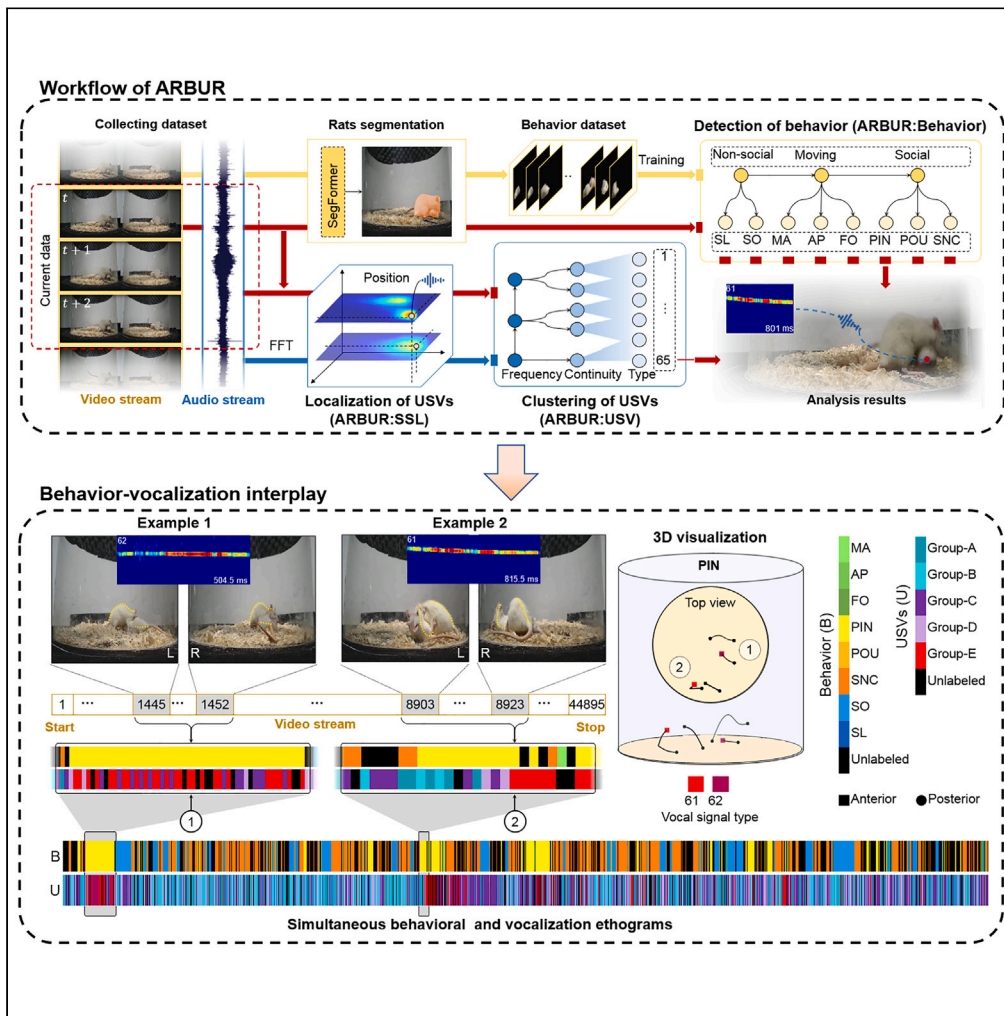


Article

ARBUR, a machine learning-based analysis system for relating behaviors and ultrasonic vocalizations of rats



Zhe Chen,
Guanglu Jia, Qijie
Zhou, ..., Toshio
Fukuda, Qiang
Huang, Qing Shi

shiqing@bit.edu.cn

Highlights

Automatic analysis system reveals underlying behavior-vocalization interplay of rats

Data-driven comprehensive clustering of ultrasonic vocalizations

Hierarchical classification can discriminate easy-to-confuse social behaviors

3D sound source localization of freely behaving rats



Article

ARBUR, a machine learning-based analysis system for relating behaviors and ultrasonic vocalizations of rats

Zhe Chen,^{1,2} Guanglu Jia,^{2,3} Qijie Zhou,^{2,3} Yulai Zhang,^{1,2} Zhenzhen Quan,⁴ Xuechao Chen,^{2,3} Toshio Fukuda,⁵ Qiang Huang,^{2,3} and Qing Shi^{2,3,6,*}

SUMMARY

Deciphering how different behaviors and ultrasonic vocalizations (USVs) of rats interact can yield insights into the neural basis of social interaction. However, the behavior-vocalization interplay of rats remains elusive because of the challenges of relating the two communication media in complex social contexts. Here, we propose a machine learning-based analysis system (ARBUR) that can cluster without bias both non-step (continuous) and step USVs, hierarchically detect eight types of behavior of two freely behaving rats with high accuracy, and locate the vocal rat in 3-D space. ARBUR reveals that rats communicate via distinct USVs during different behaviors. Moreover, we show that ARBUR can indicate findings that are long neglected by former manual analysis, especially regarding the non-continuous USVs during easy-to-confuse social behaviors. This work could help mechanistically understand the behavior-vocalization interplay of rats and highlights the potential of machine learning algorithms in automatic animal behavioral and acoustic analysis.

INTRODUCTION

Nonverbal communication and vocalization are two vital means of natural communication during social interaction across the animal kingdom and human society.^{1,2} Nonverbal communication such as facial expression,³ body posture, waggle dance,⁴ and social play⁵ can be received by nearby conspecifics and often analyzed generally by human observers as social behavior,⁶ whereas vocalization can be communicated remotely to convey information in dark and tortuous environments. Both social behavior and vocalization can be readily recorded experimentally as outer observable states that reflect the inner emotional states³ or even neural circuit activities,^{7–9} and therefore are actively investigated in animal behavior research.¹⁰ During social engagement, both the social behavior and vocalization of individuals influence each other in a continuous and interactive manner, contributing to the variable group social dynamics. Deciphering the behavior-vocalization interplay^{11,12} would thus reveal insights into the neural basis of social interaction. For such deciphering efforts, rats are widely used as social animal models because of their innate sociability¹³ and proven emotion-related ultrasonic vocalizations (USVs).¹⁴ For example, the alarm sounds (22-kHz USVs) of trapped rats may induce overwhelming distress in freely moving rats via emotional contagion and further evoke pro-social behavior in the free ones.¹⁵ However, a dedicated system for analyzing social communications of freely behaving rats is still missing to promote the mechanistic understanding of the interplay between social behavior and USVs.

The barrier is caused by the difficulty of relating the different types of behavior and USVs of rats in complex social contexts (to reveal the underlying behavior-vocalization interplay) because of the following challenges. First of all, a method for the unbiased and automatic clustering of USVs that not only covers both non-step (continuous) and step (non-continuous) signals in spectrograms¹² but also comprehensively incorporates the structure of frequency and duration,¹⁶ is still missing. Moreover, the automatic behavior detection of rats under social interaction poses a major challenge in constructing the behavior-specific features^{17,18}; therefore, it is difficult to discriminate different types of easy-to-confuse social behavior.^{19–22} In addition, allocating the recorded USVs to the vocal rat also brings an obstacle because of the proximity and top-view partial overlapping of socially engaged rats.²³

In this Article, we propose a machine learning-based Analysis system for Relating Behaviors and USVs of Rats (named ARBUR). ARBUR clusters ultrasonic syllables into user-defined or automatically calculated subgroups from a comprehensive perspective, detects eight types of behavior (Table 1) based on the lateral camera view, and locates the vocal rat in two-rat scenarios with free social interaction. Using ARBUR,

¹School of Medical Technology, Beijing Institute of Technology, Beijing, China

²Key Laboratory of Biomimetic Robots and Systems, Beijing Institute of Technology, Ministry of Education, Beijing, China

³Intelligent Robotics Institute, School of Mechatronical Engineering, Beijing Institute of Technology, Beijing, China

⁴Key Laboratory of Molecular Medicine and Biotherapy, School of Life Science, Beijing Institute of Technology, Beijing, China

⁵Institute of Innovation for Future Society, Nagoya University, Nagoya, Japan

⁶Lead contact

*Correspondence: shiqing@bit.edu.cn

<https://doi.org/10.1016/j.isci.2024.109998>



Table 1. Behavioral definitions of freely behaving rats

Behavior	Definition
SO (solitary)	Both rats engage in separate activities.
ST (still)	Both rats keep still.
MA (moving away)	A rat moves away from another rat.
AP (approaching)	A rat moves toward another rat until contact.
FO (following)	A rat follows another rat without contact.
PIN (pinning)	A rat holds down another rat lying on its back.
POU (pouncing)	A rat lies on the back of another rat.
SNC (social nose contact)	A rat touches another rat with its nose tip.

we reveal that the distinct behaviors of rats are associated with different USVs. Therefore, based on the simultaneously recorded raw video and audio streams (see [Figure S1](#)), ARBUR can relate the behavior and USVs of rats in freely interacting social contexts for downstream qualitative and quantitative analyses. For example, we show that the pinned rat produces significantly more 22-kHz aversive USVs compared with the bully one or during other social behaviors or solitary states. Moreover, ARBUR indicates several novel findings about still-associated or moving-associated USVs, which have long been neglected in the manual analysis-dominated research.

