTHs (top).

Step coding neurons could represent either a perceptual category (i.e., sweet vs salty) independently of concentration and licking direction, or the imminent licking direction independently of the stimulus. To disambiguate these options we analyzed error trials—that is, trials in which a mouse received a mixture and licked toward the wrong direction. If a step coding neuron had an average firing rate that was consistently elevated for trials where the same choice was taken, regardless of stimulus, the neuron was defined as a “step-choice” neuron (Figure 4A, right). On the contrary, if a step coding neuron had an average firing rate that was elevated for trials where the same stimulus was presented, regardless of the choice, the step coding neuron was considered a “step-perception” neuron (Figure 4A, middle).

Comparing the first 500 ms of the sampling period to the last 500 ms of the delay period, there was a similar proportion of neurons whose response profile could be significantly fit by either linear or step functions (6.5%, 41/626 for the first 500 ms vs 9.4%, 59/626 for the last 500 ms, Chi-squared $p = 0.061$) (Figure 4B). It is worth noting that neurons with no significant fits could still show significant responses to task events. For instance, in the first 500 ms 57.1% (334/585) were taste responsive and 10.1% (59/585) were taste selective (i.e., they responded differently to the mixtures), while in the last 500 ms 42.0% (238/567) showed preparatory responses in anticipation of lateral licks (63 of which were selective for a specific direction).

Within the sub-groups of neurons with significant fits there was a change in the distribution of response profile types over time: there was a significant increase in the proportion of step coding neurons from the beginning to the end of the trial (24.4%, 10/41 vs 62.7%, 37/59, Chi-squared $p < 0.001$); equivalently, there was a significant decrease in the proportion of linear coding neurons (75.6%, 31/41 vs 37.3%, 22/59, Chi-squared $p < 0.001$). The increase in step coding neurons from the beginning to the end of the trial was driven not by changes in the amount of step-perception neurons, which remained stable (24.4%, 10/41 vs 37.3%, 12/59, Chi-squared $p = 0.631$), but by an increase in the amount of step-choice neurons (0/41 vs 42.4%, 25/59, Chi-squared $p < 0.001$) (Figure 4B). Furthermore, the change in distribution over time seemed to be driven mostly by separate

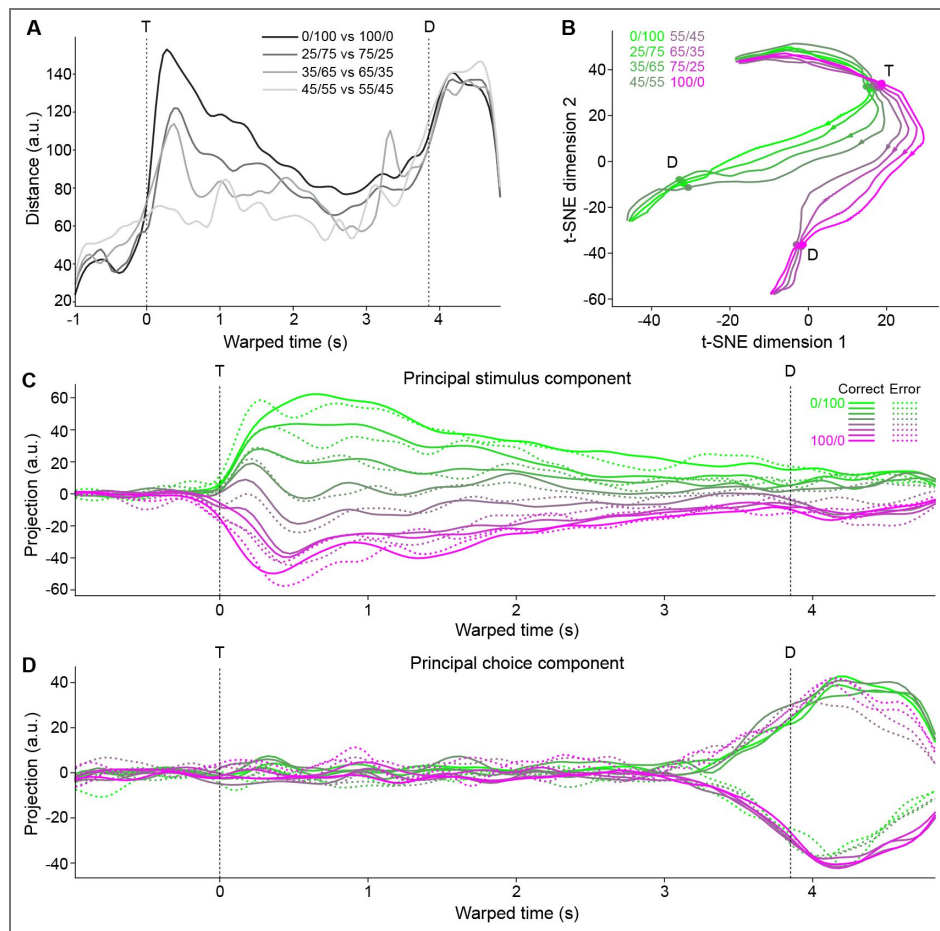


Figure 3. Low-dimensional population activity trajectories.

A: Euclidean distances between pairs of trial-averaged pseudo-population activity trajectories. **B:** t-SNE of trial-averaged pseudo-population trajectories for all stimuli (%Sucrose/%NaCl) based on pairwise Euclidean distances between activities. **C:** One-dimensional linear projections of trial-averaged pseudo-population trajectories onto the demixed principal component explaining maximum stimulus-specific variance. Solid lines are correct trial averages; dotted lines are incorrect trial averages. **D:** Same as **C**, but for the demixed principal component explaining maximum choice-specific variance. T: time of first central lick; D: time of first lateral lick.

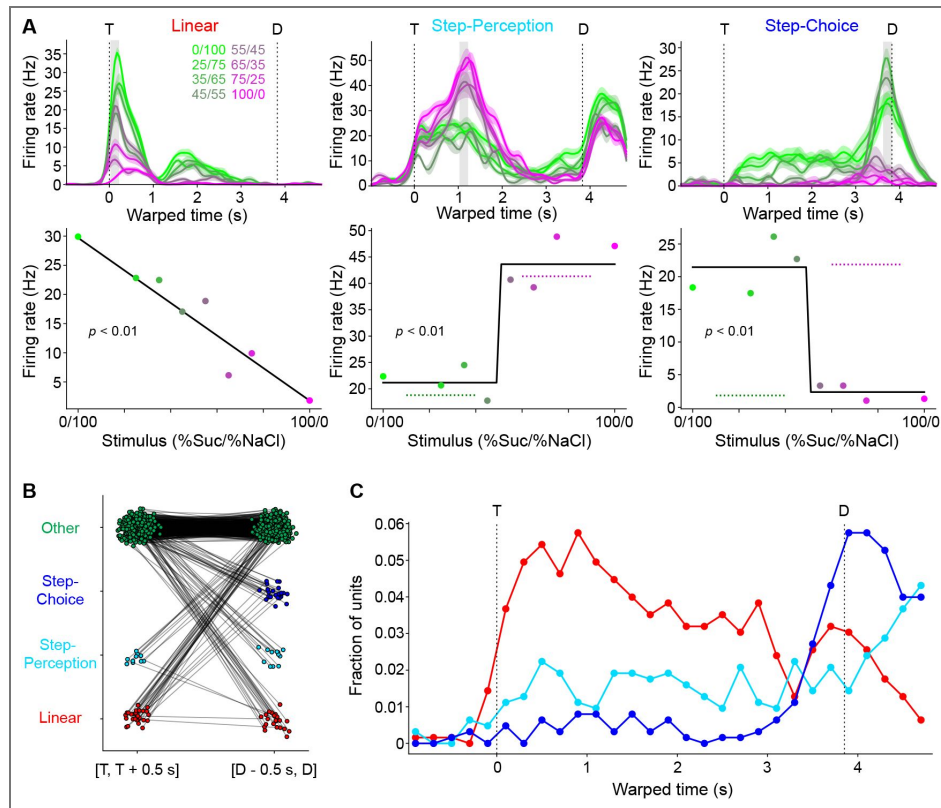


Figure 4. Classification of single unit coding types.

A: Representative single unit PSTH (top) and response profiles (bottom) exemplifying the different coding types within a time window (grey bar, top): linear (left), step-perception (middle), and step-choice (right). Step-perception (middle) and step-choice (right) types were disentangled by comparing correct trials to error trials (dashed lines in bottom plots). Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl). **B:** Visualization of each neuron's coding type label (vertical axis) between two time windows (horizontal axis). Each neuron is a point in both windows, with lines connecting the same neurons. T: time of first central lick; D: time of first lateral lick. **C:** Distribution of coding types across all neurons (pooled over all sessions) over time. For each time point (a window ~200 ms wide), the coding type classification analysis depicted in **A** was applied to each neuron.

populations of linear and step coding neurons, rather than a single population of coding neurons that switched its coding type: 77.4% of the linear coding neurons at the beginning of the trial had no significant response fit at the end of the trial; similarly, 86.5% of step coding neurons at the end of the trial had no significant fit at the beginning of the trial. The remaining neurons multiplexed across time.

To investigate the dynamics of single neuron responses, the time course of the distribution of fits was computed by running the response profile analysis with a moving window of ~200 ms. When quantified in this way, we found that 24.8% (155/626) of neurons were linear coding in at least one bin, 26.7% (167/626) were step-perception, and 18.2% (114/626) step-choice (see [Table S2](#) for session-by-session data). That said, the peak percent of coding units in any specific bin was much lower (maximum 10.7% total). Visual inspection of [Figure 4C](#) indicates the same trend of switching from mostly linear coding to mostly step-choice coding over time ([Figure 4C](#)) seen at the population level ([Figure 3](#)). In addition, the analysis reveals a small but consistently present proportion of step-perception neurons.

Altogether, single neuron analyses confirm the population results on coding dynamics and extend those findings to show that GC single units can encode in a binary fashion both the choice of licking direction (i.e., left vs right) as well as the perceptual category (i.e., sweet vs salty). These results also raise questions about the functional significance of the relatively low percentage of linear and step coding neurons and their contribution to population dynamics and task performance.

Recurrent neural networks capture experimental neural and behavioral results

To investigate the functional role of single neuron response types described above, we relied on a computational approach and, for each recording session, built a recurrent neural network (RNN) constrained by the single neuron data and capable of reproducing behavioral performance. For each experimental session, we modeled the simultaneously recorded neurons as a fraction of a larger system of units that received external stimulus input, noise input, and recurrent input from other units. The model was partitioned into units that were trained to reproduce the neural activity directly observed during the experiment, termed “constrained,” and those that were not, termed “unconstrained.” The ratio of total units to constrained units was fixed at 5.88 (thus, constrained units were ~17% of each network). An additional external unit allowed for the model to produce “choice activity” by weighting the firing rates of all internal units. Discrete choices were obtained by thresholding the average choice activity over a decision window (from $D - 100$ ms to D , for $D = 3.9$ s, the average decision time in the experimental dataset). The network’s training incorporated two processes: reproducing the experimentally observed patterns of neural activity within the constrained population, while simultaneously selecting the appropriate behavioral response to each stimulus. That is, the output of each constrained unit was trained to match the PSTHs of a corresponding neuron actually recorded from the mouse GC during behavior. This approach, like the one described in [Cohen et al. \(2020\)](#), enhanced the biological realism of the RNN trained to perform the task since its internal activity was explicitly instructed to resemble true neural activity. [Figure 5A](#) illustrates the key components of our RNN model.

The models successfully learned to perform the decision-making task and produced psychometric curves qualitatively similar to real animals’ when presented with noisy mixture stimuli ([Figure 5B](#)). Although the psychometric functions were significantly different between model and experiment (extra-sum-of-squares F-test, $p < 0.001$), with a greater slope for the model’s (0.15 vs 0.08), this was unsurprising given that we tuned the level of input noise only to match overall accuracy, which was 77.2%, comparable to the 77.2% we saw from mice (t-test, $p = 0.964$). At the same time, the models successfully reproduced the experimentally observed neural activity. [Figure 5C](#) shows an example experimental PSTH (left) and the activity of the corresponding unit in the model trained to reproduce it when presented with noiseless mixture stimuli (right). This unit has a root-mean-squared-error between model output and target PSTH of 1.86 Hz; the median value across all constrained units was 1.26 Hz. Additional examples of constrained units are provided in

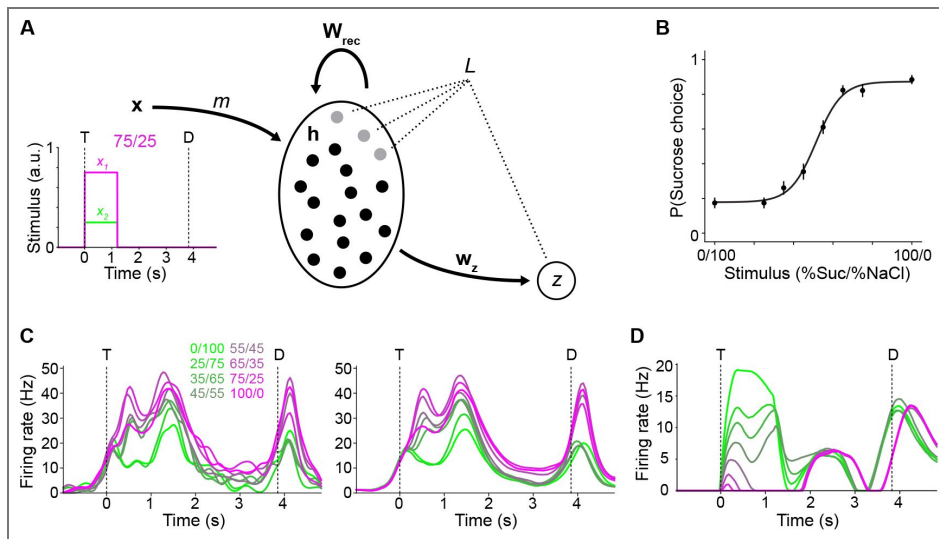


Figure 5. Recurrent neural network design and behavior.

A: Model architecture. N neurons are modeled as dynamic units with internal activity h that is influenced by the external stimulus input ($m(x)$); the time course of an example x for mixture stimulus 75/25 is shown), recurrent input (via w_{rec}), and noise input (not shown). A decision unit z measures the network's choice by taking a weighted sum of activities via w_z . The loss function L is minimized during training based on choice (z) and the activity of the constrained units (grey dots). T: time of stimulus onset; D: decision time. **B:** Psychometric curve fit to across-model means for the probability of the sucrose choice as a function of the stimulus. Circles and error bars represent mean and s.e.m. **C:** Example of experimentally observed PSTH (left) and the corresponding firing rate activity for the unit in the network trained to match it (right). **D:** Example firing rate activity for a unit in the network not explicitly trained to match any experimentally observed PSTH. Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl).

Figure S2 [↗](#). Unconstrained units learned to produce a variety of responses to stimuli, some of which resembled patterns seen in the experimental dataset. An example unconstrained unit is shown in Figure 5D [↗](#), and breakdowns of all responses are shown in Figures S3 [↗](#)-4 [↗](#). In total, 49.8% (1832/3681) of units were taste responsive and 40.8% (1502/3681) of units showed preparatory activity leading up to decisions (see Table S3 [↗](#) for across-session counts), comparable to the experimental findings.

Additional support for the overall agreement between experimental and model dynamics and coding schemes came from population analyses. 160 trials (20 per mixture stimulus) of data from each model were simulated, using the same levels of noise added to stimuli as were used to produce realistic psychometric curves. We then pooled all the units (constrained and unconstrained) across models and applied the same dPCA procedure we used on experimental data to find the principal stimulus- and choice-coding components. Population activity projected on the principal stimulus component (Figure 6A [↗](#)) showed a graded separation of trajectories according to stimulus and regardless of choice, while the projection on the principal choice component (Figure 6B [↗](#)) showed a binary separation of trajectories according to choice and regardless of stimulus. The time courses suggested a transition from stimulus-representing activity to choice-representing activity over the period between central and lateral licks. These patterns of activity are qualitatively consistent with those observed in the experimental data and demonstrate that the RNNs produce biologically realistic population dynamics.

To further validate the realism of the network and identify the salient features of neural activity on a single unit level, the same response profile classification analysis performed on experimental neurons was applied to model units. Using the results of the simulations described above, response profiles were calculated for all units (for correct and error trials, separately) in 200 ms bins, and the same procedure was used to assign a coding type label—linear, step-perception, step-choice, or “other”—to each response profile in each bin. The Venn diagram in Figure 6C [↗](#) (left) displays the breakdown in percentages of all units that exhibited each possible combination of coding types in at least one bin during the task period. In total, 24.2% (885/3681) of units were classified as linear coding, 26.8% (987/3681) of units were step-perception coding, and 25.5% (943/3681) of units were step-choice coding at some time point (see Table S4 [↗](#) for session-by-session data). As in the case of the experimental results, the model’s neurons that did not show significant fits (56.8% of units) could still be taste responsive (29.4%, 615/2092), taste selective (9.4%, 197/2092), and show preparatory responses (23.4%, 490/2092) (Figures S3 [↗](#)-4 [↗](#)).

Analyses on the time course of the distribution of coding types—linear, step-perception, and step-choice—across all units (pooled over all 23 models) showed qualitative agreement with the experimental findings (Figure 6C [↗](#), right). As in the experimental data, the peak percent of coding units in any specific bin was relatively low (maximum 14.9% total). The prevalence of linear coding units peaked soon after mixture sampling, while the step-choice coding units peaked soon before lateral licking. As in the experimental data, the model produced a lower percentage of step-perception coding units whose prevalence tiled the entire trial. The results of Figure 6 [↗](#) held even when the analyses were restricted to either the constrained (Figure S5 [↗](#)) or unconstrained (Figure S6 [↗](#)) units, with the constrained results qualitatively appearing even more similar to the experimental ones, as expected, and the unconstrained results appearing nearly indistinguishable from the full model results.

Altogether, the results show that the RNN models trained to reproduce behavioral performance and, in a fraction of the units, the observed neural activity, generate population and single unit coding patterns analogous to those observed in GC of behaving mice.

