1
2
3
4
5 **Decoding mouse behavior to explain single-trial decisions and their relationship with**
6 **neural activity.**
7
8
9
10 *Yves Weissenberger, Andrew J. King\*, Johannes C. Dahmen\**
11
12
13 Department of Physiology, Anatomy and Genetics, University of Oxford, Oxford, UK.
14
15
16 * These authors jointly supervised this work
17
18 Correspondence should be addressed to Y.W. (yves.weissenberger@dpag.ox.ac.uk) or J.C.D.
19 (johannes.dahmen@dpag.ox.ac.uk).
20
21
22
23
24
25
26
27 **Abstract**
28 Models of behavior typically focus on sparse measurements of motor output over long
29 timescales, limiting their ability to explain momentary decisions or neural activity. We developed
30 data-driven models relating experimental variables to videos of behavior. Applied to mouse
31 operant behavior, they revealed behavioral encoding of cognitive variables. Model-based
32 decoding of videos yielded an accurate account of single-trial behavior in terms of the
33 relationship between cognition, motor output and cortical activity.
34
35
36
37
38
39
40
41
42
43
44

**Main Text**

Advances in neural recording technologies have enabled activity to be measured from thousands of neurons simultaneously[1,2]. By eliminating the need for averaging activity across trials, these methods are providing unprecedented insights into neural function. But to fully realize their promise, we also require similarly comprehensive descriptions of behavior that can be used to bridge the gap between neural activity and function.

However, even in highly-controlled experimental settings, such as during a sensory decision-making task, quantitative descriptions of behavioral variability remain elusive[3,4]. Analyses of session-level choice-statistics have shown that decisions are influenced by a variety of factors[5,6] . Nevertheless, it remains extremely challenging to identify the factors underlying single-trial decisions from currently available behavioral readouts. This severely limits the functional interpretation of brain activity, which often relies on such behavioral readouts to link neural activity to cognitive processes.

The interpretation of neural activity is further complicated by correlations between experimental variables (e.g. cognitive variables or environmental stimuli) and motor output. Indeed, such correlations can confound the neural encoding of an experimental variable like a decision with the encoding of the associated motor output, i.e. the enactment of the decision.

One approach to overcoming these issues is the detailed quantitative study of behavior[4]. Classical approaches[7] focus on simple measures (e.g. aggregate choice-statistics) that are easy to relate back to experimental variables. However, these measures lack the capacity or temporal resolution that is required to robustly link neural activity to the computations underpinning trial-by-trial behavior. Although recent approaches have begun to address these shortcomings by performing unsupervised decompositions of detailed behavioral measurements[8,9] , their output can be difficult to relate to experimental variables, thereby limiting their scope.

We sought a novel and generally applicable approach to the challenge of quantifying behavior which combines the strengths of previous methods. We took a data-driven approach and developed statistical models of dense behavioral measurements. Our objective was to find representations of behavior that can account for an animal's motor output whilst remaining easily relatable to cognitive and stimulus-related variables. Crucially, we attempted to find such representations directly in the data, without *a priori* knowledge. In doing so, we aimed to extract a comprehensive and interpetable account of behavior that can support detailed analysis of neural activity.

We analyzed video data from head-fixed mice (n = 11 sessions from 6 mice) performing a sound detection task **(Fig. 1a)**, and used variational autoencoders, which are Bayesian latent-variable models (LVM)[10,11], as a starting point for modelling animals' motor output. The aim of the model was to find low-dimensional representations of the video data that enable frame-by-frame reconstructions at pixel-level resolution **(Fig. 1b i)**.

Models of behavior are useful only to the extent that they can be related to experimental variables, such as an animal's decisions or the underlying neural activity. We therefore formalized the notion of relatability as linear predictability from these variables. This yielded a novel model, which we refer to as a behavioral autoencoder (BAE), the cost function of which is augmented with an additional penalty term. This term encourages learning a representation of behavior that

2

88  is explicable in terms of *a priori* defined variables of interest **(Fig. 1b ii)** (see *Methods*). We then
89  fitted this model to videos acquired during task performance.
90       The sound detection task provided a rich set of observed and hidden variables **(Fig. 1c)**,
91  which may explain momentary variations in animals' motor output. We therefore used both sets
92  of variables (henceforth referred to collectively as experimental variables) to augment the model's
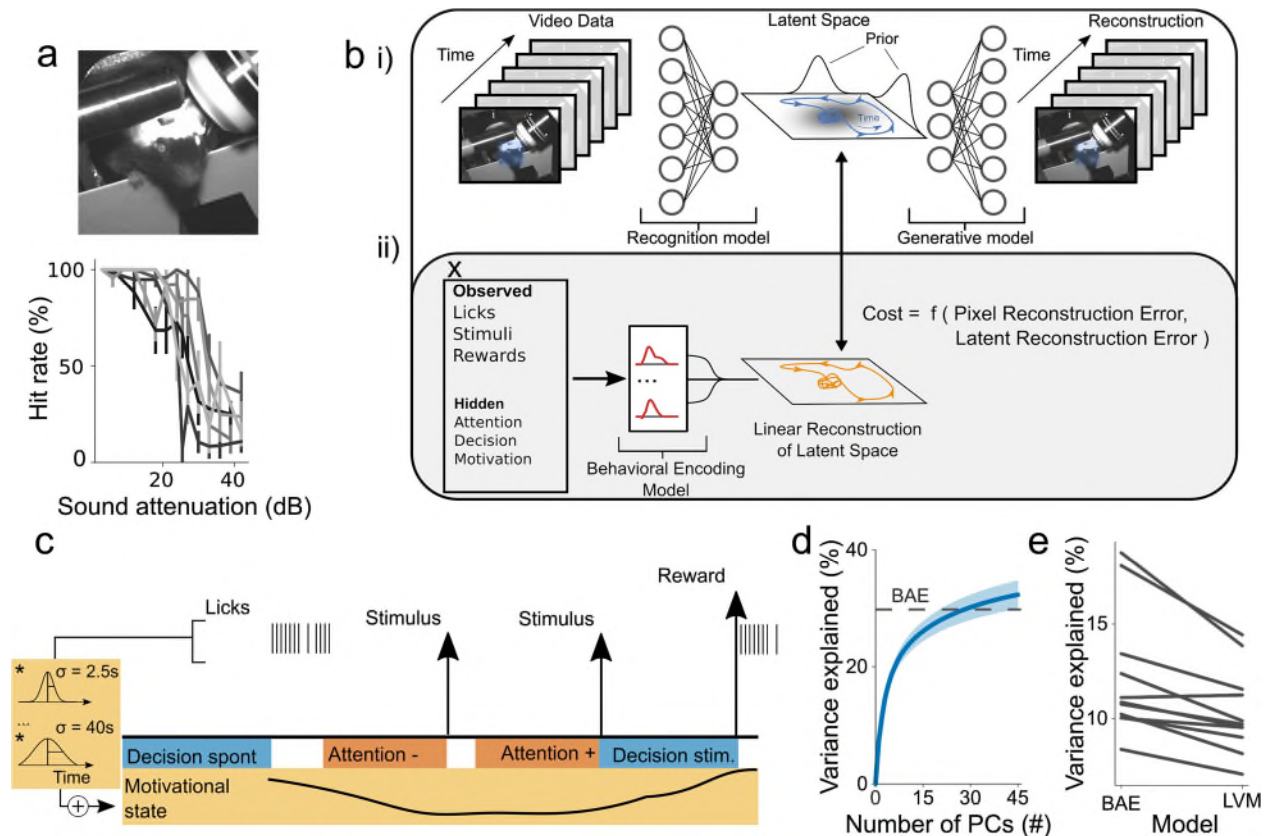93  cost function.
94



95
96
97       **Figure 1** Model Structure and performance. **a) (*top*)** Image of a mouse in the experimental
98  setup. **(*bottom*)** Example psychometric functions ($\pm$95% binomial confidence intervals)
99  illustrating performance in the sound detection task (each curve depicts performance of one
100 mouse in a single session; all curves are from different mice). **b)** Schematic of the LVM and
101 BAE. **(i)** The LVM is parameterized by two sequential deep neural networks. The first network
102 parameterizes a recognition model that maps from video data to a low-dimensional latent space.
103 The second network parameterizes a generative model which maps from the latent space back
104 into pixel space and reconstructs the video data. **(ii)** The BAE encompasses the LVM and a
105 behavioral encoding model that maps experimental variables into an approximation of the latent
106 space. This is used to encourage latent representations to be linearly predictable from
107 experimental variables **x** by an additional penalty term, which structures representations in the
108 latent space. **c)** Schematic illustrating the definition of hidden variables (see *Methods*). Briefly, an
109 animal was considered attentive on a given trial if the stimulus was of low intensity and the trial
110 was a hit-trial. It was considered inattentive on a given trial if the stimulus was of low intensity and
111 the trial was a miss-trial. An animal was considered to engage in 'stimulus-driven' licking if a

112  stimulus occurred in a 540-ms window preceding the onset of a lick bout; otherwise the licking
113  was considered to be 'spontaneous'. A high lick rate was interpreted to be indicative of reward
114  seeking and, thus, a state of high motivation. Motivational state regressors were created by
115  convolving licks with a series of Gaussian filters that were fitted individually and then summed.
116  Relative timescales across elements of the figure are not to scale. **d)** Performance of the BAE
117  (dashed line; latent states were inferred using the recognition model) compared with a principal
118  component analysis (PCA) based reconstruction (mean $\pm 2$ s.e.m) as a function of number of
119  PCs. Here, BAE reconstructions used the recognition model. **e)** Comparison of the LVM and the
120  BAE's ability to reconstruct videos using the behavioral encoding model (paired-sample t-test; *p*
121  $= 4.1 \cdot 10^{-4}$).