RESULTS

Comprehensive clustering of the vocal repertoire

We recorded 43,357 ultrasonic vocalizations (USVs) to create the vocal repertoire. To cluster the USVs in an unbiased and comprehensive way, we present the comprehensive three-step clustering algorithm ARBUR:USV (see [Figure S2](#) and [Table 2](#)). Representative clustering examples are shown in [Figure 2A](#). In the first step, ARBUR:USV divides the USVs into 2 clusters according to their mean frequency-indicated emotional states (aversive 22-kHz USVs (ranging from 15 kHz to 32 kHz) or appetitive 50-kHz USVs (above 32 kHz)). This is a special treatment for rat USV analysis since rats produce meaningful USVs with a frequency near 22 kHz while mice don't. Another special consideration in ARBUR:USV is the inclusion of step signals, which are usually overlooked in other research. Moreover, note that 50-kHz USVs are classified into four groups according to their mean frequency and duration: high-peak-frequency (HPF) 50-kHz USVs and low-peak-frequency (LPF) 50-kHz USVs, because two aggregated areas were observed in the density map ([Figure 2B](#)), indicating two kinds of 50-kHz USVs may be present despite of their similar frequency contours (shape in the spectrograms). Therefore, in the second step, the 50-kHz USVs are classified into four groups (groups A-D) according to their continuity (non-step or step) and distribution in frequency-duration space (HPF or LPF) ([Figure 2A](#)). In the third step, ARBUR:USV separates the four 50-kHz groups (groups A-D) into subgroups within each group according to their frequency contours, thus allowing users to investigate the correlation between USV contours and other biological factors (e.g., behavior types) in relevant research. In the meantime, the 22-kHz group (group E) is divided into five subgroups based on mean frequency and duration because the frequency contours of 22-kHz USVs do not vary much.

While the first two steps divide the USVs according to their biological relevance (mean frequency, natural distribution in frequency-duration space), the third clustering step separates the USVs within each group according to their contour features, which raises the question of determining the optimal number of clusters to balance over-clustering with under-clustering. Here, ARBUR:USV provides three solutions for users to choose. First, the number of clusters for groups A-E can be subjectively chosen according to their experiences. This may be suitable

Table 2. ARBUR:USV achieves comprehensive automated clustering of rodent's USVs compared with current research

	Manner of clustering	Mean frequency	Duration	Non-step contour	Step contour	Choice of optimal clusters (number)	Rodents
Wright et al. ²⁴	Manual	✓	–	✓	✓	Subjective (15)	Rats
Riede et al. ²⁵	Manual	✓	–	✓	✓	Subjective (6)	Rats
Burgdorf et al. ²⁶	Manual	✓	–	✓	✓	Subjective (3)	Rats
Fonseca et al. ²⁷	Automated	✓	–	✓	✓	Subjective (11)	Mice
Sangiamo et al. ¹²	Automated	–	–	✓	–	Progressive (22)	Mice
Takahashi et al. ¹⁶	Automated	✓	✓	–	–	Bayesian information criterion (3)	Rats
MUPET ²⁸	Automated	–	✓	✓	–	Subjective (20–200)	Mice
DeepSqueak ²⁹	Automated	–	✓	✓	–	Elbow (20)	Mice
AVA ³⁰	Automated	–	✓	✓	–	Subjective (–)	Mice
ARBUR:USV	Automated	✓	✓	✓	✓	Multiple choices(65)	Rats

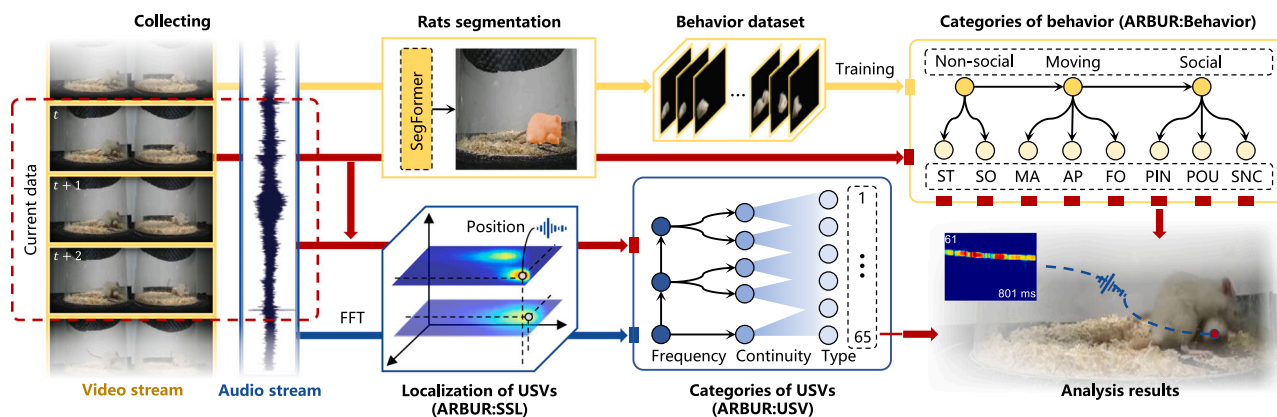


Figure 1. Illustration of the working flow of ARBUR

ARBUR takes simultaneously recorded video and audio streams as input and outputs the current behavior type, the ultrasonic vocalization (USV) in spectrogram with cluster type and duration indicated, and the labeled position of the vocal rat for each frame of the binocular images in the video frame.

for experienced animal behaviorist to investigate USV-related rat behaviors with fixed cluster size. Second, ARBUR:USV can calculate the optimal number of clusters for each group according to the elbow method (see [Method details](#) in the [STAR Methods](#)).²⁹ It basically relies on the relationship (curve) of total within-cluster error (TWCE) versus the increasing number of clusters. With increasing cluster numbers, TWCE tends to decrease exponentially. The elbow method finds the elbow point of the curve to be the optimal number. For the collected dataset, the optimal numbers for groups A-E are 13, 28, 13, 14, and 8, respectively (see [Figure S3](#)). Third, ARBUR:USV also provides the progressive method, which considers the clustering degree for each cluster. In particular, the cluster number, starting with two, increases by one if the average innerpoint percent of clustering results does not reach a user-defined threshold. For example, if the threshold is set to be 0.98, the optimal cluster number for group E is 8 (see [Figure S4](#)). It can produce a higher or lower number of clusters if stricter or looser thresholds are used for particular purposes. In short, ARBUR:USV provides three methods for choosing the number of clusters for groups A-E. Here, to balance simpler visualization and the elbow method result, we subjectively set the numbers for groups A-E to be 10, 25, 10, 15, and 5, respectively.

[Figure 2A](#) shows the clustering examples. Note that the USV clusters were re-sorted in descending order according to contour slope (for clusters 1–60) or duration (for clusters 61–65). It is intuitive that the USVs within each cluster share similar contours, and the mean contours of each cluster differ from those of the others (see [Figure S5](#)). We then quantified the clustering results. We show that ARBUR:USV achieves much higher inter-cluster distance compared with intra-cluster variance across the 65 clusters (see [Figure S6](#)). Moreover, we quantified the 65 clusters of the USVs, and showed that they vary substantially in terms of duration, frequency range, mean frequency, and signal counts across clusters (see [Figure S7](#)).

Hierarchical behavior classification of two freely behaving rats

ARBUR:Behavior is based on a hierarchical classification architecture, which we optimized for analyzing the species-specific behaviors of two freely behaving rats (see [Method details](#) in [STAR Methods](#)). We constructed different classification features based on behavioral states (non-social, moving, and social) and designed three classifiers to distinguish behavioral categories further. ARBUR:Behavior uses the binocular side-view video stream (segmented by SegFormer³¹ to remove the background) as input to discriminate eight behaviors of rats ([Figure 1B](#)). ARBUR:Behavior can also run in a single-shot mode if moving behaviors are not considered.

We annotated 1,265 randomly extracted segments from all recorded videos (including 43,357 segments) to test the performance of ARBUR:Behavior. As demonstrated in [Figure 3A](#), the module can classify the behaviors of two freely behaving rats well, especially the moving behaviors with large displacement and the social behaviors with complex postures. We show that the desired detection precision, recall, and F1-score of each behavior are achieved (average precision: 0.874, average recall: 0.862, average F1-score: 0.868, [Figure 3B](#)). The confusion matrix shows that ARBUR:Behavior can discriminate the easy-to-confuse social behaviors with high F1-scores, such as PIN (pinning) and POU (pouncing) ([Figure 3C](#)). In contrast to studies based on end-to-end architectures,^{19,32,33} ARBUR:Behavior is capable of discriminating specific social behaviors (PIN, POU, and SNC (social nose contact)) and can classify a higher number of behavior categories while achieving comparable accuracy.

In addition, we tested the performance of current mainstream deep learning classification algorithms^{34–37} for the direct classification of the eight behaviors on our dataset. ResNet achieves the best performance with an accuracy of only 0.333 using an eight-class classifier (see [Figure S8A](#)). For the two- and three-class classification of three social behaviors, the highest accuracies of these algorithms are 0.767 and 0.578 (see [Figure S8B](#)). We argue that features with behavioral specificity can not be directly extracted in the images using the end-to-end approach, resulting in poor classification performance. Therefore, some studies extracted the pose and tracking data of rats to classify behaviors in an unsupervised way and achieved high accuracies.^{38,39} However, such methods are unable to discriminate specific social behaviors. In contrast, ARBUR:Behavior detects non-social behaviors using the location feature of rats in the image and moving behaviors by estimating the optical