Model perturbations reveal behavioral significance of coding unit types

Generating a series of RNNs that aligned with experimentally observed neural activity and behavioral performance allowed for an exploration of the functional role of the different types of coding units (linear, step-perception, and step-choice). A series of virtual “ablation” experiments

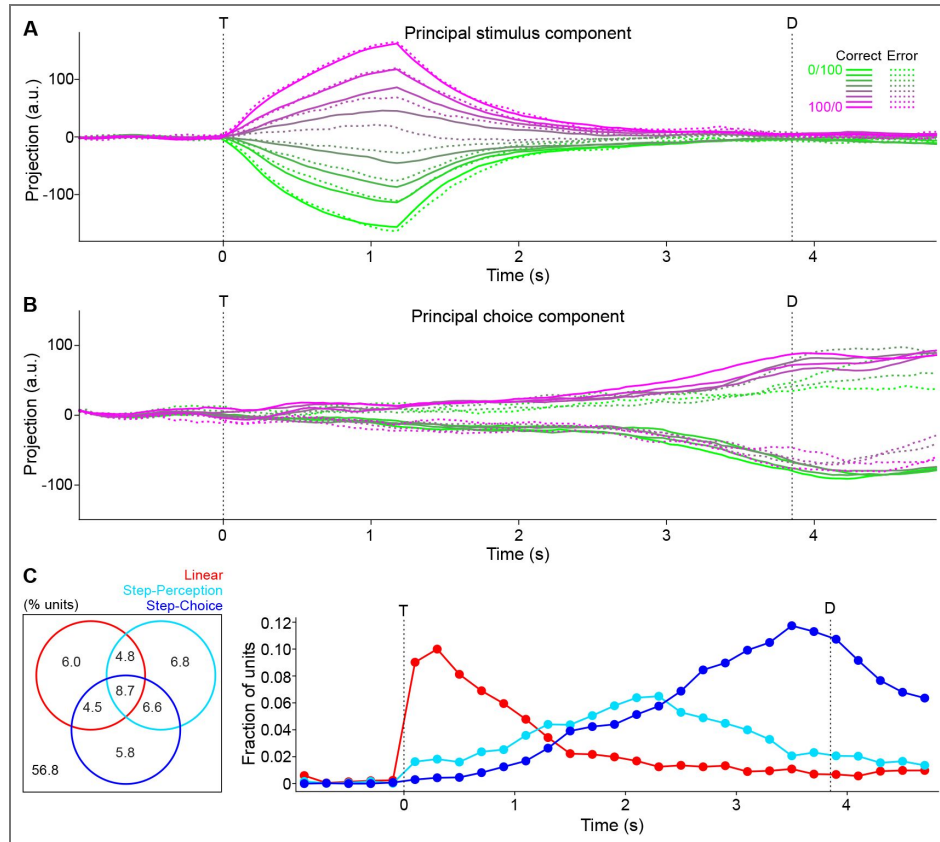


Figure 6. Modeled population activity and single unit coding properties.

A: Trial-averaged pseudo-population activity trajectories projected onto demixed principal component of maximal stimulus-specific variance. Solid lines are correct trial averages; dotted lines are incorrect trial averages. Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl). T: time of stimulus onset; D: decision time. **B:** Same as **A** but for the demixed principal component of maximal choice-specific variance. **C:** Left: Venn diagram showing percentages of neurons with all possible combinations of coding types over time. Right: Distribution of coding types across all units (pooled over all models) over time.

were conducted by re-running simulations while clamping the firing rates of specific sub-populations of units to 0. “Ablation” experiments were conducted for units that exhibited linear coding at any time point in the original simulations, those that were classified as step-perception coding at any time point, those that were labeled as step-choice at any time, and those that never followed any of these coding patterns (i.e., the “other” units).

All three coding types contributed to model dynamical activity along the stimulus- and choice-coding dimensions, as projecting the post-ablation activity onto the originally identified axes resulted in a noticeable blunting (Figure 7A). Quantitatively, mean absolute projection values were significantly reduced (post-hoc Bonferroni-corrected paired t-tests; stimulus projections: control vs linear, $p < 0.001$; control vs perception, $p < 0.001$; control vs choice, $p < 0.001$; choice projections: control vs linear, $p < 0.001$; control vs perception, $p < 0.001$; control vs choice, $p < 0.001$). Similarly, new stimulus- and choice-coding dimensions identified after ablating did not align with the old ones (Figure 7B) (absolute cosine similarities between vectors; stimulus components: control vs linear, 0.208; control vs perception, 0.244; control vs choice, 0.557; choice components: control vs linear, 0.016; control vs perception, 0.169; control vs choice, 0.021). The effects of the “ablations” could simply be the result of the removal of roughly a quarter of the units in a highly recurrent network leading to a large non-specific disruption of dynamics. However, this is not the case, as the “ablation” of the “other” units (i.e., those that do not show any significant fit), which constitute a much larger fraction of the network (56.8%, 2092/3681), left dynamics largely intact. Activity projections onto the original stimulus- and choice-coding dimensions after ablating “other” units (Figure 7A) were similar to the control condition without ablation (Figure 6A-B) (stimulus projection: $p = 0.105$; choice projection: $p > 0.999$) and newly identified coding dimensions (after the ablations) overlapped highly with the originals (i.e., before the ablations; Figure 7B) (stimulus component: 0.947; choice component: 0.959). This is relevant as neurons in this group show firing modulations to task events (see above).

In terms of behavioral impact, all three coding types were also necessary for normal model performance, as task accuracy dropped significantly upon ablating any of them, whereas task performance was unaffected by ablating the “other” units (post-hoc paired t-tests with Bonferroni correction; control vs linear, $p < 0.001$; control vs step-perception, $p < 0.001$; control vs step-choice, $p < 0.001$; control vs “other,” $p = 0.197$; Figure 7C). Furthermore, the psychometric functions were not significantly different between the control and ablated “other” conditions, while they were different between the control and all other conditions (extra-sum-of-squares F-tests with Bonferroni correction; control vs linear, $p < 0.001$; control vs perception, $p < 0.001$; control vs choice, $p < 0.001$; control vs “other,” $p = 0.235$).

To determine the relative contributions of the constrained and unconstrained units to these results, we also carried out the coding type ablation simulations while restricting the ablation to the constrained or unconstrained sub-populations. The fractions of each coding category comprised of constrained units were 194/885 for linear, 147/987 for step-perception, 129/943 for step-choice, and 353/2092 for “other” units (Table S4). In terms of dynamics, there was a significant main effect of the constrained vs unconstrained population on stimulus coding activity (2-way within-subjects ANOVA with factors coding type [linear/step-perception/step-choice/other] and constraint [constrained/unconstrained]; constraint main effect $p < 0.001$; Figure S7A) and choice coding activity (2-way within-subjects ANOVA with factors coding type [linear/step-perception/step-choice/other] and constraint [constrained/unconstrained]; constraint main effect $p < 0.001$; Figure S7B); on average, ablating the unconstrained population impaired dynamics more, most likely due to its larger population size. Stimulus coding activity was sufficiently diminished by ablating linear or step-perception units restricted to the unconstrained population, but this was not the case for step-choice units at our $\alpha = 0.01$ significance threshold (Dunnett test vs control; constrained linear, $p = 0.291$; unconstrained linear, $p < 0.001$; constrained step-perception, $p = 0.661$; unconstrained step-perception, $p < 0.001$; constrained step-choice, $p = 0.953$; unconstrained step-choice, $p = 0.045$; constrained “other,” $p > 0.999$; unconstrained “other,” $p = 0.739$; Figure S7A). Choice coding activity was impaired by ablating linear, step-perception, or step-choice units restricted to either the constrained or unconstrained populations (Dunnett test vs

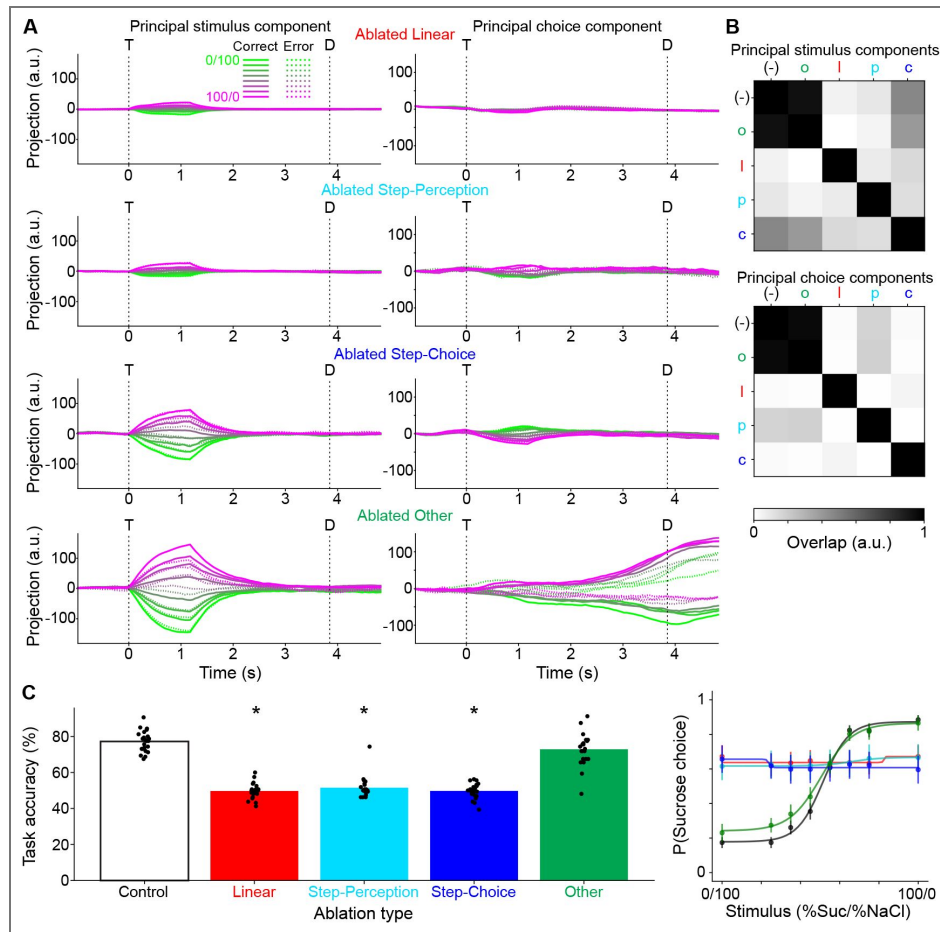


Figure 7. Effect of selective ablations on model dynamics and behavior.

A: Model dynamics after selectively ablating linear coding units, step-perception coding units, step-choice coding units, or "other" units. Post-ablation pseudo-population activity is projected onto the stimulus- (left column) and choice-coding (right column) components identified in the control condition (i.e., the same ones in Figure 6A-B). Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl); solid and dashed lines correspond to correct and error trials. T: time of stimulus onset; D: decision time. **B:** Pairwise overlaps between stimulus- (top) and choice-coding (bottom) components for control (-) and each ablation condition (o: "other," l: linear, p: step-perception, c: step-choice). **C:** Behavioral performance of the model after selectively ablating categories of coding units. Left: across-model distributions of task accuracy vs ablation condition. Bars represent means. * indicates significant difference vs control condition (post-hoc paired t-test Bonferroni-adjusted $p < 0.01$). Right: psychometric functions fit to across-model mean probability of sucrose choice for different ablation conditions. Circles and error bars represent mean and s.e.m.

control; constrained linear, $p < 0.001$; unconstrained linear, $p < 0.001$; constrained step-perception, $p < 0.001$; unconstrained step-perception, $p < 0.001$; constrained step-choice, $p < 0.001$; unconstrained step-choice, $p < 0.001$; constrained “other,” $p = 0.499$; unconstrained “other,” $p = 0.357$; [Figure S7B](#)). In terms of behavioral performance, again there was a significant main effect of the constrained vs unconstrained population (2-way within-subjects ANOVA with factors coding type [linear/step-perception/step-choice/other] and constraint [constrained/unconstrained]; constraint main effect $p < 0.001$; [Figure S7C](#)), with un-constrained ablations having larger impact, and task accuracy was still reduced by ablating linear, step-perception, or step-choice units, regardless of their restriction to the constrained or unconstrained populations (Dunnett test vs control; constrained linear, $p < 0.001$; unconstrained linear, $p < 0.001$; constrained step-perception, $p < 0.001$; unconstrained step-perception, $p < 0.001$; constrained choice, $p = 0.002$; unconstrained choice, $p < 0.001$; constrained “other,” $p = 0.649$; unconstrained “other,” $p > 0.999$; [Figure S7C](#)).

Finally, to analyze the temporal aspect of coding unit relevance to dynamics and behavior, we conducted additional ablation simulations with refined targets and time windows, based on two periods of interest: the beginning, i.e., the first 1.2 s after stimulus onset, and the end, i.e., the last 1.2 s prior to decision. For each period, we silenced the units with a particular coding type label at some point within that period for the entirety of the period. In terms of dynamics, stimulus coding was significantly reduced only by silencing the linear units in the beginning period (Dunnett test vs control; linear beginning, $p < 0.001$; linear end, $p = 0.990$; step-perception beginning, $p = 0.016$; step-perception end, $p = 0.999$; step-choice beginning, $p > 0.999$; choice end, $p = 0.994$; [Figure S8A](#)), whereas choice coding was blunted by silencing step-choice units in the end period or by silencing linear units in either the beginning or the end periods (Dunnett test vs control; linear beginning, $p < 0.001$; linear end, $p < 0.001$; step-perception beginning, $p = 0.508$; step-perception end, $p = 0.515$; step-choice beginning, $p = 0.969$; step-choice end, $p < 0.001$; [Figure S8B](#)). In terms of behavioral performance, task accuracy was impaired by silencing linear units in the beginning period, step-perception units in the end period, or step-choice units in the end period (Dunnett test vs control; linear beginning, $p < 0.001$; linear end, $p = 0.433$; step-perception beginning, $p = 0.242$; step-perception end, $p < 0.001$; step-choice beginning, $p > 0.999$; step-choice end, $p < 0.001$; [Figure S8C](#)). Overall, this suggests dynamics and behavior are mostly driven by linear units in the beginning period and categorical units in the end period of the trial.

In summary, the effect of these ablations on dynamics and behavior suggests that the model was quite robust to perturbation in general, but sensitive to manipulations that targeted linear and step coding units.

Discussion

GC plays a fundamental role in representing multiple sensory, affective, cognitive, and motor processes associated with a gustatory experience ([Yamamoto et al., 1985](#); [Samuelsen et al., 2012](#); [Gardner and Fontanini, 2014](#); [Vincis and Fontanini, 2016](#)). This function is performed through time-varying patterns of neural activity that sequentially encode for different variables associated with a taste-related task. Neural dynamics in GC have been extensively studied both in single neurons and at the population level ([Sadacca et al., 2016](#); [Mahmood et al., 2023](#); [Mazzucato et al., 2015](#); [Mendoza et al., 2024](#); [Livneh and Andermann, 2021](#)), and their role in producing behavior is beginning to be elucidated ([Kusumoto-Yoshida et al., 2015](#); [Mukherjee et al., 2019](#); [Vincis et al., 2020](#)). Yet, the relationship among single neuron firing patterns, population dynamics, and behavior is not completely understood. Here we investigate how sub-populations of GC neurons defined by their single unit response profiles influence population dynamics and behavioral performance in the context of a taste mixture-based 2AC task.

By recording neuronal activity with high-density probes in the GC of behaving mice, we unveiled population and single neuron dynamics associated with a taste mixture 2AC task. We found that both population and single neuron activities go through a phase in which mixtures are linearly coded by their components' concentration to a phase where stimuli are binarily coded. Additional analyses of single neuron activity show that units with binary coding could further be divided as

either representing the predominant mixture component—that is, its overall taste quality (sweet vs salty, “step-perception”)—or the directional licking decisions cued by the stimulus (left vs right, “step-choice”). While some neurons showed only one coding pattern, others could have different response profiles at different times during the trial. Overall, the neurons whose tuning curves could be significantly fit by a linear or step function at any point in time were less than 50% of the total number of recorded neurons, with each specific type constituting no more than 27%.

To study the functional role of such groups of single units, we built RNN models (Valente et al., 2022 [↗](#); Cohen et al., 2020 [↗](#)) of GC (one per session) and trained them to perform the taste mixture 2AC task while also reflecting the experimentally observed single neuron firing. The RNNs matched experimentally observed single neuron PSTHs, GC population dynamics, and rodent behavioral performance. Perturbing the model by removing different groups of neurons with distinct coding properties showed that linear and step coding neurons are necessary for neural population dynamics as well as behavioral performance. Ablation of the neurons that did not fit into the above categories had no impact on dynamics and performance, highlighting the functional importance of the single neuron firing patterns identified in this study.

Altogether the results presented in this study explain the role of single neuron firing patterns in a decision-making task and validate a data-driven, machine learning-based approach to modeling and generating hypotheses about the functional significance of system components (Song et al., 2016 [↗](#); Barak, 2017 [↗](#); Yang and Wang, 2020 [↗](#); Valente et al., 2022 [↗](#)).

Coding of task-related variables in GC of mice engaged in a taste mixture 2AC task

Since the pioneering work of Katz et al. (2001) [↗](#), it has been known that GC population dynamics can sequentially encode different aspects of a gustatory experience, from somatosensation to chemosensation to hedonic evaluation. More recent work has extended these findings to include GC population-level dynamics coding for taste-predictive cues, expectation, licking preparation, and abstract decision-making in the context of taste-based behavioral tasks (Stapleton, 2007 [↗](#); Samuelsen et al., 2012 [↗](#); Livneh et al., 2017 [↗](#); Fonseca et al., 2018 [↗](#); Vincis et al., 2020 [↗](#); Lang et al., 2023 [↗](#)). For instance, Vincis et al. (2020) [↗](#) showed population activity trajectories separating according to sensory quality (sweet vs bitter) before licking decision (lick left vs lick right) in GC of mice performing a taste-based 2AC task.

The population-level analyses presented here add to the growing body of evidence for GC dynamics' involvement in taste-based decision-making (Miller and Katz, 2010 [↗](#); Fonseca et al., 2018 [↗](#); Vincis et al., 2020 [↗](#); Lang et al., 2023 [↗](#); Jezzini and Padoa-Schioppa, 2024 [↗](#); Kogan and Fontanini, 2024 [↗](#); Zheng et al., 2025 [↗](#)). In the context of the taste mixture-based 2AC task presented here, we found neurons that discriminate between predominantly-sucrose and predominantly-NaCl mixtures throughout the trial time course, with population firing rates featuring two peaks, one related to the sampling of taste and one preceding lateral licking. Decoding of population activity showed that the representation of taste mixtures peaks early while choice-related coding peaks just before lateral licking. We applied dimensionality reduction techniques, t-SNE and demixed PCA (Kobak et al., 2016 [↗](#)), to extract the dominant trajectories coding for the task components in a neural space. The low-dimensional activity trajectories revealed a graded, linear separation based on mixture stimuli early on with sampling, and a binary separation based on selected choice later, prior to the decision. These population-level results, along with previous studies (Kogan and Fontanini, 2024 [↗](#); Maier and Katz, 2013 [↗](#)), led us to search for single neurons with linear and step coding responses. Indeed, a subset of single units whose tuning curves represent mixture stimuli as a linear or step function were identified. Based on the analysis of correct and error trials, the step coding units could be further divided as either representing the predominant mixture component (step-perception) or the directional licking decisions cued by the stimulus (step-choice). It is worth mentioning that while some neurons displayed only one coding pattern, others showed coding patterns that could vary in time, providing evidence for multiplexing at the single neuron level. Regardless, the time course of the prevalence of linear and step-choice responses was consistent with the dynamics observed at the

population level, with the former peaking during sampling and the latter before lateral licking. While the existence of these coding types suggests a role in driving population dynamics, the relatively limited presence of specific coding types in any time bin may cast doubts on their functional significance. In particular, step-perception units appeared in a uniformly low proportion throughout most of the trial and, surprisingly, peaked after the decision, perhaps serving as a feedback signal for learning by encoding the percept to be compared with the choice and outcome that just occurred. In this way, the percept could guide future adjustments of the choice following errors. The rarity of step-perception units and their post-decision peak therefore raised additional questions about the functional role of single unit coding patterns.