122  _____
123
124
125

126          To assess the model's performance, we quantified the reconstruction quality and capacity
127  of the experimental variables to explain behavioral latent states. Qualitative and quantitative
128  analyses revealed accurate reconstruction of the video data (mean $r^2$ = 30%, s.e.m = 3%)
129  **(Supplementary Fig. 1, Supplementary Video 1)**. Quantitatively, a 10-dimensional BAE
130  outperformed optimal linear methods, which required three-fold greater dimensionality to account
131  for the same variance **(Fig. 1d, Supplementary Fig. 2a)**. Importantly, learned representations
132  were highly interpretable, as assayed by measuring their predictability from experimental
133  variables **(Supplementary Fig. 2b)**. Furthermore, augmentation of the cost function in the BAE
134  significantly improved this predictability over that provided by the LVM **(Fig. 1e, Supplementary
135  Fig. 2b)**. Together, these findings suggest that the model learned comprehensive and
136  interpretable representations of the animals' behavior.
137          We then asked which experimental variables were encoded (i.e. expressed) in the
138  animals' behavior by quantifying the capacity of individual variables to explain behavioral latent
139  states. Although we found that all variables are encoded in behavior **(Fig. 2a),** this may arise
140  simply because many of them are correlated. We therefore quantified the effect of excluding
141  subsets of regression parameters, relating to a single experimental variable, on model-fit quality
142  (*see Methods*). This revealed that only a subset of variables uniquely accounted for variance in
143  the data **(Fig. 2b)**. Time into session accounted for most variance, reflecting the fact that the
144  animals' resting posture gradually changed over the course of the session. Additionally, we
145  consistently found that the animals' motivational state (operationalized as a smoothed lick time
146  series, **Fig. 1c;** see *Methods*) was explicitly encoded in behavior **(Supplementary Fig. 3a,b)**. By
147  contrast, we found no evidence that trial-by-trial variations in attention or stimulus presentation
148  were expressed in behavior **(Fig. 2a,b, Supplementary Fig. 3c)**. The latter result suggests that
149  the animals' behavioral response to the stimulus is largely embodied by its decision to lick.
150          Given the importance of single-trial analyses in decision-making paradigms[12,13], we next
151  investigated the behavioral correlates of decision-making processes. The non-zero false alarm
152  rates observed in our data suggest that multiple processes drive mouse licking. We therefore
153  sought to test whether distinct causes of licking (i.e. spontaneous vs. stimulus-driven) were
154  differentially encoded in behavior **(Fig. 1c, Fig. 2a,b)**. To do so, we attempted to decode the
155  causes of licking on a lick-by-lick basis.
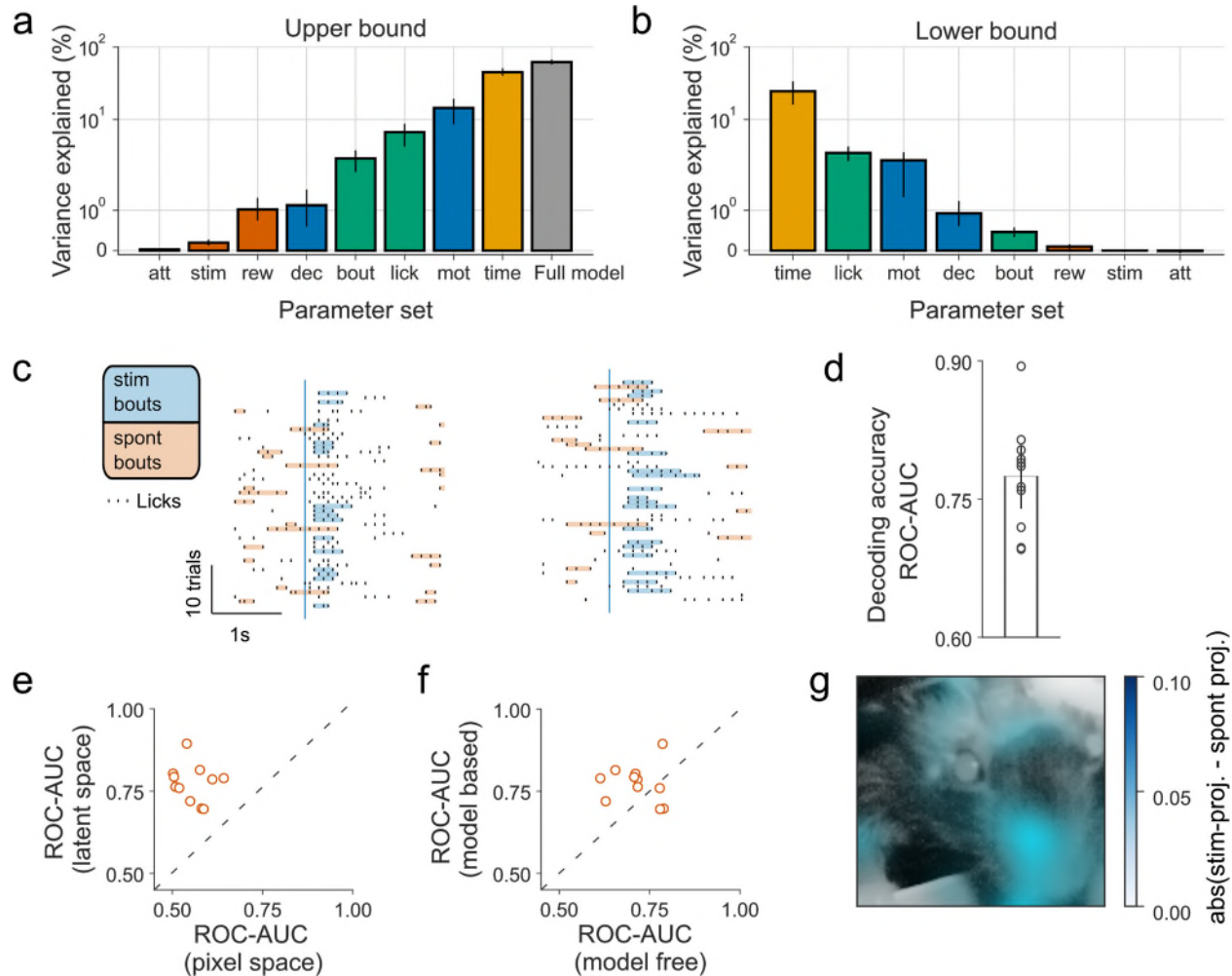
156
157



158
159

160 **Figure 2** Encoding and decoding behavior. **a)** Estimation of upper bounds on extent of encoding
161 by only regressing parameter sets belonging to a single variable. Variables are sorted according
162 to their ability to predict latent states. **b)** Estimation of lower bound on extent of encoding by
163 removing regressors relating to a single variable, one at a time, and subtracting cross-validated
164 $r^2$ for full model performance from $r^2$ for models with individual components removed. Error bars
165 show bootstrapped 95% confidence intervals. **c)** Excerpts from two example sessions showing
166 lick-bouts defined as either stimulus-driven (blue) or spontaneous (orange) depending on their
167 timing relative to the stimulus onset (blue vertical line). **d)**  Decoding of intention (i.e.
168 classification of bout type) by inverting behavioral encoding models reveals accurate decoding
169 (mean ROC-AUC = 0.78; s.e.m = 0.01). Error-bars show $\pm$ 2 s.e.m. Circles are individual data-
170 points. **e)** Decoding in latent space is more accurate than decoding in pixel space (paired
171 samples t-test; $p = 3.9 \cdot 10^{-6}$). **f)** Model-based decoding performs better than model-free (SVM)
172 decoding (paired samples t-test; $p = 0.0086$).  **g)** Difference between the BAE's estimate of a
173 stimulus and a spontaneous bout overlayed on an image of a mouse. Estimates were created
174 by projecting linear predictions of stimulus-driven and spontaneous bouts into pixel space. In

175     this case, informative pixels are clustered around the snout. (att=attention; stim=stimulus
176     presentation; rew=reward delivery; dec=decision basis (spontaneous vs stimulus-driven licking);
177     bout=lick-bout initiation; mot=motivational state)

178     _____

179

180

181

182

183     We grouped licks into bouts **(Fig. 2c, Supplementary Fig. 4)** and selected a
184     counterbalanced set (*see Methods*) of stimulus-driven (fast response times on trials with loud
185     stimuli) and spontaneous (outside of the peri-stimulus period) lick-bouts. We then decoded (i.e.
186     predicted) the causes of these bouts using the latent states within the ~500 ms preceding the first
187     lick of each bout. Previous work has demonstrated that the inversion of encoding models offers a
188     powerful and parsimonious approach to decoding[14,15]. We therefore constructed model-based
189     decoders based on the inversion of the behavioral-encoding models **(Fig. 1b)**. Consistent with
190     results from the encoding perspective, we were able to decode, on a bout-by-bout basis, whether
191     a stimulus preceded a bout or not **(Fig. 2d)**. Thus, the animals' behavior preceding a lick bout
192     allowed us to infer whether a stimulus drove that bout.

193     Further analysis demonstrated that decoding accuracy was higher in the latent-space than
194     in pixel-space **(Fig. 2e)** and that model-based decoding out-performed comparable model-free
195     support vector machines (SVM) **(Fig. 2f)**. Importantly, decoding is unlikely to be driven by motor
196     preparation **(Supplementary Fig. 5a-d)**. Finally, the generative capabilities of the BAE enabled
197     us to project linear approximations of stimulus-driven and spontaneous lick bouts back into pixel
198     space. This visual account of the basis of their classification revealed that idiosyncratic behaviors
199     associated with lick bouts formed the basis for classification **(Fig. 2g, Supplementary Videos**
200     **2,3)**.

201     Model-based decoding thus offers a data-driven alternative to *a priori* analysis of behavior.
202     In doing so, it both provides a way of automatically identifying behavioral correlates of
203     experimental variables and a means of classifying behavior based on these correlates. In turn,
204     this yields an interpretable account of momentary behavior that can readily be employed to
205     improve our understanding of neural activity.

206     To demonstrate this, we sought to explicitly benchmark model-based and *a priori*
207     classifications of trial-by-trial decisions against neural activity. Previous work has demonstrated
208     that behavioral choice correlates with the activity of neurons in primary auditory cortex (A1)[16-18].
209     We reasoned that by comparing the behavioral categorization of bout-by-bout intent with neural
210     activity, we would be able to compare the two classification approaches.

