

1 **Rodent ultrasonic vocal interaction resolved with millimeter precision using hybrid**
2 **beamforming**

3 M. L. Sterling¹, B. Englitz^{1*}

4 ¹Computational Neuroscience Lab, Donders Institute for Brain, Cognition and Behaviour, Radboud
5 University, Nijmegen, The Netherlands

6
7 ***Correspondence:** englitz@science.ru.nl

8
9 **Author contributions:**

10 MLS: software coding, data analysis, manuscript writing, manuscript editing, figure preparation

11 BE: experiment design, conduct experiments, data analysis, manuscript writing, manuscript
12 editing, figure preparation.

13

14 **Competing Financial Interests:** The authors declare that they do not have any competing
15 financial interests associated with the present work.

16

17 **Keywords:**

18 mouse social interaction, ultrasonic vocalization, USV, vocal communication, sound localization,
19 beamforming, automatic tracking

20 **Data & Code Availability**

21
22 During the review process, reviewers can access all Data and Code via the link below:

23 <https://data.donders.ru.nl/login/reviewer-208072048/jJ4c-oRNlIp3yArkjYQ0lAW9FjpiHL8foxSzw1FDA>

24 Upon acceptance, these materials will be made available to the public.

25

26 **Abstract**

27

28 Ultrasonic vocalizations (USVs) fulfill an important role in communication and navigation in many
29 species. Because of their social and affective significance, rodent USVs are increasingly used as
30 a behavioral measure in neurodevelopmental and neurolinguistic research. Reliably attributing
31 USVs to their emitter during close interactions has emerged as a difficult, key challenge. If
32 addressed, all subsequent analyses gain substantial confidence.

33 We present a hybrid ultrasonic tracking system, HyVL, that synergistically integrates a
34 high-resolution acoustic camera with high-quality ultrasonic microphones. HyVL is the first to
35 achieve millimeter precision (~3.4-4.8mm, 91% assigned) in localizing USVs, ~3x better than
36 other systems, approaching the physical limits (mouse snout ~ 10mm).

37 We analyze mouse courtship interactions and demonstrate that males and females
38 vocalize in starkly different relative spatial positions, and that the fraction of female vocalizations
39 has likely been overestimated previously due to imprecise localization. Further, we find that male
40 mice vocalize more intensely when interacting with two mice, an effect mostly driven by the
41 dominant male.

42 HyVL substantially improves the precision with which social communication between
43 rodents can be studied. It is also affordable, open-source, easy to set up, can be integrated with
44 existing setups, and reduces the required number of experiments and animals.

45 Introduction

46
47 Ultrasonic vocalizations (USVs) fulfill an important role in animal ecology as means of
48 communication or navigation in many rodents¹⁻⁵, bats⁶, frogs⁷, cetaceans⁸, and even some
49 primates.^{9,10} In many of these species, USVs have been shown to be present innately and to have
50 significance at multiple stages of life, from neonates¹¹ to adults¹, often with diverse functions as
51 distress/alarm calls^{11,12}, courtship signals¹³, territorial defense signals¹⁴, private communication¹⁰,
52 and echolocation⁶. USVs have been extensively studied in mice, where their communicative
53 significance has been widely demonstrated by their influence on conspecific behavior¹⁵⁻²⁰ (also in
54 line with observational studies²¹⁻²⁴). USVs can be grouped into different types that are highly
55 context-dependent^{17,18,22,25-37}, and USV syntax itself is predictive of USV sequence.³⁸ Taken
56 together, the current literature suggests USVs convey affective and social information in different
57 behavioral contexts. This is further supported by the modulatory effect that testosterone and
58 oxytocin have on USV production.³⁹⁻⁴⁵ Importantly, the neuronal circuitry underlying USVs has
59 recently been identified and is being studied extensively.^{20,25,46-53}

60 Because of their social and affective significance and our growing mechanistic
61 understanding, mouse USVs are increasingly being used as a behavioral measure in
62 neurodevelopmental and neurolinguistic translational research.^{26,32,49,54-59} Their manipulation and
63 precise measurement not only provide the basis for tackling many fundamental questions but also
64 pave the way, via advanced animal models, for the discovery of essential, novel drug targets for
65 many debilitating conditions such as autism-spectrum disorder^{58,60}, Parkinson's disease⁶¹, stroke-
66 induced aphasia⁶², epilepsy aphasia syndromes⁶³, progressive language disorders⁶⁴, chronic
67 pain⁶⁵, and depression/anxiety disorders⁶⁶, where ultrasonic vocalizations serve as a biomarker
68 for animal well-being and normal development. Consequently, we expect the scientific
69 importance of mouse USVs to continue to increase in the coming years, highlighting the necessity
70 to advance the methods required for their study. In recent years, substantial advances have been
71 made in USV detection⁶⁷⁻⁷¹, classification^{67,68,70,72}, and localization.⁷³⁻⁷⁶

72 Localization is of particular importance during social interactions, when most USVs are
73 emitted and any meaningful analysis of USV properties rests on a reliable assignment of each
74 USV to its emitter. This task is complex for multiple reasons: (i) most USVs are emitted at close
75 range, (ii) social behavior often requires free movement of the animals, and (iii) USV production
76 is invisible.^{37,77} With reliable assignment, all subsequent analyses can be conducted with
77 substantial confidence concerning each USV's emitter. Although USVs could in theory be
78 classified and assigned based on their shape^{13,78-81}, this approach will depend strongly on

79 different behavioral contexts and strains. Recent advances in acoustic localization^{74–76} have
80 improved the localization accuracy to 11-14 mm, however, close-up snout-snout interactions -
81 which is when a large fraction of USVs are emitted - require an even higher precision.

82 We have developed an advanced localization system for USVs in which a high resolution
83 'acoustic camera' consisting of 64 ultrasound microphones with an array of 4 high-quality
84 ultrasound microphones. Both systems can individually localize USVs but exhibit rather
85 complementary patterns of localization errors. We fuse them into a hybrid system that exploits
86 their respective advantages in sensitivity, detection, and localization accuracy. We achieve a
87 median absolute localization error of 3.4-4.8 mm, translating to an assignment rate of ~91%.
88 Compared to the previous state of the art^{73,75}, the accuracy represents a three-fold improvement
89 that halves the proportion of previously unassigned USVs. Given the physical dimensions of the
90 mouse snout ($\varnothing \sim 10$ mm), this likely approaches the physical limit of localizability for USVs. We
91 successfully apply it to and analyze dyadic and triadic courtship interactions between male and
92 female mice. We demonstrate that the fraction of female vocalizations has likely been
93 overestimated in previous analyses, due to a lack of precision in sound localization.