To address these questions on the functional role of the single unit activity described above, we relied on a modeling approach that would allow us to reproduce the experimentally recorded single unit activity as well as the behavioral performance for each session.

Using RNNs to investigate the role of single neuron coding patterns

Determining the functional role of a particular brain region's dynamics typically relies on optogenetic manipulations, which have become the gold standard due to their high temporal precision and the availability of specific genetically- and/or anatomically-targeted viral constructs (Li et al., 2019; Emiliani et al., 2022). Optogenetic silencing of GC and its inputs at different times during a trial has revealed the dynamic role of GC in encoding of palatability information (Lin et al., 2021), gaping behaviors (Mukherjee et al., 2019), expectation (Kusumoto-Yoshida et al., 2015), and taste-based decision-making (Vincis et al., 2020). As powerful as the approach is, it cannot be applied to selectively perturb neurons that are characterized by a specific coding pattern with no known genetic or connectivity signature. In other words, the specific coding populations we explored here cannot be manipulated by traditional optogenetic techniques. To address this gap, we relied on simulated manipulations in computational models where there is full control over all individual units.

Previous modeling efforts on GC have largely focused on replicating population-level phenomena (e.g., sequences of metastable states) starting from *a priori* architectures and tuning parameters by hand to reproduce population dynamics and behavioral performance (Miller and Katz, 2010; Mazzucato et al., 2015, 2016, 2019; Lang et al., 2023). While advancing our knowledge of the properties and origin of GC population dynamics, these efforts have not directly furthered our understanding of the role of single unit firing patterns. Here we turned to RNNs as an unbiased method for reproducing single unit firing activity (Barak, 2017; Yang and Wang, 2020). In particular, RNNs do not assume a specific synaptic matrix ahead of training. We took advantage of the relative ease of training RNNs (Song et al., 2016; Paszke et al., 2019) to learn single neuron activity and behavioral performance, capturing population dynamics in the process (Perich et al., 2020; Cohen et al., 2020; Valente et al., 2022; Rajan et al., 2016). Our RNN was composed of “constrained” and “unconstrained” units. Each constrained unit was paired with a target neuron from the experimental dataset, and the mismatch between model and experimental PSTHs was included in the loss function. Unconstrained units were included to represent unobserved neurons in the experimental recordings, and to add degrees of freedom to the network tasked with matching the constrained units' activities to targets. Each unit was driven by stimulus input, recurrent input, and noise input. The external decision unit was trained to select the correct choice given the stimulus; thus, we modeled an ideal scenario in which stimuli have matched intensities and no motor biases exist. Unlike previous models (Lang et al., 2023; Cisek et al., 2009), an external preparatory input prior to decision was not necessary. While this is not proof that GC lacks preparatory inputs, it demonstrates that these dynamics can be produced internally.

The tractability of training RNNs via automatic differentiation typically comes at the expense of realism and mechanistic understanding when considering the network as a model of the brain. In contrast, spiking neural networks offer much increased biophysical plausibility, but are much harder to train (Bohte, 2011; Li et al., 2021; DePasquale et al., 2016) and often require *a priori* decisions about their architecture. Here we took additional measures to alleviate the trade-

off between these two approaches by training a subset of neurons to reproduce the experimental PSTHs. This enhances biological realism and allows for RNN predictions to have more meaningful interpretations (Cohen et al., 2020 [↗](#); Rajan et al., 2016 [↗](#)). For our ablation studies, we identified the sub-populations of coding units based on their response patterns in the simulations. We then ran new simulations while clamping the firing rates of units in specific sub-populations to 0. We found that all three coding types—linear, step-perception, and step-choice—were required for normal population dynamics and behavior. We showed the impact of these unit types on population-level trajectories measured with demixed PCA and on behavioral performance measured with a psychometric function. The removal of linear coding units made the network unable to produce appropriate stimulus-related and choice-related population activity as well as flattened the psychometric function. Surprisingly, the same effect was obtained by the ablation of step-perception units. Ablation of step-choice neurons had a less dramatic effect on stimulus-related population activity but flattened choice-related trajectories. Crucially, ablating all other neurons (which were much more plentiful) left dynamics and behavior intact. This lack of effect from ablation was not because the activity of those neurons was unrelated to the task. Indeed, many represented stimuli and upcoming choices, only in ways outside of our defined coding types. This demonstrates the disproportionate importance of linear, step-perception, and step-choice coding types to the model compared to alternative coding strategies, and highlights their relevance. Interestingly, our results emphasize the importance of step-perception neurons, which appear to be responsible for bridging early stimulus-evoked dynamics with successive choice-related activity.

The work presented here leaves some open questions that present promising avenues for future research. Accumulating evidence indicates that GC's ongoing, taste-evoked, and decision-making activity is supported by metastable dynamics (Jones et al., 2007 [↗](#); Mazzucato et al., 2015 [↗](#); Sadacca et al., 2016 [↗](#); Lang et al., 2023 [↗](#)), activity manifested at the ensemble level as coordinated changes between internal (hidden) states. Although this work focuses on single-unit activity, it does so in the context of population dynamics. We therefore believe that our framework is compatible with alternative approaches centered on ensemble activity. For instance, neurons with the coding patterns identified here could coordinate with others to participate in the formation of different metastable states. However, our model was not specifically constructed to exhibit metastability as we fit it to single neuron activities and behavioral performance without enforcing the structural connectivity constraints required to produce metastable dynamics (Mazzucato et al., 2015 [↗](#)). Pre-umably, though, the most faithful model of GC would capture both single unit coding types and metastability, and future work could explore the implementation of a constrained fitting procedure that achieves this. Recent evidence also indicates that discrimination learning on this mixture task (that is, learning to improve difficult mixture pair discriminations via additional training) is associated with an increase in choice selective cells during the later delay period (Kogan and Fontanini, 2024 [↗](#)). The mechanistic origins of this finding are unclear, and a modified version of the model presented here may prove to be a useful tool for clarifying them.

Conclusions

In conclusion, our work shows that GC, a well-studied model for understanding cortical dynamics in sensory areas, can encode taste mixtures linearly or categorically, and choices categorically, during a mixture-based decision-making task. These phenomena are observed at the population level as well as at the single neuron level, and our findings are consistent with a dynamic progression of coding from representation of stimulus information to decision-making, with coding of perceptual category providing a bridging signal. The different types of coding sub-populations all make essential contributions to population dynamics and behavioral performance in our model of GC, underscoring the relevance of even small groups of neurons encoding task-relevant variables in very specific ways. It is worth noting that the neurons outside of these groups, whose activity was not necessary for normal population dynamics or behavioral performance in this particular task, may be critical for other taste-based decision-making tasks. Indeed, these “other” neurons may constitute a reservoir from which functionally-significant units

emerge during learning according to task-specific demands. We believe the modeling approach taken here is a powerful means of analyzing the impact of single units dynamics on network activity and performance and makes a case for data-driven, interpretable RNNs as useful tools in neuroscience that compromise between realistic biophysical models and “black box” machine learning models.

While the research presented here focuses on GC, its implications go beyond taste. Both the findings and the approach are relevant for understanding population and single neuron dynamics in areas where sensory, cognitive, and motor activity are jointly encoded.

Methods

Stereotaxic surgeries

Mice were anesthetized using a cocktail of ketamine (70 mg/kg) and dexmedetomidine (1 mg/kg) via intraperitoneal injection. After the animal was fully anesthetized, the head was shaved and cleaned with iodine and 70% ethanol. The animal was then transferred onto a stereotaxic apparatus. During the surgery, the depth of anesthesia was monitored via visual inspection of breathing rate, toe pinch reflex, and whisking. A heating pad was used to maintain body temperature. After the skin was excised, the skull was exposed and cleaned with saline, dry swabs, and 70% ethanol. A small amount of Vetbond (3M) was used to secure the skin around the edge of the incision. A pencil was used to trace the coronal, interfrontal, sagittal, and lambdoid sutures. GC craniotomy sites (AP: +1.2 mm, ML: ± 3.7 mm relative to Bregma) were marked with a permanent marker and covered with Kwik-Sil (World Precision Instruments). A craniotomy site above the cerebellum was drilled and a ground wire soldered to a male pin was placed (A-M system, Cat.No.786000). The midline of a custom head bar was aligned with the interfrontal and sagittal suture markings and positioned 1 mm posterior to Bregma. It was then secured with dental acrylic (C&B Metabond), covering both the skull and the top of the head bar. DV coordinates of Bregma and Lambda were remeasured, and a calibration point was marked on the head bar for stereotaxic reference.

Immunohistochemistry

Mice were deeply anesthetized with 220 mg/kg pentobarbital sodium (390 mg/ml) and were perfused with phosphate buffer saline (PBS) followed by 4% paraformaldehyde (PFA) in PBS. After post-fixing overnight in 4% PFA, the brains were sliced at 50 μ m with a vibratome (Leica VT-1000S). The brain slices were counterstained with Hoechst 33342 (1:5000 dilution, H3570, Thermo Fisher, Waltham, MA). Sections were mounted, cover-slipped, and imaged using a fluorescent microscope (Olympus BX51WI).

A Python-based GUI for Histological E-data Registration in Brain Space (HERBS) was used to register the slice images onto Allen CCF mouse atlas based on anatomical features for 2D and 3D visualization (Fuglstad et al., 2023 [↗](#); github.com/Whitlock-Group/HERBS [↗](#)). Reconstruction and visualization of electrode track trajectories was performed with open-source Allen CCF Tools (Shamash et al., 2018 [↗](#); github.com/cortex-lab/allenCCF [↗](#)) in a custom MATLAB script.

Statistical tests

For simple distribution comparisons, we used Mann-Whitney U tests (i.e., rank-sum tests) or t-tests (Python: *Scipy stats*). For comparing proportions, we used Chi-squared tests or 2-tailed binomial tests (Python: *Scipy stats*). For within-subjects comparisons with more than two levels, we used 1-way repeated measures ANOVAs (Python: *AnovaRM* from *statsmodels*) and followed up significant results with post-hoc paired t-tests (Python: *Scipy stats*). Bonferroni corrections were implemented by multiplying *p* values by *K*-choose-2, where *K* is the number of levels. For within-subjects comparisons with two predictors, we used 2-way repeated measures ANOVAs (MATLAB: *fitrm* and *ranova*; MathWorks) to analyze interactions and main effects. To compare all groups to a single

control group post hoc, we used Dunnett's test (Python: *Scipy stats*). Additional statistical tests (such as extra-sum-of-squares F-test; see below) were implemented with custom code in Python or MAT-LAB. Significance level was taken as $\alpha = 0.01$.

Behavioral data

All psychometric curves (Figures 1B, 5B, 7C) are least-squares fitted 4-parameter logistic functions of the form:

$$y(s) = (p_1 - p_4) [1 + \exp(-p_2(s - p_3))]^{-1} + p_4$$

where y is the probability of a sucrose choice, $s \in (0, 100)$ is the %Sucrose of the stimulus, p_1 is the upper asymptote, p_2 is the slope parameter, p_3 is the inflection point, and p_4 is the lower asymptote. Rather than averaging psychometric curves over sessions/models, a single psychometric was always fitted to the session-/model-averaged response as a function of the stimulus. Fitting was performed with Python's *scipy* library, using the *curve_fit* method from the *optimize* module, with restrictions on the parameters: $0 \leq p_1 \leq 1$; $15 \leq p_3 \leq 85$; $0 \leq p_4 \leq 1$.

Psychometric curves were compared to each other using an extra-sum-of-squares F-test. The F statistic was calculated as (Motulsky and Christopoulos, 2004; Maxwell et al., 2017):

$$F_{\text{stat}} = \frac{\text{SSE}_1 - \text{SSE}_2}{\text{SSE}_2} \cdot \frac{\text{df}_2}{\text{df}_1 - \text{df}_2}$$

where SSE_1 is the sum of squared errors when fitting all data points with a single curve, SSE_2 is the sum of squared errors when fitting the separate data points with two separate curves, $\text{df}_1 = N_{\text{points}} - 4$, and $\text{df}_2 = N_{\text{points}} - 8$ (N_{points} is the total number of data points). The p value was calculated as the area under the $F(\text{df}_1 - \text{df}_2, \text{df}_2)$ distribution to the right of F_{stat} and evaluated at the $\alpha = 0.01$ significance level. Bonferroni corrections were applied by multiplying p by K -choose-2, where K is the total number of psychometrics being considered.

Electrophysiological data acquisition

Prior to each recording session, Neuropixel 1.0 probes were coated with Vybrant™ DiI (Thermo-Fischer) or neuro-DiO (Biotium). Before recording, the animal was placed in an induction chamber under 2.5% isoflurane for 2 to 3 minutes and then placed on a stereotaxic frame where anesthesia was maintained with 1 to 2% isoflurane. Dental acrylic over the Kwik-Sil was removed to expose the craniotomy site. A craniotomy was drilled based on the marker and cleaned with gel foam and saline. The animal was then transferred to the behavior platform. A multi-probe motorized manipulator system was used for recordings (New Scale Technologies). Bregma and Lambda were calibrated based on their relative distances from the calibration point. Neuropixels trajectory explorer with the Allen CCF mouse atlas was used to visualize the probe location in the brain (https://github.com/petersaj/neuropixels_trajectory_explorer). Craniotomies were kept moist with frequent application of saline during recordings and were sealed with Kwik-Sil and covered with a thin layer of dental cement after recording. The open-source software package SpikeGLX (<https://billkarsh.github.io/SpikeGLX/>) was used for data acquisition.

Spike sorting

Spike sorting was performed with Kilosort 2 and Kilosort 4 (Pachitariu et al., 2024; github.com/MouseLand/Kilosort) and sorted clusters were manually curated using Phy (Cyrille Rossant, International Brain Laboratory) and custom MATLAB scripts. Units were identified with distinct clusters in waveform principal component space and a clear refractory period (> 1 ms) in auto-correlation histograms.

Firing rate data

To analyze firing rate dynamics with respect to two discrete trial events—the time of the first central lick, T, and the time of the first lateral lick, D—simultaneously, we used a warped time scale. Spike trains were aligned to T, and a fixed number of bins (77) was used to calculate firing

rates between T and D. This inter-event interval, $IEI = D - T$, varied from trial to trial, with an average duration of 3.85 s across the entire dataset (thus, the mean bin duration was ~50 ms). The fixed number of bins ensured firing rate timeseries could be aligned to both T and D across all trials from all sessions. Firing rates before T and after D were calculated using 50 ms bins. After averaging over trials to construct PSTHs, firing rate activities were smoothed using acausal Gaussian kernels 11 bins wide.

auROC

Area under the receiver operating characteristic curve (auROC) was used to measure each neuron's average difference in firing rate between %Sucrose < %NaCl and %Sucrose > %NaCl trials over time. Each ROC curve was constructed as $\Pr(R_1 > \theta)$ vs $\Pr(R_2 > \theta)$ for R_1 the firing rate in %Sucrose < %NaCl trials, R_2 the firing rate in %Sucrose > %NaCl trials, and θ the threshold parameter. Thus, an auROC value of 0 indicated firing rates in predominantly-sucrose trials were always greater, a value of 1 indicated firing rates in predominantly-NaCl trials were always greater, and a value of 0.5 represented complete overlap between distributions of firing rates in both trial types. Peak auROC values were identified as those with greatest absolute difference from 0.5, and neurons were labeled as sucrose- or NaCl-"preferring" based on whether this peak value was closer to 0 or 1.

Responsivity and selectivity analyses

In line with previous work (Kogan and Fontanini, 2024 [↗](#)), we classified neurons as responsive in the sampling (T to T + 0.5 s) or delay (D - 0.5 s to D) periods if their firing rate distributions were significantly different from baseline (T - 3 s to T - 2.5s for sampling baseline; D - 5.5 s to D - 5 s for delay baseline). We then checked responsive neurons to see if their firing rate distributions within sampling or delay periods were significantly different between predominantly-sucrose and predominantly-NaCl trials—if so, they were considered selective. All statistical comparisons were done via Mann-Whitney U tests at the $\alpha = 0.01$ significance level. Only correct trials were used. We applied the same analyses to RNN units but with T - 0.5 s to T as the baseline for both sampling and delay periods.

t-SNE

A t-distributed stochastic neighbor embedding (t-SNE) was used as a non-linear dimensionality reduction approach for visualizing pseudo-population firing rate trajectories. This method aims to preserve the true distance structure among points in the original, high-dimensional space when mapping to the low-dimensional embedding space. We concatenated correct trial PSTHs to construct an $(N_{\text{stim}} \cdot N_{\text{time}}) \times N_{\text{neu}}$ pseudo-population matrix for $N_{\text{stim}} = 8$ the number of unique stimuli, $N_{\text{time}} = 117$ the total number of time points, and $N_{\text{neu}} = 626$ the total number of neurons in the pseudo-population. The number of columns in the matrix was then reduced to 2 via the embedding—we used the *TSNE* class from the *manifold* module of Python's *scikit-learn* library with default parameters.

dPCA

A demixed principal component analysis (dPCA) was performed as a linear, supervised alternative to t-SNE for dimensionality reduction. Detailed methods are described in Kobak et al. (2016) [↗](#). Briefly, this analysis began by organizing the firing rate data in a multi-dimensional array X of shape $N_{\text{trials}} \times N_{\text{neurons}} \times S \times Q \times T$, where $T = 117$ is the number of time bins, $Q = 2$ is the number of choices, $S = 8$ is the number of stimuli, $N_{\text{neurons}} = 626$ is the total number of recorded neurons, and N_{trials} is the maximum number of trials across all sessions for any stimulus and choice combination (correct and error trials are included). Whenever a session had no trials for a particular stimulus/choice combination, we fit simple linear models to each time bin's available data and used them to predict the missing data:

$$y(t) = \beta_0(t) + \beta_1(t) \cdot \text{stim} + \beta_2(t) \cdot \text{choice} + \beta_3(t) \cdot \text{stim} \cdot \text{choice},$$

where stim (in %Sucrose) is the stimulus value, $\text{choice} \in \{0, 1\}$ is the animal's choice (0 for “NaCl choice” and 1 for “sucrose choice”), $y(t)$ is the predicted firing rate in time bin t , and the $\beta(t)$ s are the fitted regression weights. At least 2 trials per session and stimulus/choice combination are required for optimizing the regularization hyperparameter in dPCA (if only 1 existed, it was copied), though the main dPCA analysis operates on the trial-averaged data matrix $\tilde{\mathbf{X}}$ of shape $N_{\text{neurons}} \times SQT$. The analysis partitions this matrix into marginalized contributions from each variable, $\tilde{\mathbf{X}} = \sum_{\phi} \tilde{\mathbf{X}}_{\phi}$, where $\phi \in \{s, q, t, sq, st, qt, sqt\}$ is a variable combination (s for stimulus, q for choice, t for time; we joined s with st , q with qt , and sq with sqt), then minimizes the loss L_{ϕ} for each:

$$L_{\phi} = \|\tilde{\mathbf{X}}_{\phi} - \mathbf{F}_{\phi} \mathbf{D}_{\phi} \tilde{\mathbf{X}}\|_2^2 + SQT \|\mathbf{F}_{\phi} \mathbf{D}_{\phi} \tilde{\mathbf{C}}^{1/2}\|_2^2 + \mu \|\mathbf{F}_{\phi} \mathbf{D}_{\phi}\|_2^2$$

where \mathbf{F}_{ϕ} and \mathbf{D}_{ϕ} are encoders and decoders, respectively, $\tilde{\mathbf{C}} = \frac{1}{SQT} \langle \tilde{\mathbf{X}}_{\phi} \tilde{\mathbf{X}}_{\phi}^T \rangle_{\phi}$ is the average covariance matrix over all parameter combinations ($\langle \cdot \rangle_{\phi}$ denotes averaging over ϕ and $(\cdot)^T$ denotes matrix/vector transposition), μ is the regularization hyperparameter, and $\|\cdot\|_2$ denotes the Frobenius norm for matrices and the 2-norm for vectors. After fitting, one-dimensional projections of pseudo-population activity were obtained by calculating the inner product of \mathbf{d}_{ϕ} and the N_{neurons} -dimensional vector over time, where \mathbf{d}_{ϕ} is the row of \mathbf{D}_{ϕ} that maximized projected variance.