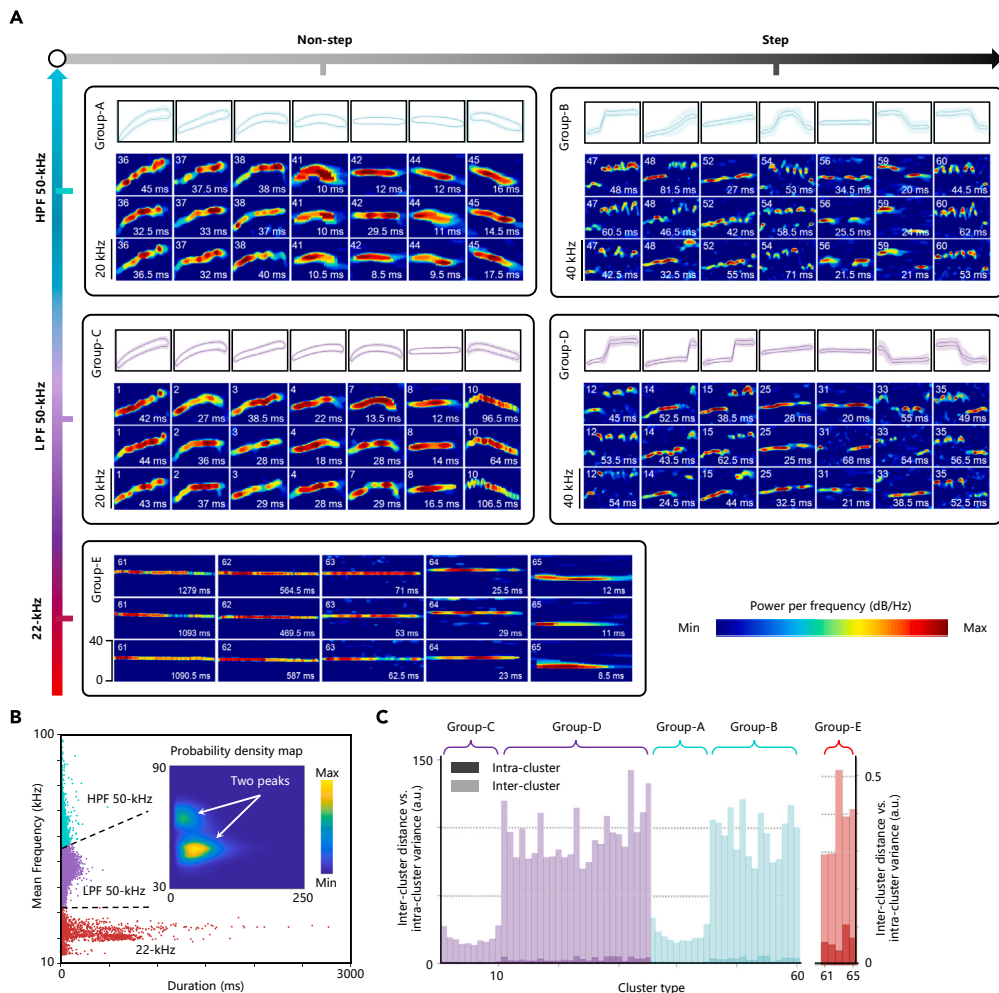


Figure 2. ARBUR:USV clusters the vocal repertoire in an unbiased and comprehensive manner

(A) Representative examples show that USVs within the same clusters share similar contours. Groups A-D, top panels, mean value of contour \pm s.e.m. Bottom three panels, three selected USVs spectrograms with cluster type (top left) and duration (bottom right) are indicated. Each column represents one cluster. Intensity is indicated by color (red, maximum; blue, minimum). Groups A-E represent high-frequency non-step 50-kHz USVs (clusters 36–45), high-frequency step 50-kHz USVs (clusters 46–60), low-frequency non-step 50-kHz USVs (clusters 1–10), low-frequency step 50-kHz USVs (clusters 11–35) and 22-kHz USVs (clusters 61–65), respectively. The frequency range of spectrograms: 20 kHz for group A and C; 40 kHz for group B and D; 0–40 kHz for group E. (B) Intuitive demonstration of the classification within the 50-kHz USVs, which was made because two independent peaks were observed in the probability density map (inset). (C) Quantification of the clustering results, which shows high inter-cluster distances and small intra-cluster variances.

flow of the rat as well as its relative orientation. To discriminate social behaviors, it extracts “histogram of oriented gradients” (HOG)⁴⁰ descriptors to represent the statistics of social posture features in a standardized way.

We also tested the performance of mainstream machine learning-based classifiers for two- and three-class classifiers on three social behaviors. The support vector machine (SVM) achieved the highest classification accuracy (0.843) in binary classification, and the Receiver Operating Characteristic (ROC) curve shows that it has the best overall classification performance (see Figure S8C). Therefore, we designed an SVM-based decision tree classifier for discriminating social behaviors (see Method details). ARBUR:Behavior has a greater advantage in terms of F1-score/accuracy and the number of behavior categories, compared with hand-crafted features used in existing studies.^{20,22,41} On balance, ARBUR:Behavior outperforms existing studies with respect to comprehensive detection performance (Table 3). ARBUR:Behavior successfully detected 33,428 behaviors (with a confidence of over 0.8) out of 43,357 segments of video frames corresponding to recorded USVs (Figure 3D; Table 3).

Locating the vocal rat in 3-D space

To locate the vocal rat during free-behaving scenarios, we developed a novel algorithm (ARBUR:SSL) that incorporates the lateral binocular view (for rat nose reconstruction in the Cartesian space), behavior classification results, height-sensitive distributed microphone

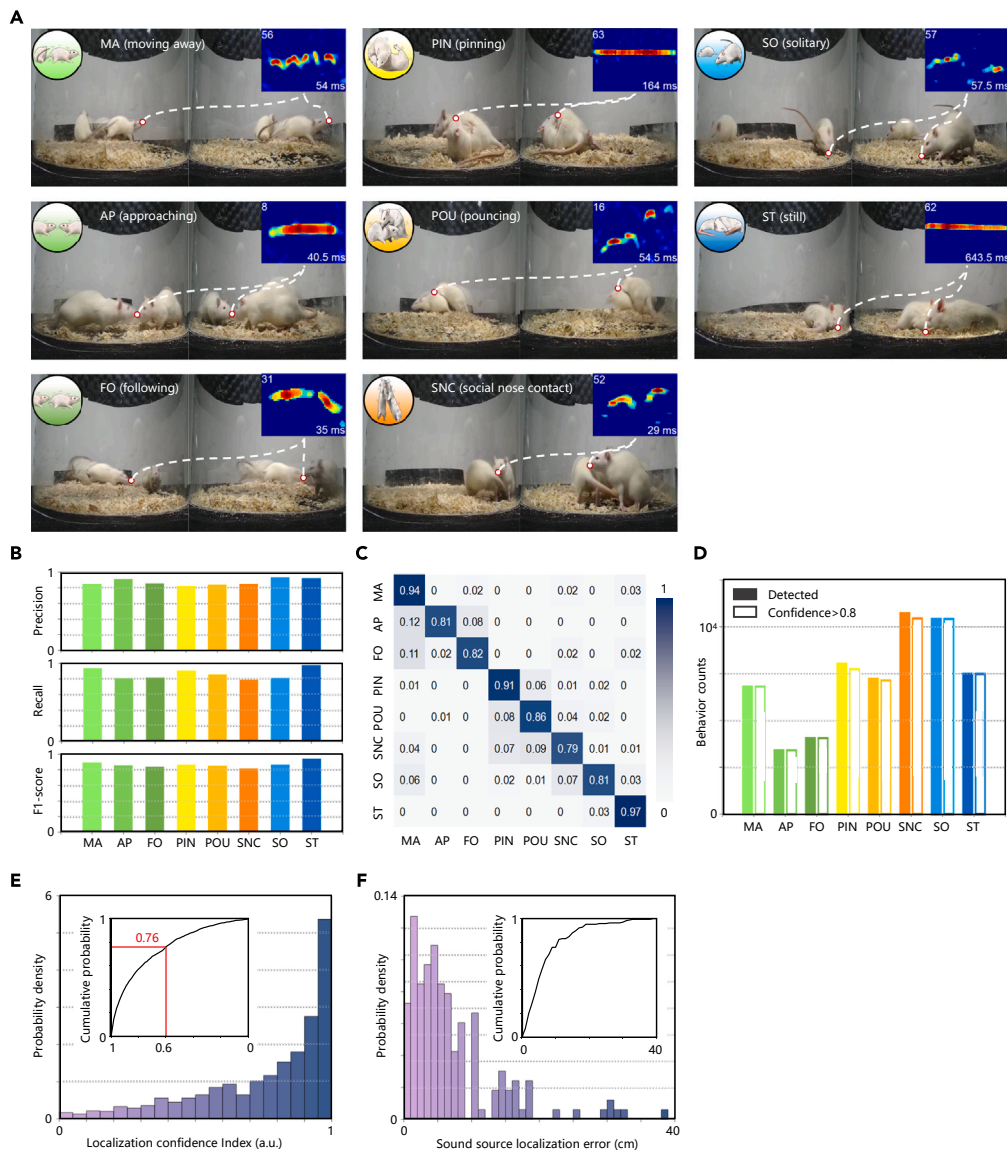


Figure 3. ARBUR simultaneously achieves accurate behavior detection and sound source localization in three dimensions

(A) Examples of the ARBUR outputs for eight behaviors: left and right camera view, behavior type (top left of each panel), spectrograms of USVs with cluster (inset, top left), and duration (inset, bottom right) and location of the vocal rat in the two camera views (red circle). Frequency range of spectrograms: 20 kHz for cluster 8; 0–40 kHz for clusters 62 and 63; 40 kHz for all other clusters.

(B) Quantification of the performance of ARBUR:Behavior.

(C) Confusion matrix showing the identification rates of ARBUR:Behavior.

(D) Counts of different types of behavior detected (logarithmic) across behavior types.

(E) Distribution of the localization confidence index (LCI) obtained by ARBUR:SSL; inset, cumulative probability vs. LCI.

(F) Distribution of the sound source localization error by ARBUR:SSL; inset, cumulative probability vs. error.

configuration, and 3-D probabilistic triangulation of sound source to locate the vocal rat in 3-D space. ARBUR:SSL takes the 3D reconstructed nose position, the current behavior type, and four-channel audio streams as input. It can comprehensively leverage these information and output the pixel position of the vocal rat in both binocular images, which is useful for follow-up expert quantitative and qualitative analysis.