94

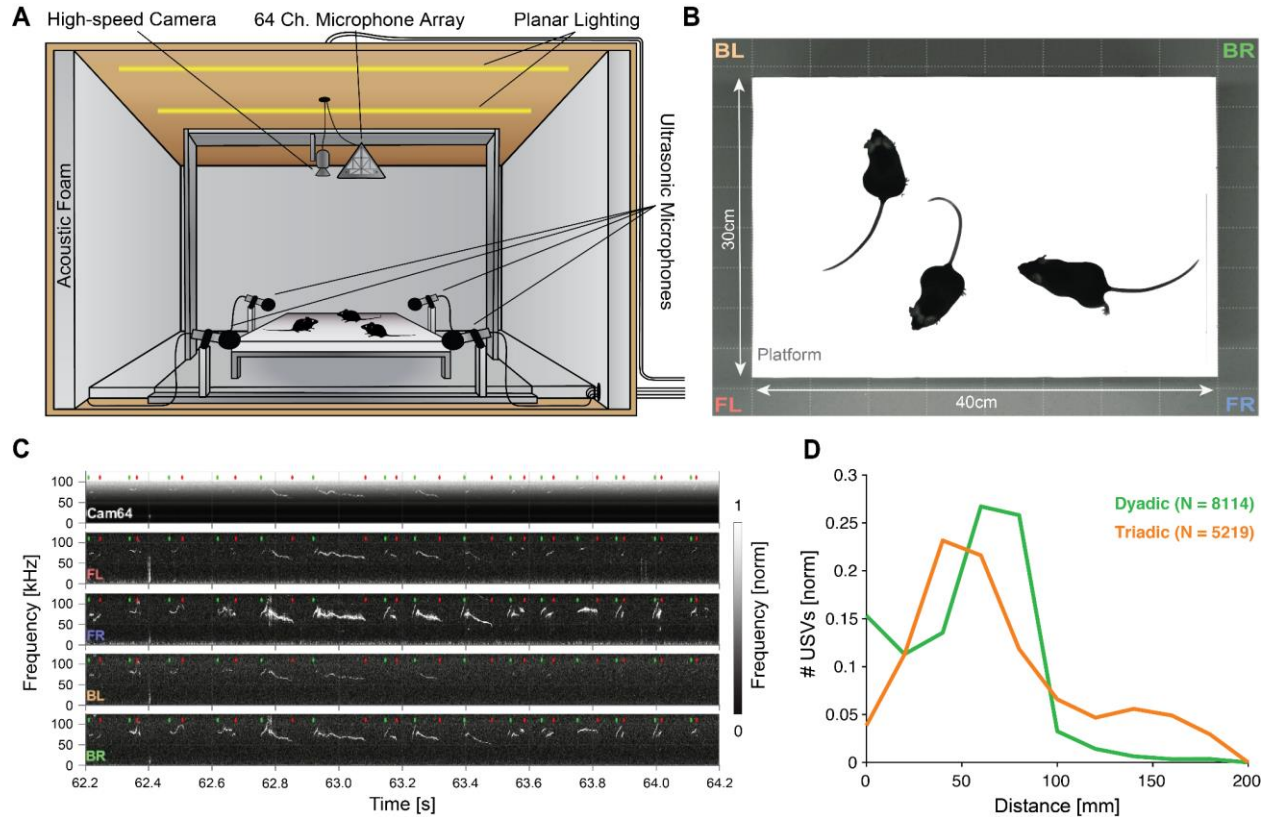
95

Results

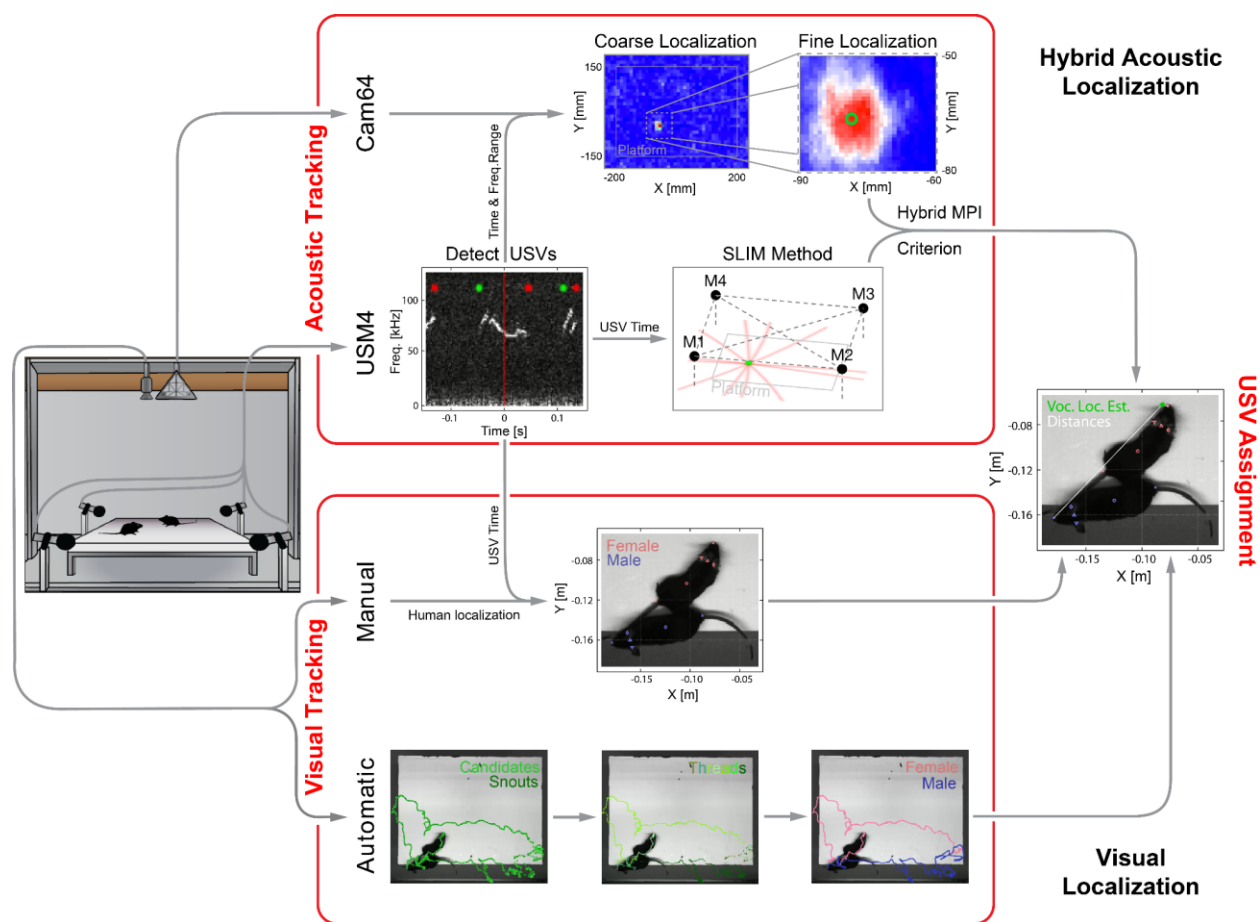
96

97 We analyzed courtship interactions of mice in dyadic and triadic pairings. The mice interacted on
98 an elevated platform inside an anechoic booth (see *Fig. 1A*, for details see *Materials & Methods:*
99 *Recording Setup*). Each trial consisted of 8 minutes of free interaction while movements were
100 tracked with a high-speed camera (see *Fig. 1B*), and ultrasonic vocalizations (USVs) were
101 recorded with a hybrid acoustic system composed of 4 high-quality microphones (i.e., USM4) as
102 well as a 64-channel microphone array (*Cam64*, often referred to as an acoustic camera; see *Fig.*
103 *1C* for raw data samples, green and red dots mark the start and stop times of USVs).

104 Most USVs were emitted in close proximity in dyadic and triadic pairings (see *Fig. 1D*).
105 Reliably assigning most USVs to their emitter therefore requires a highly precise acousto-optical
106 localization system. The presently developed Hybrid Vocalization Localizer (HyVL) system is the
107 first to achieve sub-centimeter precision, i.e. ~3.4-4.8 mm (see *Fig. 2* for an overview). This
108 accuracy on the acoustic side is achieved by combining the complementary strengths of the



109 **Figure 1:** Mice emit ultrasonic vocalizations (USVs) in close proximity during courtship behavior.
110 **A** Two or three mice of different sexes were allowed to interact freely on an elevated platform. Vocalizations
111 were recorded with 4 high-quality ultrasonic microphones in a rectangular arrangement around the platform
112 and a 64-channel microphone array ('Cam64') mounted above the platform. The spatial location of the pair
113 was recorded visually with a high-speed camera. The platform was located in an ultrasonically sound-proof
114 and anechoic box and illuminated uniformly using an array of LEDs.
115 **B** Sample image from the camera that shows the high contrast between the mice and the interaction
116 platform. The two-letter abbreviations indicate the locations of the 4 high-quality microphones (F = front,
117 B = back, L = left, R = right).
118 **C** Sample spectrograms from the four ultrasonic microphones and the average of all Cam64 microphones
119 for a bout of vocalizations (start/end times marked by green/red dots). The Cam64 microphones are of
120 lower quality than the USM4 microphones, evidenced by the rising noise floor for higher frequencies,
121 affecting very high frequency USVs.
122 **D** Most USVs in the present paradigm were emitted in close proximity to the interaction partners, with the
123 vast majority within 10 cm snout-snout distance (i.e., ~93 and 72% for dyadic and triadic, respectively).



124

125 **Figure 2:** Overview of the combined acoustic and visual tracking pipeline.

126 (Top) Acoustic tracking of animal vocalizations was enabled by a hybrid acoustic system, which recorded
 127 the sounds in the booth using a 64-channel ultrasonic microphone array ('Cam64') and 4 high-quality
 128 ultrasonic microphones ('USM4'). Vocalizations were automatically detected using USM4 data (start/end
 129 times marked by green/red dots) and then localized on the platform using both the SLIM algorithm on USM4
 130 data and delay-and-sum beamforming on the corresponding Cam64 data. The Cam64 localization
 131 proceeded in two steps: first coarse (10 mm resolution), then fine centered around the coarse peak at 1 mm
 132 resolution (30 x 30 mm local window). The local, weighted average (green circle) was then used as the
 133 USV origin localized by Cam64. For each USV, the Cam64 localization was chosen if its SNR > 5, otherwise
 134 the USM4/SLIM estimate was used (for details, see *Materials & Methods: Localization of Ultrasonic*
 135 *Vocalizations*).