Overlaps between components obtained via dPCA were calculated as an absolute cosine similarity between vectors:

$$\text{overlap}(\mathbf{u}, \mathbf{v}) = \frac{|\mathbf{u}^T \mathbf{v}|}{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2}$$

Minimum-distance decoding

For each session, we trained minimum-distance classifiers to decode several task-relevant variables from population firing rate activity over time. These custom classifiers may be thought of as variants of 1-nearest neighbor classifiers where neighbors are centroids (class averages). At each point in time, all trials were represented as population firing rate vectors and their associated labels $\{(\mathbf{r}_i, l_i)\}$ for i indexing trials and $l_i \in U$, the set of unique labels (e.g., for decoding Choice, $U = \{\text{Left}, \text{Right}\}$). Each trial i was held out, and mean population firing rate vectors were calculated for each label $u \in U$:

$$\mathbf{c}_u = \langle \mathbf{r}_j \rangle_{j:(j \neq i) \wedge (l_j = u)}$$

where \wedge denotes logical “and.”

The held-out trial was then assigned the class label corresponding to the nearest mean firing rate vector:

$$\text{class}(\mathbf{r}_i) = \hat{l}_i = \text{argmin}_u d(\mathbf{r}_i, \mathbf{c}_u)$$

where d is a distance metric. We chose the Euclidean distance metric, $d(\mathbf{r}_i, \mathbf{c}_u) = \|\mathbf{r}_i - \mathbf{c}_u\|_2$.

The class-balanced leave-one-out test accuracy was calculated as:

$$\frac{1}{|U|} \sum_u \frac{\sum_i I((\hat{l}_i = l_i) \wedge (l_i = u))}{\sum_i I(l_i = u)}$$

where $|U|$ is the cardinality of U , i.e., the number of unique labels, and I is the indicator function:

$$I(z) = \begin{cases} 0, & z \text{ is false} \\ 1, & z \text{ is true} \end{cases}$$

A one-tailed α -significance threshold for the session-averaged decoding accuracy was calculated from the binomial distribution as k/N , where N is the average number of trials and k is the minimum number of hits such that the tail probability to the right of k is still below α . That is:

$$k = \min \left\{ x : \sum_{y: y \geq x} \text{Pr}(y) < \alpha \right\}$$

where $\Pr(y)$ is the probability mass function of a binomial random variable parameterized by N and p . We used $\alpha = 0.01$, $N = 137$ (the average number of trials across all sessions), and $p \in \{1/8, 1/2\}$ (the theoretical chance levels of decoding each variable).

We used a 2-way within-subjects ANOVA with factors decoding (stimulus or choice) and time (sampling or delay) to compare the time courses of decoding. Each decoding time course had its theoretical chance level subtracted prior to the test, and decoding time courses were averaged within the first 10 bins after T (T to $\sim T + 0.5$ s) and last 10 bins before D ($\sim D - 0.5$ s to D) for sampling and delay, respectively.

Response profiles

Single unit response profiles were analyzed with a least-squares regression-based pipeline similar to that of Maier and Katz (2013). Let $x_{k,t}$ represent a single neuron's pre-processed firing rate data (i.e., already time-warped and smoothed) for k indexing trials and t indexing time. Let $\sigma(k)$ be the stimulus administered on trial k (in terms of %Sucrose, for simplicity) and $o(k)$ be the outcome of the trial (i.e., correct or error). For each neuron from each session, the response profile r was calculated as a function of the stimulus s (again, in terms of %Sucrose) for a given time window w :

$$r(s) = \langle \langle x_{k,t} \rangle_{t:t \in w} \rangle_{k:(o(k)=\text{correct}) \wedge (\sigma(k)=s)}$$

where again $\langle \cdot \rangle_x$ denotes averaging over x and \wedge denotes logical "and". The shape of each response profile was then labeled "linear," "step," or "other" by comparing to template shapes: a 2-parameter line,

$$f_{\text{line}}(s) = p_1 s + p_2$$

and a 3-parameter step function,

$$f_{\text{step}}(s) = \begin{cases} p_1, & s < p_3 \\ p_2, & s \geq p_3 \end{cases}$$

The comparison was carried out by finding the best fit (in the least-squares sense) for each template, and then comparing the resulting F values. The F_{stat} value for each template was calculated from the extra-sum-of-squares F formula above (**Methods: Behavioral data**), now applied to a fitted vs null model comparison rather than a nested model comparison (Motulsky and Christopoulos, 2004; Maxwell et al., 2017):

$$F_{\text{stat}} = \frac{\text{SSE}_1 - \text{SSE}_2}{\text{SSE}_2} \cdot \frac{\text{df}_2}{\text{df}_1 - \text{df}_2} = \frac{\sum_s (r(s) - \bar{r})^2 - \sum_s (r(s) - f(s))^2}{\sum_s (r(s) - f(s))^2} \cdot \frac{N_d - N_p}{N_p - 1}$$

where SSE_1 is the sum of squared error for the null model (the mean of the data), SSE_2 is the sum of squared error for the fitted model, df_1 and df_2 are the corresponding degrees of freedom, \bar{r} is the average of $r(s)$, N_d is the number of data points, and N_p is the number of parameters. The shape of the response profile was then assigned as the template with the largest associated F_{stat} value, as long as this F_{stat} value's corresponding p value (area under the $F(\text{df}_1 - \text{df}_2, \text{df}_2)$ distribution to the right of F_{stat}) was less than $\alpha = 0.005$ (0.01,

Bonferroni-corrected for the number of tests). If this was not the case, the response profile's shape was assigned as "other." Least-squares fitting for f_{line} was performed as a constrained fit, subject to $f_{\text{line}}(s) \geq 0$ for all s , using Python's *scipy* library (*minimize* method of the *optimize* module). Least-squares fitting for f_{step} was performed manually by varying $p_3 \in \{40, 50, 60\}$ and, for each, finding the optimal remaining parameters as $p_1 = \langle r(s) \rangle_{s:s < p_3}$, $p_2 = \langle r(s) \rangle_{s:s \geq p_3}$.

To sub-classify neurons labeled "step" into "step-perception" and "step-choice," we incorporated error trials. To be considered a "step-choice" neuron, two criteria had to be met: (i) the inflection point of the step response profile (p_3) had to occur at $s = 50$; (ii) the average firing rate of the neuron had to be consistently (between correct and error trials) greater for trials where the same direction was chosen. For example, if the firing rate averaged over correct trials of $s > 50$ was greater than the average over correct trials of $s < 50$, then we required its firing rate averaged over error trials of $s < 50$ to be greater than the average over error trials of $s > 50$ (since that would

indicate a consistent “preference” for trials of one lick direction). More explicitly, we calculated $\delta_x = (r_x(s))_{s: s>50} - (r_x(s))_{s: s<50}$ for $x \in \{\text{correct, error}\}$ (where $r_{\text{correct}}(s) = r(s)$ and $r_{\text{error}}(s)$ is obtained by replacing the condition $\sigma(k) = \text{correct}$ with $\sigma(k) = \text{error}$ in the above definition of $r(s)$) and required $\delta_{\text{correct}} \delta_{\text{error}} < 0$ for the “step-choice” label. If no error trials existed for stimulus values on either side of $s = 50$ (relevant for some simulated data), “step” neurons were not sub-classified.

We examined the proportions of all neurons pooled over all sessions and subjects that were classified as each shape at the beginning and end of the trial using windows $w_1 = [T, T + 0.5 \text{ s}]$ and $w_2 = [D - 0.5 \text{ s}, D]$ on the un-warped time scale. Chi-squared tests of proportions were done without Yates’ correction, using *chi2_contingency* from the *stats* module of Python’s *scipy* library. We also visualized the proportions over the full trial time course by using a moving window 4 bins wide (~200 ms) on the warped time scale, stepped by 4 bins (the final window was only 1 bin wide since the total number of time points was not divisible by 4).

RNN: Components and dynamics

Recurrent neural network models were inspired by Cohen et al. (2020) [Cohen et al. \(2020\)](#) and adapted from opensource code (Valente et al., 2022 [Valente et al., 2022](#); github.com/adrian-valente/lowrank_inference). Each RNN used here is comprised of N_c constrained and N_u unconstrained artificial units, for N_c equal to the number of simultaneously recorded neurons in the corresponding experimental session and the total number set at $N = N_c + N_u = \text{round}(5.88 \cdot N_c)$. The additional decision unit, described later, is not included in the total count.

The internal activity of all units in the network, $\mathbf{h} \in \mathbb{R}^N$, is influenced by external stimulus input, noise input, and recurrent input. The external stimulus input, $m(\mathbf{x})$, is the mixture stimulus modeled as a 2-dimensional vector—e.g., 25/75 (%Sucrose/%NaCl) is $\mathbf{x} = [0.25, 0.75]^T$ —passed through the non-linear mapping $m: \mathbb{R}^2 \rightarrow \mathbb{R}^N$ defined by $m(\mathbf{x}) = \mathbf{A}^{(2)} \tanh(\mathbf{A}^{(1)} \mathbf{x})$ for matrices $\mathbf{A}^{(1)} \in \mathbb{R}^{100 \times 2}$, $\mathbf{A}^{(2)} \in \mathbb{R}^{N \times 100}$, and \tanh applied element-wise to vectors. The noise input is $\boldsymbol{\eta}$ for $\eta_i \sim \mathcal{N}(0, \sigma_\eta)$ resampled at each time step from a Gaussian distribution with mean 0 and standard deviation $\sigma_\eta = 0.05/\alpha$ (α defined below). The recurrent input is $\mathbf{W}_{\text{rec}} f(\mathbf{h} + \mathbf{b})$ for the synaptic matrix $\mathbf{W}_{\text{rec}} \in \mathbb{R}^{N \times N}$, the input bias $\mathbf{b} \in \mathbb{R}^N$, and the transfer function f , applied element-wise to vectors, is a rectified linear function with a maximum value of 80, i.e., $f(z) = \min(\max(z, 0), 80)$. The output of the transfer function is interpreted as a firing rate, i.e., $\mathbf{r} = f(\mathbf{h} + \mathbf{b})$. Network activity evolves according to:

$$\tau \frac{d\mathbf{h}}{dt} = -\mathbf{h} + m(\mathbf{x}) + \mathbf{W}_{\text{rec}} f(\mathbf{h} + \mathbf{b}) + \boldsymbol{\eta}$$

for time constant τ . We integrated this equation in discrete time using the forward Euler algorithm with step size $\alpha = dt/\tau = 0.2$. The model’s decision over time is governed by an additional decision unit that is functionally external from the network (i.e., it does not appear in the equation for \mathbf{h} above). This unit is denoted by z and its internal activity is driven by the firing rates of all N units in the network:

$$\tau \frac{dz}{dt} = -z + \mathbf{w}_z^T f(\mathbf{h} + \mathbf{b})$$

for weight vector $\mathbf{w}_z \in \mathbb{R}^N$. This equation is integrated in parallel with the one for \mathbf{h} and the model’s binary choice is determined by interpreting $c = \tanh(z)$ as “sucrose choice” if $c > 0$ and as “NaCl choice” if $c < 0$.

RNN: Training

The goal during training is for the model to respond to any particular stimulus input by producing the correct choice (i.e., value of c) during a pre-defined decision window as well as constrained firing rate activities (i.e., first N_c components of \mathbf{r}) that match the correct trial PSTHs of the neurons in the corresponding session. Thus, we define an overall loss as:

$$L = \lambda_{\text{beh}} L_{\text{beh}} + \lambda_{\text{neu}} L_{\text{neu}}$$

for behavioral loss L_{beh} , neural loss L_{neu} , and associated weights $\lambda_{\text{beh}} = 150$ and $\lambda_{\text{neu}} = 1$ to counter-balance the different error scales. The behavioral loss is:

$$L_{\text{beh}} = \frac{1}{\beta_{\text{beh}}} \left(\sum_k \sum_{t \in \text{win}_0} c_{k,t}^2 + \sum_k \sum_{t \in \text{win}_D} (c_{k,t} - \gamma_k)^2 \right)$$

where k indexes stimuli, t indexes time, $\text{win}_0 = [T - 1 \text{ s}, T]$ is the pre-stimulus window for T the stimulus delivery time, $\text{win}_D = [D - 0.1 \text{ s}, D]$ is the decision window for D the decision time, $c_{k,t}$ is the network's value of c in response to stimulus k at time t , γ_k is the correct decision in response to stimulus k (e.g., for $k = 0/100$, $\gamma_k = -1$; for $k = 100/0$, $\gamma_k = +1$), and β_{beh} is the total number of terms in the sums. The neural loss is:

$$L_{\text{neu}} = \frac{1}{\beta_{\text{neu}}} \sum_{k,t,n} (r_{k,t,n} - \rho_{k,t,n})^2$$

where n indexes constrained neurons, $r_{k,t,n}$ is the network's firing rate output for neuron n in response to stimulus k at time t , ρ is the experimentally observed correct trial PSTHs, and β_{neu} is the total number of terms in the sum.

Training was carried out in PyTorch (Paszke et al., 2019): L was minimized with respect to the network's trainable parameters— $\mathbf{A}^{(1)}$, $\mathbf{A}^{(2)}$, \mathbf{W}_{rec} , \mathbf{b} , \mathbf{w}_z , and the initial value of \mathbf{h} —via backpropagation, using the Adam optimizer with a learning rate of 0.01 and gradient clipping above 1.0. Elements of $\mathbf{A}^{(1)}$ and $\mathbf{A}^{(2)}$ were randomly initialized as $\mathcal{N}(0, 1)$. Elements of \mathbf{W}_{rec} were randomly initialized as $\mathcal{N}(0, 0.1/\sqrt{N})$, with self-connections prohibited. Elements of \mathbf{w}_z were randomly initialized as $\mathcal{N}(0, 1/N)$. Elements of \mathbf{b} and the initial values of \mathbf{h} were initialized at 0. Training proceeds for 2000 iterations or until $L < 1$.

RNN: Simulations

Single trial simulations (both during and after training) were conducted for 5.9 s each, matching the time-warped window of analysis used for experimental data. The stimulus \mathbf{x} was “on” 1 s after trial start and lasted for 1.2 s. When off, $\mathbf{x} = [0, 0]^T$. After training, simulations were carried out with persistent additional noise, changing the external input current to $m(\mathbf{x} + \mathbf{e})$ with $e_i \sim \mathcal{N}(0, \sigma)$ resampled at each time step. The level of external noise, σ , was tailored to each model such that its overall task accuracy was within 5% of the animal's accuracy from the corresponding experimental session (σ ranged from 0.10 to 1.15). We simulated 20 trials per stimulus, and choices were obtained from the sign of the average value of c over the decision window, which covers 0.1 s before the decision time (which occurs 3.9 s after stimulus start) up to the decision time.

For ablation simulations, elements of $\mathbf{r} = f(\mathbf{h} + \mathbf{b})$ corresponding to firing rates of units belonging to specific sub-populations of interest (identified based on results from simulations without ablation) were clamped to 0 for all time, except in Figure S8, where the firing rates were clamped to 0 only in specific temporal windows. Ablations can result in the model becoming highly biased toward one choice direction; to limit inclusion of inferred missing data for the dPCA of model dynamics across ablation conditions, we excluded models that did not have at least 2 trials for each chosen direction under all conditions. Out of 23 models, the number of models passing criterion depended on the conditions considered and were 14 for Figures 6A-B, 7A-B, S5A-B, and S6A-B; 13 for Figure S8A-B; and 8 for Figure S7A-B.

Supplementary materials

Session	Total	Taste Responsive	Taste Selective	Delay Responsive	Delay Selective
1	17	6	3	1	1
2	44	28	6	19	7
3	21	12	5	8	3
4	48	30	9	23	12
5	29	14	3	12	5
6	68	43	10	35	7
7	16	8	1	3	3
8	41	26	9	25	12
9	17	8	0	6	2
10	38	23	4	18	7
11	14	8	4	7	4
12	7	4	2	1	1
13	11	9	2	3	1
14	27	21	7	12	3
15	9	4	1	6	2
16	8	3	1	2	2
17	30	15	1	14	6
18	26	19	8	15	7
19	31	17	5	8	1
20	28	11	2	12	3
21	40	22	2	15	2
22	19	15	5	7	4
23	37	24	3	23	4
<i>Total</i>	<i>626</i>	<i>370</i>	<i>93</i>	<i>275</i>	<i>99</i>

Supplementary Table 1. Session-by-session responsive and selective neuron counts for experimental data.

Responsivity indicates a difference in firing rate distributions between baseline and a window of interest (from the first central lick to 500 ms after it for taste; from 500 ms before the first lateral lick to the first lateral lick for delay). Selectivity indicates a difference in firing rate distributions between categories within the window of interest (predominantly-sucrose vs predominantly-NaCl for taste; left vs right for delay).

Session	Total	Linear	Step-Perception	Step-Choice	Other
1	17	4	6	1	8
2	44	13	9	8	27
3	21	3	11	1	9
4	48	22	11	9	23
5	29	8	5	6	18
6	68	16	16	10	38
7	16	3	6	2	9
8	41	8	12	17	17
9	17	1	2	4	10
10	38	6	4	12	22
11	14	6	4	2	4
12	7	3	1	1	4
13	11	4	2	0	5
14	27	5	8	5	14
15	9	4	0	0	5
16	8	3	2	2	4
17	30	6	9	4	17
18	26	8	14	12	7
19	31	5	11	8	14
20	28	6	7	2	16
21	40	7	6	4	25
22	19	6	9	2	5
23	37	8	12	2	21
<i>Total</i>	626	155	167	114	322

Supplementary Table 2. Session-by-session neuron coding type counts for experimental data.