The rat-rat social interaction scenario makes the microphone array far above the rats (50 mm above ground) compared with the experimental sets for mouse interaction, which makes the recorded USVs signal intensity (therefore signal-to-noise ratio) much weaker than others. This poses a new challenge to us because the mainstream sound source localization (SSL) algorithms suitable for mouse interaction experiments work much worse in our dataset and, therefore undesirable for our purposes. To address this challenge, we optimized a classic

Table 3. ARBUR:Behavior achieves accurate behavior classification

	Input modality	Algorithm characteristics	Background subtraction	Social behavior	Behavior categories	F1-score/ Accuracy
DeepEthogram ¹⁹	RGB (top view)	End-to-end (supervised)	–	✓	6	0.638/-
SIPEC:BehaveNet ³³	RGB (top view)	End-to-end (supervised)	✓	–	3	0.720/-
DeepAction ²⁰	RGB (top view)	Hand-crafted (supervised)	–	✓	12	-/0.739
VSAMBR ⁴¹	RGB (side view)	Hand-crafted (supervised)	✓	–	8	-/0.771
DeepCaT-z ³²	RGB-D (top view)	End-to-end (supervised)	✓	–	4	0.822/-
OpenLabCluster ³⁹	Pose and tracking data	Hand-crafted (supervised)	–	–	8	-/0.842
MARS ²²	RGB (top view)	Hand-crafted (supervised)	–	✓	3	0.878/-
B-SOID ³⁸	Pose and tracking data	Hand-crafted (supervised)	–	–	11	-/0.915
ARBUR:Behavior	RGB (side view)	Hand-crafted (supervised)	✓	✓	8	0.868/ 0.870

method²³ and proposed a new index (localization confidence index, LCI), which measures the consistency between 3-mic SSL and 4-mic SSL results, and hence can rule out those low-SNR USVs. Using the proposed method, we achieved a median SSL error of 5.01 cm in the test dataset, which is suitable for rat-rat social scenarios. Since the large size of adult rats (>18 cm), snout-to-snout distance in rat-rat interaction scenarios is much bigger than in mouse-mouse interactions (see Figure S9). Such an SSL error can be usable for safely assigning the USVs to the vocal rat in rat-rat interaction scenarios.

Another concern that haunted the experts working on SSL is the lack of accurate test datasets. Mainstream studies reportedly induced mouse USVs using heterosexual urine,^{42,43} but the visual localization of rodent snout would bring unnecessary errors, and lack the ability of snout height detection. In this article, we constructed a rat USV test dataset for evaluating SSL errors using USVs collected from still rats, which brings two advantages: 1) all the USVs are produced at the same spatial source position because the rat kept still throughout the process; 2) 3D position of the vocal rat nose/snout is provided, therefore allowing evaluating 3D SSL algorithms.

We further show that a higher LCI is associated with smaller sound source localization (SSL) error (see Figure S10). Therefore, only USV signals with an LCI over 0.6 are accepted, which consist of 76% of the valid vocal repertoire (29772 USVs, intersection of the valid outputs of ARBUR:USV and ARBUR:Behavior, see Figure S11). ARBUR:SSL shows decent horizontal-plane localization precision, with less than 10 cm error for over 78% of the data, which is considerable, considering the much weaker USVs in our dataset, compared to existing algorithms^{12,23} (see Figures 3F and S10A), and more importantly, suitable for our purposes. Moreover, by evaluating the LCI at different heights (reconstructed from the binocular view), ARBUR:SSL can pinpoint the vocal rat even if the two rats overlap in the top-down projection plane (which is more frequent in, for example, pinning and pouncing). ARBUR:SSL assigned 17,235 USVs out of 29,772 USVs to the vocal rat. Figure 3A shows examples. In contrast to other sound source localization methods for rodents' USVs, ARBUR:SSL features the ability to locate the vocal rat in 3-D space while preserving the localization precision with weaker USV signals (Table 4).

Revealing the behavior-vocalization interplay

By seamlessly incorporating these three modules, ARBUR can build connections between different types of rat behavior and vocalizations. Figure 4A shows examples (also see Video S1). It is intuitive that when the rats are engaged in non-aggressive play behavior, such as social nose contact and pouncing, more appetitive 50-kHz USVs are communicated. On the contrary, the submissive one emits more aversive 22-kHz USVs when engaged in aggressive play behavior, i.e., pinning. We investigated whether these two hypotheses would hold in the whole

Table 4. Comparison of ARBUR:SSL with mainstream sound source localization methods for rodents' USVs

	Dimen. of SSL	Number of mics	Principle of SSL	Median error (cm)	Mic-to-source distance (cm)	Estimation area (cm ²)	Rodents
Matsumoto et al. ⁴⁴	2-D	4	Direction of arrival	1.58	47	20 × 20	Mice
Heckman et al. ⁴⁵	1-D	2	Hyperbolic	0.85	46	7 × 7.5	Mice
Warren et al. ⁴⁶	2-D	8	Hyperbolic	1.80	38	66 × 66	Mice
Neunuebel et al. ²³	2-D	4	Hyperbolic	3.87	33	66 × 66	Mice
Oliveira-Stahl et al. ⁴²	2-D	4	Hyperbolic	1.31	27.8	30 × 40	Mice
Sterling et al. ⁴³	2-D	4 + 1(64-channel)	Hyperbolic	0.48	27.8	30 × 40	Mice
ARBUR:SSL	3-D	4	Hyperbolic	5.01	55	D50	Rats

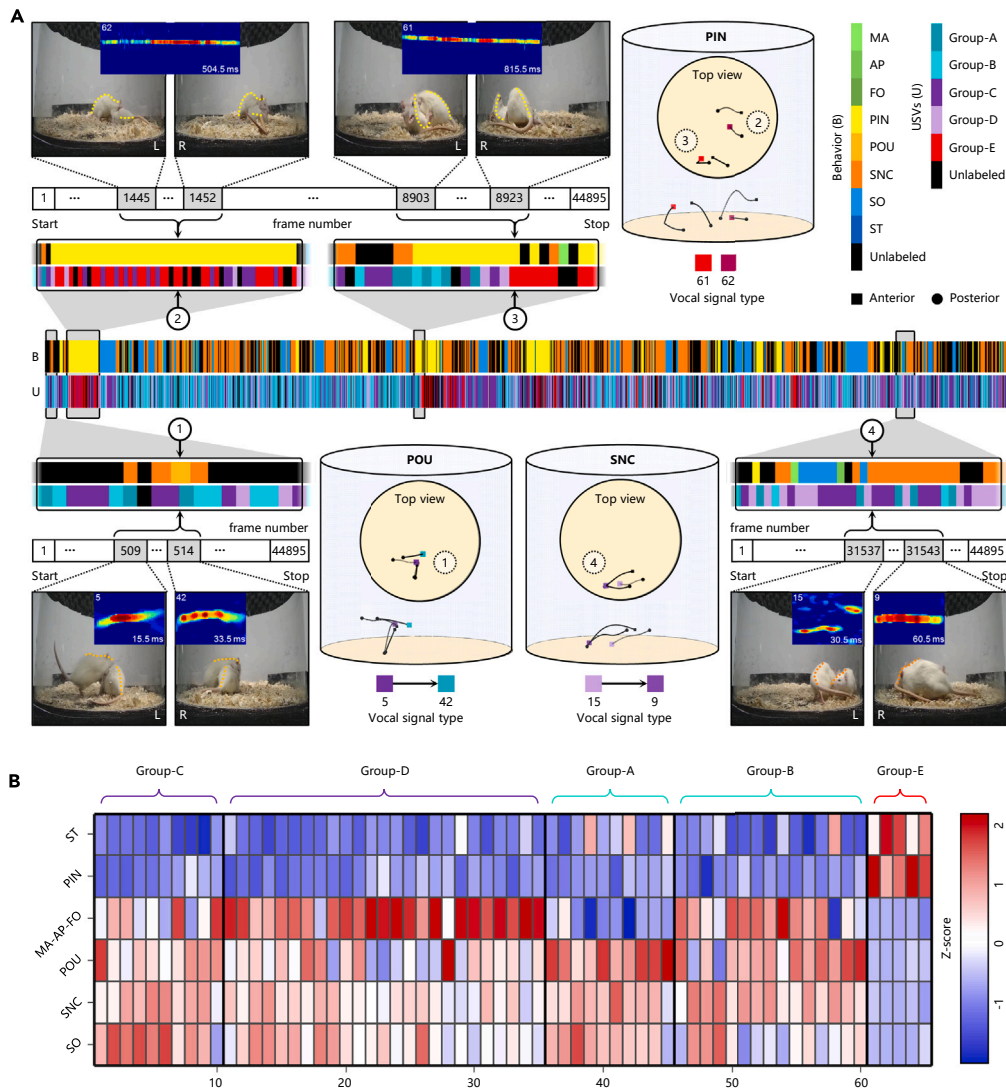


Figure 4. Revealing the behavior-vocalization interplay using ARBUR

(A) Examples showing how 22-kHz USVs are associated with pinning (top), and appetitive 50-kHz USVs are associated with pouncing (bottom left) and social nose contact (bottom right). Frequency range of spectrograms: 20 kHz for clusters 5, 9, and 42; 0–40 kHz for clusters 61 and 62; and 40 kHz for cluster 15.

(B) Quantification of behavior-associated USV distributions, which indicates whether a certain cluster of USVs is emitted during a specific type of behavior above chance (red), at chance (white), or below chance (blue). Groups A-E, the same as in Figure 2.

repertoire. Therefore, we quantified the USV proportions of rats in different types of behavior (Figure 4B). Results show that the submissive rat produces significantly more aversive 22-kHz USVs during the pinning process. Moreover, in social nose contact and pouncing, both rats produce more 50-kHz USVs.

The finding that during an aggressive interaction (i.e., pinning), the submissive one produces significantly more aversive 22-kHz USVs, is consistent with former research.^{26,47,48} We also show that even though pouncing and pinning are similar types of behavior, they actually have distinct USV distributions. In particular, during pinning, the submissive rat rarely emits appetitive 50-kHz USVs, whereas, during pouncing, the submissive one produces mainly 50-kHz USVs. This may indicate that the submissive rat has distinct feelings and emotional states when engaging in the two types of behavior. ARBUR can reveal such a distinction because it can not only accurately discriminate these behaviors, but locate the submissive rat spatially. These findings are in line with former studies and therefore validate the effectiveness of ARBUR.

In addition, ARBUR reveals that several USVs, mostly 22-kHz USVs, are recorded when both rats are in the still state which resembles sleeping (Figure 4B). The still state detected by ARBUR is characterized by both rats lying flat for several minutes or longer, with their bodies immobilized and eyes closed. This indicates that some USVs may be emitted by rats unconsciously (possibly during sleeping), and these 22-kHz USVs may neither convey emotional information nor trigger behavioral changes of the other rat. However, the still-associated (or sleeping-associated) USVs have rarely been focused on because of their independence from social interaction. Furthermore, ARBUR shows

that rats produce more step 50-kHz USVs than non-step ones during motion (Figure 4B), indicating a part of step 50-kHz USVs may be caused by rats' locomotion. Recent research has demonstrated the tight correlation between respiratory activity and frequency features of USVs,²⁵ which indicated that locomotion may induce more step 50-kHz USVs by affecting rat respiratory activity.