136 (Bottom) Animals were tracked visually on the basis of concurrently acquired videos (50 FPS). Two tracking
 137 strategies were employed: (i) manual tracking in the video frames corresponding to the mid-point of USVs
 138 in all recordings and (ii) automatic tracking for all frames in dyadic recordings.

139 (i) *Manual visual tracking*: the observer was presented with a combined display of the vocalization
 140 spectrogram and the concurrent video image at the temporal midpoint of each USV and annotated the
 141 snout and head center (i.e., midpoint between the ears).

142 (ii) *Automatic visual tracking*: Started with finding the optimal locations of each marker based on marker
143 estimate clouds produced by *DeepLabCut* (DLC; see Mathis et al., 2018) for all frames. Next, these marker
144 positions were assembled into spatiotemporal threads with the same, unknown identity based on a
145 combination of spatial and temporal analysis. Finally, the thread ends still loose were connected based on
146 quadratic spatial trajectory estimates for each marker, yielding the complete track for both mice (see
147 *Materials & Methods: Automatic Visual Animal Tracking and Suppl. Fig. 1*).

148
149 USM4 and Cam64 data. The Cam64 data is processed using acoustic beamforming⁸² which
150 delivers highly precise estimates (MAE = ~4-5 mm), but is not sensitive enough for very high-
151 frequency USVs (see Fig. S2). The USM4 data is analyzed using our SLIM algorithm,⁷³ which
152 delivers accurate (MAE = ~11-14 mm) and less frequency-limited estimates. The methods exhibit
153 a complementary pattern of localization errors, which predestines them for high synergy when
154 combined (see below).

155 For each USV, a choice is made between the USM4/SLIM and Cam64/Beamforming
156 estimates based on a comparison of each method's USV-specific certainty and the relative
157 position of the mice to the estimates, using an extended, hybrid Mouse Probability Index (MPI⁷⁶).
158 HyVL is the first system of its kind that exploits a hybrid microphone array to overcome the
159 limitations of each subarray. The positions of the mice are obtained via manual and automatic
160 video tracking using *DeepLabCut*,⁸³ each of which achieve millimeter precision for localizing the
161 snout .

162 Overall, 83 recordings were collected from 14 male and 4 female mice. Of these
163 recordings, 55 were dyadic and 28 triadic. In all trials combined, 13406 USVs were detected.

164 165 *Precision of USV Localization*

166 Assigning USVs to individual mice required combining high-speed video imaging with the HyVL
167 location estimates at the times of vocalization. We manually tracked the animal snouts at the
168 temporal midpoint of each USV to obtain near-optimal ground truth position estimates (see Fig. 2).
169 We first assessed the relative structure of the localization errors between both methods,
170 USM4/SLIM (Fig. 3A, green) and Cam64/Beamforming (red, each dot is a USV). While most
171 errors were small, and clustered close to the origin of the graph (evidenced by the small Median
172 Absolute Errors (MAE), shown as horizontal and vertical lines, respectively), the less frequent,
173 larger errors exhibited an L-shape. This error pattern is an optimal situation for combining
174 estimates from the two methods, to compensate for each other's limitations. While the Cam64
175 data can compensate for single microphone noise through the large number of microphones, the
176 nature of its micro-electromechanical systems (MEMS) microphones deteriorates for very high

177 frequencies (see *Fig. S2B*). Conversely, the USM4 microphones show an excellent noise level
178 across frequencies (see *Fig. S2A*) but can produce erroneous estimates if there is noise in a
179 single microphone and have an intrinsic limitation in spatial accuracy due to the physical size of
180 their receptive membrane ($\varnothing \sim 20$ mm).

181 We therefore designed an analytical strategy to combine the estimates of both systems to
182 optimize the number of reliably assignable USVs, while evaluating the resulting spatial accuracy
183 alongside. Briefly, the location estimates of both methods each come with an estimate of
184 localization uncertainty. First, we assess for each method's estimate how reliably it can be
185 assigned to one of the mice, taking into account the positions of the other mice. This is quantified
186 using the Mouse Probability Index (MPI),⁷⁶ which compares the probability of assignment to a
187 particular mouse to the sum of probabilities for all mice, weighted by the estimate's uncertainty. If
188 the largest MPI exceeds 0.95, it is considered a reliable assignment to the corresponding mouse.
189 If both methods allowed reliable assignments, the one with smaller residual distance was chosen.
190 If only one method was reliable for a particular USV, its estimate was used. If neither method
191 allowed for reliable assignment, the USV was not used for further analysis. This typically happens
192 if the snouts are extremely close or the USV is very quiet. This approach outperformed many
193 other combination approaches in accuracy and assignment percentage, e.g. Maximum Likelihood
194 (see *Materials & Methods: Assigning USVs and Discussion* for details).

195 HyVL performed significantly better than either method alone, allowing a total of 91.1% of
196 USVs to be assigned at a spatial accuracy of 4.8 mm (MAE). This constitutes a substantial, 2.9-
197 fold improvement in accuracy over the previous state of the art, the SLIM algorithm.⁷³ On the full
198 set of USVs where both microphone arrays were recording (N = 7982), HyVL outperformed both
199 USM4/SLIM and Cam64/Beamforming significantly, both in residual error (SLIM: 14.8 mm;
200 Cam64: 5.33 mm; HyVL: 5.07 mm; $p < 10^{-10}$ for all comparisons, Wilcoxon rank sum test) and
201 percentage of reliably assigned USVs (SLIM: 74.4%; Cam64: 80%; HyVL: 91.1%).
202 Cam64/Beamforming performed even more precisely on its reliably assignable subset (4.56 mm),
203 which was, however, smaller than the HyVL set. This difference emphasizes the complementarity
204 of the two methods and thus the synergy through their combination. There was no significant
205 difference between tracking on dyadic and triadic recordings (HyVL: 5.0 mm vs. 5.1 mm, $p = 0.71$,
206 Wilcoxon rank sum test) with correspondingly similar selection percentages (92 vs. 90%,
207 respectively).

208 The accuracies above are an average over localization performance at any distance. In
209 particular during close interaction, USVs will often be reflected or obstructed, complicating
210 localization. While this constitutes the realistic challenge during mouse social interactions, we

211 also investigated the 'ideal', unobstructed performance of HyVL by comparing the performance
212 on USVs emitted when all animals were 'far' (>100 mm) apart, i.e. more than ~20 times the
213 average accuracy of HyVL, as well as for a single male mouse on the platform. For the *far* USVs
214 the reliably assignable fraction increased to 97.9%, and the accuracy significantly improved to
215 3.8 mm (*Fig. 3C* gray, $p=8.6 \times 10^{-7}$, Wilcoxon rank sum test). For the *single animal* USVs, the
216 accuracy was even better at 3.4 mm with 98.4% reliably assigned (*Fig. 3C*, blue).

217 Next, we inspected separate localization along the X and Y axis to check for anisotropies
218 of localization (*Fig. 3D/E*, histograms normalized to maximum). The position of the closest animal
219 aligned precisely with the estimated position in both dimensions, indicated by the high density
220 along the diagonal (*Pearson* $r > 0.99$ for both dimensions) and the MAE's along the X and Y
221 direction separately ($X = 3.1$ mm, $Y = 2.8$ mm). These one-dimensional accuracies might be of
222 relevance for interactions where movement is restricted.