Neurons are assigned coding type labels if they exhibit the response profile pattern in any time bin (as per analysis in [Figure 4C](#)) and, thus, the labels are not mutually exclusive.

Session	Total		Taste Responsive		Taste Selective		Delay Responsive		Delay Selective	
	Con.	Unc.	Con.	Unc.	Con.	Unc.	Con.	Unc.	Con.	Unc.
1	17	83	10	55	10	42	8	48	5	38
2	44	215	9	25	5	15	10	14	7	10
3	21	102	11	62	8	38	13	51	8	42
4	48	234	10	79	9	52	9	59	10	49
5	29	142	20	97	8	77	19	74	13	79
6	68	332	25	144	20	109	30	115	18	115
7	16	78	10	45	7	27	12	36	13	42
8	41	200	19	79	13	64	17	62	10	21
9	17	83	8	34	3	19	6	29	4	13
10	38	185	10	75	5	35	15	67	10	29
11	14	68	6	42	3	31	9	37	9	40
12	7	34	3	23	2	18	5	22	3	24
13	11	54	8	39	5	31	2	25	6	34
14	27	132	13	75	8	56	16	98	10	71
15	9	44	6	34	8	38	9	29	3	38
16	8	39	7	29	4	25	6	26	4	27
17	30	146	16	86	11	60	10	55	7	49
18	26	127	19	74	10	46	17	55	14	37
19	31	151	25	101	12	86	18	61	11	69
20	28	137	14	91	8	59	10	57	10	67
21	40	195	20	106	13	87	24	84	16	71
22	19	93	13	60	5	43	10	43	7	29
23	37	181	17	78	12	65	17	63	12	59
<i>Total</i>	<i>626</i>	<i>3055</i>	<i>299</i>	<i>1533</i>	<i>189</i>	<i>1123</i>	<i>292</i>	<i>1210</i>	<i>210</i>	<i>1053</i>

Supplementary Table 3. Session-by-session responsive and selective unit counts for model data.

Responsivity indicates a difference in firing rate distributions between baseline and a window of interest (from stimulus onset to 500 ms after it for taste; from 500 ms before the decision to the decision for delay). Selectivity indicates a difference in firing rate distributions between categories within the window of interest (predominantly-sucrose vs predominantly-NaCl for taste; left vs right for delay).

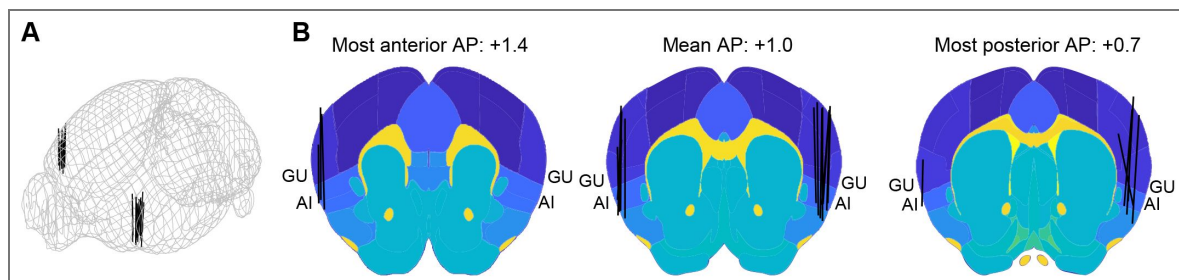
Supplementary Table 4. Session-by-session unit coding type counts for model data.

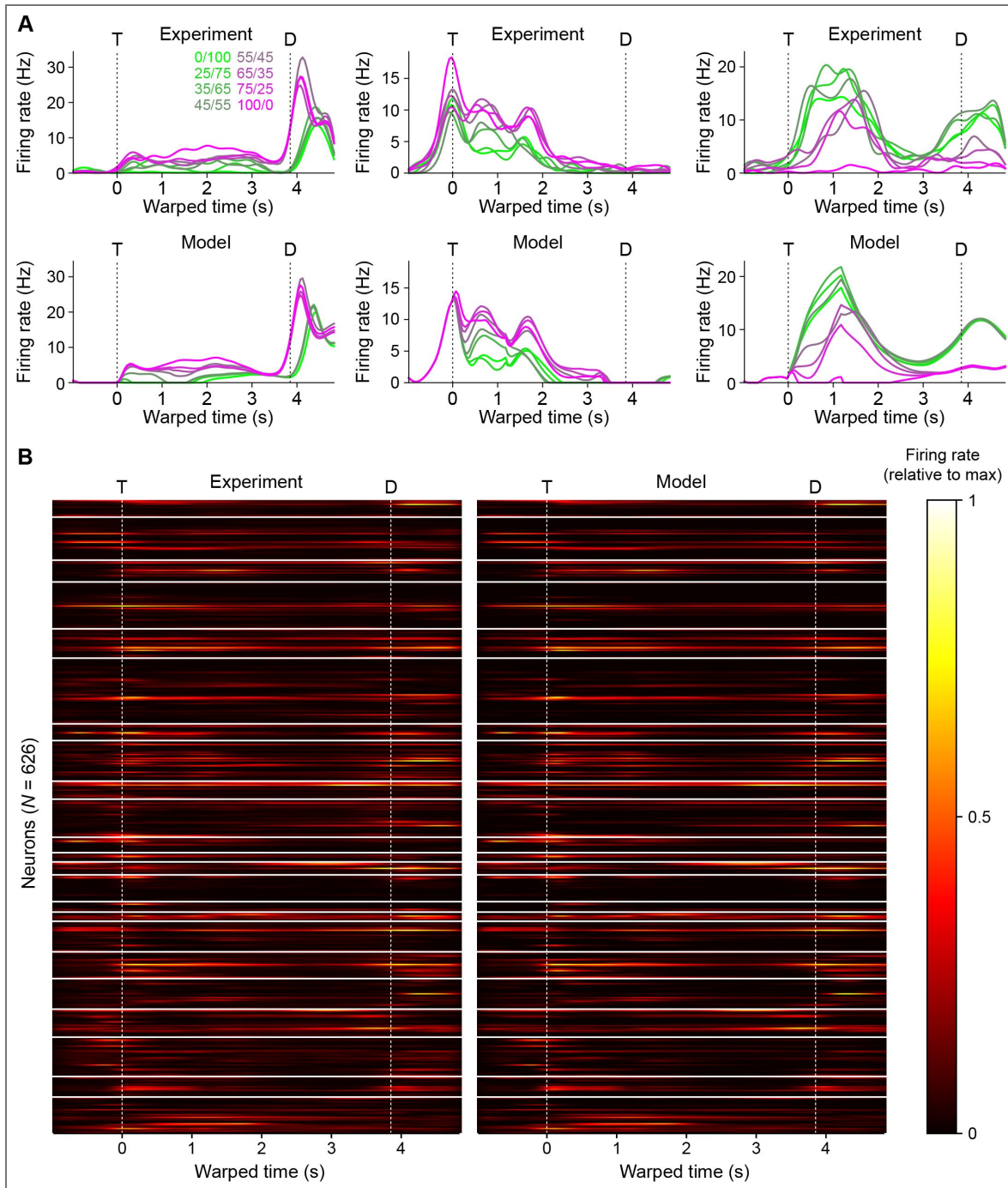
Units are assigned coding type labels if they exhibit the response profile pattern in any time bin (as per analysis in Figure 6C) and, thus, the labels are not mutually exclusive. Con.: constrained units; Unc.: unconstrained units.

Session	Total		Linear		Step-Perception		Step-Choice		Other	
	Con.	Unc.	Con.	Unc.	Con.	Unc.	Con.	Unc.	Con.	Unc.
1	17	83	5	17	6	32	6	42	10	28
2	44	215	8	9	5	18	9	14	33	191
3	21	102	6	13	4	29	2	36	14	56
4	48	234	9	37	8	46	3	25	37	175
5	29	142	11	57	12	61	3	57	11	52
6	68	332	13	53	8	70	11	80	47	204
7	16	78	11	17	12	34	9	32	3	39
8	41	200	12	21	4	28	1	4	27	163
9	17	83	4	12	7	23	3	9	10	56
10	38	185	12	28	3	19	8	27	26	137
11	14	68	7	33	8	40	4	32	4	21
12	7	34	2	15	2	22	0	19	3	6
13	11	54	4	32	2	38	4	35	4	10
14	27	132	10	39	9	50	8	44	14	52
15	9	44	3	19	0	3	4	39	4	4
16	8	39	4	15	5	17	3	29	0	6
17	30	146	4	22	4	38	7	46	20	88
18	26	127	16	31	17	39	5	27	6	67
19	31	151	14	38	6	50	14	58	12	63
20	28	137	9	53	3	49	4	41	18	62
21	40	195	16	48	10	62	11	57	19	102
22	19	93	5	31	4	25	6	28	9	48
23	37	181	9	51	8	47	4	33	22	109
<i>Total</i>	<i>626</i>	<i>3055</i>	<i>194</i>	<i>691</i>	<i>147</i>	<i>840</i>	<i>129</i>	<i>814</i>	<i>353</i>	<i>1739</i>

Supplementary Figure 1. Neuropixels probe trajectory reconstruction.

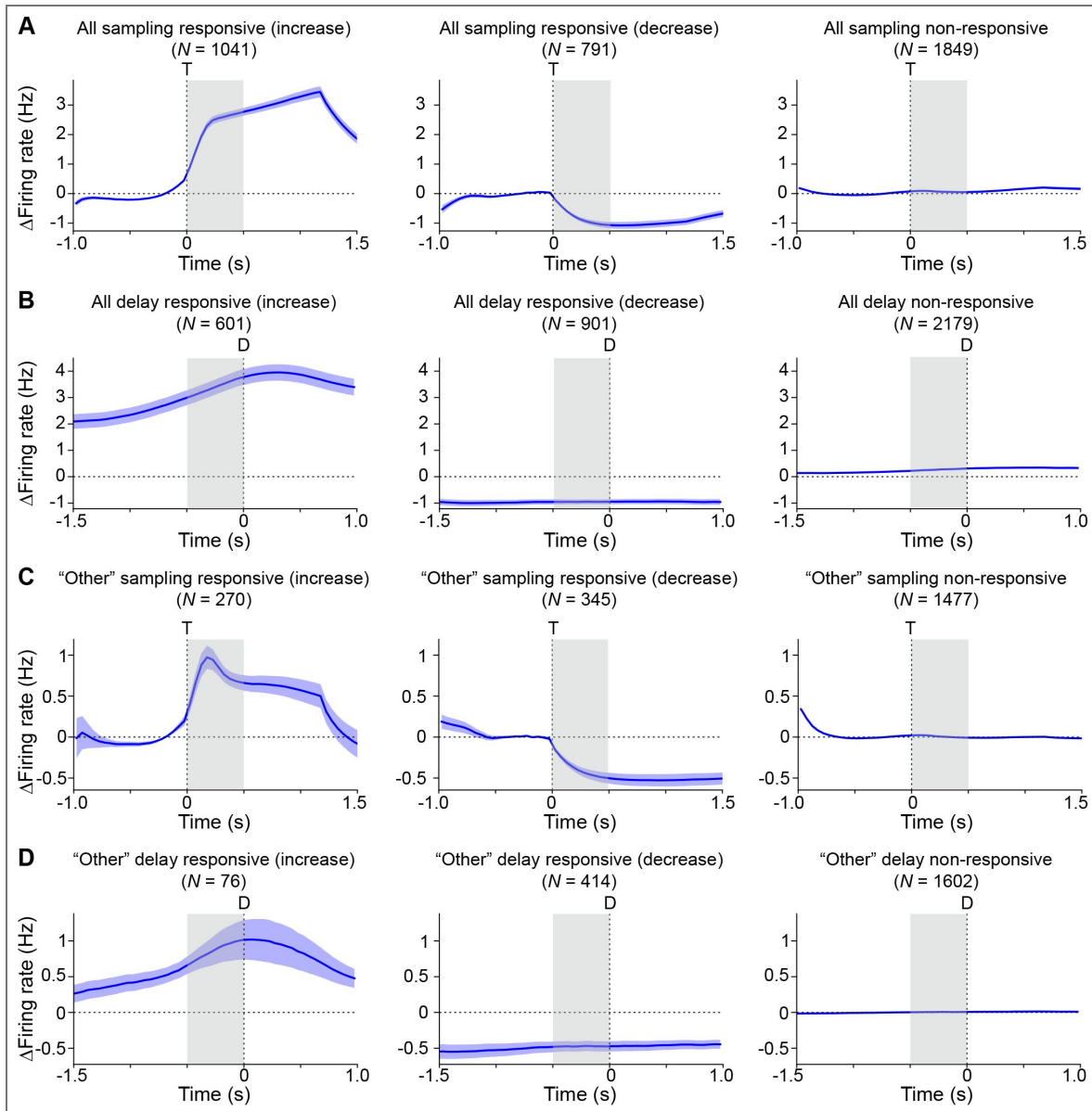
A: 3D reconstruction of the 23 probe trajectories from the experimental dataset. **B:** 2D reconstruction of the same 23 probe trajectories, overlaid on the Allen Brain Atlas at varying anteroposterior (AP) distances (relative to Bregma in mm) around GC. At these coordinates, both GU (gustatory areas) and AI (anterior insular areas) account for GC. Reconstructions performed with open-source Allen CCF Tools (Shamash et al., 2018; github.com/cortex-lab/allenCCF).





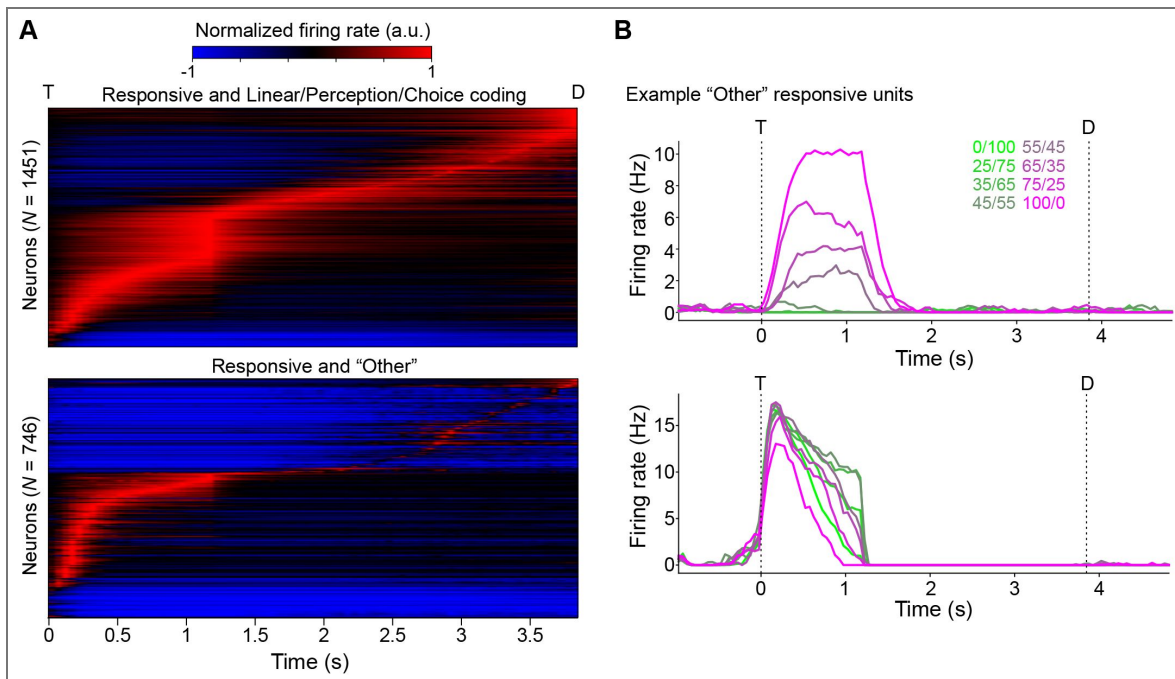
Supplementary Figure 2. Model constrained unit activity.

A: Three examples (columns) of experimentally-observed PSTHs (top) and corresponding model unit firing rate activities trained to match them (bottom). Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl). **B:** Comparison of firing rate activities (stimulus-averaged PSTHs) between all experimental neurons and their corresponding model constrained units. Vertical whitespace separates individual sessions/models. Firing rates are normalized to the maximum within each session. T: time of first central lick/stimulus onset; D: time of first lateral lick/decision time.



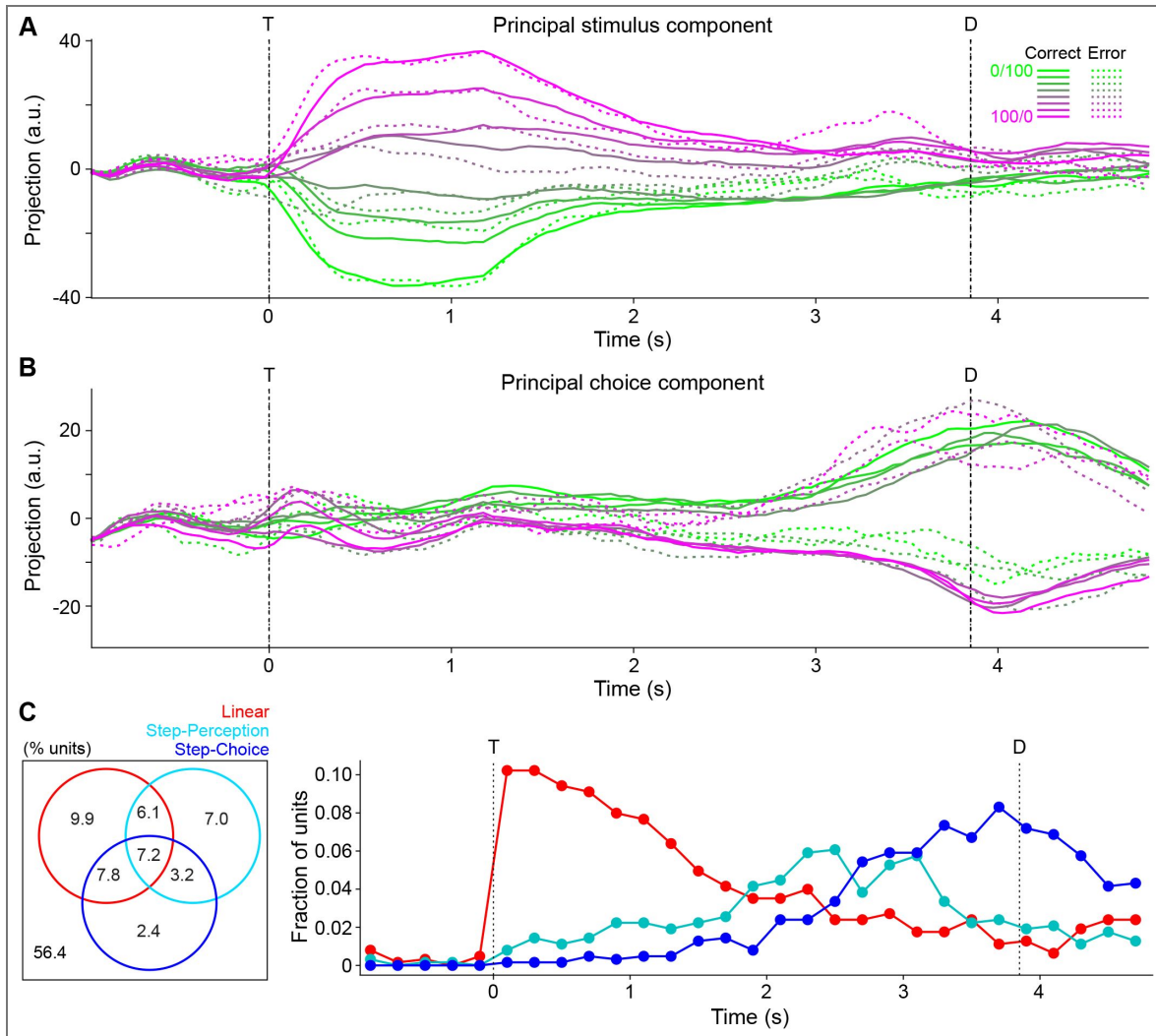
Supplementary Figure 3. RNN unit responsiveness.