211     We therefore performed two-photon calcium imaging of excitatory layer 2/3 neurons in A1
212     of three mice **(Fig. 3a-c)**. To assess whether neural activity covaries with behavioral choice, we
213     computed choice probabilities[12] (CPs) , and identified a subpopulation of L2/3 neurons with
214     significant CPs **(Fig. 3d,e; Supplementary Fig. 6)**. CPs calculated by comparing hit-trials and
215     miss-trials were both significantly correlated with **(Fig. 3f)** and not systematically different from
216     **(Supplementary Fig. 7a)** those calculated by comparing hit-trials with level-matched hit-trials in
217     which animals responded prematurely (i.e. with a latency of <120 ms, which is faster than mouse

6

218    reaction times). These results argue that CPs reflected sensorimotor coupling, rather than licking
219    or reward consumption, and were thus used as a benchmark measure of behavioral classification.
220         Given the non-zero false-alarm rates observed in our data, a subset of hit-trials likely
221    occurred as a result of spontaneous behavior, rather than the learned stimulus-response
222    association. In light of the robust choice encoding in A1, we reasoned that, neurally, these trials
223    should more closely resemble miss-trials than hit-trials. If our decoder is able to correctly reclassify
224    those hit-trials on which licking was spontaneous, we should observe larger CPs. Consistent with
225    this expectation, we found that CPs were indeed larger when calculated based on decoded
226    causes of behavior (mean = 0.71; s.e.m=0.005), than on *a priori* criteria (mean = 0.67; s.e.m =
227    0.0034), i.e. defining all trials with licking in a window 150-600 ms after the stimulus and no pre-
228    stimulus licking as hit trials **(Fig 3g., Supplementary Fig. 7b)**. This suggests that model-based
229    decoding of video data can provide a more accurate readout of behavior than readouts based on
230    *a priori* definitions imposed by the task structure.
231         Finally, we sought to use the behavioral models to further clarify the relationship between
232    neural encoding of movement-related and choice-related variables. To relate neural activity to
233    these variables, we fitted a linear model that attempts to explain neurons' frame-by-frame activity
234    using experimental variables as well as behavioral latent-states. This approach allowed us to
235    dissociate movement- and decision-related influences on neural activity, as during the inter-trial
236    interval movement and decisions are decoupled. Fitting these models to the activity of each
237    neuron thus yielded parameters quantifying how the activity of a given neuron covaries with the
238    animal's behavior. To further examine whether movement-related influences on neural activity
239    underlie CPs, we attempted to predict neurons' CPs from these parameters. We found that the
240    relationship between a neuron's activity and behavioral latent states was poorly predictive of its
241    CP **(Fig. 3h)**. Together with the behavioral controls **(Fig. 3f)**, these findings strongly suggest that
242    neural tuning to motor variables does not underlie choice-related activity in A1.
243         Recent work has demonstrated that animals' movements are predictive of neural activity
244    across cortical regions, including sensory cortex[19]. Consistent with this result, we were better
245    able to predict neural activity using both behavioral latent states and experimental variables as
246    regressors, than experimental variables alone **(Fig. 3i)**. However, this could either reflect
247    genuine neural tuning to motor output or be mediated via effects of internal variables on both
248    neural activity and motor output. The comprehensive representations learned by the BAE
249    allowed us to differentiate these two possibilities by quantifying how well A1 population activity
250    predicts animals' movements. If neurons in A1 are truly tuned to motor output, we should be
251    able to accurately reconstruct behavioral latent states from the measured neural activity.
252    Contrary to this prediction, we were poorly able to predict behavioral latent states from neural
253    activity (mean $r^2$ = 3%; range 1%-5%). These findings strongly argue that motor output has, at
254    most, a small effect on auditory cortical activity and that correlations between the two are likely
255    mediated by variables such as an animal's decision that affect both movement and neural
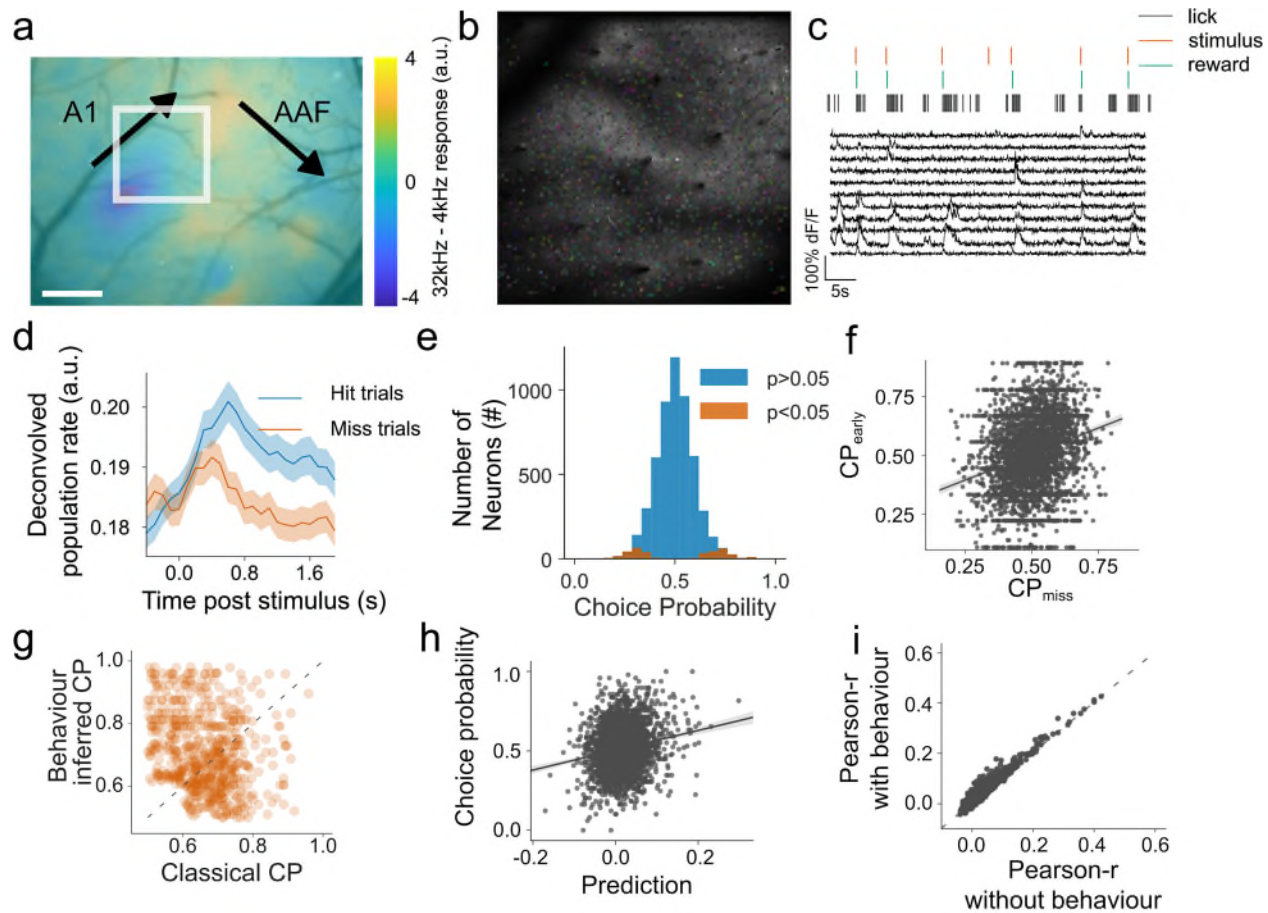256    activity.
257

**Figure 3** Behaviorally-decoded choices reflect neural activity**. a)** Functional localization of auditory cortical fields using wide-field single photon imaging. Scale bar shows $500\mu m$. **b)** Example imaging field (~$900\mu m^2$; region in white square in **a** with regions of interest (n = 976) randomly colored. **c)** Activity of ten neurons from **b. d)** Across the entire population of recorded neurons, we observed significant choice-related activity that emerged shortly after stimulus onset. Shaded regions are $\pm 2$ s.e.m. **e)** Distribution of choice probabilities (CPs). Significant CPs ($p <$ 0.05, permutation-test 500 shuffles) were measured in 378 of 5339 neurons (7.1 %). This is a larger subpopulation than would be expected by chance (binomial-test $p = 2.1 \cdot 10^{-119}$). **f)** CPs calculated by comparing hit and miss trials and CPs calculated from hit and 'early hit' trials are correlated (r = 0.26; $p = 1.3 \cdot 10^{-69}$) across neurons. **g)** CPs, plotted here as distance from 0.5, are greater when trial classification is based on model-based decoding rather than *a priori* criteria (paired sample t-test; $p = 3.6 \cdot 10^{-44}$). See **supplementary Fig. 6b** for raw CPs. **(h)** CPs are poorly predicted (mean$r^2$ = 1%), on a neuron-by-neuron basis, from neural tuning to behavioral latent states as assessed by fitting a multi-linear regression model. **(i)** Including behavioral latent states into a linear regression model to predict neural activity significantly improves fit quality (paired sample t-test; $p < 1 \cdot 10^{-80}$).