DISCUSSION

Understanding the behavior-vocalization interplay of rats is inhibited by the difficulty of relating the behaviors and the USVs of freely behaving rats in complex social contexts. In this study, we propose a machine learning-based analysis system (named ARBUR) to relate rat vocalizations to their free behaviors. ARBUR contains three modules: 1) an comprehensive three-step algorithm (ARBUR:USV) that clusters rat USVs in an unbiased manner by considering their mean frequency, duration, and both non-step (continuous) and step (non-continuous) contours in the spectrograms; 2) a machine learning-based hierarchical framework (ARBUR:Behavior) that detects the type of social behavior from a laterally binocular view; and 3) a sound source localization algorithm (ARBUR:SSL) that spatially allocates the USVs to the vocal rats.

Rats produce diverse USVs in terms of duration, mean frequency, and frequency contours.²⁴ Existing studies have only clustered the USVs in a way that considers some of these factors in an automated^{12,16} or manual^{14,25} manner (Table 2). Among them, a recent study leveraged unsupervised learning to cluster mouse non-step USVs into 22 categories based on their contours, which showed potential for unbiased clustering of a broader spectrum of USVs.¹² However, an automated clustering algorithm that comprehensively takes all these three factors and step signals under consideration remains missing. To our knowledge, ARBUR:USV is the first automated clustering method that considers the step USVs of rodents. Although it is designed and optimized to cluster rat USVs, it can be readily generalized for mouse USV applications. Compared with algorithms from existing studies, the proposed clustering algorithm features: 1) the comprehensive consideration of frequency, duration, and contour information; and 2) the inclusion of both step and non-step USVs. Moreover, ARBUR:USV provides multiple choices for users to determine the optimal clusters, including subjective setting, the progressive method, and the elbow method. It should therefore be adaptive to a wide range of applications regarding rodent acoustic clustering and analysis.

ARBUR:Behavior is a hierarchically supervised learning algorithm that archives state-of-the-art performance in classifying eight types of common rat behaviors, including three easy-to-confuse social behaviors (pinning, pouncing, and social nose contact). We show that the mainstream deep learning classification algorithms^{34–37} are not competent for the end-to-end classification of complex rat behaviors. For example, ResNet achieves the best performance with an accuracy of only 0.333 using an eight-class classifier (see Figure S8A). This highlights the necessity of the hierarchical classification of complex rat behaviors. Therefore, ARBUR:Behavior hierarchically classifies rats' non-social, moving, and social state behaviors sequentially. It detects non-social behaviors using the location feature of rats in the image and moving behaviors by estimating the optical flow of the rat as well as its relative orientation. To discriminate tricky social behaviors, ARBUR:Behavior extracts "histogram of oriented gradients" (HOG)⁴⁰ descriptors to represent the statistics of social posture features in a standardized way. This empowers ARBUR:Behavior with high-accuracy discrimination of rat social behaviors. ARBUR:Behavior outperforms existing studies with respect to comprehensive detection performance (Table 3).

ARBUR:SSL can locate the vocal rat during two-rat free-behaving scenarios by incorporating the lateral binocular view (for rat nose reconstruction in the Cartesian space), behavior classification results, height-sensitive distributed microphone configuration, and 3-D probabilistic triangulation of sound source. The large rat-to-mic distance in rat-rat interacting scenarios makes the recorded USV signal intensity (therefore signal-to-noise ratio) much weaker than others. ARBUR:SSL addressed this challenge by measuring the consistency between 3-mic SSL and 4-mic SSL results, hence ruling out those low-SNR USVs. It achieves desirable SSL precision suitable for rat-rat social scenarios. Another concern that haunted the experts working on SSL is the lack of accurate test datasets. Maintream studies reportedly induced mouse USVs using heterosexual urine,^{42,43} but the visual localization of rodent snout would bring unnecessary errors, and lack the ability of snout height detection. In this work, we constructed a rat USV test dataset for evaluating SSL errors using USVs collected from still rats, which brings two advantages: 1) all the USVs are produced at the same spatial source position because the rat kept still throughout the process; 2) 3D position of the vocal rat nose/snout is provided, therefore allowing evaluating 3D SSL algorithms. ARBUR:SSL can therefore contribute to SSL research by highlighting the potential of refining low-SNR vocal signals and stressing the value of USVs emitted by still or sleeping rats for evaluating SSL algorithms.

By seamlessly combining these three modules, ARBUR features the advantages of comprehensive and unbiased USV clustering, hierarchical high-accuracy rat behavior detection, and 3-D sound source localization to reveal the latent behavior-vocalization interplay of rats. For example, during a socio-aggressive interaction (i.e., pinning), the submissive one produces significantly more aversive 22-kHz USVs, which is consistent with former research.^{26,47,48} We also show that even though pouncing and pinning are similar types of behavior, they actually have distinct USV distributions. Moreover, ARBUR indicates two novel findings. First, rats may unconsciously emit 22-kHz USVs during the still state. The still state resembles sleeping and should not be confused with freezing because it features both rats lying flat on the floor for several minutes or longer, with their eyes closed. In contrast, freezing follows the cessation of an ongoing behavior (moving, grooming, and so forth),^{49,50} typically not occurring while calmly lying with the eyes closed. Second, rats possibly produce more discontinuous (step) 50-kHz USVs during moving. It's possible that the movement of rats can cause changes in their respiratory behavior, thereby affecting USV vocalization. But the direct evidence to support this hypothesis is missing, and more solid examinations are required to reveal how movement behavior or specific actions affect breathing patterns, and correspondingly affect the continuity of emitted USVs. These two directions have long been neglected in the manual analysis-dominated research, and call for more stringent verification to bring substantial knowledge. We also note that although ARBUR is designed and optimized to investigate the behavior-vocalization interplay of rats, it can be potentially generalized to other rodents that also communicate through both non-verbal signals and USVs (for example, mice and mole rats).

In summary, we proposed a machine learning-based analysis system, which can not only automatically reveal the well-understood behavior-associated vocalizations that were carefully concluded by other behavioral researchers, but also hold the promise to indicate novel findings that can be hardly found by manual analysis, especially regarding step USVs and the active/passive rat-associated USVs during easy-to-confuse social behaviors. This work highlights the potential of machine learning algorithms in automatic animal behavioral and acoustic analysis and could help mechanistically understand the interactive influence between the behaviors and USVs of rats.

Limitations of the study

Despite the above advancements, ARBUR can be further improved in the following aspects. First, ARBUR:USV rejects the concurrent USVs (simultaneously recorded USVs at different frequencies), which may be important during social communication and could be separated by intelligent segmentation algorithms combined with high-precision SSL, further enriching the vocal repertoire. Moreover, ARBUR:Behavior lacks the ability to detect potential behavioral changes in freely behaving rats, which is essential for discovering the underlying mechanisms of novel behavior-vocalization and could be enhanced by incorporating unsupervised learning to uncover latent structures and categories in behavioral space.^{51,52} In addition, ARBUR:SSL could be improved in terms of the limited number of visual reconstructions of rat noses by incorporating more camera views in the future, fulfilling the need for a more reliable three-dimensional sound source localization. Using ARBUR, we revealed that rat USV distribution is biased by rat behavior, indicating that the USVs of rats are associated with their behaviors. For example, a submissive rat significantly up-regulates the aversive 22-kHz USVs during pinning. Moreover, during pouncing and social nose contact, both rats produce substantially more appetitive 50-kHz USVs than they do during pinning.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE
- RESOURCE AVAILABILITY
 - Lead contact
 - Materials availability
 - Data and code availability
- METHOD DETAILS
 - Animals and experiment preparation
 - Audio segmentation
 - Vocal signals clustering
 - Behavior classification
 - Sound source localization
- QUANTIFICATION AND STATISTICAL ANALYSIS

SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.isci.2024.109998>.

ACKNOWLEDGMENTS

This study was funded by the National Natural Science Foundation of China under Grant 62088101, the Science and Technology Innovation Program of Beijing Institute of Technology under Grant 2022CX01010, and the National Science and Technology Key major projects (STI2030-Major Projects 2022ZD02068000).

AUTHOR CONTRIBUTIONS

Q.S. conceived and supervised the project. Z.C. wrote the article and analyzed the results. Z.C. and G.J. implemented the methods, conducted the experiments, and processed the data. Q.Z. and Y.Z. assisted in conducting the experiments. All authors contributed to discussions. All authors edited and approved the article.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: March 8, 2024