A: Activity of all RNN units grouped by responsiveness during the sampling period. If the unit's firing rate distribution during the sampling period (T to $T + 0.5$ s for T the stimulus onset time) was significantly different from its baseline ($T - 0.5$ s to T) firing rate distribution, it was sampling responsive and grouped by whether its mean firing rate increased (left) or decreased (middle); otherwise it was non-responsive (right). **B:** Activity of all RNN units grouped by responsiveness during the delay period. Same as **A** except the firing rate distribution of interest is calculated over $D - 0.5$ s to D for D the decision time. **C** and **D:** Same as **A** and **B**, respectively, except that the only units considered are those labeled "other" by the response profile analysis of Figure 6C. Firing rates are expressed relative to baseline, and traces are population mean \pm s.e.m.



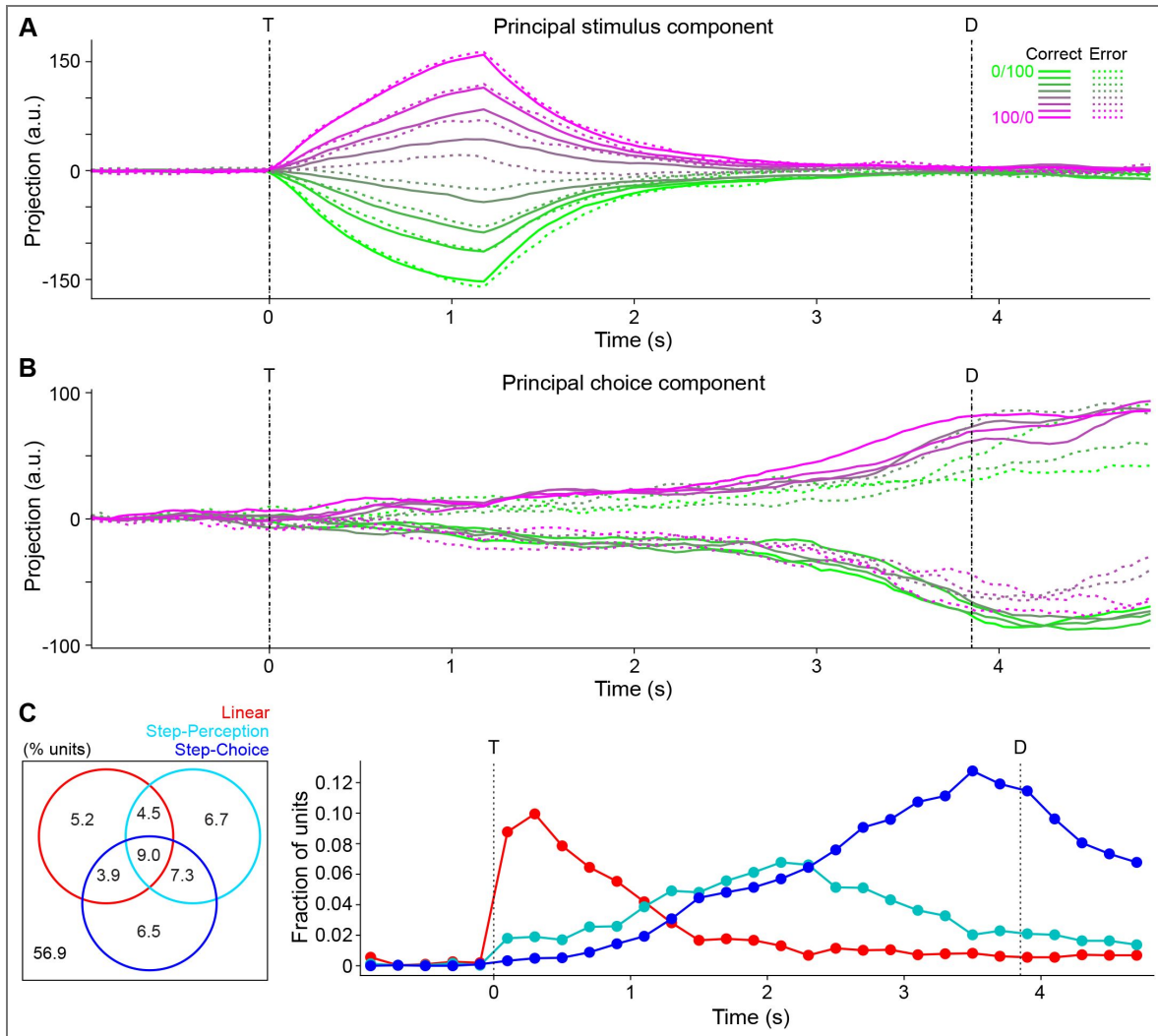
Supplementary Figure 4. RNN unit activity patterns.

A: Heatmaps of firing rate activities for units that responded significantly during the sampling and/or delay periods, broken down into coding units (linear, step-perception, and/or step-choice) (top) and "other" units (not linear, not step-perception, and not step-choice) (bottom). Firing rates are expressed relative to baseline and normalized to the maximum absolute value. T: time of stimulus onset; D: decision time. **B:** Two example "other" unit responses. Both respond significantly during the sampling period, but neither response pattern matches the linear or step templates. Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl).



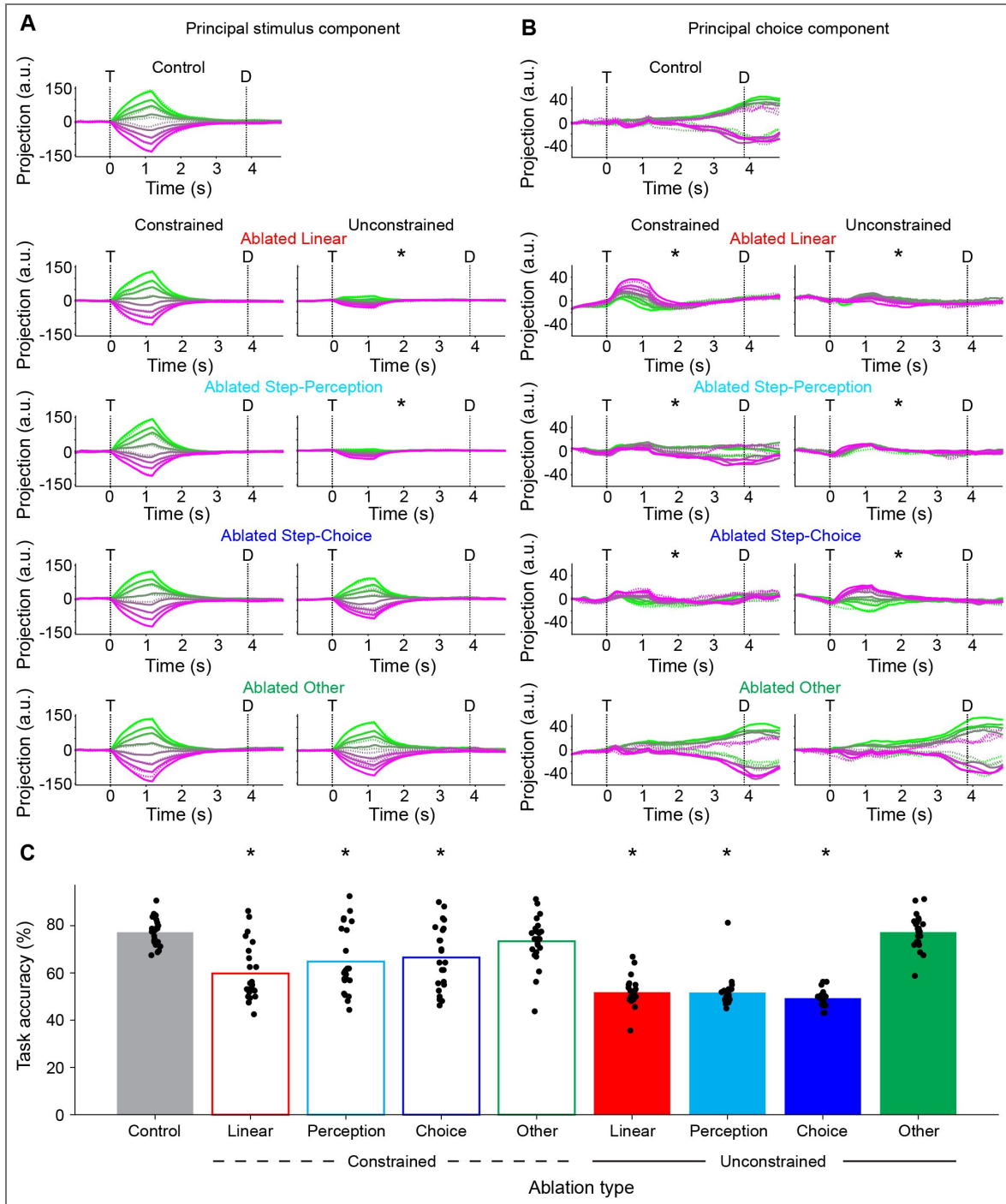
Supplementary Figure 5. Modeled population activity and single unit coding properties: constrained units only.

Compare with Figures 6 and S6. **A:** Trial-averaged constrained pseudo-population activity projected onto demixed principal component of maximal stimulus-specific variance. Solid lines are correct trial averages; dotted lines are incorrect trial averages. Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl). T: time of stimulus onset; D: decision time. **B:** Same as **A** but for the demixed principal component of maximal choice-specific variance. **C:** Left: Venn diagram showing percentages of constrained units (pooled over all models) with all possible combinations of coding types over time. Right: Distribution of coding types across constrained units over time.



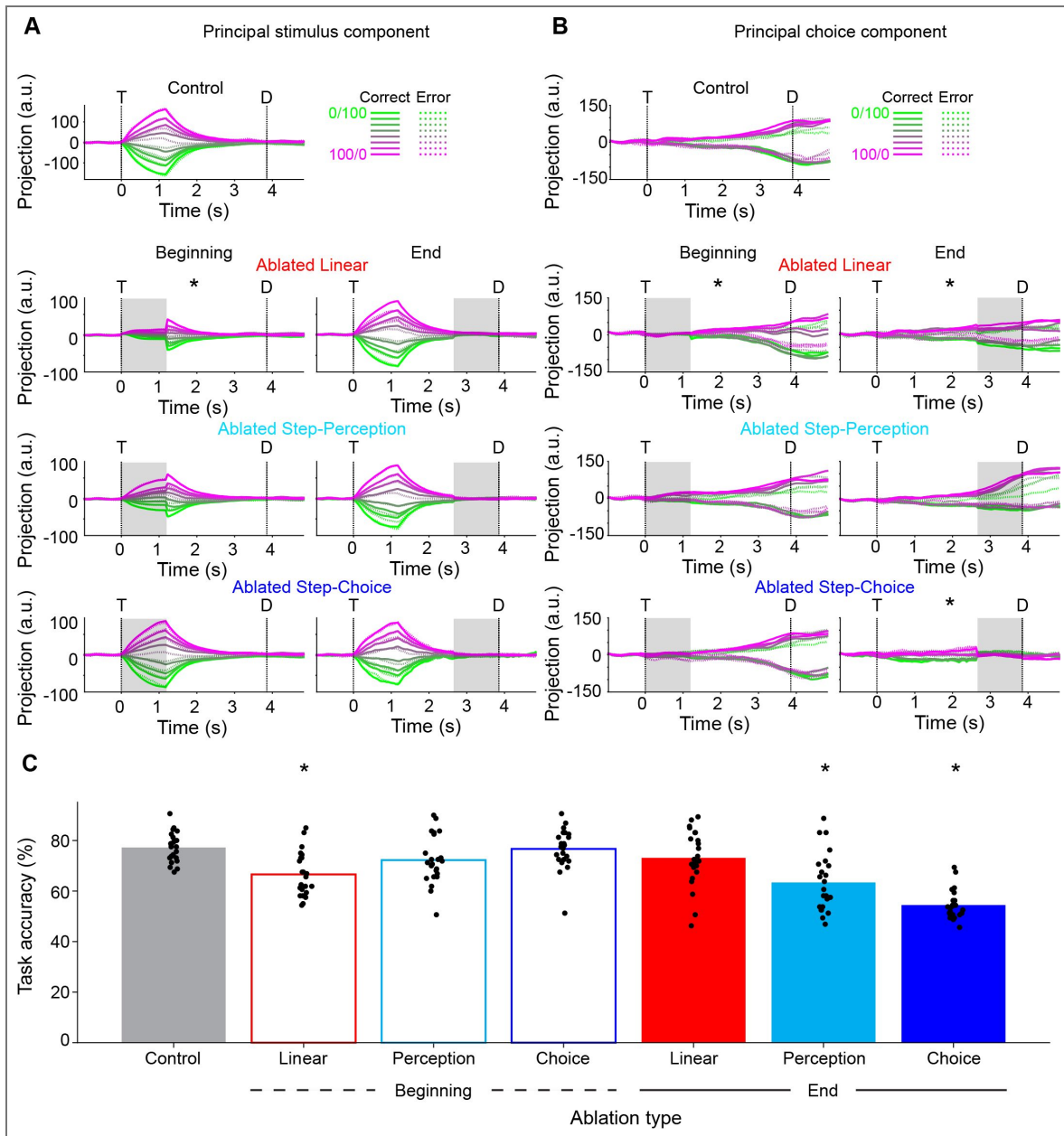
Supplementary Figure 6. Modeled population activity and single unit coding properties: unconstrained units only.

Compare with Figures 6 and S5. **A:** Trial-averaged unconstrained pseudo-population activity projected onto demixed principal component of maximal stimulus-specific variance. Solid lines are correct trial averages; dotted lines are incorrect trial averages. Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl). T: time of stimulus onset; D: decision time. **B:** Same as **A** but for the demixed principal component of maximal choice-specific variance. **C:** Left: Venn diagram showing percentages of unconstrained units (pooled over all models) with all possible combinations of coding types over time. Right: Distribution of coding types across unconstrained units over time.



Supplementary Figure 7. Effect of selective ablations on model dynamics and behavior: constrained vs unconstrained.

A-B: Model dynamics after selectively ablating linear coding units, step-perception coding units, step-choice coding units, or “other” units in the constrained (left columns) or unconstrained (right columns) populations. Post-ablation pseudo-population activity is projected onto the stimulus (**A**) and choice-coding (**B**) components identified in the control condition (top). Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl); solid and dashed lines correspond to correct and error trials. * indicates significant difference in mean absolute projections vs corresponding control condition (Dunnett’s test $p < 0.01$). T: time of stimulus onset; D: decision time. **C:** Behavioral performance of all models after selectively ablating categories of coding units. Bars represent means. * indicates significant difference in task accuracy vs control condition (Dunnett’s test $p < 0.01$).



Supplementary Figure 8. Effect of temporally restricted selective ablations on model dynamics and behavior: beginning vs end.

A-B: Model dynamics after selectively ablating linear coding units, step-perception coding units, step-choice coding units, or “other” units at the beginning of the trial (left columns) or the end of the trial (right columns). Post-ablation pseudo-population activity is projected onto the stimulus (**A**) and choice-coding (**B**) components identified in the control condition (top). Color scale corresponds to different mixture stimuli (%Sucrose/%NaCl); solid and dashed lines correspond to correct and error trials. * indicates significant difference in mean absolute projections vs corresponding control condition (Dunnett’s test $p < 0.01$). T: time of stimulus onset; D: decision time. The beginning is [T, T + 1.2 s]; the end is [D - 1.2 s, D]. **C:** Behavioral performance of all models after selectively ablating categories of coding units in the beginning or end of the trial. Bars represent means. * indicates significant difference in task accuracy vs control condition (Dunnett’s test $p < 0.01$).

Data availability

Experimental dataset available by request. Modeling dataset and code for all analyses available at: <https://github.com/llang6/linear-categorical> [↗](#).

Acknowledgements

Work supported by R01DC018227 from NIH/NIDCD (A.F.), 1UF1NS115779 from NIH/NINDS Brain Initiative (AF and GLC), American Association of University Women (AAUW) International Fellowship (C.Y.Z), K12GM102778 from NIH/NGM (J.M.B).