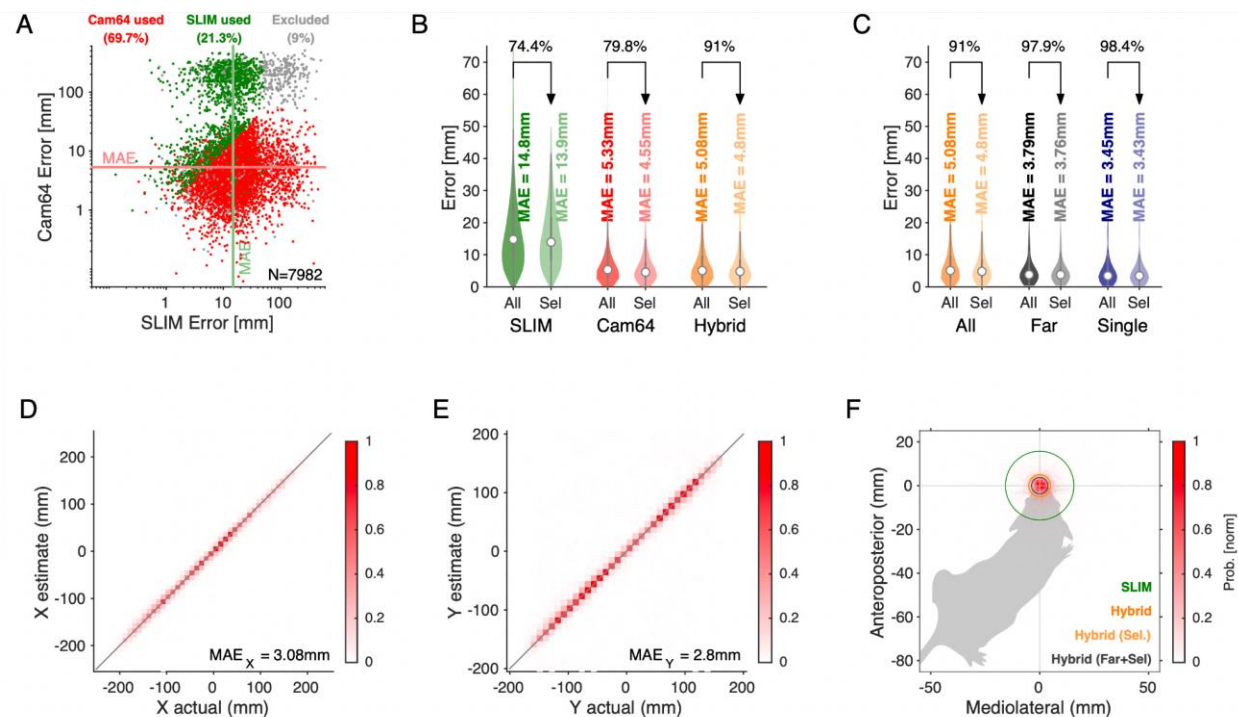
223 Lastly, we visualized the localization density relative to the mouse that the vocalization
224 was assigned to (*Fig. 3F*). Combining both dimensions and appropriately rotating them, the
225 estimated position of the USVs is shown relative to the mouth. The density is narrowly centered
226 on the snout of the mouse (circle radius = MAE: green: SLIM method; orange: HyVL; light orange:
227 HyVL assigned USVs, gray: Far assigned USVs).

228 In summary, the HyVL system provides a substantial improvement in the localization
229 precision. In comparison to other methods, its precision also allows a larger fraction of
230 vocalizations to be reliably assigned and retained for later analysis, which enables a near
231 complete analysis of vocal communication between mice or other vocal animals.

232

233 *Sex Distribution of Vocalizations During Social Interaction*

234 Courtship interactions between mice lead to high rates of vocal production, but are challenging
235 due to the relative proximity, including facial contact. Previous studies using a single microphone
236 have often assumed that only the male mouse vocalized,^{84–87} while more recent research has



237
 238 **Figure 3:** Spatial accuracy of localizing USVs during mouse social interaction improves ~3-fold over the
 239 state of the art.⁷³

240 **A** The vast majority of USVs is localized with very small errors for both methods, concentrated close to the
 241 axes and thus hardly visible, evidenced by the median average errors (MAE) for Cam64 (light red line) and
 242 SLIM (light green line). The fewer larger errors form an L-shape, emphasizing the synergy of a hybrid
 243 approach that compensates for the weaknesses of each method. Location estimates were excluded (gray)
 244 if they were >50mm from either mouse, or the hybrid MPI<0.95.

245 **B** The hybrid localization system HyVL (orange) combines the virtues of SLIM and Cam64 enabling the
 246 localization of 91.1% of all USVs (light orange), achieving an MAE = 4.8 mm. Cam64 localization (red)
 247 alone only includes 74.4% of all USVs, but at an MAE = 4.5 mm (light red). SLIM-based localization (green)
 248 only includes 80.0% of all USVs, at an MAE = 14.6 mm (light green, see *Materials & Methods: USV*
 249 *Assignment* for details on the relation between accuracy and selection criteria).

250 **C** USVs emitted when all animals were >100 mm apart and a single mouse condition was used to assess
 251 the ideal accuracy of HyVL. For the far condition, virtually all USVs (332/339, 97.9%) were assigned at an
 252 MAE = 3.8 mm, similarly to the single animal condition (MAE = 3.4 mm, 251/255, 98.4%).

253 **D/E** Comparison of actual with estimated snout locations along the X (horizontal; **D**) and Y (vertical; **E**)
 254 dimensions indicating strong agreement. Colors indicate peak-normalized occurrence rates.

255 **F** Centered overlay of USV localizations relative to emitter snout. Precision is depicted as a circle with a
 256 radius equivalent to the median absolute error (green: SLIM; orange: HyVL, all USVs; light orange: HyVL,
 257 selected USVs, dark gray: HyVL, when mice >100 mm apart).

258 concluded that female mice vocalize as well.^{76,88} Female vocalizations were typically less
259 frequent, but constituted a substantial fraction of the vocalizations (11-18%). Below, we
260 demonstrate that the accuracy of the localization system can be an important factor for
261 conclusions about the contribution of different sexes to the vocal interaction.

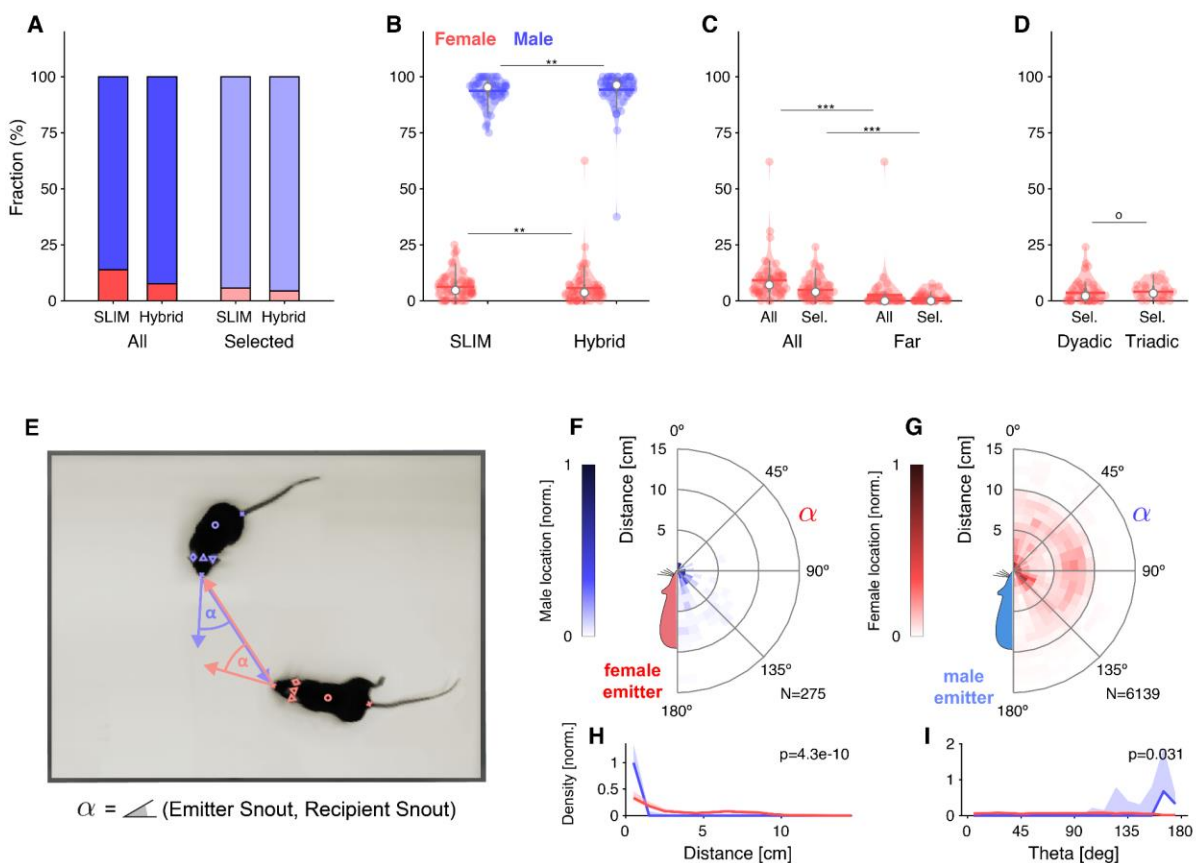
262 Overall, female vocalizations constitute the minority of vocalizations overall. Naive
263 estimation without MPI selection using SLIM estimates ~14%, while HyVL tallies it at just 7%
264 (Fig. 4A). Applying MPI selection, SLIM estimates only 5.5%, while HyVL arrives at significantly
265 less, just 4.4% ($p=0.002$, paired Wilcoxon signed rank test, Fig. 4A/B), while reliably classifying
266 91.1% of all vocalizations.