_____

8

279     In summary, our novel class of Bayesian model enables comprehensive and interpretable
280     quantification of momentary behavior. Application of this model demonstrated robust encoding of
281     cognitive variables in animals' behavior and enabled us to disentangle neural encoding of
282     cognitive and motor variables. We constructed model-based decoders whose application
283     provided sub-second accounts of behavior which more accurately reflected neural activity than
284     behavioral readouts imposed by task structure. Combined with recent methods for pose
285     estimation[20] , we envision that our approach will be able to extract simple readouts of complex
286     behavior . Finally, while we have deployed our model in the context of a sensory decision-making
287     task, these methods should be broadly applicable to both basic and clinical research seeking to
288     relate neural activity, computation and behavior.
289
290
291
292     **References**
293

294     1.      Jun, J. J. *et al.* Fully integrated silicon probes for high-density recording of neural
295     activity. *Nature* **551,** 232–236 (2017).
296     2.      Sofroniew, N. J., Flickinger, D., King, J. & Svoboda, K. A large field of view two-photon
297     mesoscope with subcellular resolution for in vivo imaging. *Elife* **5,** e14472 (2016).
298     3.      Gomez-Marin, A. & Mainen, Z. F. Expanding perspectives on cognition in humans,
299     animals, and machines. *Current Opinion in Neurobiology* **37,** 85–91 (2016).
300     4.      Gomez-Marin, A., Paton, J. J., Kampff, A. R., Costa, R. M. & Mainen, Z. F. Big
301     behavioral data: psychology, ethology and the foundations of neuroscience. *Nat. Neurosci.* **17,**
302     1455–1462 (2014).
303     5.      Busse, L. *et al.* The Detection of Visual Contrast in the Behaving Mouse. *J. Neurosci.* **31,**
304     11351–11361 (2011).
305     6.      Bak, J. H., Choi, J., Witten, I., Akrami, A. & Pillow, J. W. Adaptive optimal training of
306     animal behavior. in *NIPS* 1939–1947 (2016).
307     7.      Wichmann, F. A. & Hill, N. J. The psychometric function: I. Fitting, sampling, and
308     goodness of fit. *Percept. Psychophys.* **63,** 1293–313 (2001).
309     8.      Wiltschko, A. B. *et al.* Mapping Sub-Second Structure in Mouse Behavior. *Neuron* **88,**
310     1121–1135 (2015).
311     9.      Egnor, S. E. R. & Branson, K. Computational Analysis of Behavior. *Annu. Rev. Neurosci.*
312     **39,** 217–236 (2016).
313     10.     Kingma, D. P. & Welling, M. Auto-Encoding Variational Bayes. *Iclr* 1–14 (2014).
314     doi:10.1051/0004-6361/201527329
315     11.     Rezende, D. J., Mohamed, S. & Wierstra, D. Stochastic Backpropagation and
316     Approximate Inference in Deep Generative Models. (2014). doi:10.1051/0004-6361/201527329
317     12.     Britten, K. H., Newsome, W. T., Shadlen, M. N., Celebrini, S. & Movshon, J. A. A
318     relationship between behavioral choice and the visual responses of neurons in macaque MT.
319     *Vis. Neurosci.* **13,** 87–100 (1996).
320     13.     Parker, A. J. & Newsome, W. T. Sense and the single neuron: Probing the Physiology of
321     Perception. *Annu. Rev. Neurosci.* **21,** 227–277 (1998).

322  14.  Pillow, J. W., Ahmadian, Y. & Paninski, L. Model-Based Decoding, Information
323  Estimation, and Change-Point Detection Techniques for Multineuron Spike Trains. *Neural*
324  *Comput.* **45,** 1–45 (2010).
325  15.  Zhang, K. *et al.* Interpreting Neuronal Population Activity by Reconstruction : Unified
326  Framework With Application to Hippocampal Place Cells Interpreting Neuronal Population
327  Activity by Reconstruction : Unified Framework With Application to Hippocampal Place Cells.
328  *Am. Physiol. Soc.* **79,** 1017–1044 (2014).
329  16.  Bizley, J. K., Walker, K. M. M., Nodal, F. R., King, A. J. & Schnupp, J. W. H. Auditory
330  cortex represents both pitch judgments and the corresponding acoustic cues. *Curr. Biol.* **23,**
331  620–625 (2013).
332  17.  Jaramillo, S., Borges, K. & Zador, A. M. Auditory Thalamus and Auditory Cortex Are
333  Equally Modulated by Context during Flexible Categorization of Sounds. *J. Neurosci.* **34,** 5291–
334  5301 (2014).
335  18.  Francis, N. A. *et al.* Small Networks Encode Decision-Making in Primary Auditory
336  Cortex. *Neuron* **97,** 885–897.e6 (2018).
337  19.  Najafi, F. & Churchland, A. K. Perceptual Decision-Making: A Field in the Midst of a
338  Transformation. *Neuron* **100,** 453–462 (2018).
339  20.  Abe, T. *et al.* DeepLabCut: markerless pose estimation of user-defined body parts with
340  deep learning. *Nat. Neurosci.* **21,** 1281–1289 (2018).
341
342
343  **Acknowledgements**
344
352
353
354  **Author Contributions**
355
356  Y.W. conceived the study and the model. Y.W. and J.C.D. designed the experiments. Y.W. and
357  J.C.D. performed surgeries. Y.W. performed experiments. Y.W. analysed the data. A.J.K.
358  provided infrastructure and resources. A.J.K. and J.C.D. supervised the project. Y.W., A.J.K.
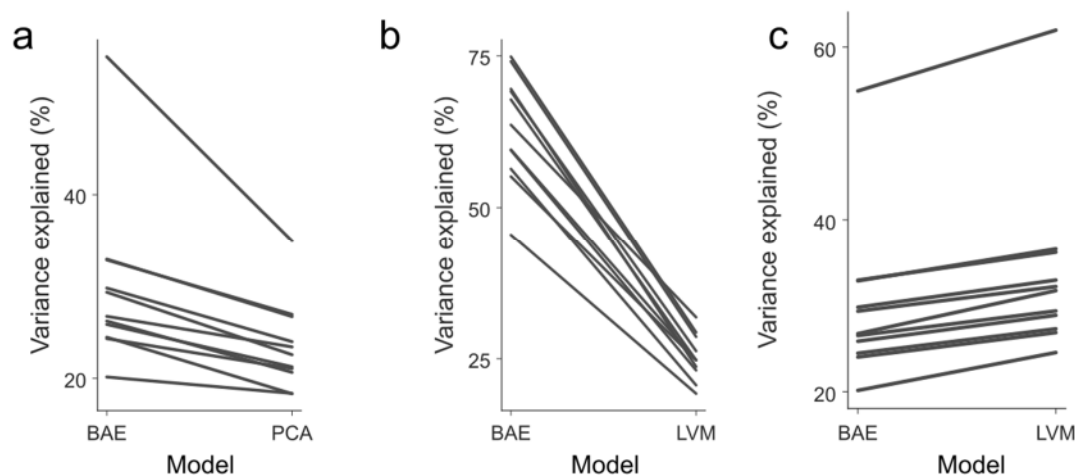359  and J.C.D. wrote the manuscript.
360
361
362  **Competing Interests**
363
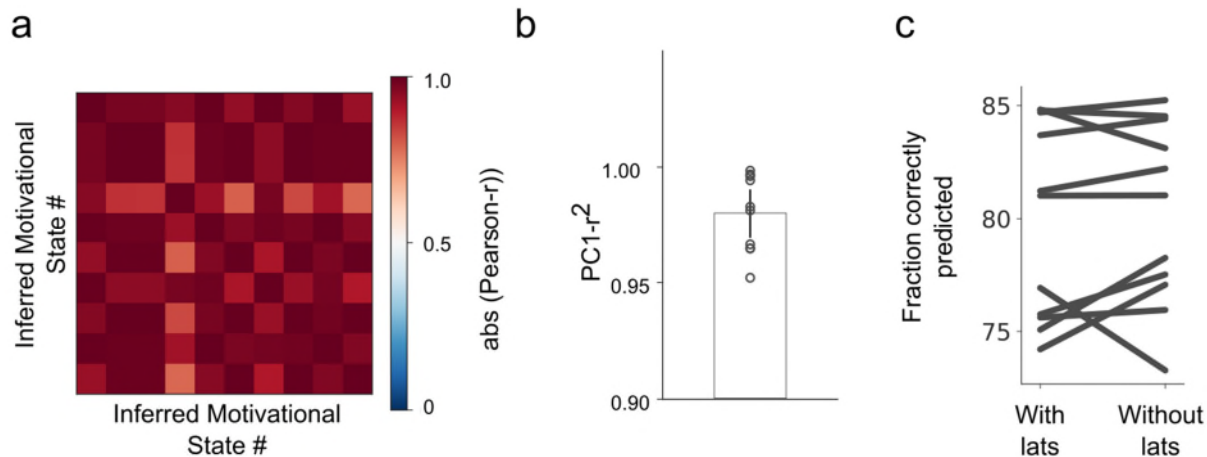364  The authors declare no competing interests.
365

## Supplementary Figures



**Supplementary Figure 1.** Visualization of reconstructions from the latent space. Example of a video frame in its raw and preprocessed form as well as its reconstruction. In the preprocessing step, each pixel of video data had its mean subtracted and was divided by its variance.
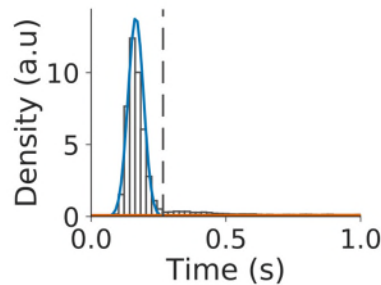


**Supplementary Figure 2** Quantitative analysis of pixel-space reconstructions of video data by various models. **a)** Pairwise comparison of reconstructions of the video data by BAE and PCA. For BAE reconstructions shown here, we performed one full pass through the model, using the recognition model to obtain latent-states and the generative model to obtain pixel-space reconstructions. Each line represents a single session. In all cases, BAE outperforms PCA (paired t-test; p=0.0002). **b)** To assess how well latent states can be predicted from experimental variables we compared the ability of the BAE and LVM **(Fig 1b)** to predict behavioral latent states. The BAE out performed the LVM in all sessions (paired t-test; $p=3.5 \cdot 10^{-10}$), demonstrating enhanced, linear predictability of latent-states as a result of the augmentation of the model's cost function. **c)** Pixel-space reconstructions, created by a full pass of the video data through the LVM (i.e. video data are passed through the encoder network to generate latent variables, which are then passed to the decoder network, reconstructing the full images) are better than BAE (Paired t-test; $p=1.7 \cdot 10^{-7}$).
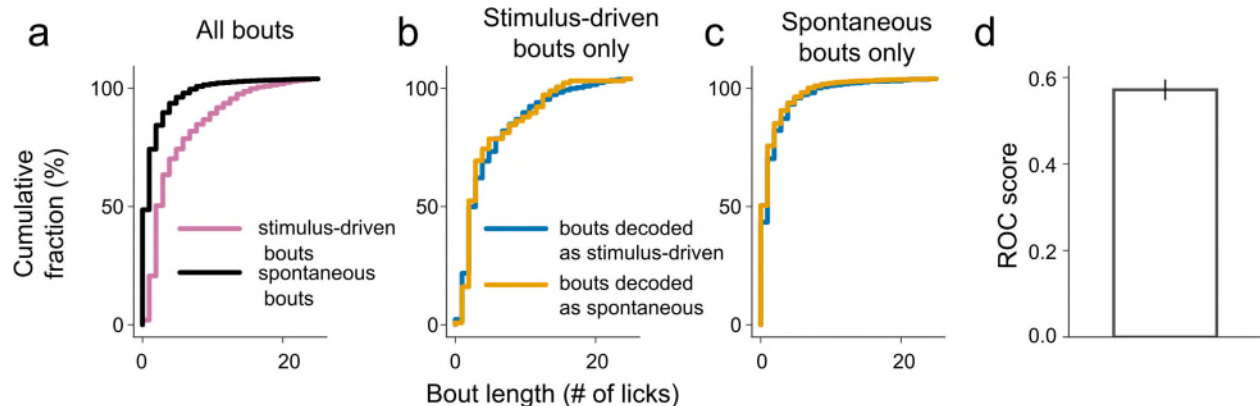
11