Revised: April 1, 2024

Accepted: May 14, 2024

Published: May 18, 2024

REFERENCES

- Barker, A.J., Veviuurko, G., Bennett, N.C., Hart, D.W., Mograby, L., and Lewin, G.R. (2021). Cultural transmission of vocal dialect in the naked mole-rat. *Science* 371, 503–507. <https://doi.org/10.1126/science.abc6588>.
- Padilla-Coreano, N., Batra, K., Patarino, M., Chen, Z., Rock, R.R., Zhang, R., Hausmann, S.B., Weddington, J.C., Patel, R., Zhang, Y.E., et al. (2022). Cortical ensembles orchestrate social competition through hypothalamic outputs. *Nature* 603, 667–671. <https://doi.org/10.1038/s41586-022-04507-5>.
- Dolensek, N., Gehrlach, D.A., Klein, A.S., and Gogolla, N. (2020). Facial expressions of emotion states and their neuronal correlates in mice. *Science* 368, 89–94. <https://doi.org/10.1126/science.aaz9468>.
- Dong, S., Lin, T., Nieh, J.C., and Tan, K. (2023). Social signal learning of the waggle dance in honey bees. *Science* 379, 1015–1018. <https://doi.org/10.1126/science.ade1702>.
- VanRyzin, J.W., Marquardt, A.E., Argue, K.J., Vecchiarelli, H.A., Ashton, S.E., Arambula, S.E., Hill, M.N., and McCarthy, M.M. (2019). Microglial Phagocytosis of Newborn Cells Is Induced by Endocannabinoids and Sculpted Sex Differences in Juvenile Rat Social Play. *Neuron* 102, 435–449.e6. <https://doi.org/10.1016/j.neuron.2019.02.006>.
- Robinson, G.E., Fernald, R.D., and Clayton, D.F. (2008). Genes and social behavior. *Science* 322, 896–900. <https://doi.org/10.1126/science.1159277>.
- Li, S.W., Zeliger, O., Strahs, L., Báez-Mendoza, R., Johnson, L.M., McDonald Wojciechowski, A., and Williams, Z.M. (2022). Frontal neurons driving competitive behaviour and ecology of social groups. *Nature* 603, 661–666. <https://doi.org/10.1038/s41586-021-04000-5>.
- Schneider, S., Lee, J.H., and Mathis, M.W. (2023). Learnable latent embeddings for joint behavioural and neural analysis. *Nature* 617, 360–368. <https://doi.org/10.1038/s41586-023-06031-6>.
- Wei, D., Talwar, V., and Lin, D. (2021). Neural circuits of social behaviors: Innate yet flexible. *Neuron* 109, 1600–1620. <https://doi.org/10.1016/j.neuron.2021.02.012>.
- Murugan, M., Jang, H.J., Park, M., Miller, E.M., Cox, J., Taliaferro, J.P., Parker, N.F., Bhawe, V., Hur, H., Liang, Y., et al. (2017). Combined Social and Spatial Coding in a Descending Projection from the Prefrontal Cortex. *Cell* 171, 1663–1677.e16. <https://doi.org/10.1016/j.cell.2017.11.002>.
- Knutson, B., Burgdorf, J., and Panksepp, J. (1998). Anticipation of play elicits high-frequency ultrasonic vocalizations in young rats. *J. Comp. Psychol.* 112, 65–73. <https://doi.org/10.1037/0735-7036.112.1.65>.
- Sangiameo, D.T., Warren, M.R., and Neunuebel, J.P. (2020). Ultrasonic signals associated with different types of social behavior of mice. *Nat. Neurosci.* 23, 411–422. <https://doi.org/10.1038/s41593-020-0584-z>.
- Venniro, M., and Shaham, Y. (2020). An operant social self-administration and choice model in rats. *Nat. Protoc.* 15, 1542–1559. <https://doi.org/10.1038/s41596-020-0296-6>.
- Brudzynski, S.M. (2013). Ethotransmission: communication of emotional states through ultrasonic vocalization in rats. *Curr. Opin. Neurobiol.* 23, 310–317. <https://doi.org/10.1016/j.conb.2013.01.014>.
- Ben-Ami Bartal, I., Decety, J., and Mason, P. (2011). Empathy and pro-social behavior in rats. *Science* 334, 1427–1430. <https://doi.org/10.1126/science.1210789>.
- Takahashi, N., Kashino, M., and Hironaka, N. (2010). Structure of rat ultrasonic vocalizations and its relevance to behavior. *PLoS One* 5, e14115. <https://doi.org/10.1371/journal.pone.0014115>.
- Lorbach, M., Kyriakou, E.I., Poppe, R., van Dam, E.A., Noldus, L.P.J.J., and Veltkamp, R.C. (2018). Learning to recognize rat social behavior: Novel dataset and cross-dataset application. *J. Neurosci. Methods* 300, 166–172. <https://doi.org/10.1016/j.jneumeth.2017.05.006>.
- Vogt, N. (2021). Automated behavioral analysis. *Nat. Methods* 18, 29. <https://doi.org/10.1038/s41592-020-01030-1>.
- Bohnslav, J.P., Wimalasena, N.K., Clousing, K.J., Dai, Y.Y., Yarmolinsky, D.A., Cruz, T., Kashlan, A.D., Chiappe, M.E., Orefice, L.L., Woolf, C.J., and Harvey, C.D. (2021). DeepEthogram, a machine learning pipeline for supervised behavior classification from raw pixels. *Elife* 10, e63377. <https://doi.org/10.7554/eLife.63377>.
- Harris, C., Finn, K.R., Kieseler, M.L., Maechler, M.R., and Tse, P.U. (2023). DeepAction: a MATLAB toolbox for automated classification of animal behavior in video. *Sci. Rep.* 13, 2688. <https://doi.org/10.1038/s41598-023-29574-0>.
- Kabra, M., Robie, A.A., Rivera-Alba, M., Branson, S., and Branson, K. (2013). JAABA: interactive machine learning for automatic annotation of animal behavior. *Nat. Methods* 10, 64–67. <https://doi.org/10.1038/nmeth.2281>.
- Segalin, C., Williams, J., Karigo, T., Hui, M., Zelikowsky, M., Sun, J.J., Perona, P., Anderson, D.J., and Kennedy, A. (2021). The Mouse Action Recognition System (MARS) software pipeline for automated analysis of social behaviors in mice. *Elife* 10, e63720. <https://doi.org/10.7554/eLife.63720>.
- Neunuebel, J.P., Taylor, A.L., Arthur, B.J., and Egnor, S.E.R. (2015). Female mice ultrasonically interact with males during courtship displays. *Elife* 4, e06203. <https://doi.org/10.7554/eLife.06203>.
- Wright, J.M., Gourdon, J.C., and Clarke, P.B.S. (2010). Identification of multiple call categories within the rich repertoire of adult rat 50-kHz ultrasonic vocalizations: effects of amphetamine and social context. *Psychopharmacology (Berl)* 211, 1–13. <https://doi.org/10.1007/s00213-010-1859-y>.
- Riede, T. (2013). Stereotypic laryngeal and respiratory motor patterns generate different call types in rat ultrasound vocalization. *J. Exp. Zool. A Ecol. Genet. Physiol.* 319, 213–224. <https://doi.org/10.1002/jez.1785>.
- Burgdorf, J., Kroes, R.A., Moskal, J.R., Pfau, J.G., Brudzynski, S.M., and Panksepp, J. (2008). Ultrasonic vocalizations of rats (*Rattus norvegicus*) during mating, play, and aggression: Behavioral concomitants, relationship to reward, and self-administration of playback. *J. Comp. Psychol.* 122, 357–367. <https://doi.org/10.1037/a0012889>.
- Fonseca, A.H., Santana, G.M., Bosque, O.G., Bampi, S., and Dietrich, M.O. (2021). Analysis of ultrasonic vocalizations from mice using computer vision and machine learning. *Elife* 10, e59161. <https://doi.org/10.7554/eLife.59161>.
- Van Segbroeck, M., Knoll, A.T., Levitt, P., and Narayanan, S. (2017). MUPET-Mouse Ultrasonic Profile ExTraction: A Signal Processing Tool for Rapid and Unsupervised Analysis of Ultrasonic Vocalizations. *Neuron* 94, 465–485.e5. <https://doi.org/10.1016/j.neuron.2017.04.005>.
- Coffey, K.R., Marx, R.E., and Neumaier, J.F. (2019). DeepSqueak: a deep learning-based system for detection and analysis of ultrasonic vocalizations. *Neuropsychopharmacology* 44, 859–868. <https://doi.org/10.1038/s41386-018-0303-6>.
- Goffinet, J., Brudner, S., Mooney, R., and Pearson, J. (2021). Low-dimensional learned feature spaces quantify individual and group differences in vocal repertoires. *Elife* 10, e67855. <https://doi.org/10.7554/eLife.67855>.
- Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J.M., and Luo, P. (2021). SegFormer: Simple and efficient design for semantic segmentation with transformers. *Adv. Neural Inf. Process. Syst.* 34, 12077–12090.
- Gerós, A., Cruz, R., de Chaumont, F., Cardoso, J.S., and Aguiar, P. (2022). Deep learning-based system for real-time behavior recognition and closed-loop control of behavioral mazes using depth sensing. Preprint at bioRxiv. <https://doi.org/10.1101/2022.02.22.481410>.
- Marks, M., Qiuhuan, J., Sturman, O., von Ziegler, L., Kollmorgen, S., von der Behrens, W., Mante, V., Bohacek, J., and Yanik, M.F. (2022). Deep-learning based identification, tracking, pose estimation, and behavior classification of interacting primates and mice in complex environments. *Nat. Mach. Intell.* 4, 331–340. <https://doi.org/10.1038/s42256-022-00477-5>.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022.
- Liu, Z., Mao, H., Wu, C., Feichtenhofer, C., Darrell, T., and Xie, S. (2022). A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11976–11986.
- Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., and Jégou, H. (2021). Training data-efficient image transformers & distillation through attention. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 10347–10357.
- Hsu, A.I., and Yttri, E.A. (2021). B-SOid, an open-source unsupervised algorithm for identification and fast prediction of behaviors. *Nat. Commun.* 12, 5188. <https://doi.org/10.1038/s41467-021-25420-x>.
- Li, J., Keselman, M., and Shlizerman, E. (2022). OpenLabCluster: Active learning based clustering and classification of animal behaviors in videos based on automatically extracted kinematic body keypoints. Preprint at bioRxiv. <https://doi.org/10.1101/2022.10.10.511660>.