267 Using HyVL instead of SLIM significantly reduces the fraction of female vocalizations,
268 suggesting that less accurate algorithms overestimate the female fraction (only results for MPI-
269 selected USVs shown, Fig. 4B). Considering only vocalizations that are emitted when the snouts
270 are >50 mm apart, further significantly reduces the fraction to female USVs to 1.1% after MPI
271 selection ($p=5.2 \times 10^{-8}$, Wilcoxon Rank Sum test). Comparing the percentage of female
272 vocalizations between dyadic and triadic trials, no significant differences were found ($p=0.22$,
273 Wilcoxon Rank Sum test, Fig. 4D).

274 Beyond the absolute distance between the mouths of the mice, high-accuracy localization
275 of USVs allows one to position the bodies of the animals relative to one another at the times of
276 vocalization by combining acoustic data with multiple concurrently tracked visual markers. This
277 provides an occurrence density of other mice relative to the emitter (Fig. 4E).

278 Female mice appear to emit vocalizations in very close snout-snout contact, with a small
279 fraction of vocalizations occurring when the male snout is around the hind-paws/anogenital region
280 (Fig. 4F). Male mice emit vocalizations both in snout-snout-contact, but also at greater distances,
281 which dominantly correspond to a close approach of the male's snout to the female anogenital
282 region (Fig. 4G). This was verified separately with a corresponding analysis, where the recipient's
283 tail-onset was used instead (not shown).

284 In summary, the combination of high-precision localization and selection using the MPI
285 indicate that female vocalizations may be even less frequent than previously thought. When they
286 vocalize, the mice appear to almost exclusively be in close snout-snout contact. As this is
287 incidentally also the condition which has the highest chance of mis-assignments, even the
288 remaining female vocalizations need to be treated with caution.



289

290 **Figure 4:** Analysis of Sex-dependent Vocalizations can depend on Localization Accuracy

291 **A** Female vocalizations constitute a small fraction of the total set of vocalizations. The female fraction further

292 reduces with increased precision and when selecting vocalizations based on the MPI.

293 **B** Using the hybrid method instead of SLIM significantly reduces the fraction of female vocalizations,

294 suggesting that less accurate algorithms overestimate the female fraction (only results for MPI-selected

295 USVs shown).

296 **C** The fraction of female vocalizations further reduces if only USVs are considered that are emitted while

297 all animal snouts were >50 mm apart from each other. This indicates a preference of female mice to

298 vocalize in close snout-snout contact, however, this entails that female vocalizations are more prone to

299 confusion with male vocalizations due to their relative spatial occurrence.

300 **D** There was no difference in the female fraction of USVs between dyadic and triadic pairings (2 male and

301 2 female conditions combined here).

302 **E** High-accuracy localization of USVs allows one to analyze the relative spatial vocalization preferences of

303 the mice, i.e. their occurrence density in relation to the relative position of other mice to the emitter. We

304 quantified this by collecting the position of the non-vocalizing mice at the times of vocalization, in relation

305 to the vocalizing mouse. α corresponds to the angle between the emitter's snout and the snout of other

306 mice.

307 **F** Female mice appear to emit vocalizations in very close snout-snout contact, with a small fraction of
308 vocalizations also occurring when the male mouse around the hind-paws/anogenital region.

309 **G** Male mice emit vocalizations both in snout-snout-contact, but also at greater distances, which dominantly
310 correspond to a close approach of the male's snout to the female anogenital region. This was verified
311 separately with a corresponding analysis, where the recipient's tail-onset was used instead (not shown).

312 **H** Radial distance density of receiver animals, marginalized over directions, shows a significant difference,
313 with females vocalizing mostly when males (blue) are in close proximity of the snout, while males vocalize
314 when the female mouse's snout is very close (corresponding to snout-snout contact), but also when the
315 female's snout is about 1 body length away (snout-anogenital interaction). Plot show medians and
316 percentile-based confidence bounds.

317 **I** Direction density of receiver animals, marginalized over distances, shows borderline significance with
318 female mice vocalizing mostly when the male snout is located behind their own snout.

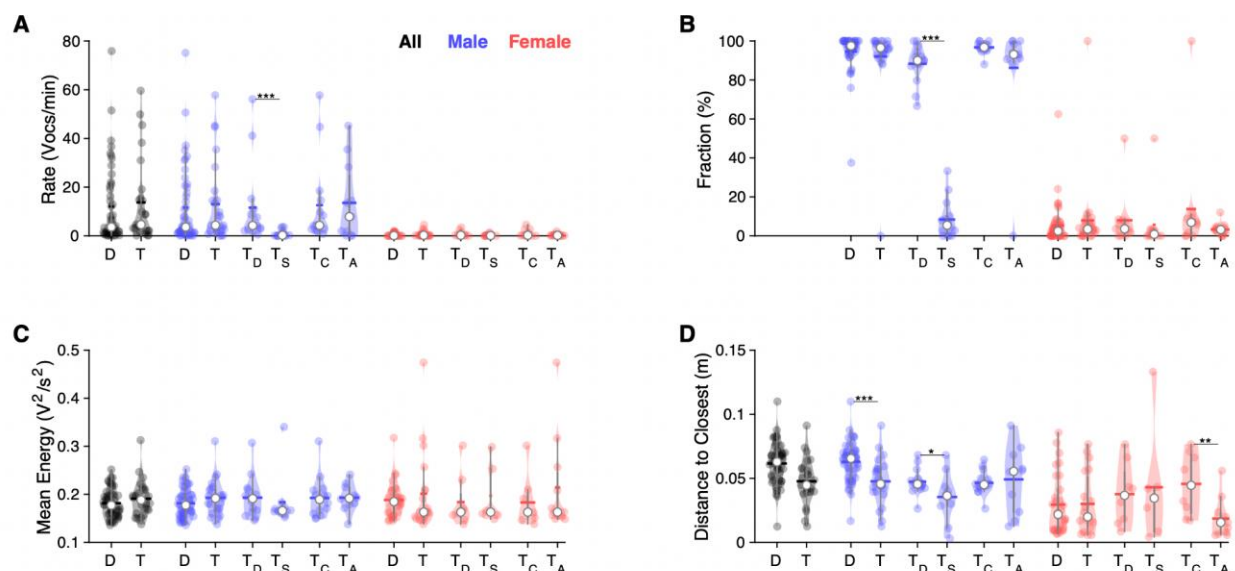
319

320 *Vocalization Rate Analysis*

321 In dyadic trials, one female and one male mouse interacted, whereas in triadic trials either two
322 males and one female or two females and one male mouse interacted. In the analysis of triadic
323 interactions, we separate competitive and alternative contexts depending on whether a mouse
324 had to compete with another same sex mouse or could interact with two opposite sex mice,
325 respectively. For triadic trials we further separate the same-sex mice into dominant and
326 subordinate, based on who vocalized more.

327 Overall vocalization rates did not differ significantly between dyadic and triadic conditions
328 for male and female mice (see *Fig. 5A*). However, in competitive interactions between males, one
329 male mouse significantly and strongly dominated the 'conversation', with on average 9-fold more
330 vocalizations than the other male mouse (T_D vs T_s , both comparisons: $p < 0.0001$, Wilcoxon Sum
331 of Ranks test; retains $p < 0.05$ after Bonferroni correction for multiple testing, *Fig. 5A,B*). While the
332 present division into dominant and subordinate mouse based on a higher vocalization rate within
333 a recording will always lead to a significant difference, the quantitative difference between them
334 is the striking aspect in this comparison. Overall male vocalization rates were similar in
335 competitive and alternative triadic trials. Female vocalization rates were similar across all
336 compared conditions.