391
392

**Supplementary Figure 3** Further analysis of behavioral correlates of cognitive variables. **a)** Analysis of the encoding model from an example session, which shows that motivational state explains variance not accounted for by licking, suggested that an animal's motivational state is externalized in behavior **(Fig 2a,b)**. However, there is a chance that the encoded quantity may not actually reflect motivation, but changes in posture that are unrelated to the animal's motivational state. Motivation, in the context of our behavioral task, may be measured along a one-dimensional continuum, that is to say that at each point in time animals have a certain level of motivation. Therefore, if the measured quantity truly reflects motivation, we reasoned that different parts of the animal's posture, reflected in the ten behavioral latent-states, should change in a coordinated fashion. In contrast to this, if the measured quantity is just related to slow changes in posture, there is no *a priori* reason that the different behavioral latent states should change in a correlated fashion. To distinguish these possibilities we calculated the weighted sum of motivation regressors for each latent variable. Regressors were weighted by the values of fitted regression parameters for each latent variable. We refer to this sum as the inferred motivational state. We then measured the correlation between the inferred motivational states fitted to each latent state. Shown is an example correlation matrix, constructed by cross-correlating the inferred motivational states for each latent variable. This example illustrates that inferred motivational states, fitted to each behavioral latent-state independently, are highly correlated, consistent with the hypothesis that the extracted variable is related to the animals' motivational state rather than arising from spurious changes in posture. **b)** To quantify the extent to which the motivational state variables may be described by a one-dimensional quantity, we performed principal component analysis and quantified the variance explained by the first principal component. We found that in all sessions a single principal component captured more than 95% of the variance across motivational variables. **c)** Analysis of encoding model parameters suggested that attention was not expressed in animal's behavior. To further test this, we performed a logistic regression analysis and tried to predict trial-by-trial decisions, asking whether knowledge of latent-states preceding stimulus onset helped us in doing so. We compared performance of a baseline model to performance of an extended model that included the latent-states preceding stimulus onset. The baseline model included the intensity of the presented stimulus and whether the previous trial was a hit- or miss-trial. Expanding this model by including behavioral latent states preceding stimulus presentation did not improve the

12

424    model's ability to predict whether a given trial is a hit- or miss-trial (paired sample t-test; $p =$
425    0.32). These results bolster the conclusion that attention is not encoded in the animals' behavior
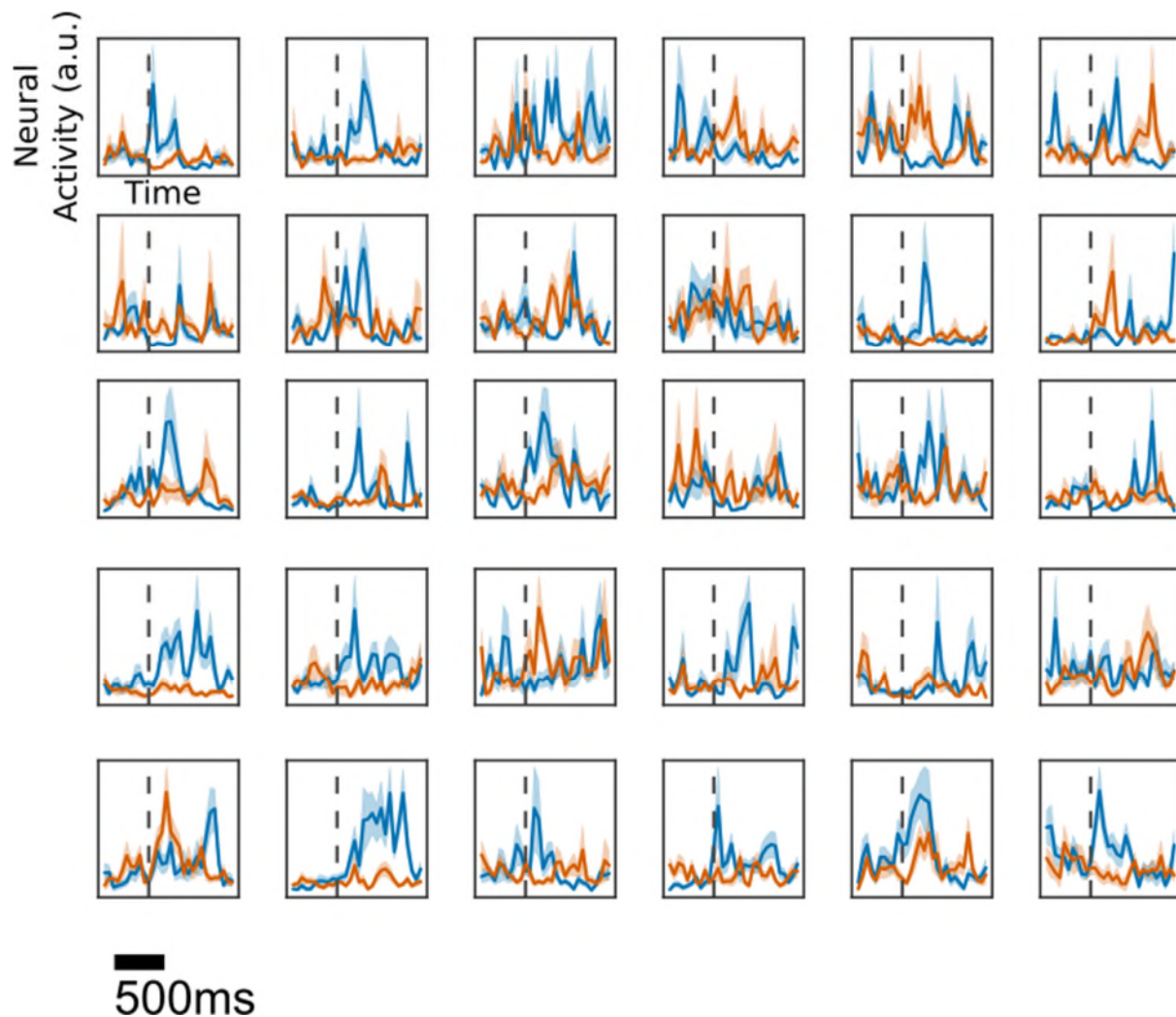426    preceding stimulus onset.

427

428

429



430

431    **Supplementary Figure 4** Mouse licking behavior is organized into bouts. Distribution of inter-
432    lick intervals across all sessions and animals (white histogram bars). Gaussians fitted to intra-
433    bout inter-lick intervals (blue curve) and between-bout inter-lick intervals (orange curve)
434    overlaid, together with the optimal separation boundary (dashed vertical gray line).

435

436

437

438



439

440

441    **Supplementary Figure 5** Excluding motor preparation and time as bases for classifying
442    behavior. **a)** Significant differences in bout lengths (quantified in terms of number of licks in a
443    bout) exist between stimulus-driven and spontaneous bouts. Therefore, stimulus-driven and
444    spontaneous bouts could be associated with differences in motor preparation that the decoder
445    might be able to exploit for its classification. **b)** Partitioning of only stimulus-evoked bouts
446    according to decoder classification reveals no differences in bout length as a function of the
447    decoder's classification. **c)** Partitioning of only spontaneous bouts according to decoder
448    classification also revealed no difference in bout length as a function of the decoder's
449    classification. This suggests that decoder performance is not driven by potential differences in
450    motor preparation between short and long lick bouts. **d)** To estimate the extent to which the

13

451   decoder relies on differences in bout length to perform classification, we measured how well
452   bout length could predict decoding performance. To do so, we computed the area under the
453   receiver operating characteristic curve (mean=0.56; s.e.m=0.01) and found that bout length was
454   a poor predictor of the decoder's decision.
455
456
457
458



459
460   **Supplementary Figure 6** Representative examples of neurons with significant choice
461   probabilities. Each panel shows the average activity (mean±s.e.m) of a single neuron in a
462   window surrounding stimulus onset (dashed vertical line). The y-axis of each panel is
463   normalized to show the full dynamic range of each neuron. Blue curves show mean activity
464   during hit-trials; orange curves show mean activity during miss-trials. Examples shown are
465   taken from all animals.
466
467

468
469 **Supplementary Figure 7** Further analysis of choice-related activity. **a)** Choice probabilities
470 calculated by comparing hit and miss trials ($CP_{miss}$) and choice probabilities computed by
471 comparing hit vs early hit trials ($CP_{early}$) are not significantly different in magnitude (paired-
472 sample t-test $p = 0.68$). **b)** Full distribution of classical versus behavior inferred choice
473 probabilities.
474 _____
475
476
477
478 **Supplementary Video 1.** Example pre-processed video and associated reconstructions using
479 the BAE. Latent states were estimated using the recognition model.
480
481 **Supplementary Video 2.** Estimation, via the BAE, of the mean video sequence preceding
482 stimulus-driven and spontaneous lick bouts, respectively. Estimation is based on data from one
483 example session. These pre-lick bout sequences were estimated by reconstructing latent states
484 using the behavioral encoding model and projecting these latent states into pixel space using
485 the generative model.
486
487 **Supplementary Video 3.** Example sets of video sequences preceding stimulus-driven and
488 spontaneous lick bouts from a single session. Data shown in video are temporally
489 counterbalanced such that simultaneously shown clips are close in time. Data are from the
490 same session as Supplementary Video 2.
491
492
493
494
495
496
497
498
499

**Methods**

501

502

*Animals*

504

505        All experiments were approved by the local ethical review committee at the University of
506 Oxford and licensed by the UK Home Office. One female C57BL/6NTac.*Cdh23753A>G* (Harlan
507 Laboratories, UK) mice[23], 3 female (C57B6.129S-Gt(ROSA)26Sortm95.1(CAG-GCaMP6f)Hze
508 [Jax: 024105] x C57B6.Cg-Tg(Camk2a-cre)T29-1Stl/J [Jax:005359]), one male
509 (*Igs7tm93.1(tetO-GCaMP6f)Hze* Tg(Camk2a-tTA)1Mmay/J [ Jax: 024108] x Rbp4_KL100-Cre,
510 MMRRC: 037128; Gerfen et al., 2013)  and one male Rbp4-cre mouse were used for behavioral
511 experiments. Neural data were obtained from the three (C57B6.129S-
512 Gt(ROSA)26Sortm95.1(CAG-GCaMP6f)Hze x C57B6.Cg-Tg(Camk2a-cre)T29-1Stl/J) mice. All
513 experiments were performed before mice reached 12 weeks of age, preceding the onset of age-
514 related sensorineural hearing loss in C57BL/6J strains[21,22].