40. Dalal, N., and Triggs, B. (2005). Histograms of oriented gradients for human detection. In 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 886–893.
41. Jhuang, H., Garrote, E., Mutch, J., Yu, X., Khilnani, V., Poggio, T., Steele, A.D., and Serre, T. (2010). Automated home-cage behavioural phenotyping of mice. *Nat. Commun.* 1, 68. <https://doi.org/10.1038/ncomms1064>.
42. Oliveira-Stahl, G., Farboud, S., Sterling, M.L., Heckman, J.J., van Raalte, B., Lenferink, D., van der Stam, A., Smeets, C.J.L.M., Fisher, S.E., and Englitz, B. (2023). High-precision spatial analysis of mouse courtship vocalization behavior reveals sex and strain differences. *Sci. Rep.* 13, 5219. <https://doi.org/10.1038/s41598-023-31554-3>.
43. Sterling, M.L., Teunisse, R., and Englitz, B. (2023). Rodent ultrasonic vocal interaction resolved with millimeter precision using hybrid beamforming. *Elife* 12, e86126. <https://doi.org/10.7554/eLife.86126>.
44. Matsumoto, J., Kanno, K., Kato, M., Nishimaru, H., Setogawa, T., Chinzorig, C., Shibata, T., and Nishijo, H. (2022). Acoustic camera system for measuring ultrasound communication in mice. *iScience* 25, 104812. <https://doi.org/10.1016/j.isci.2022.104812>.
45. Heckman, J.J., Proville, R., Heckman, G.J., Azarfar, A., Celikel, T., and Englitz, B. (2017). High-precision spatial localization of mouse vocalizations during social interaction. *Sci. Rep.* 7, 3017. <https://doi.org/10.1038/s41598-017-02954-z>.
46. Warren, M.R., Sangiamo, D.T., and Neunuebel, J.P. (2018). High channel count microphone array accurately and precisely localizes ultrasonic signals from freely-moving mice. *J. Neurosci. Methods* 297, 44–60. <https://doi.org/10.1016/j.jneumeth.2017.12.013>.
47. Litvin, Y., Blanchard, D.C., and Blanchard, R.J. (2007). Rat 22kHz ultrasonic vocalizations as alarm cries. *Behav. Brain Res.* 182, 166–172. <https://doi.org/10.1016/j.bbr.2006.11.038>.
48. Lukas, M., and Wöhr, M. (2015). Endogenous vasopressin, innate anxiety, and the emission of pro-social 50-kHz ultrasonic vocalizations during social play behavior in juvenile rats. *Psychoneuroendocrinology* 56, 35–44. <https://doi.org/10.1016/j.psyneuen.2015.03.005>.
49. Fanselow, M.S., and Bolles, R.C. (1979). Naloxone and shock-elicited freezing in the rat. *J. Comp. Physiol. Psychol.* 93, 736–744. <https://doi.org/10.1037/h0077609>.
50. Hashimoto, S., Inoue, T., and Koyama, T. (1996). Serotonin reuptake inhibitors reduce conditioned fear stress-induced freezing behavior in rats. *Psychopharmacology (Berl)* 123, 182–186. <https://doi.org/10.1007/BF02246175>.
51. Brattoli, B., Büchler, U., Dorkenwald, M., Reiser, P., Füll, L., Helmchen, F., Wahl, A.S., and Ommer, B. (2021). Unsupervised behaviour analysis and magnification (uBAM) using deep learning. *Nat. Mach. Intell.* 3, 495–506. <https://doi.org/10.1038/s42256-021-00326-x>.
52. Huang, K., Han, Y., Chen, K., Pan, H., Zhao, G., Yi, W., Li, X., Liu, S., Wei, P., and Wang, L. (2021). A hierarchical 3D-motion learning framework for animal spontaneous behavior mapping. *Nat. Commun.* 12, 2784. <https://doi.org/10.1038/s41467-021-22970-y>.
53. Jia, G., Chen, Z., Zhou, Q., Zhang, Y., Chen, X., Fukuda, T., Huang, Q., and Shi, Q. (2023). ARBUR: A machine learning-based Analysis system for Relating the Behavior and USVs of Rats (Zenodo). <https://doi.org/10.5281/zenodo.8081539>.
54. Tachibana, R.O., Kanno, K., Okabe, S., Kobayashi, K.I., and Okanoya, K. (2020). USVSEG: A robust method for segmentation of ultrasonic vocalizations in rodents. *PLoS One* 15, e0228907. <https://doi.org/10.1371/journal.pone.0228907>.
55. Percival, D.B., and Walden, A.T. (1993). *Spectral Analysis for Physical Applications* (Cambridge University Press).

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Experimental models: Organisms/strains		
Sprague-Dawley Rat	SPF biotechnology, Beijing	https://www.spfbiotech.com
Software and algorithms		
ARBUR	This paper	https://github.com/Guanglu-Jia/ARBUR
MATLAB	The Mathworks, Inc.	https://www.mathworks.com/products/matlab.html

RESOURCE AVAILABILITY

Lead contact

Further information and requests for resources should be directed to and will be fulfilled by the lead contact, Qing Shi (shiqing@bit.edu.cn).

Materials availability

This study did not generate new unique reagents.

Data and code availability

All data needed to evaluate the conclusions in the paper are present in the Article. The datasets generated and/or analyzed during the current study are available from the [lead contact](#) upon reasonable request.

The open-source software ARBUR and its tutorials are freely available at GitHub (<https://github.com/Guanglu-Jia/ARBUR>) and Zenodo (<https://doi.org/10.5281/zenodo.8081539>).⁵³

Any additional information required to reanalyze the data reported in this paper is available from the [lead contact](#) upon request.

METHOD DETAILS

Animals and experiment preparation

Animals

Adult (8-10 weeks) male (n = 4) and female (n = 4) Sprague-Dawley rats (stock number 198499, 212822, SPF biotechnology, Beijing, China) were used. Rats were kept on a 12/12 h light/dark cycle at a consistent ambient temperature (26 ± 1 °C) and humidity ($50 \pm 5\%$), and all experiments were performed during the light cycle. Food and water were accessed *ad libitum*. All experimental procedures were approved by the Institutional Animal Care and Use Committee of the Beijing Institute of Technology, Beijing, China. Before all recording experiments, rats were singly housed for at least 3 days to minimize group housing effects on the social behavior of the rats.

Video and audio stream recording

The video and audio simultaneous recording experimental set-up is shown in [Figure S1](#). Rats were able to behave freely in a circular open-field arena made of a transparent acrylic barrel with a base diameter of 0.5 m and a height of 0.5 m. A black velvet cloth was placed underneath the recording area as well as wood chips to avoid reflecting light. Encircling the recording area with sound-absorbing material above can reduce sound reflections from the walls. Two synchronized cameras were mounted at an angle of 60 degrees on either side of the acquisition area. Four ultrasonic microphones (CM16/CMPA, Avisoft Bioacoustics, Glienicke, Germany) are fixed evenly to the edges of the transparent acrylic barrel. The video stream was recorded at 25 frames per second. The four channels of the audio stream from the four microphones were simultaneously recorded through data acquisition equipment (UltraSoundGate 416H, Avisoft Bioacoustics). The video and audio streams were synchronized by aligning the recording onset timing. Both the streams were stored on a high-performance computer to avoid recording time shifting (32G RAM). The video and audio stream recording experiments were conducted over 4 days (12 hours of continuous recording per day) with no external human interference.

Audio segmentation

Audio streams were segmented automatically with multi-taper spectral analysis. The vocal signals from microphone 1 were segmented using the USVSEG algorithm,^{44,54} and the parameters of the algorithm were optimized as follows. First, multi-taper spectrograms (time-frequency matrices) were generated using six discrete prolate spheroidal sequences of length 512 as windowing functions (NW=3).⁵⁵ Then the

spectrograms were flattened by replacing the first three cepstral coefficients with zero to reduce transient broadband noises and by subtracting the median spectrum. After flattening, we thresholded the flattened spectrograms with a threshold of 4.5 to further reduce the environmental noise. These spectrograms were band-pass filtered between 10 and 100 kHz. In addition, the maximum duration of an audio segment was set to 5000 ms and the minimum duration was set to 5 ms. Sound elements with a more than 30 ms gap were judged as two individual syllables and segmented as two independent spectrograms accordingly. Finally, these segmentation timings from microphone 1 were used to segment vocal signals from the other three microphones. A total of 43357 ultrasonic signals were segmented from all the audio data collected.

Vocal signals clustering

Extraction of frequency contours

The segmented spectrograms (time-frequency matrices) were filtered sequentially by the median, Gaussian, and threshold filters to enhance the signal-to-noise ratio. The filtered spectrograms were then binarized for detecting the boundaries of the frequency contours (edges in the spectrograms) by the modified Moore's neighbor tracking algorithm with the Jacobian's termination condition. Those contours with less than 22 points and a maximum intensity inside the contour of less than 15 were detected as isolated islands and eliminated. For the discontinuous USVs with multiple contours, only the longest four contours were selected as the final contours to recapitulate the original complex signals. Harmonics were considered as those frequency contours with 90% overlap in time, and only the frequency contour of the lowest mean frequency was preserved whereas others were eliminated. Other than harmonics, those frequency contours that share an overlap of duration over 3 ms were considered overlapping vocalizations emitted by two rats and also abolished.

Coarse clustering

Coarse clustering consists of two steps. In the first step, the USVs were clustered coarsely into 2 groups according to their mean frequency: aversive 22-kHz and appetitive 50-kHz USVs. In the second step, the 2 groups were further divided into 5 groups (A-D) according to their distribution of duration and mean frequency of frequency contours: so-called 22-kHz, low-peak-frequency/high-peak-frequency step/non-step 50-kHz USVs (see [Figure S2A](#)). Specifically, those USVs with a mean frequency of less than 32 kHz were classified as 22-kHz signals. The rest of the USVs (50-kHz) were clustered into 2 groups (low-peak-frequency 50-kHz and high-peak-frequency 50-kHz) by k-means clustering according to two features: duration and mean frequency. These USVs were further classified into 4 groups according to their continuity of frequency contours: non-step (continuous) USVs containing only one contour and step USVs containing more than one contour.

Determining the optimal number of clusters

In the third step, groups A were further divided using unsupervised learning within their groups. For the appropriate setting of the cluster number, ARBUR:USV provides users with three choices: subjective, automatic, and progressive methods. The subjective method allows users to set an integer number for each cluster. For example, in [Figure 2](#), cluster numbers for groups A-E are set as 10, 25, 10, 15, and 5, respectively. The automatic method relies on the well-known elbow method to avoid over-clustering, which works as follows. For each cluster number, the total within-cluster error (TWCE) is calculated. By sliding from 2 clusters to 100 clusters (200 for group B for its complexity), forming a TWCE-clusters curve. Then, for each point on this curve, its left and right sides are used for linear fitting, and the total fitting error is stored, which forms a second curve: fitting error versus cluster number (see [Figure S3](#) as an example). The global bottom of the second curve is called the inflection point, which indicates the optimal number of clusters. The progressive method takes a user-defined intra-cluster variance parameter to find the appropriate number. It calculates the innerpoint percentage (percent of points falling within the mean $\pm 2.5SD$ range of each cluster). Theoretically, with increasing clusters, the intra-cluster variance decreases. This method starts with 2 clusters, calculates the average innerpoint percent, and stops when the user-defined threshold is first reached. For example, if the threshold is set as 98%, that is, an average of 98% points should be the inner points for all clusters and all features. The optimal cluster should be 8 for group E (see [Figure S4](#)).