Additional information

Funding

Funder	Grant reference number
NIH Common Fund	R01DC018227
NIH Common Fund	1UF1NS115779
NIH Common Fund	K12GM102778
American Association of University Women (AAUW)	International Fellowship

Author ORCID iDs

Liam Lang: <https://orcid.org/0000-0002-5606-6370>

Giancarlo La Camera: <https://orcid.org/0000-0001-7834-6472>

Alfredo Fontanini: <https://orcid.org/0000-0003-4561-9563>

References

- Arieli E, Younis N, Moran A.** (2022) Distinct Progressions of Neuronal Activity Changes Underlie the Formation and Consolidation of a Gustatory Associative Memory. *Journal of Neuroscience* **42**:909-921 <https://doi.org/10.1523/JNEUROSCI.1599-21.2021> | PubMed
- Barak O.** (2017) Recurrent Neural Networks as Versatile Tools of Neuroscience Research. *Current Opinion in Neurobiology* **46**:1-6 <https://doi.org/10.1016/j.conb.2017.06.003> | PubMed
- Bohte SM** (2011) Error-Backpropagation in Networks of Fractionally Predictive Spiking Neurons. In: Honkela T, Duch W, Girolami M, Kaski S (Eds). *Artificial Neural Networks and Machine Learning – ICANN 2011* **6791** Berlin, Heidelberg: Springer Berlin Heidelberg. pp. 60-68 https://doi.org/10.1007/978-3-642-21735-7_8
- Bouaichi CG, Vincis R.** (2020) Cortical Processing of Chemosensory and Hedonic Features of Taste in Active Licking Mice. *Journal of Neurophysiology* **123**:1995-2009 <https://doi.org/10.1152/jn.00069.2020> | PubMed
- Buetfering C, Zhang Z, Pitsiani M, Smallridge J, Boven E, McElligott S, Häusser M.** (2022) Behaviorally Relevant Decision Coding in Primary Somatosensory Cortex Neurons. *Nature Neuroscience* **25**:1225-1236 <https://doi.org/10.1038/s41593-022-01151-0> | PubMed
- Churchland AK, Ditterich J.** (2012) New Advances in Understanding Decisions among Multiple Alternatives. *Current Opinion in Neurobiology* **22**:920-926 <https://doi.org/10.1016/j.conb.2012.04.009> | PubMed
- Cisek P, Puskas GA, El-Murr S.** (2009) Decisions in Changing Conditions: The Urgency-Gating Model. *The Journal of Neuroscience* **29**:11560-11571 <https://doi.org/10.1523/JNEUROSCI.1844-09.2009> | PubMed
- Cohen Z, DePasquale B, Aoi MC, Pillow JW** (2020) Recurrent Dynamics of Prefrontal Cortex during Context-Dependent Decision-Making. *bioRxiv* <https://doi.org/10.1101/2020.11.27.401539>

- DePasquale B, Churchland MM, Abbott LF (2016) Using Firing-Rate Dynamics to Train Recurrent Networks of Spiking Model Neurons. *arXiv* <https://doi.org/10.48550/ARXIV.1601.07620>
- Emiliani V, Entcheva E, Hedrich R, Hegemann P, Konrad KR, Lüscher C, Mahn M, Pan ZH, Sims RR, Vierock J, et al. (2022) Optogenetics for Light Control of Biological Systems. *Nature Reviews Methods Primers* **2**:55 <https://doi.org/10.1038/s43586-022-00136-4> | PubMed
- Fonseca E, De Lafuente V, Simon SA, Gutierrez R. (2018) Sucrose Intensity Coding and Decision-Making in Rat Gustatory Cortices. *eLife* **7**:e41152 <https://doi.org/10.7554/eLife.41152> | PubMed
- Fuglstad JG, Saldanha P, Paglia J, Whitlock JR (2023) Histological E-data Registration in Rodent Brain Spaces. *eLife* **12**:e83496 <https://doi.org/10.7554/eLife.83496> | PubMed
- Gardner MPH, Fontanini A. (2014) Encoding and Tracking of Outcome-Specific Expectancy in the Gustatory Cortex of Alert Rats. *Journal of Neuroscience* **34**:13000-13017 <https://doi.org/10.1523/JNEUROSCI.1820-14.2014> | PubMed
- Goltstein PM, Reinert S, Bonhoeffer T, Hübener M. (2021) Mouse Visual Cortex Areas Represent Perceptual and Semantic Features of Learned Visual Categories. *Nature Neuroscience* **24**:1441-1451 <https://doi.org/10.1038/s41593-021-00914-5> | PubMed
- Guo ZV, Hires SA, Li N, O'Connor DH, Komiyama T, Ophir E, Huber D, Bonardi C, Morandell K, Gutnisky D, et al. (2014a) Procedures for Behavioral Experiments in Head-Fixed Mice. *PLOS One* **9**:e88678 <https://doi.org/10.1371/journal.pone.0088678> | PubMed
- Guo ZV, Li N, Huber D, Ophir E, Gutnisky D, Ting JT, Feng G, Svoboda K. (2014b) Flow of Cortical Activity Underlying a Tactile Decision in Mice. *Neuron* **81**:179-194 <https://doi.org/10.1016/j.neuron.2013.10.020> | PubMed
- Jezzini A, Padoa-Schioppa C. (2024) Neuronal Activity in the Gustatory Cortex during Economic Choice. *The Journal of Neuroscience* **44**:e2150232024 <https://doi.org/10.1523/JNEUROSCI.2150-23.2024> | PubMed
- Jones LM, Fontanini A, Sadacca BF, Miller P, Katz DB (2007) Natural Stimuli Evoke Dynamic Sequences of States in Sensory Cortical Ensembles. *Proceedings of the National Academy of Sciences* **104**:18772-18777 <https://doi.org/10.1073/pnas.0705546104> | PubMed
- Jun JJ, Steinmetz NA, Siegle JH, Denman DJ, Bauza M, Barbarits B, Lee AK, Anastassiou CA, Andrei A, Aydin Ç, et al. (2017) Fully Integrated Silicon Probes for High-Density Recording of Neural Activity. *Nature* **551**:232-236 <https://doi.org/10.1038/nature24636> | PubMed
- Katz DB, Simon SA, Nicolelis MAL (2001) Dynamic and Multimodal Responses of Gustatory Cortical Neurons in Awake Rats. *The Journal of Neuroscience* **21**:4478-4489 <https://doi.org/10.1523/JNEUROSCI.21-12-04478.2001> | PubMed
- Kobak D, Brendel W, Constantinidis C, Feierstein CE, Kepecs A, Mainen ZF, Qi XL, Romo R, Uchida N, Machens CK (2016) Demixed Principal Component Analysis of Neural Population Data. *eLife* **5**:e10989 <https://doi.org/10.7554/eLife.10989> | PubMed
- Kogan JF, Fontanini A. (2024) Learning Enhances Representations of Taste-Guided Decisions in the Mouse Gustatory Insular Cortex. *Current Biology* **34**:1880-1892.e5 <https://doi.org/10.1016/j.cub.2024.03.034> | PubMed
- Kusumoto-Yoshida I, Liu H, Chen BT, Fontanini A, Bonci A. (2015) Central Role for the Insular Cortex in Mediating Conditioned Responses to Anticipatory Cues. *Proceedings of the National Academy of Sciences* **112**:1190-1195 <https://doi.org/10.1073/pnas.1416573112> | PubMed
- Lang L, La Camera G, Fontanini A. (2023) Temporal Progression along Discrete Coding States during Decision-Making in the Mouse Gustatory Cortex. *PLOS Computational Biology* **19**:e1010865 <https://doi.org/10.1371/journal.pcbi.1010865> | PubMed
- Li N, Chen S, Guo ZV, Chen H, Huo Y, Inagaki HK, Chen G, Davis C, Hansel D, Guo C, et al. (2019) Spatiotemporal Constraints on Optogenetic Inactivation in Cortical Circuits. *eLife* **8**:e48622 <https://doi.org/10.7554/eLife.48622> | PubMed

- Li Y, Guo Y, Zhang S, Deng S, Hai Y, Gu S. (2021) Differentiable Spike: Rethinking Gradient-Descent for Training Spiking Neural Networks. In: Advances in Neural Information Processing Systems 34. pp. 23426-23439
- Lin JY, Mukherjee N, Bernstein MJ, Katz DB (2021) Perturbation of Amygdala-Cortical Projections Reduces Ensemble Coherence of Palatability Coding in Gustatory Cortex. *eLife* **10**:e65766 <https://doi.org/10.7554/eLife.65766> | PubMed
- Livneh Y, Andermann ML (2021) Cellular Activity in Insular Cortex across Seconds to Hours: Sensations and Predictions of Bodily States. *Neuron* **109**:3576-3593 <https://doi.org/10.1016/j.neuron.2021.08.036> | PubMed
- Livneh Y, Ramesh RN, Burgess CR, Levandowski KM, Madara JC, Fenselau H, Goldey GJ, Diaz VE, Jikomes N, Resch JM, et al. (2017) Homeostatic Circuits Selectively Gate Food Cue Responses in Insular Cortex. *Nature* **546**:611-616 <https://doi.org/10.1038/nature22375> | PubMed
- Livneh Y, Sugden AU, Madara JC, Essner RA, Flores VI, Sugden LA, Resch JM, Lowell BB, Andermann ML (2020) Estimation of Current and Future Physiological States in Insular Cortex. *Neuron* **105**:1094-1111.e10 <https://doi.org/10.1016/j.neuron.2019.12.027> | PubMed
- Mahmood A, Steindler JR, Katz DB (2025) Perceptual Processing of Tastes Is Performed by the Amygdala-Cortical Loop. *bioRxiv* <https://doi.org/10.1101/2025.07.01.662567>
- Mahmood A, Steindler J, Germaine H, Miller P, Katz DB (2023) Coupled Dynamics of Stimulus-Evoked Gustatory Cortical and Basolateral Amygdalar Activity. *Journal of Neuroscience* **43**:386-404 <https://doi.org/10.1523/JNEUROSCI.1412-22.2022> | PubMed
- Maier JX, Katz DB (2013) Neural Dynamics in Response to Binary Taste Mixtures. *Journal of Neurophysiology* **109**:2108-2117 <https://doi.org/10.1152/jn.00917.2012> | PubMed
- Maxwell SE, Delaney HD, Kelley K. (2017) *Designing Experiments and Analyzing Data: A Model Comparison Perspective* (Third 3) New York: Routledge. <https://doi.org/10.4324/9781315642956>
- Mazzucato L, La Camera G, Fontanini A. (2019) Expectation-Induced Modulation of Metastable Activity Underlies Faster Coding of Sensory Stimuli. *Nature Neuroscience* **22**:787-796 <https://doi.org/10.1038/s41593-019-0364-9> | PubMed
- Mazzucato L, Fontanini A, La Camera G. (2015) Dynamics of Multistable States during Ongoing and Evoked Cortical Activity. *The Journal of Neuroscience* **35**:8214-8231 <https://doi.org/10.1523/JNEUROSCI.4819-14.2015> | PubMed
- Mazzucato L, Fontanini A, La Camera G. (2016) Stimuli Reduce the Dimensionality of Cortical Activity. *Frontiers in Systems Neuroscience* **10** <https://doi.org/10.3389/fnsys.2016.00011> | PubMed
- Mendoza G, Fonseca E, Merchant H, Gutierrez R. (2024) Neuronal Sequences and Dynamic Coding of Water-Sucrose Categorization in Rat Gustatory Cortices. *iScience* **27** <https://doi.org/10.1016/j.isci.2024.111287> | PubMed
- Miller P, Katz DB (2010) Stochastic Transitions between Neural States in Taste Processing and Decision-Making. *The Journal of Neuroscience* **30**:2559-2570 <https://doi.org/10.1523/JNEUROSCI.3047-09.2010> | PubMed
- Motulsky H, Christopoulos A. (2004) *Fitting Models to Biological Data Using Linear and Nonlinear Regression: A Practical Guide to Curve Fitting* Oxford University Press.
- Mukherjee N, Wachutka J, Katz DB (2019) Impact of Precisely-Timed Inhibition of Gustatory Cortex on Taste Behavior Depends on Single-Trial Ensemble Dynamics. *eLife* **8**:e45968 <https://doi.org/10.7554/eLife.45968> | PubMed
- Niessing J, Friedrich RW (2010) Olfactory Pattern Classification by Discrete Neuronal Network States. *Nature* **465**:47-52 <https://doi.org/10.1038/nature08961> | PubMed
- Pachitariu M, Sridhar S, Pennington J, Stringer C. (2024) Spike Sorting with Kilosort4. *Nature Methods* **21**:914-921 <https://doi.org/10.1038/s41592-024-02232-7> | PubMed

- Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, *et al.* (2019) PyTorch: An Imperative Style, High-Performance Deep Learning Library. In: Advances in Neural Information Processing Systems 32. <https://doi.org/10.48550/arxiv.1912.01703>
- Perich MG, Arlt C, Soares S, Young ME, Mosher CP, Minxha J, Carter E, Rutishauser U, Rudebeck PH, Harvey CD, *et al.* (2020) Inferring Brain-Wide Interactions Using Data-Constrained Recurrent Neural Network Models. *bioRxiv* <https://doi.org/10.1101/2020.12.18.423348>
- Rajan K, Harvey CD, Tank DW (2016) Recurrent Network Models of Sequence Generation and Memory. *Neuron* **90**:128-142 <https://doi.org/10.1016/j.neuron.2016.02.009> | PubMed
- Reinert S, Hübener M, Bonhoeffer T, Goltstein PM (2021) Mouse Prefrontal Cortex Represents Learned Rules for Categorization. *Nature* **593**:411-417 <https://doi.org/10.1038/s41586-021-03452-z> | PubMed
- Sadacca BF, Mukherjee N, Vladusich T, Li JX, Katz DB, Miller P. (2016) The Behavioral Relevance of Cortical Neural Ensemble Responses Emerges Suddenly. *The Journal of Neuroscience* **36**:655-669 <https://doi.org/10.1523/JNEUROSCI.2265-15.2016> | PubMed
- Samuelsen CL, Gardner MPH, Fontanini A. (2012) Effects of Cue-Triggered Expectation on Cortical Processing of Taste. *Neuron* **74**:410-422 <https://doi.org/10.1016/j.neuron.2012.02.031> | PubMed
- Shamash P, Carandini M, Harris KD, Steinmetz NA (2018) A Tool for Analyzing Electrode Tracks from Slice Histology. *bioRxiv* <https://doi.org/10.1101/447995>
- Shuler MG, Bear MF (2006) Reward Timing in the Primary Visual Cortex. *Science* **311**:1606-1609 <https://doi.org/10.1126/science.1123513> | PubMed
- Song HF, Yang GR, Wang XJ (2016) Training Excitatory-Inhibitory Recurrent Neural Networks for Cognitive Tasks: A Simple and Flexible Framework. *PLOS Computational Biology* **12**:e1004792 <https://doi.org/10.1371/journal.pcbi.1004792> | PubMed
- Stapleton JR (2007) Ensembles of Gustatory Cortical Neurons Anticipate and Discriminate between Tastants in a Single Lick. *Frontiers in Neuroscience* **1**:161-174 <https://doi.org/10.3389/neuro.01.1.1.012.2007> | PubMed
- Valente A, Pillow JW, Ostojic S. (2022) Extracting Computational Mechanisms from Neural Data Using Low-Rank RNNs. *Advances in Neural Information Processing Systems* **35**:24072-24086 <https://doi.org/10.52202/068431-1748>
- Vincis R, Chen K, Czarnecki L, Chen J, Fontanini A. (2020) Dynamic Representation of Taste-Related Decisions in the Gustatory Insular Cortex of Mice. *Current Biology* **30**:1834-1844.e5 <https://doi.org/10.1016/j.cub.2020.03.012> | PubMed
- Vincis R, Fontanini A. (2016) A Gustocentric Perspective to Understanding Primary Sensory Cortices. *Current Opinion in Neurobiology* **40**:118-124 <https://doi.org/10.1016/j.conb.2016.06.008> | PubMed
- Yamamoto T, Yuyama N, Kato T, Kawamura Y. (1985) Gustatory Responses of Cortical Neurons in Rats. II. Information Processing of Taste Quality. *Journal of Neurophysiology* **53**:1356-1369 <https://doi.org/10.1152/jn.1985.53.6.1356> | PubMed
- Yang GR, Wang XJ (2020) Artificial Neural Networks for Neuroscientists: A Primer. *Neuron* **107**:1048-1070 <https://doi.org/10.1016/j.neuron.2020.09.005> | PubMed
- Zheng CY, Blackwell JM, Fontanini A. (2025) Deficits in Taste-Guided Behaviors and Central Processing of Taste in the Transgenic TDP-43Q331K Mouse Model of Frontotemporal Dementia. *Neurobiology of Disease* **207**:106850 <https://doi.org/10.1016/j.nbd.2025.106850> | PubMed

Peer reviews

Reviewer #1 (Public review):

The manuscript provides several important findings that advance our current knowledge about the function of the gustatory cortex (GC). The authors used high density electrophysiology to record neural activity during a sucrose/NaCl mixture discrimination task. They observed population-based activity capable of representing different mixtures in a

linear fashion during the initial stimulus sampling period as well as representing the behavioral decision (i.e., lick left or right) at a later time point. Analyzing this data at the single neuron level, they observed functional subpopulations capable of encoding the specific mixture (e.g., 45/55), tastant (e.g., sucrose), and behavioral choice (e.g., lick left). To test the functional consequences of these subpopulations, they built a recurrent neural network model in order to "silence" specific functional subpopulations of GC neurons. The virtual ablation of these functional subpopulations altered virtual behavioral performance in a manner predicted by the subpopulation's presumed contribution.

Strengths:

Building a recurrent neural network model of the gustatory cortex allows the impact of the temporal sequence of functionally identifiable populations of neurons to be tested in a manner not otherwise possible. Specifically, the author's model links neural activity at the single neuron and population level with perceptual ability. The electrophysiology methods and analyses used to shape the network model are appropriate. Overall, the conclusions of the manuscript are well supported.

Weaknesses:

One minor weakness is the mismatch between the neural analyses and behavioral data. Neural analyses (i.e. population activity trajectories) indicate a separation of the neural activity associated with each mixture. Given this analysis, one might expect the psychometric curve to have a significantly steeper slope. One potential explanation is the concentration of the stimuli utilized in the mixture discrimination task. The authors utilize equivalent concentrations, rather than intensity matched concentrations. In this case, a single stimulus can (theoretically) dominant the perception of a mixture resulting in a biased behavioral response despite accurate concentration coding. Given the difficulty of iso-intensity matching concentrations, this concern is not paramount.

<https://doi.org/10.7554/eLife.109313.2.sa3>

Reviewer #2 (Public review):

Lang et al. investigate the contribution of individual neuronal encoding of specific task features to population dynamics and behavior. Using a taste based decision-making behavioral task with electrophysiology from the mouse gustatory cortex and computational modeling, the authors reveal that neurons encoding sensory, perceptual, and decision-related information with linear and categorical patterns are essential for driving neural population dynamics and behavioral performance. Their findings suggest that individual linear and categorical coding units have a significant role in cortical dynamics and perceptual decision-making behavior.

Overall, the experimental and analytical work is of very high quality, and the findings are of great interest to the taste coding field, as well as to the broader systems neuroscience field.

I initially had some suggestions for further analyses to clarify the contribution of constrained and unconstrained units. In the revised version, the authors have performed all the suggested analyses, further strengthening their conclusions.

<https://doi.org/10.7554/eLife.109313.2.sa2>

Reviewer #3 (Public review):

Primary taste cortex neurons show a variety of dynamic response profiles during taste decision making tasks, reflecting both sensory and decision variables. In the present study,

Lang et al., set out to determine how neurons with distinct response profiles contribute to perceptual decisions about taste stimuli.

The methods with regard to the behavioral task and electrophysiological recordings/data analysis are straightforward, solid and appropriate. The computational model is presented in a clear and conceptually intuitive manner, although the details are outside of my area of expertise.

The experimental design features a simple 2-alternative forced choice task that yielded clear psychometric curves across a range of stimuli. In vivo recordings were performed using neuropixels and yielded an appropriate sample of single neuron responses. The strength of the model lies in the fact that it consists of single neurons whose response profiles mimic those recorded in vivo, and allows neuron-selective manipulation.

By virtually lesioning specific subsets of neurons in the network, the authors demonstrate that a relatively small populations of neurons with specific tuning profiles were sufficient to produce the observed neural dynamics and behavioral responses. This effect was selective as lesioning other responsive neurons did not affect overall response dynamics or performance.

These findings provide new insight into the relation between the response profiles of single neurons in sensory cortex, their population-level activity dynamics, and the perceptual decisions they inform.

The approach is particularly innovative as it uses computational modeling to target functionally-defined "cell types", which cannot necessarily be targeted by more conventional genetic approaches.