515

516

517

*Click detection task*

519

520        Three days before mice commenced behavioral training, we started restricting their
521 access to water and  acclimatising them  to handling and head-fixation. Throughout the training
522 and testing period the mice' body weight remained above 80% of their pre-restriction body
523 weight.  Mice were trained daily to lick in response to a 0.05-ms biphasic click stimulus
524 presented at 80 dB SPL. There were two types of trials: stimulus trials (80 dB SPL click; water
525 reward for licking) and catch trials (no stimulus; no reward for licking). These were randomly
526 interleaved at an inter-trial interval drawn from a uniform distribution between 6s and 12s. If
527 mice licked during a 1.5 s window following onset of the stimulus, a water drop (2 µl) was
528 delivered immediately. Once mice reached high performance levels (> 80 % correct on stimulus
529 trials), which took 2-5 sessions, they were moved to the testing phase in which stimuli were
530 presented at different intensities. Stimuli were randomly interleaved and presented over a
531 maximum range of 38 dB SPL to 80 dB SPL (3-dB steps). The range of stimulus levels
532 presented in a given session was, in some cases, adjusted according to the animals' sensitivity.
533 Behavioral data were acquired in blocks lasting between 7 and 30 minutes. Typical sessions
534 lasted approximately forty minutes during which mice performed approximately 250 trials.
535        Data were excluded, in a block-wise manner according to several criteria. Firstly, mice
536 needed to have undergone at least two testing sessions prior to the sessions considered for
537 inclusion. Secondly, to be able to reliably identify stimulus-driven bouts, we required hit-rates for
538 the loudest stimuli to exceed 95%. Finally, to be able to reliably identify hit-trials as being
539 stimulus driven, we required false-alarm rates to be below 45%. Of the 12 sessions (two per
540 mouse) passing these criteria, one had to be excluded because of video frames missing as a
541 result of camera failure.

542

543

544 *Apparatus*
545
546       The behavioral apparatus was controlled from a computer running Windows 7 using
547 MATLAB (Mathworks) interfaced with a National Instruments board (NI- DAQ USB-6008) for
548 data acquisition. Stimuli were presented using MATLAB 2016a (Mathworks) running
549 psychtoolbox. Stimuli were digital-to analog converted using a commercial soundcard (ASUS
550 Xonar-U7), amplified (Portable Ultrasonic Power Amplified; Avisoft Bioacoustics) and played
551 through a free-field electrostatic speaker (Vifa; Avisoft Bioacoustics), positioned approximately
552 15 cm in front of the mouse's snout.
553       Stimuli were calibrated using an M500 microphone (Pettersson), which was itself
554 referenced to a sound-level calibrator (Iso-Tech SLC-1356). Click volumes were calibrated by
555 integrating the recorded RMS of clicks over the mouse hearing range (1-100kHz) and
556 comparing it to the RMS of stimuli from the reference sound-level calibrator.
557       Video frame acquisition was triggered by the frame clock of the two-photon microscope,
558 such that one video frame was acquired for every two microscope frames, resulting in an
559 acquisition rate of ~13 Hz at a resolution of 640 x 480 pixels. The camera, a DMK23UV024 (The
560 Imaging Source) mounted with a M5018-MP2 (Computar) lens, was positioned approximately
561 30 cm in front of and 30 cm above the behavior apparatus, aligned to have the mouse's face
562 and most of its body in the field of view. Regions of interest showing the mouse's face
563 (**Supplementary Fig. 1**) were drawn manually (approximately 150 x 150 pixels in size) on each
564 dataset. These regions of interest were used for further analysis.
565
566
567
568 *Widefield calcium imaging*
569
570 The widefield imaging system consisted of a 470nm LED (M470L3, Thorlabs), a digital camera
571 (340M-GE, Thorlabs) and a 2X objective (TL2X-SAP, Thorlabs) mounted on a Thorlabs
572 Bergamo II microscope body. Images were acquired at a rate of 10 Hz and a resolution of 96 by
573 128 pixels using ThorCam (Thorlabs) software. Sound waveforms were generated in LabView
574 (National Instruments) and presented on the same hardware as described above. For the
575 frequency mapping of auditory cortical fields we presented 500 ms long sinusoidally amplitude
576 modulated (SAM) tones with a modulation frequency and depth of 10 Hz and 100%,
577 respectively. Each map was based on the responses to 15 repeats of one low carrier frequency
578 (4 kHz or 5.04 kHz) and 15 repeats of one high carrier frequency (25.4 kHz or 32 kHz) SAM
579 tone, presented at either 55 dB SPL or 65 dB SPL and at a rate of 0.33Hz. Frequency maps
580 (**Fig. 3a**) were generated by calculating the average response (mean signal intensity in a 1-s
581 window following sound onset minus mean signal in a 1-s window preceding sound onset) to the
582 low-frequency and high-frequency stimulus, subtracting one from the other, color-coding the
583 resulting image and superimposing it on a grayscale image of the bloodvessel pattern.
584
585
586
587

17

588    *Two-photon data acquisition*

589

590    Two photon imaging was performed as described previously[24]. Briefly, image acquisition
591    was carried out using a commercially available two-photon laser-scanning system (B-Scope;
592    Thorlabs). A SpectraPhysics Mai-Tai eHP laser fitted with a DeepSee prechirp unit (70fs pulse
593    width, 80MHz repetition rate) provided the laser beam for two photon excitation. The beam was
594    directed into a Conoptics modulator and then through the objective (16x 0.8NA water immersion
595    objective; Nikon). The beam was scanned across the brain using an 8-kHz resonance scanner
596    (X) and a galvanometric mirror (Y). The resonance scanner was used in bidirectional mode,
597    enabling acquisition of 512 x 512 pixels at a frame-rate of approximately 26 Hz. Emitted photons
598    were filtered (525/50) and collected and amplified by GaAsP photomultiplier tubes
599    (Hamamatsu). ScanImage was used to acquire data and control the microscope. All imaging
600    was done between 150 and $250 \mu m$ below the cortical surface.

601

602

603

604    *Latent variable model*

605

606    The mathematics underlying variational autoencoders[10,11], on which our models are
607    based, has been covered in great detail elsewhere (see e.g. Doersch, 2016[25] for a tutorial) so
608    we will give only a brief summary here. Given some observed high-dimensional series of pixel
609    intensities (i.e. video data) $X$, we seek to explain variation in $X$ by assuming that some low-
610    dimensional underlying latent variables, $z$, give rise to the data. Ideally, the quantity we would
611    seek to maximize when fitting the model is thus $P(X)$, the probability of the data. We can relate
612    $z$ to $P(X)$ mathematically by conditioning:

613
$$P(X) = \int p(X|z) P(z) \, dz \approx \frac{1}{n} \sum_{i=1}^{n} P(X|z_i) \tag{1}$$

614

615    where we note that any integral can be approximated by a finite sum over samples of $z_i$. This
616    formulation has the important property that by specifying the functional form of $p(X|z)$ and a
617    method of sampling $z_i$ we can evaluate $P(X)$ and hence quantify the performance of the model.
618    For analytical tractability and ease of sampling, we assert that $P(z)$ is a Gaussian distribution
619    with 0 mean and diagonal, unit covariance.

620

621
$$P(z) = N(0 \,|\, I) \tag{2}$$

622

623    Based on the continuous values of pixel intensities, we further specify $P(X_i|z_i)$ to be a normal
624    distribution:

625

626
$$P(X_i|z_i) = N(\mu = f_\phi(z_i); \Sigma = I) \tag{3}$$

627

628    where $f_\phi(z)$ is a deterministic function, with parameters $\phi$, that map latent variables, $z$, into pixel
629    space. In practice, we implement $f_\phi(z)$ as a multi-layer neural network.