337 The mean vocalization energy of dominant males in triadic pairings tends to be higher
338 than those of submissive males in triadic pairings, however, this result did not reach significance
339 in the present dataset (see *Fig. 5C*). No effects of vocalization energy were found in females.



340
 341 **Figure 5:** In triadic interaction, one male vocalizes dominantly and males vocalize even closer to females
 342 **A** Overall, vocalization rates were comparable between dyadic (D) and triadic (T) conditions. Male mice
 343 (blue) vocalized at higher rates than female mice (red). However, this was restricted to the dominant male
 344 mouse (T_D : dominant = emitted more USVs within same-sex) in triadic, competitive (2m/1f) conditions (see
 345 text for all p-values.). Male vocalization rates were similar in competitive (T_C : with same-sex competitors;)
 346 and alternative (T_A : no same-sex competitor, i.e. for male vocs: 2f/1m) pairings. Female vocalization rates
 347 remained low and similar across all conditions. T_S : submissive mouse = emitted fewer USVs within same
 348 sex during competitive trial; white dot: median; horizontal bar: mean.

349 **B** While the fraction of USVs emitted by males was overall comparable between D and T pairings, the
 350 dominant male (T_D) emitted a substantially larger fraction than their submissive counterpart (T_S), roughly a
 351 factor of 9. In competitive pairings, male mice tended to emit an overall larger fraction of all USVs than in
 352 alternative pairings (T_C vs. T_A), but this is unsurprising as both males vocalize. In female mice, the overall
 353 fraction of USVs in D and T pairings was also similar (see details in Results for potential caveats of the
 354 dominant/subordinate classification).

355 **C** In triadic pairings, dominant male mice tended to vocalize more intensely than in dyadic pairings,
 356 however, this difference was not significant at the current sample size. No significant differences were found
 357 for female mice.

358 **D** Male mice emitted USVs in closer proximity to the closest female mouse in triadic compared to
 359 dyadic interactions. Female mice generally emitted USVs at closer distances (see also Fig. 4F/H), in
 360 particular for alternative vs. competitive pairings.

361
 362 The distance to the closest animal of opposite sex was found to be even closer during
 363 triadic trials (see Fig. 5D), driven purely by male vocalizers ($p=0.0003$, Wilcoxon Sum of Ranks
 364 test): the distance to the closest animal does not change between conditions for vocalizing

365 females ($p=0.975$, Wilcoxon Sum of Ranks test). Interestingly, the distance to the closest animal
366 was larger for females at the time of vocalization when they had a same-sex competitor on the
367 interaction platform with them than when they were the only female (T_c vs. T_a , $p=0.0068$, Wilcoxon
368 Sum of Ranks test).

369 In summary, in competitive triadic interactions, one of the male mice took a strongly
370 dominant role and approached the female even more closely when vocalizing than in dyadic
371 pairings. The latter suggests that more vocalizations occurred during snout-snout interactions
372 than in dyadic interactions, where the bulk was in snout-to-anogenital interactions.
373 Correspondingly, female mice also vocalized more snout-snout in alternative interactions from
374 their perspective, pointing to overall more vocalizations during snout-snout interactions when two
375 male mice competed for a female mouse.

376

Discussion

377

378 We have developed and evaluated a novel, hybrid sound localization system (HyVL) for ultrasonic
379 vocalizations (USVs) emitted by mice and other rodents. USVs are innately used by rodents to
380 communicate social and affective information and are increasingly being used in neuroscience as
381 a behavioral measure in neurodevelopmental and neurolinguistic research. In the context of
382 dyadic and triadic social interactions between mice, we demonstrate that HyVL achieves a
383 groundbreaking increase in localization accuracy down to ~3.4-4.8 mm, enabling the reliable
384 assignment of >90% of all USVs to their emitter. Further, we demonstrate that this can be
385 combined with automatic tracking, enabling a near-complete and automated analysis of vocal
386 interaction between rodents. The showcased analyses demonstrate the advantages obtained
387 through more precise localization, further discussed below. HyVL is based on an array of high-
388 quality microphones in combination with a commercially available, affordable acoustic camera.
389 With our freely available code, this system can be readily reproduced by other researchers and
390 has the potential to revolutionize the study of natural interactions of mice.

391

Comparison with previous approaches for localizing vocalizations

392 Localization accuracy was first systematically reported by Neunuebel et al.⁷⁶ using a 4-
393 microphone setup and a maximum likelihood approach,⁸⁹ who attained an MAE of ~38 mm that
394 conferred an assignment rate of 14.6-18.1% (their Table 1, *assigned* relative to *detected* or
395 *localized*). Originating from the same research group, Warren et al.⁷⁵ employed both a 4 and 8
396 microphone setup in a follow-up study, achieving an MAE of ~30 mm for 4 microphones (~52%
397 assignment rate) and of ~20 mm with 8 microphones (~62% assignment rate), both using a
398 jackknife approach to increase robustness of localization. Stahl et al.⁷³ introduced the SLIM
399 algorithm, reaching an MAE of ~11-14 mm (~80-85% assignment rate depending on the dataset)
400 using 4 microphones. Presently, we advance the state-of-art in multiple ways: we use 68
401 microphones, combining a 64 channel 'acoustic camera' with 4 high quality ultrasonic
402 microphones. While the acoustic camera has relatively basic micro-electromechanical systems
403 (MEMS) microphones, it is inexpensive, features a high degree of integration and correspondingly
404 easy operation. Combining the complementary strengths of the two arrays is the key advantage
405 of the present approach over previous approaches, as it allows for a quantum leap in accuracy
406 (3.4-4.8 mm, 91% assignment rate), while keeping the complexity of the system manageable. A
407 comparable alternative might be a 16-channel array from high-quality microphones, which would,
408 however, be substantially more expensive (~€40,000) as well as cumbersome to build and refine.
409

410 A future generation of MEMS microphones might make the use of the high-quality microphones
411 unnecessary and thus further simplify the system setup, allowing for inexpensive, small-form
412 factor deployment (see below).

413

414 *Expected impact for future research*

415 Mice and rats are social animals,^{90,91} and isolated housing⁹² or testing⁹³ can affect subsequent
416 research outcomes. Social isolation also has direct effects on the number and characteristics of
417 USVs, at least in males.^{94,95} Sangiamo et al.⁸⁸ demonstrated that distinct USV patterns can be
418 linked to specific social actions and the latter that locomotion and USVs influence each other in a
419 context-dependent way. Using HyVL, such analyses could be extended to more close-range
420 behaviors, when a substantial fraction of the vocalizations are emitted (see *Fig. 1D*). The
421 development of more unrestricted behavioral paradigms, made viable by increased localization
422 precision, will thus also likely prove valuable to the fields of human language impairment and
423 animal behavior. As an added benefit, better USV localization will also likely increase lab animal
424 well-being via (i) more social contact in specific cases where they spend much time with their
425 conspecifics in the testing environment, or when the home environment is the testing environment
426 (e.g., PhenoTyper; Noldus Information Technologies), and (ii) a reduced need for (non-)invasive
427 markers.

428 Here, we conducted a limited set of showcase analyses on the spatial characteristics of
429 vocalization behavior. As expected, the system was accurate enough to assign vocalizations
430 during many snout-snout interactions as well as other, slightly more distant interactions, e.g. snout
431 contact with the ano-genital region of the dyadic partner. We found the male mice to vocalize
432 most while making snout contact with the abdomen and ano-genital region of the female wild-
433 type. Females vocalized predominantly during snout-snout contact, with the male approaching
434 from behind.

435 This highlights an example of how localization accuracy can shape our understanding of
436 roles in social interaction between mice: A recent, pivotal study⁷⁶ demonstrated that female mice
437 vocalize during courtship interactions. Research from our group⁷³ concluded further that mice
438 primarily vocalize in snout-snout interactions, incidentally the condition that makes assignment
439 the most difficult. While the present results maintain that female mice vocalize, the fraction
440 appears to be lower than previously thought. We, however, emphasize that this conclusion still
441 requires further study under different social contexts, e.g. interaction of more mice as in some of
442 the previous studies.^{34,88}

443 The compact form factor of the HyVL microphone arrays, in particular the Cam64, enables
444 studies of social interaction in home cages. There, rodents are less stressed and likely to exhibit
445 more natural behavior, in particular if the home cage includes enrichments. The relatively low
446 hardware costs for HyVL allows deployment of multiple systems to cover larger and more natural
447 environments.