<https://doi.org/10.7554/eLife.109313.2.sa1>

Author response:

The following is the authors' response to the original reviews.

Reviewer #1 (Public review):

This manuscript provides several important findings that advance our current knowledge about the function of the gustatory cortex (GC). The authors used high-density electrophysiology to record neural activity during a sucrose/NaCl mixture discrimination task. They observed population-based activity capable of representing different mixtures in a linear fashion during the initial stimulus sampling period, as well as representing the behavioral decision (i.e., lick left or right) at a later time point. Analyzing this data at the single neuron level, they observed functional subpopulations capable of encoding the specific mixture (e.g., 45/55), tastant (e.g., sucrose), and behavioral choice (e.g., lick left). To test the functional consequences of these subpopulations, they built a recurrent neural network model in order to "silence" specific functional subpopulations of GC neurons. The virtual ablation of these functional subpopulations altered virtual behavioral performance in a manner predicted by the subpopulation's presumed contribution.

Strengths:

Building a recurrent neural network model of the gustatory cortex allows the impact of the temporal sequence of functionally identifiable populations of neurons to be tested in a manner not otherwise possible. Specifically, the author's model links neural activity at the single neuron and population level with perceptual ability. The electrophysiology methods and analyses used to shape the network model are appropriate. Overall, the conclusions of the manuscript are well supported.

Weaknesses:

One potential concern is the apparent mismatch between the neural and behavioral data. Neural analyses indicate a clear separation of the activity associated with each mixture that is independent of the animal's ultimate choice. This would seemingly indicate that the animals are making errors despite correctly encoding the stimulus. Based solely on the neural data, one would expect the psychometric curve to be more "step-like" with a significantly steeper slope. One potential explanation for this observation is the concentration of the stimuli utilized in the mixture discrimination task. The authors utilize equivalent concentrations, rather than intensity-matched concentrations. In this case, a single stimulus can (theoretically) dominate the perception of a mixture, resulting in a biased behavioral response despite accurate concentration coding at the single neuron level. Given the difficulty of iso-intensity matching concentrations, this concern is not paramount. However, the apparent mismatch between the neural and behavioral data should be acknowledged/addressed in the text.

We thank the Reviewer for the insightful comments and thoughtful suggestions. Our electrophysiological recordings show that GC dynamically encodes stimulus concentration of mixture elements, dominant perceptual quality, and decisions of directional lick. With regard to the encoding of mixtures, the clear separation of activity associated with each mixture (Figure 3) is present at a trial-averaged pseudo-population level, and average activities associated with more similar, intermediate mixtures are closer to each other in this space. At a single trial level activities evoked by similar, intermediate mixtures are much harder to separate. This increased similarity can lead to behavioral errors resulting from either incorrect encoding of the stimulus or from the inability to interpret the stimulus to guide the correct decision. The psychometric function, which shows that more distinct stimuli (100/0 vs 0/100) lead to fewer mistakes than more ambiguous, intermediate mixtures (55/45 vs 55/45), is consistent with the increased ambiguity of responses to intermediate mixtures.

The Reviewer is correct that there could be a slight mismatch in the perceived intensity of the mixture components. This mismatch could be the reason for the slight asymmetry in our psychometric function (Figure 1B). However, it is not uncommon for mice in these 2AC tasks to also have a motor laterality bias in their responses that manifests itself for the more ambiguous stimuli. We chose not to model this bias given its subtlety and its unknown origin. Rather, we chose to model an ideal scenario in which stimuli have matched intensity and no motor bias exists. In the revised manuscript we discuss this issue.

Reviewer #1 (Recommendations for the authors):

(1) The apparent mismatch between neural and behavioral data. I am providing more details in this section to hopefully better illustrate my concern.

(a) Based on the author's psychometric curve, sucrose appears to be a more salient signal causing the behavior to be shifted (e.g., a 50/50 mixture results in a >60% predicted behavioral performance). If both sucrose and salt were intensity-matched, a 50/50 mixture should result in a behavioral performance near 50%. The increased salience of sucrose could cause the animals to have lower overall performance despite accurate neural encoding. Alternatively, certain animals could display a strong side bias, skewing the data slightly. These issues have seemingly been fixed in the model data, which displays a more balanced psychometric curve. Accordingly, the model data seemingly displays a larger shift in error trials as compared to correct trials (Figure 6A).

The reviewer is correct in observing that the average experimental psychometric curve in Figure 1B shows a slight shift in favor of the sucrose side with a 50/50 mixture. We fit psychometric curves to each session and the mean value of $P(\text{Sucrose choice} \mid \text{Stimulus} =$

50/50) across sessions was significantly different from 0.5 (one-sample t-test, $p = 0.003$), with 5 probabilities below 0.5 and 18 above it.

This slight bias could be attributed to a slight mismatch in the perceived intensity of the mixture components and/or lateral motor biases. In any case, it is subtle and its origins were not a focus of this study.

Models were not trained to match the animals' psychometric curves, but rather to choose correctly in an ideal scenario where stimuli have matched intensities. This explains why the model simulations lack the bias observed in animal behavior data.

We do not believe that there is a mismatch between the experimental behavioral and neural data, as trial-averaged pseudo-population trajectories are farther in neural space for more discriminable stimuli and closer in neural space for more similar stimuli, consistent with behavioral performance that is high for more discriminable stimuli and low for more similar stimuli. Moreover, as the model also shows, a clear separation of trial-averaged trajectories still results in a sigmoidal performance function for trial-to-trial behavior.

Finally, subtle behavioral biases would not necessarily be expected to appear in our dPCA analyses since we used this technique to find a single axis that best separates all stimuli conditions regardless of choice when the pseudo-population data are projected upon it. Additional modes of activity that explain less overall variance might better reflect biases.

(b) Although I am not an expert at these analyses, I wonder whether the elevated bump (i.e., >0) in Figure 3C of the 55/45 mixture that occurs early in the stimulus presentation further supports the hypothesis mentioned above and could indicate an early signal of salience/increased intensity?

The reviewer is correct that the 55/45 trajectory features a brief positive wave right after stimulus delivery before going negative. While this may be related to stimuli not being explicitly balanced for intensity, it could also reflect a signal related to ambiguity or balanced mixtures. We are hesitant to interpret this positive deflection as conclusive evidence of a bias in neural activity, given its short duration and the natural variability of neural signals.

(2) The increase in step-perception neurons after the decision period is confusing (Figure 4C). The text states (line 246) "the analysis reveals a small and time-invariant proportion of step-perception neurons". However, the proportion doubles after the decision-making process, which is seemingly a significant change. Why does this occur? This observation is noticeably missing from the network data. Could it be attributed to a mislabeling of "step-choice" neurons, given the correlation between the left/right decision and sweet/salty? Either way, it is very noticeable and should be addressed.

We cannot be sure of the reason for the increase in step-perception neurons after decisions. One possibility is that they are acting as feedback for learning, encoding the percept to compare with choice and outcome to improve performance. The model, which presumably learns the task differently from the animals, does not seem to leverage this signal for its own learning. We have modified the text, now referring to a "small but consistently present proportion" of step-perception neurons, and included this proposed explanation in the Discussion.

(3) *Optional: I think the authors are missing an opportunity to analyze the temporal aspect of this multiplex code using their network-based modeling approach. A significant proportion of neurons fall into different categories (i.e., step-perception/linear, etc.) at different time points. However, the virtual ablation experiments remove any neuron that falls into one of these categories at any time. By limiting the cell-specific virtual ablation to specific time windows, you could (I think) provide stronger evidence for the temporal sequence of the encoding of these perceptual aspects.*

This was an excellent suggestion for an additional modeling experiment, so we performed it. A new supplemental figure (Figure S8) and additional text in the revised manuscript showcase the results. In summary:

In terms of behavioral results, ablating the linear coding units in the beginning (that is, silencing all units that are labeled linear in any bin within the first 1.2 s after stimulus onset for the entirety of the 1.2 s) significantly reduces performance, as does ablating the step-perception or step-choice coding units at the end (1.2 s prior to choice). The remaining combinations of coding type and timing of the ablation do not affect performance.

Regarding the dynamics of coding types (compare Figure 7A), stimulus coding activity was significantly blunted only by ablating the linear coding units in the beginning, whereas choice coding activity was diminished by ablating the choice coding units at the end or by ablating the linear coding units at either the beginning or the end.

Reviewer #2 (Public review):

Lang et al. investigate the contribution of individual neuronal encoding of specific task features to population dynamics and behavior. Using a taste-based decision-making behavioral task with electrophysiology from the mouse gustatory cortex and computational modeling, the authors reveal that neurons encoding sensory, perceptual, and decision-related information with linear and categorical patterns are essential for driving neural population dynamics and behavioral performance. Their findings suggest that individual linear and categorical coding units have a significant role in cortical dynamics and perceptual decision-making behavior.

Overall, the experimental and analytical work is of very high quality, and the findings are of great interest to the taste coding field, as well as to the broader systems neuroscience field.

I have a couple of suggestions to further enhance the authors' important conclusions:

My main comment is the distinction between constrained and unconstrained units. The authors train a small percentage of units to match the real neural data (constrained units), and then find some unconstrained units that are similar to the real neural data and some that are not. As far as I could tell, the relative fraction of constrained and unconstrained units in the trained RNN is not reported; I assume the constrained ones are a much smaller population, but this is unclear. The selection of different groups of neurons for the RNN ablation experiments appears to be based on their response profiles only. Therefore, if I understood correctly, both constrained and unconstrained units are ablated together for a given response category (e.g., linear or step-perception). It would be useful, therefore, to separately compare the effects of constrained vs. unconstrained RNN units.

We thank the Reviewer for the constructive feedback. The Reviewer is correct that ablations were carried out with respect to response categories only and included both constrained and unconstrained units.

The ratio of total units to constrained units was fixed at 5.88, thus constrained units were ~17% of the network and unconstrained units were ~83%. This value is specified in the Methods (RNN: Components and dynamics), but we have reported it in the Results of the revised manuscript for clarity.

We have also edited the Methods because they wrongly stated that the ratio of unconstrained (rather than total) units to constrained units was 5.88.

Specifically:

(1) For the analyses in the initial version of the manuscript, the authors should specify how many units in each ablation category are constrained and unconstrained.

In the revised manuscript, we have specified the fractions of constrained and unconstrained units within each response category. For convenience, they are reported here: linear = 194 constrained and 691 unconstrained units; step-perception = 147 constrained and 840 unconstrained units; step-choice = 129 constrained and 814 unconstrained units; “other” = 353 constrained and 1739 unconstrained units.

(2) The authors should repeat Figure 6, but only for unconstrained units to test how much of the effects in the initial version of Figure 6 are driven by constrained vs. unconstrained RNN units.

In the revised version we have included two additional supplemental figures (Figures S5-6) where the analyses of Figure 6 are carried out separately for constrained and unconstrained units. In short, the results for the constrained units strongly resemble those for the experimental data, while the results for the unconstrained units strongly resemble those for all model units.

(3) The authors should repeat Figure 7, but performing ablations separately on the constrained and unconstrained units to examine how the network behaves in each case and the resulting “behavioral” effect.

The revised version includes a supplemental figure (Figure S7) with the results of these additional ablation simulations.

In summary:

In terms of behavioral performance, the prior results showing that ablating linear, step-perception, or step-choice units significantly impairs performance, while ablating “other” has no significant effect, hold even if ablation is restricted to only constrained or only unconstrained units. There is a significant main effect of constrained vs unconstrained; on average, ablating the unconstrained population impairs performance more, most likely due to their larger population size.

In terms of dynamics, to impair stimulus coding by ablating step-choice units, you must ablate them all; to impair stimulus coding by ablating linear or step-perception units, however, ablating just the unconstrained ones suffices. As before, ablating linear, step-perception, or step-choice units significantly impairs choice coding activity, while ablating “other” units does not; these results hold even if ablation is restricted to only constrained or only unconstrained units. Finally, there is again a significant main effect of constrained vs unconstrained; on average, ablating the unconstrained population impairs dynamics more, most likely due to the larger population size.

Reviewer #2 (Recommendations for the authors):

(1) In addition to panel 5B, it would be informative to show data from individual mice and the corresponding RNNs trained on each mouse, to assess how closely they match. If available, including one representative example of a good match and one of a less accurate match would help the reader get a better sense of the data.

Figure 5B shows the average behavioral performance of the model. Individual models were not trained directly on the psychometric curves of experimental sessions; they were trained to perform the task correctly. After successful training, model simulations were run with input noise to be able to produce a sigmoidal psychometric curve. However, although the input noise was tuned to capture the overall correct rate of the corresponding experimental session, we did not attempt to match the details of the psychometric curve. See also the next reply.

(2) In addition to panel 5C, it would be useful to add examples of experimentally observed PSTHs and the corresponding activity trajectory for the units in the RNN trained to match them, for all the other coding patterns (step-perception and step-choice).

We note that the PSTH in 5C is not an example of a linear coding unit as the Reviewer implies, but simply one with a good fit, and here the model's output was produced in the absence of input noise. In order to classify step-perception and step-choice responses one needs error trials, but the model was trained without this input noise that induces errors (and produces a sigmoidal psychometric function) to match experimental PSTHs from correct trials only. Post-training simulations were then run with input noise to induce error trials, and model unit response profiles were classified based on this. However, there is no guarantee that error trials in the model match the error trials in the experiment; therefore, step-perception and step-choice units in the model may or may not be step-perception and step-choice units in the data. Despite this limitation, the revised manuscript includes additional examples, in Figure S2, of experimentally observed PSTHs and their corresponding model activity, to supplement Figure 5C and provide a better sense of the goodness-of-fit.

(3) Electrophysiological data in Figure 2 - It would be helpful to provide statistics on how many neurons change their activity in each session.

In the revised manuscript we have included across-session statistics for proportions of neurons that are taste-responsive and that show decision preparatory activity. We have also included tables (Tables S1 and S3) with the numbers of neurons that are taste-responsive and that show preparatory activity for each session in the experimental and model data.

(4) Peak auROC selection - How was the peak auROC selected? Selecting only one bin for the peak could be potentially problematic and may result in the incorrect identification of an outlier that does not faithfully represent the neuron's overall activity. The peak selection could instead be based on several consecutive bins showing a consistent trend. If this approach was already implemented, the authors should explicitly describe it in the Methods section.

Peak auROC was selected from a single bin (with average duration about 50ms). While it is true that this may result in outlier neurons that transiently prefer one stimulus strongly but more consistently prefer the other, we opted for a simple criterion to sort the neurons into two categories for visualization. Adopting more stringent criteria that consider multiple bins may result in neurons that cannot be placed in either category, and we wanted a way to examine the entire pseudo-population. Also, the entire auROC trace is visualized in the heatmap, so potential outliers are not hidden and can be assessed by eye.

Reviewer #3 (Public review):

Primary taste cortex neurons show a variety of dynamic response profiles during taste decision-making tasks, reflecting both sensory and decision variables. In the present study, Lang et al. set out to determine how neurons with distinct response profiles contribute to perceptual decisions about taste stimuli.

The methods, with reference to the behavioral task and electrophysiological recordings/data analysis, are straightforward, solid, and appropriate. The computational model is presented in a clear and conceptually intuitive manner, although the details are outside of my area of expertise.

The experimental design features a simple 2-alternative forced-choice design that yielded clear psychometric curves across a range of stimuli. In vivo recordings were performed using Neuropixels and yielded an appropriate sample of single neuron responses. The strength of the model lies in the fact that it consists of single neurons whose response profiles mimic those recorded in vivo, and allows neuron-selective manipulation.

By virtually lesioning specific subsets of neurons in the network, the authors demonstrate that a relatively small population of neurons with specific tuning profiles was sufficient to produce the observed neural dynamics and behavioral responses. This effect was selective as lesioning other responsive neurons did not affect overall response dynamics or performance.

These findings provide new insight into the relation between the response profiles of single neurons in sensory cortex, their population-level activity dynamics, and the perceptual decisions they inform.

The approach is particularly innovative as it uses computational modeling to target functionally-defined "cell types", which cannot necessarily be targeted by more conventional genetic approaches.

We thank the Reviewer for the positive assessment of our study.

Reviewer #3 (Recommendations for the authors):

(1) Introduction: I'm missing a clearly stated specific hypothesis and what is predicted on the basis of that hypothesis. What is the alternative?

The null hypothesis is that single neuron activity patterns, even when clearly structured, do not matter for population activity or behavior. Alternatively, they do matter for these phenomena, and our model supports the alternative hypothesis. We have made this hypothesis clearer in the Introduction.

(2) Discussion: Much of the text is a recap of the Introduction and Results sections. Please elaborate on the specific insights gained from the findings. The idea that tuned neurons in the sensory cortex are the basis for perception and perceptual decisions concerning the features being represented by those neurons is generally accepted. What the present study adds to this insight could be described more explicitly. On the other hand, the idea that small populations of tuned neurons are responsible for perception of taste/perceptual decisions about taste appears in contrast with previous accounts where stimulus features/decisions are reflected in correlated changes in activity across distributed populations of taste cortical neurons, including ones that are not necessarily tuned or even overtly responsive. How do the present findings relate to this idea?

This is a very good point about reconciling these findings with past ones that have focused on coordinated changes across ensembles of neurons, i.e., metastable dynamics of internal (hidden) states. There is a brief mention of metastability toward the end of the Discussion, but we agree it deserves elaboration.

This work does emphasize single unit activity, but in the context of, and as relevant to, population activity. We believe that the findings and frameworks of previous studies and those presented here are compatible rather than mutually exclusive. There is no reason why neurons with the coding patterns we studied here cannot coordinate with others to participate in the formation of different metastable states. The question of which—neurons with specific response profiles, or ensemble activity patterns that may involve these neurons?—is necessary and sufficient for producing perception and behavior during the mixture-based decision-making task is interesting but rather difficult to answer because of the single units' contribution to both alternatives. One would need to utilize a manipulation that disrupts ensemble coordination without disrupting single unit activity to differentiate between them. We have made these points clearer in the Discussion.

(3) Results: RNNs were based on data from single sessions -- how many neurons of each tuning type were observed in each session? In particular, there were 23 sessions but only 25 neurons total tuned to choice, suggesting that modelled choice neurons were based on ~1 neuron.

The revised manuscript includes the session-by-session breakdown of response types for both experiment and model in two supplementary tables (Tables S2 and S4). We note that there are 25 neurons tuned to choice during the last 500 ms of the trial prior to decision, but 114 out of 626 neurons in total are tuned to choice in some time bin in the experimental data.

(4) Minor: Indicate the time windows used for analysis of stimulus sampling, delay, and choice on the figures.

The revised manuscript now includes the illustration of sampling and delay windows in Figure 2C-D, since we averaged the values over these windows for use in a 2-way ANOVA. All other figures either are associated with bin-by-bin analyses and have the first central and lateral licks (T and D) indicated, or have the time windows specified (e.g., Figure 4B, which uses [T, T + 0.5 s] and [D - 0.5 s, D]).

<https://doi.org/10.7554/eLife.109313.2.sa0>