630   However, with high-dimensional data, naive sampling approaches are inefficient to the
631   point of intractability because for most values of $z_i$, $p(X_i|z_i) \approx 0$. To enable efficient sampling,
632   allowing us to tractably approximate the above integral, we construct an auxiliary distribution
633   $Q(z_i|X_i)$which enables us to draw samples from $P(z_i)$ such that the sampled $z_i$ are likely to give
634   rise to $X_i$. In practice, we assume that
635
636
$$Q(z_i|X_i) = N(z_i|\mu = g_\theta(X_i); \Sigma = h_\theta(X_i)) \qquad (4)$$

637
638   where $g$and $h$ are deterministic functions of $X$, parameterised by $\theta$, which are implemented by a
639   deep neural network. However, naively sampling $Q(z|X)$,rather than$P(z)$, to evaluate $P(X)$ will
640   result in biased estimates. To circumvent these issues we apply standard identities from the
641   Variational Bayesian literature[7] to derive:
642
643
$$L(\theta,\phi) = \log P(X) - D_{z\sim Q(z|X)}(Q(z|X) || P(z|X)) = -E_{z\sim Q(z|X)}[\log p_\phi(X|z)] +$$
644
$$D(Q_\theta(z|X) || P(z)) \qquad (5)$$

645
646   where $D(p||q)$ denotes the KL-Divergence (a measure of difference between probability
647   distributions) between $p$ and $q$. The left hand side of this equation is the quantity we seek to
648   maximize. Doing so maximizes the likelihood of the data $P(X)$while minimizing the difference
649   between our approximation of $Q(z|X)$ and the true $P(z|X)$. Since both $Q_\theta(z|X)$ and $P(z)$ are
650   Gaussian, this divergence has a closed form solution. Similarly, we can arrive at a
651   computationally tractable form of the expectation $E_{z\sim Q(z|X)}[\cdot]$ by using a single sample from
652   $Q(z|X)$to make the approximation. Furthermore, tractable derivatives of this cost function are
653   available[10,11] .
654   We extend this model to encourage learning of interpretable latent representations. We
655   achieved this by adding an additional term to the cost function. Specifically, we fitted a
656   behavioral encoding model (see *Behavioral encoding model* for details), mapping from task
657   variables to the latent variables $z$ using a linear regression model with parameters $\beta$. We
658   augment the cost function with the error term of this regression model to obtain a more
659   interpretable model in which the values of latent variables $z$are linearly predictable from
660   variables of interest.
661
662
$$L(\theta,\phi;\beta) = -E_{z\sim Q(z|X)}[\log p_\phi(X|z)] + D(Q_\theta(z|X) || P(z))$$
663
$$-E_{z\sim Q(z|X)}[\log p_\beta(z|V)] \qquad (6)$$

664
665   Importantly, the prior on the latent space acts to regularize the latent parameters
666   preventing overfitting. Additionally, our behavioral encoding model only biases the learning of
667   weights, it does not bias the inferred latent representation.
668
669
670
671
672

673 *Data analysis*

674

675 Data were analysed in Matlab and Python 3.6.2 augmented with standard libraries for scientific

676 computing[26-31]. Unless stated otherwise, standard algorithms (e.g. principal component analysis)

677 are implemented using reference implementations from these libraries. A reference

678 implementation of the behavioral autoencoder, together with an example video dataset is

679 available for use and alteration at www.github.com/yves-weissenberger/bae.

680         All statistical tests were, unless otherwise stated, implemented using reference

681 implementations in standard libraries for scientific computing in Python. All statistical tests were

682 two-tailed.

683

684

685

686 *Model implementation*

687

688         The hierarchical Bayesian model is implemented using the Python library *Tensorflow*[32].

689 The model is comprised of two sequential networks termed recognition model and generative

690 model, respectively. All neural activation functions were rectified-linear unless otherwise stated.

691         The recognition model is a four-layer network. The first two layers are comprised of

692 convolutional units (256 and 128 units), and kernel sizes three and five pixels, respectively. In

693 both cases, the stride of kernels was set to two pixels. These layers were followed by a fully

694 connected layer with 100 units and a final bipartite layer comprised of 10 linear and 10 softplus

695 units, mapping to the mean and covariance of the latent space, respectively.

696         The decoder network consisted of two fully connected layers with 100 and 500 units,

697 respectively, followed by a final fully connected linear layer mapping the previous layers' output

698 into pixel space. Our network was trained using a 60/20/20 train/validation/test split. To optimize

699 the cost function we used AdamOptimizer[33] with the learning rate set to 0.005. Hyperparameters

700 were, once heuristically optimized using a separate dataset not included in this report, held fixed

701 for all analyses reported here.

702

703

704

705 *Lick bout analysis*

706

707         To separate licks into bouts, we fitted a two-component Gaussian Mixture Model

708 (implemented by the *GaussianMixture* class of the scikit-learn library) to the inter-lick interval

709 (ILI) distribution of all mice. We thereby separated the ILI distribution into two components which

710 we interpreted as corresponding to within bout ILIs and across bout ILIs. In doing so, we

711 determined the optimal separation window for dividing licks into bouts as the point at which the

712 probability of the fitted Gaussian with the larger mean exceeded that of the smaller one. Doing

713 so, we found that a window of ~266ms provided the optimal separation window for

714 differentiating within-bout licks from across-bout licks.

715

716

717   *Behavioral encoding model*

718

719        Our behavioral encoding model was a linear-regression model mapping from the set of
720   observed and hidden variables $V$ to inferred latent-states $z_i$ using parameters $\beta$. The set of
721   observed variables we used comprised licks, rewards, lick-bout initiations (defined as the first
722   lick in a bout of licks) and sound stimuli. The timestamps of each of these observed event types
723   were discretized to construct a set of $T \times 1$ vectors (where $T$ is the length of the session), either
724   set to 1 on the camera frame at which the event occurred (click, reward) or two frames
725   preceding an event (lick-bout initiation, lick), as these movements will be initiated before a lick is
726   completed, and 0 everywhere else. In the case of the clicks, we also analyzed the data after
727   scaling entries in the vector according to sound level, but this made no qualitative or quantitative
728   difference (data not shown).
729        The set of hidden variables was comprised of decision basis, attention and motivational
730   state. Decision basis was a $T \times 2$ binary vector whose first and second columns signified
731   whether a stimulus-driven or spontaneous lick-bout occurred, respectively. An entry in the first
732   column was set to a value of 1 at five frames (~380 ms) preceding the onset of a lick-bout if a
733   stimulus preceded the lick-bout within a ~600 ms window (this window represents the 70th
734   percentile of the across-animal reaction time distribution). Analogously, an element was set to 1
735   in the second column if no stimulus preceded the bout and the bout was initiated outside the
736   peri-stimulus period. This period was defined as the period from ~150 ms prior to onset of the
737   stimulus to ~1.5 s following the onset of the stimulus.
738        Attention was  a $T \times 2$ binary vector whose first column signified that the animal was
739   attentive. We reasoned that detection of particularly loud stimuli was not affected by attention
740   and therefore did not include these in this analysis. An element in the first column was set to 1
741   at five frames preceding the onset of a stimulus if that stimulus was presented at a low intensity
742   (average hit-rate at that intensity <75%) and the trial was a hit trial. Analogously, an element in
743   the second column was set to 1 on miss trials.
744        Motivational state was a $T \times 5$ continuously valued vector approximating the extent of
745   reward seeking. We constructed each row of this matrix by convolving the vector of licks with a
746   Gaussian distribution. We derived this definition of motivational state based on recent work
747   demonstrating that in head-fixed mice, increased motivation is associated with increased
748   baseline lick rates[34]. The Gaussian for each row had a different standard deviation reflecting our
749   *a priori* uncertainty about the timescales of motivational fluctuations. The standard deviations
750   ranged from ~2.5 s  to ~40 s multiplied in powers of two.
751        We additionally included a set of time regressors, a $T \times 10$ vector, where each row is a
752   continuous low frequency oscillation, to account for slow drifts in posture over time. The period
753   of these oscillations ranged from ~1450 s to ~2150 s. To enable events to affect latent-states at
754   future time points, all the above vectors (with the exception of motivational-state and time) were
755   multiplied with a Toeplitz matrix giving rise to a series of lagging regressors extending 5 frames
756   into the future.
757        The Design Matrix $\hat{V}$ was then constructed by concatenating these vectors together with
758   an offset term yielding the following regression model

759

760

21

$$p(z|\,\beta;\hat{V}) \ = \ N(\beta \cdot \hat{V}|\,I) \tag{7}$$

where

$$\hat{V} = [v^{offset}, v^{time} \ v^{lick} \cdot K, v^{bout} \cdot K, \ v^{rew} \cdot K, v^{stim} \cdot K,$$

$$v_1^{att} \ \cdot K, \ v_2^{att} \ \cdot K, v_1^{dec} \ \cdot K, \ v_2^{dec} \cdot K, v^{mot}\,] \tag{8}$$

Linear models were regularized using an L2-penalty term. Fitting, as well as regularization parameter selection was implemented using the scikit-learn function *RidgeCV*. Fit quality estimation was performed using repeated, nested K-fold cross validation (five folds; four repeats). In the inner K-fold loop (five folds), the training data were used for fitting and hyperparameter selection, while in an outer loop fit quality was assessed using the held-out data.

*Analysis of behavioral-encoding model parameters*

To determine the importance of each regressor in the behavioral encoding model, we performed two complementary analyses to bound the extent of their encoding. This was required because of the collinearity of regressors. To obtain a lower bound on strength of encoding, we quantified the effect of excluding subsets of regression parameters, relating to a single experimental variable (e.g. $v^{bout}$), on cross-validated fit quality. Secondly, to obtain an upper bound, we included only parameters relating to a single experimental variable in the regression model. Each of these models was fitted to latent-states extracted after the initial, global fitting process. Model performance was estimated, as during initial fitting, using repeated, nested K-fold cross validation (six folds; four repeats). In the inner K-fold loop (five folds), we determined the optimal regularization parameter. In the outer loop, we attempted to assign hit or miss labels to a held-out test set of trials based on fit parameters.

*Logistic-regression analysis of attentional state*

To determine whether trial-by-trial attentional states were externalized in behavior, we attempted to use behavioral latent-states preceding stimulus onset to predict whether a given trial was a hit or miss trial. To do so, following fitting of our latent variable model and the determination of behavioral latent-states, we fitted a logistic regression model to subjects' trial-by-trial choices. Logistic regression was implemented using the *sklearn* function *LogisticRegression* using the Newton Conjugate Gradient solver and an L2 penalty. A reference model included as regressors the level of the presented stimulus and a variable indicating whether the previous trial was a hit- or miss-trial. To determine whether some correlate of attention was externalized in behavior, we compared performance of the reference model to a model which additionally included the behavioral latent states on the ten video frames preceding each stimulus onset as regressors. Model performance was estimated using a repeated, nested

22

805  K-fold cross validation (six folds; four repeats). Regularization parameters were optimized in an
806  inner K-fold loop (five folds).