448

449 *Current limitations and future improvements of the presented system*

450 The millimeter accuracy by HyVL enables the assignment of USVs even during close interaction,
451 certainly including all snout-anogenital interactions, and many snout-snout interactions. However,
452 certain snout-snout interactions are still too close to reliably assign co-occurring USVs. While the
453 MPI criterion maintains reliability even then, subsequent analysis will be partially biased due to
454 the exclusion of these USVs during the closest interactions. While a further improvement of
455 accuracy may be possible, close inspection of the sound density maps available via beamforming
456 from the Cam64 recordings suggests that the mouse's snout acts as a distributed source: the
457 sound density is rather evenly distributed on it, without a clear internal peak. During free
458 interaction, we noticed that the sound density was co-elongated with the head-direction of the
459 mouse and could thus be used as an additional feature to identify the vocalizer. However, this
460 proved unreliable during close interaction, likely due to absorption and reflection of sounds based
461 on the mice's bodies. More advanced modeling of the local acoustics or deep learning might be
462 able to resolve these issues, but would require ground truth recordings, which could be obtained
463 in interactions where one mouse is known to be silent, e.g. by cutting the vocal cords.

464 The present strategy for combining the estimates from Cam64/Beamforming and
465 USM4/SLIM was chosen as it optimized the reliably assigned percentage of USVs, while
466 minimizing the residual distance. We also tested alternative approaches, e.g. using direct
467 beamforming on the combined data from Cam64 and USM4 (unreliable estimates, due to
468 mismatch of number of microphones, not further pursued), maximum likelihood combination of
469 estimates (MAE=7.1 mm),⁹⁶ and making the selection solely depend on the MPI (MAE=5.2 mm).
470 While each of these approaches have certain, theoretically attractive features, the results were
471 worse in each case, likely due to particular idiosyncrasies of the MPI computation, the different
472 microphone characteristics, and the estimation of single-estimate uncertainty.

473 Lastly, for the purpose of online feedback during experiments and to reduce data
474 warehousing, it would be advantageous to perform the localization of USVs in real-time. This
475 would be enabled by streaming the data to a GPU, performing localization immediately and
476 keeping only a single channel, beamformed estimate of each USV. Ideally, the same device could

477 run visual tracking simultaneously, which would remove all temporal limitations on the recordings
478 in terms of data size and enable continuous audiovisual tracking.

479

480 *Conclusion*

481 HyVL delivers breakthrough accuracy and assignment rates, likely approaching the physical limits
482 of assignment. The low system costs (<€10k) in relation to its performance make HyVL an
483 excellent choice for labs studying rodent social interaction. Many recent questions regarding the
484 sequencing of vocalizations during social interactions become addressable with HyVL without
485 intrusive interventions. Its use can both refine the precision and reliability of the analysis, while
486 reducing the number of animals required to complete the research due to a larger fraction of
487 assigned USVs per animal.

488

Materials & Methods

489

490 All experimental procedures were approved by the animal welfare body of the Radboud University
491 under the protocol DEC-2017-0041-002 and conducted according to the Guidelines of the
492 National Institutes of Health.

493

Animals

494 In our experiment, 4 female C57Bl/6J-WT, 6 male C57Bl/6J-WT and 8 male C57Bl/6J-
495 *Foxp2.flox/flox;L7-cre* mice (bred locally at the animal facility) were studied. For subsequent
496 analyses, WT and KO mice were combined (see beginning of Results for reasoning). The mice
497 were 8 weeks old at the start of the experiments. After 1 week of acclimation in the animal facility,
498 the experiments were started. Mice of the same sex were housed socially (2 to 5 mice per cage)
499 on a 12-hour light/dark cycle with ad libitum access to food and water in individually ventilated,
500 conventional EU type II mouse cages at 20°C with paper strip bedding and a plastic shelter for
501 basic enrichment. Upon completion of the experiments, the animals were anesthetized using
502 isoflurane and sacrificed using CO₂.

503
504 The current experiment was performed as an add-on to an existing set of experiments,
505 whose focus included a region-specific knockout of *Foxp2* in the cerebellar Purkinje cells of the
506 male mice, denoted as C57Bl/6J-*Foxp2.flox/flox;L7-cre*. Neither previous work nor our own work
507 has detected any differences in USV production between WT and KO animals⁹⁷, so - given the
508 mostly methodological focus of the present work - we considered it acceptable to pool them in the
509 current analysis, reducing the number of animals needed, thus treating all males as WT C57Bl/6J,
510 the genotype of the female mice.

511

Recording Setup

512 The behavioral setup consisted of an elevated interaction platform in the middle of an anechoic
513 booth together with 4 circumjacent ultrasonic microphones as well as an overhanging 64-channel
514 microphone array and high-speed video camera (see *Fig. 1A*).

515
516 The booth had internal dimensions of 70 x 130 x 120 cm (L x W x H). The walls and floor
517 were covered with acoustic foam on the inside (thickness: 5 cm, black surface Basotect Plan50,
518 BASF). The acoustic foam shields against external noises above ~1 kHz with a sound absorption
519 coefficient > 0.95 (N.B., defined as the ratio between absorbed and incident sound intensity),
520 which corresponds to >26 dB of shielding apart from the shielding provided by the booth itself. In
521 addition, the foam strongly attenuates internal reflections of high-frequency sounds like USVs.

522 Illumination was provided via 3 dimmable LED strips mounted to the ceiling, providing light from
523 multiple angles to minimize shadows.

524 The support structure for the interaction platform and all recording devices was a common
525 frame constructed from slotted aluminum (30x30 mm) mounted to the floor of the anechoic booth,
526 guaranteeing precise relative positioning throughout the entire experiment. The interaction
527 platform itself was a 40x30 cm rectangle of laminated, white acoustic foam (thickness 5 cm;
528 Basotect Plan50) chosen to maximize the visual contrast with the mice and simplify the cleaning
529 of excreta. The interaction platform had no walls to avoid acoustic reflections and was located
530 centrally in the booth. Its surface was elevated 25 cm above the floor (i.e. 20 cm above the foam
531 on the booth floor), which was generally sufficient in preventing animals from leaving the platform.
532 If a mouse left the platform, data was excluded from further analysis (<5% of frames).

533 Sounds inside the booth were recorded with 2 sets of microphones: (i) 4 high-quality
534 microphones (USM4) and (ii) a 64-channel microphone array (Cam64), both recording at a
535 sampling rate of 250 kHz at 16 bits. (i) The 4 high-quality microphones (CM16/CMPA48AAF-5V,
536 AviSoft, Berlin) were placed in a rectangle that contained the platform (see *Fig. 1A*) at a height
537 exceeding the platform by 12.1 cm to minimize the amount of sound blocked by the mice during
538 interaction. The position of a microphone was defined as the center of the recording membrane.
539 Considering the directional receptivity of the microphones (~25 dB attenuation at 45°), the
540 microphones were placed a short distance away from the corners of the platform to maximize
541 sound capture (5 cm in the long direction and 6 cm in the short direction of the platform). The
542 rotation of each microphone was chosen to be such that it aimed at the platform center. The
543 microphones produce a flat (± 5 dB) frequency response within 7-150 kHz that was low-pass
544 filtered at 120 kHz to prevent aliasing (using the analog, 16th order filter, which is part of the
545 microphone amplifier). Recorded data was digitized using a data acquisition card (PCIe-6351,
546 National Instruments). (ii) In addition, a 64-channel microphone array (Cam64 custom ultrasonic
547 version, Sorama B.V.) was mounted above the platform with a relative height of 46.5 cm
548 measured to the bottom of the Cam64 and a relative lateral shift of 6.52 cm to the right of the
549 platform midpoint. The Cam64 utilizes 64 MEMS microphones (Knowles, Digital Zero-Height
550 SiSonic, SPH0641LU4H-1) for acoustic data collection that are positioned in a Fermat's spiral
551 over a circle with a ~16 cm diameter. Raw microphone data was streamed to an m.2 SSD for later
552 analysis. Synchronization between the samples acquired by the Cam64 and the ultrasonic
553 microphones was performed by presenting 2 brief acoustic clicks (realized by stepping a digital
554 output from 0 to 5V) close to one of the microphones on the Cam64 at the start and end of each

555 trial using a headphone driver (IE 800, Sennheiser). The recorded pulses were automatically
556 retrieved and used to temporally align the recording sources.