Refined clustering of non-step USVs

The non-step (continuous) signals were clustered further according to their contour shapes. The detected boundary points were centered by subtracting the mean frequency. To normalize these centered contours with different durations, we constructed a fixed-length feature vector containing 100 features (see [Figure S2B](#)). The first 50 and last 50 features represent the upper and the bottom boundaries of the centered contour, respectively. Each contour was linearly mapped to a 100-feature vector. By doing so, frequency contours with different durations or mean frequencies can be clustered appropriately according to their shapes. These USVs were categorized using k-means clustering based on the 100-feature vectors. The number of clusters for low-frequency and high-frequency USVs were both subjectively set as 10. Alternatively, the number of clusters can be determined progressively as in the study of Sangioamo et al.¹²

Refined clustering of step signals

For the step (discontinuous) signals with multiple contours, a similar 100-feature vector was constructed based on these contours. All contours of each USV were centered by subtracting their mean frequencies. Then, the upper boundaries of each contour were sorted in ascending order in time and mapped linearly to a 50-feature vector. By doing so, those inter-contour intervals were deleted, and contours were placed compactly. Similarly, the bottom boundaries were mapped to the second 50-feature vector. These step USVs were then categorized using

k-means clustering based on the 100-feature vectors concatenating these two vectors together. The numbers of clusters for low-frequency and high-frequency USVs were subjectively set as 15 and 25, respectively.

Refined clustering of 22-kHz USVs

22-kHz USVs vary a lot in duration, ranging from several microseconds to several seconds, such a feature should not be ignored in clustering. In contrast, the contour shapes varied little across USVs. Therefore, we constructed a 2-feature vector for each 22-kHz signal containing the mean frequency and duration. These 22-kHz signals were then categorized using k-means clustering based on the 2-feature vectors. The number of clusters was subjectively set to 5.

Validating clustering

To validate clustering results, dissimilarities within each cluster (intra-cluster variance) and between different clusters (inter-cluster distance) were evaluated, respectively. They were calculated based on the constructed 100-feature vectors for both step and non-step USVs (see Figure S6A). The inter-cluster distance of the i th cluster D_i was defined as the average difference between each pair of frequency contours from two different clusters within each group (one of 5 groups by coarse clustering):

$$D_i = \frac{\sum_{j_l \neq i}^{N_x} \|\bar{\mathbf{V}}_i - \bar{\mathbf{V}}_j\|_2}{(n_x - 1)},$$

where $\|\cdot\|_2$ denotes l_2 -norm, N_x is the index integer set of one of five particular groups, $\bar{\mathbf{V}}_i \in R^{100}$ is the average feature vector of the i th cluster in that group and n_x is the number of individuals in N_x . The intra-cluster variance of i th cluster VAR_i was calculated as:

$$\text{VAR}_i = \frac{\sum_m^{N_i} \|\mathbf{V}_m - \bar{\mathbf{V}}_i\|_2}{n_i},$$

where N_i is the index integer set of the i th cluster, \mathbf{V}_m is the feature vector of the m th USV in this cluster, and n_i is the number of individuals in N_i .

Behavior classification

ARBUR:Behavior uses the binocular side-view video stream (segmented by SegFormer to remove the background) as input to discriminate eight behaviors of rats. First, the rat regions within each input image were labeled by edge detection. Detected regions were recognized as outlier regions if their areas were below the maximum error range:

$$L_a^e = L_f L_a^{max},$$

where L_a^{max} is the size of the maximum area with labels in the current image, and L_f is the filter threshold. The segmented images and labels were then input into a hierarchical algorithm to classify the non-social, moving, and social state behaviors of rats as follows.

Non-social state behaviors

The movement variation (in pixels) M_p of the center position of each label was calculated to determine whether rats have moved:

$$M_p = L_{center}^N - L_{center}^1,$$

where L_{center}^N and L_{center}^1 are the center pixel positions of the labels in the last and the first image. The value of M_p for sleeping rats hardly changes. For individually moving rats, we consider the total number of labels in each image to reflect whether the rats come into contact with each other (solitary) within that sequence of images.

Moving state behaviors

Based on the behavioral definition in Table 1, we used the number of labels in the first and last images of the image sequence to initially determine the moving state behaviors. Then, we estimated the motion optical flow of the rat in the images and defined its centroid motion vector as:

$$\vec{V}_c^{label} = \sum \sum_{u,v} \vec{V}_{u,v}^{label},$$

where u and v are the pixel coordinates of the image. We determined the specific behavior by judging the vector \vec{V}_c^{label} relative direction of the two labels.

Social state behaviors

We designed a decision tree algorithm based on multiple binary SVM classifiers to detect social state behaviors. A total of 860 images were labeled (PIN: 280, POU: 300, SNC: 280), and the HOG features were extracted as the training dataset. The SVM models were trained using a

third-degree polynomial kernel function, and parameter optimization was performed for each model. Specifically, for each image, we deployed three SVM models simultaneously (PIN-POU, PIN-SNC, POU-SNC) to determine the category of the image with maximum probability. We defined the total confidence of the behavioral category for all images:

$$P_{label} = \frac{1}{2N} \sum_{n=1}^N (S_{label}(I_{Rc}^n) + S_{label}(I_{Lc}^n)),$$

where

$$S_{label}(I) = \begin{cases} 1, & P(I = label) \geq I_l \\ 0, & P(I = label) < I_l \end{cases}$$

where $P(I = label)$ is the probability that the current image category is *label* (PIN, POU, SNC), and I_l is the likelihood of a specific image belonging to the category. Behavioral category I_c was mapped through the maximum probability in each P_{label} :

$$\max(P_{PIN}, P_{POU}, P_{SNC}) \rightarrow I_c.$$

Image sequences that do not meet the aforementioned criteria were labeled as undetected.

Evaluation metrics

We use Precision, Recall and F1-score to evaluate the detection performance of ARBUR:Behavior. The number distribution of each behavior in the test set is as follows (It should be noted that the low number of AP and FO is due to the fact that our algorithm only extracted these data in all 43,357 segments). AP:24, FO:44, MA:181, PIN:185, POU:196, SNC:216, SO:230, ST:189. The test dataset is only used to evaluate the performance of the classifier and is not used for training.

Sound source localization

3D keypoints detection of rats

To obtain the spatial position of the two rats, we designed a deep learning-based keypoint detection network capable of detecting seven keypoints distributed across the rats' heads, spines, and tails. We then calculated the 3D coordinates of the keypoints via triangulation by utilizing the calibrated camera parameters and the detected key points in two image spaces. It is worth noting that for rats engaged in social behaviors, there can be occlusions of certain body parts, resulting in only a part of keypoints being detectable in each frame. Using these 3D coordinates, a curve that represents the posture of the rat in the current frame could be fitted (Figure 4A), which could be useful to determine the submissive rat and the dominant one during PIN (pinning) and POU (pouncing). For other cases, only the 3D coordinates of the rats' noses were used for the downstream SSL processes.

Estimating sound sources

The sources of USVs were estimated using an algorithm modified from Neunuebel et al.²³ For each extracted USV from four microphones, a data augmentation process was performed by slicing the recognized frequency contours into m segments. Then, four groups with 3 channels (mics) and a group with 4 mics were input to a steered response power-based sound source localization method. That is, $m \times 5$ estimates were calculated. The estimates with 4 mics ($m \times 4$) estimates and those with 3 mics ($m \times 5$) were used to calculate two 2-D probability density functions ($f_3(x, y)$, $f_4(x, y)$, $x, y \in \mathbb{R}$), respectively. Then, the localization confidence index (LCI) P_{LCI} that represents the reliability of estimation results was produced as follows:

$$P_{LCI} = \frac{f_3(x, y) | [(x, y), f_4(x, y) = \max(f_4(x, y))] |}{\max(f_3(x, y))}.$$

If the LCI is above the pre-set threshold (0.6), the SSL result is accepted as the center of estimates with maximum probability densities in $f_3(x, y)$ and $f_4(x, y)$. The SSL results are rejected otherwise.

Combining visual and SSL estimations

Because rats have a body length as long as 25 cm, the height of their sound source cannot be ignored in some cases. For example, if a rat presents a high-rearing pose, the nose can be as high as over 20 cm. Under these circumstances, the height should be considered to avoid incorrect SSL prediction.

For those behaviors with no apparent height difference between the noses of two rats (besides POU and PIN), the SSL results were assigned to rat noses based on horizontal information only. That is, the USV was assigned to the nearer rat if two noses were detected. If only one rat nose was detected, the USV was assigned to the detected one if it lies within the SSL estimation area (10 cm from the estimated point), and the USV was assigned to the undetected one otherwise. If no rat noses were detected, the USV cannot be assigned.

For pouncing and pinning behaviors, the SSL results were assigned based on not only horizontal information but also 3-D information. If two rat noses were detected, SSL results were calculated again with the nose height added. The USV was allocated to the rat with a higher LCI.

If only one rat nose was detected and an obvious height difference was observed and labeled, the USV can be allocated to the rat with a higher LCI by assuming the height of another rat. This 3D SSL process is useful if two rats are close but have obvious height differences, which is common in pinning and pouncing.

QUANTIFICATION AND STATISTICAL ANALYSIS

All statistical analyses were performed in MATLAB (version 2021b, MathWorks) using two-sample single-tailed Welch's t-tests with a significance level of $\alpha = 0.05$. Data distributions were assumed to be normal, but this was not formally tested.