807
808
809
810  *Behavioral decoding dataset*

811
812      The window for decoding extended 5 video frames backwards from the onset of the lick-
813  bouts. To ensure that lick history did not form the basis of our behavioral decoding, we only
814  selected lick-bouts in which no licks occurred in a ~610 ms window preceding bout-onset.
815  Additionally, to ensure that long-timescale covariation in posture and spontaneous bout-rates do
816  not drive decoder performance (spontaneous bout-rates are typically higher at the beginning of
817  behavioral sessions), spontaneous and stimulus-driven lick-bouts were selected in a temporally
818  counterbalanced fashion. Specifically, for each session, we counted the number of stimulus-
819  driven and spontaneous bouts. We denote the smaller of these two sets the reference set $R_1$.
820  For each bout in the reference set, we selected the bout in the larger set that was its nearest
821  neighbour, yielding a second set of bouts $R_2$. The union of these sets ($R_1 \cup R_2$) then comprised
822  the decoding dataset. This led to an unbiased selection of spontaneous and stimulus-driven
823  bouts. Decoding performance was similar when the bout distributions were not counterbalanced
824  in this fashion (data not shown). Decoding performance was estimated on a test-set held out
825  during fitting, using repeated, nested K-fold cross validation (five folds; four repeats).

826
827  *Model free decoding*

828
829      Model free decoding was performed using a linear support vector machine whose
830  regularization parameter $C$ was determined in an inner cross validation loop, as described
831  above. In addition to determining the optimal regularization parameter, variable selection was
832  performed in the inner loop, whereby the optimal set of timepoints to use for classification was
833  determined by optimizing prediction accuracy on the training set. Classification was
834  implemented by the *sklearn* function *SVC*.

835
836  *Model-based decoding*

837
838      Decoding was performed using log-likelihood ratios ($LLR$) similarly to Pillow et al[14].
839  Specifically, for each lick-bout we compared the log-likelihood of the behavioral latent-states
840  preceding the onset of a bout under the assumption that this bout was stimulus-driven, with the
841  log likelihood that the bout was spontaneous:

842
843  $$LLR = log \frac{p(V_{stim}|\ \beta;\ z)}{p(V_{spont}|\ \beta\ ;z)} = log \frac{p(z|\ \beta;V_{stim})}{p(z|\ \beta\ ;V_{spont})} + K \propto \sum_{t=1}^{H} \{(z_t\ -\ \beta\ \cdot V^t_{stim})^2\ -\ (z -$$
844  $$\beta \cdot V^t_{spont})^2\ \}$$      (9)

845
846  Where $V^t_{stim}$ is the design matrix constructed by setting the relevant entry (i.e. five frames
847  preceding bout onset) for stimulus-driven bout to 1 and the entry for spontaneous bout to 0,

848  $V_{spont}$ is the reverse, $H$ is the analysis horizon and $K$ are terms independent of $V$. A log
849  likelihood ratio greater than 0 corresponds to a lick bout that is decoded as being stimulus-
850  driven.
851         To quantify the accuracy of the decoder we performed a repeated nested, stratified K-
852  fold (six folds; four repeats) cross validation. In an inner K-fold loop (five folds), we determined
853  the optimal regularization parameter for the behavioral encoding model. This means that
854  regularization parameters were only explicitly optimized for encoding, and only implicitly
855  optimized for decoding. Decoding performance was then estimated on the held-out cross
856  validation set comprising equal numbers of stimulus-driven and spontaneous lick-bouts.
857         Pixel space decoding was performed by projecting latent-space estimates of stimulus-
858  driven (i.e. $\beta \cdot V^t{}_{stim}$) and spontaneous lick bouts (i.e. $\beta \cdot V^t{}_{spont}$) back into pixel space using
859  the trained generative model and calculating log likelihood ratios in pixel space.
860
861     $$LLR \propto \sum_{t=1}^{H} \quad \{(X_t - f_\phi(\beta \cdot V^t{}_{stim}))^2 - (X_t - f_\phi(\beta \cdot V^t{}_{spont}))^2\}$$
862     (10)
863
864  Where $f_\phi(\cdot)$ (see equation (3) ) is a neural network implementing the generative model,
865  returning the posterior mean in pixel space from some latent value.
866
867
868  *Two-photon data preprocessing*
869
870         Data preprocessing was performed in Python using the Two-Photon Analysis Toolbox:
871  twoptb (https://yves-weissenberger.github.io/twoptb/). Briefly, data were motion registered using
872  the efficient subpixel registration algorithm. Next, regions of interest (ROIs) were automatically
873  segmented (then manually curated) using a pre-trained supervised algorithm, included in the
874  toolbox, which uses the mean image to identify ROIs. Segmentation was performed in a two-
875  step process where the initial step involved finding seed regions for ROIs using a random-
876  forests classifier. In a second step, a region-growing algorithm was applied to construct ROIs.
877  Traces were extracted as an unweighted average of fluorescence within each region of interest.
878  All traces were neuropil corrected using the fluorescence averaged in a 20 x 20$\mu m$ square
879  surrounding the ROI (empirically determined correction factor: ~0.5). Traces were then baseline
880  corrected using a Kalman-filter based estimate of baseline fluorescence. Finally, spike inference
881  was performed on neuropil corrected traces using the c2s toolbox[35]. To improve temporal
882  resolution, all neural analyses were performed on inferred spike rates.
883
884  *Choice probability estimation*
885
886         For analysis of choice probabilities[12] , we selected equal numbers of hit and miss trials
887  from each stimulus level with hit-rates between 25% and 75%. This was done to maximise data
888  inclusion while preventing variation in sound-evoked activity from dominating the influence of
889  choice. To calculate choice probabilities, we measured the neural response (average neural
890  activity in a 300ms window following stimulus onset) for each trial. We then used the resulting
891  hit and miss trial response distributions to calculate the area under the receiver operating

892  characteristic curve using the *roc_auc_score* function in the *sklearn* package. P-values for
893  choice probabilities were determined by permutation testing using 2000 shuffles.
894       When calculating choice probability based on behavioral decoding, the subset of hit-
895  trials that were behaviorally decoded as spontaneous were moved from the hit-trial to the miss-
896  trial group. To avoid biased estimates as a result of class imbalances, we calculated choice
897  probability by averaging the mean accuracy for each class (hit and miss). Calculating choice
898  probabilities without such counterbalancing did not qualitatively affect conclusions (data not
899  shown).
900
901  *Neural regression model*
902
903       Regression models fitted to neural activity were identical in implementation to those
904  used in the behavioral encoding model (see above), except for the inclusion of instantaneous
905  (i.e. no time lagged regressors were used) behavioral latent-states as regressors. When neural
906  regression models were fit only to behavioral latent-states and did not include the design matrix
907  used in the behavioral encoding model, results with respect to choice encoding were
908  qualitatively similar (data not shown).
909
910
911  *Choice probability prediction*
912
913       To assess whether neural choice probabilities (CPs) were related to the covariation of
914  neural activity and movements, we analyzed the parameters of fitted neural regression models.
915  Following the fitting of neural regression models, parameters relating to behavioral latent-states
916  were extracted. We then fitted a multi-linear model, separately to each session, which
917  attempted, on a neuron-by-neuron basis, to predict the neuron's choice probability from that
918  neuron's regression model parameters related to behavioral latent-states. We reasoned that if
919  choice probability was explained by neural tuning to motor output, or indeed motion artifacts
920  unaccounted for by image registration, then, across neurons, choice probability should be
921  predictable from neurons' tuning to behavioral latent states. The multi-linear model was
922  implemented by the *OLS* class from the *statsmodels* library.
923
924
925
926  **Methods References**
927
928

929  21.    Willott, J. F., Aitkin, L. M. & McFadden, S. L. Plasticity of auditory cortex associated with
930  sensorineural hearing loss in adult C57BL/6J mice. *J. Comp. Neurol.* **329,** 402–411 (1993).
931  22.    Ison, J. R., Allen, P. D. & O'Neill, W. E. Age-related hearing loss in C57BL/6J mice has
932  both frequency-specific and non-frequency-specific components that produce a hyperacusis-like
933  exaggeration of the acoustic startle reflex. *JARO - J. Assoc. Res. Otolaryngol.* **8,** 539–550
934  (2007).

935  23.    1. Mianné, J. *et al.* Correction of the auditory phenotype in C57BL/6N mice via
936  CRISPR/Cas9-mediated homology directed repair. *Genome Med.* **8,** 16 (2016).
937  24.    Vasquez-Lopez, S. A. *et al.* Thalamic input to auditory cortex is locally heterogeneous
938  but globally tonotopic. *Elife* **6,** e25141 (2017).
939  25.    Doersch, C. Tutorial on Variational Autoencoders. (2016). doi:10.3389/fphys.2016.00108
940  26.    Rossum, G. Van & Drake, F. L. *Python Tutorial, Technical Report CS-R9526. Centrum*
941  *voor Wiskunde en Informatica (CWI)* (1995). doi:16/j.histeuroideas.2011.02.001
942  27.    Hunter, J. D. Matplotlib: A 2D graphics environment. *Comput. Sci. Eng.* **9,** 99–104
943  (2007).
944  28.    Oliphant, T. E. SciPy: Open source scientific tools for Python. *Comput. Sci. Eng.* **9,** 10–
945  20 (2007).
946  29.    Seabold, S. & Perktold, J. Statsmodels: econometric and statistical modeling with
947  Python. in *Proceedings of the 9th Python in Science Conference* 57–61 (2010).
948  30.    Pedregosa, F. *et al. Scikit-learn: Machine Learning in Python. Journal of Machine*
949  *Learning Research* **12,** (2011).
950  31.    Van Der Walt, S., Colbert, S. C. & Varoquaux, G. The NumPy array: A structure for
951  efficient numerical computation. *Comput. Sci. Eng.* **13,** 22–30 (2011).
952  32.    Abadi, M. *et al.* TensorFlow : A System for Large-Scale Machine Learning. *Proc 12th*
953  *USENIX Conf. Oper. Syst. Des. Implement.* (2016). doi:10.1126/science.aab4113.4
954  33.    Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. (2014).
955  doi:10.1063/1.4902458
956  34.    Berditchevskaia, A., Cazé, R. D. & Schultz, S. R. Performance in a GO/NOGO
957  perceptual task reflects a balance between impulsive and instrumental components of behavior.
958  *Sci. Rep.* **6,** (2016).
959  35.    Theis, L. *et al.* Benchmarking Spike Rate Inference in Population Calcium Imaging.
960  *Neuron* **90,** 471–482 (2016).
961
962
963
964