557 A high-speed camera (PointGrey Flea3 FL3-U3-13Y3M-C, Monochrome, USB3.0) was
558 mounted above the platform with a relative height of 46.5 cm measured to the bottom of the front
559 end of the lens (6 mm, Thorlabs, part number: MVL6WA) and a relative lateral shift of 4.48 cm to
560 the left of the platform midpoint. Video was recorded with a field of view of 52.2 x 41.7 cm at
561 ~50 fps and digitized at 640 x 512 pixels (producing an effective resolution of ~0.815 mm/pixel).
562 The shutter time was set to 10 ms to guarantee good exposure while keeping the illumination
563 rather dim. The frame triggers from the camera were recorded on an analog channel in the PCIe-
564 6531 card for subsequent temporal alignment with the acoustic data.

565

566 *Experimental Procedures*

567 The experiment had 3 conditions: dyadic (with 2 mice; 57 trials), triadic (with 3 mice; 28 trials), as
568 well as monadic (single male mouse, ground truth data; 8 trials). Two dyadic trials were excluded
569 from further analysis due to repeated but required experimenter interference during the
570 recordings, leaving 55 dyadic trials. Each trial consisted of 8 minutes of free interaction between
571 at least 1 female and at least 1 male mouse on the platform. Females were always placed on the
572 platform first, and males were added shortly thereafter. In the monadic case, fresh female urine
573 was placed on the platform to prompt the male mouse to vocalize. The high-speed camera and 4
574 high-quality microphones started recording after all mice had been placed on the platform and
575 continued for 8 minutes. Data points where one mouse had left the platform or the hand of the
576 experimenter was visible 10 seconds before or after (e.g., to pick up a mouse) were discarded
577 (<5% of frames). Due to the rate of data generation of the Cam64 recordings (32 MB/s), their
578 duration and timing was optimized manually. The experimenter had access to the live
579 spectrogram from the USM4 microphones, and upon the start of USVs, triggered a new Cam64
580 recording (of fixed 2 min duration). If additional USVs occurred after that point, the experimenter
581 could trigger additional recordings.

582

583 *Data Analysis*

584 The analysis of the raw data involved multiple stages (see *Fig. 2*): From the audio data, the
585 presence and origin of USVs was estimated automatically. From the video data, mice were
586 carefully tracked by hand at the temporal midpoint of each USV as a ground truth comparison for
587 their acoustically localized origin. To estimate what proportion of our precision would be lost when
588 using a faster and more scalable visual tracking method, we also tracked the mice automatically

589 during dyadic trials. The estimated locations of the mice and USVs were then used to attribute
590 the USVs to their emitter. All these steps are described in detail below.

591

592 *Audio Preprocessing:* Prior to further analysis, acoustic recordings were filtered at different
593 frequencies. USM4 data was band-pass filtered between 30-110 kHz before further analysis using
594 an inverse impulse response filter of order 20 in Matlab (function: `designfilt`, type: `bandpassiir`).
595 Cam64 data was band-pass filtered with a frequency range adapted to the frequency content of
596 each USV. Specifically, first the frequency range of the USV was estimated as the 10th to 90th
597 percentile of the set of most intense frequencies at each time point. Next, this range was
598 broadened by 5 kHz at both ends, and then limited at the top end to 95 kHz. If this range exceeded
599 50 kHz, the lower end was set to 45 kHz. This ensured that beamforming was conducted over the
600 relevant frequencies for each USV and avoided the high-frequency regions where the Cam64
601 microphones are dominated by noise (see *Fig. 1C*).

602

603 *Video Preprocessing:* The high-speed camera lens failed to produce perfect rectilinear mapping
604 and was placed off-center with respect to the interaction platform, thereby producing a nonlinear
605 radial-tangential visual distortion. We corrected for the radial distortion with:

$$606 \quad x_{rd} = x_1 + \frac{\text{atan}(r_d / \lambda)}{r_d / \lambda} * (x_{ru} - x_1) * Z_x$$

$$607 \quad y_{rd} = y_1 + \frac{\text{atan}(r_d / \lambda)}{r_d / \lambda} * (y_{ru} - y_1) * Z_y$$

608 where $[x_{rd}, y_{rd}]$ represent the radially distorted image coordinates, $[x_1, y_1]$ the coordinates of the
609 image center, r_d the Euclidean distance to the radial distortion center, λ the distortion strength,
610 $[x_{ru}, y_{ru}]$ the radially undistorted coordinates, and Z_x, Z_y axis-specific zoom factors. The tangential
611 distortion, on the other hand, we corrected with:

$$612 \quad x_{td} = x_{tu} - \frac{(x_{tu} - a_x)}{|x - a_x|} * \frac{\kappa_x}{p_y} * (y_{tu} - \Delta p_y) * Z_x$$

$$613 \quad y_{td} = y_{tu} - \frac{(y_{tu} - a_y)}{|y - a_y|} * \frac{\kappa_y}{p_x} * (x_{tu} - \Delta p_x) * Z_y$$

614 where $[x_{td}, y_{td}]$ represent the tangentially distorted image coordinates, $[x_{tu}, y_{tu}]$ the tangentially
615 undistorted coordinates, $[a_x, a_y]$ the coordinates of the tangential distortion center, $[x, y]$ the size
616 of the image, $[\kappa_x, \kappa_y]$ the tangential distortion strengths, $[p_x, p_y]$ the size of the interaction platform
617 in the undistorted image, and $[\Delta p_x, \Delta p_y]$ the offset of the platform with respect to the top-left corner
618 of the undistorted image.

619
620 *Detection of Ultrasonic Vocalizations:* USVs were detected automatically using a set of custom
621 algorithms (see `VocCollector.m`) described elsewhere.⁷² Detection was only performed on the
622 USM4 data, as their sensitivity and frequency range was generally better than for the Cam64 (see
623 *Fig. 1C*). A vocalization only had to be detected on 1 of the 4 high-quality microphones to be
624 included into the set. In total, we collected 13406 USVs, out of which 8424 occurred when the
625 Cam64 recordings were active.

626
627 *Automatic Visual Animal Tracking:* To assess whether we could reliably assign USVs to their
628 emitter in a fast and scalable way, we automatically tracked multiple body parts of interacting mice
629 — most importantly the snout and head center — for all dyadic trials (see *Fig 2*). The tracking task
630 can be separated into 2 steps: (i) identifying all candidate locations for different body parts in the
631 video and (ii) linking these body part locations over time without switching identities between the
632 animals. We performed the first step by analyzing the video data offline in the XY-plane with
633 *DeepLabCut*, a toolbox that uses a convolutional neural network to visually identify animal
634 features (DLC; see Mathis et al., 2018). To train the network, a training set was used (~1400
635 frames) containing manually placed markers for both mice (i.e., snout, tail base, head center,
636 left/right ear, body center). The training set was constructed in 2 iterations, with the problematic
637 aspects of training using the first iteration (i.e., ~800 randomly drawn frames) prompting frame
638 choice in the second iteration (i.e., ~600 manually selected frames). Subsequently, DLC was
639 trained on the basis of this data (DLC v.2.2b8, running on a GTX 1070 GPU with NVIDIA driver
640 version 390.77 on Ubuntu 18.04.1 LTS; see *Supplementary Material* for a comprehensive
641 overview of all DLC parameters). The resulting convolutional neural network was then used to
642 predict the marker locations in all frames and all trials and provided quite high spatial precision,
643 although not all visible markers were represented by a location estimate in every frame. In
644 addition, while more recent versions of DLC (v.2.2+) should also be able to track identities
645 consistently over time, we did not manage to achieve high reliability in that respect using DLC.
646 Taken together, every trial would have required hundreds of manual corrections to achieve
647 reliable tracking over the entire period without identity switches.

648 To solve these 2 problems, we developed a dedicated solution in MATLAB (see
649 `C_trackMiceDLC.m`) that used spatial location estimates generated by DLC but disregarded
650 their unreliable identity. The general strategy was as follows: first, we requested in DLC not just
651 the estimate with the highest probability rating for each marker but rather an estimate cloud of the
652 50 most likely candidate locations for each marker and each mouse (see *Suppl. Fig. 1A*). Within

653 each marker class, all estimates were treated as identical at this point. Optimal marker locations
654 were then generated by within-frame k-means clustering of the estimate clouds (N.B., with k being
655 varied in a data-driven manner), or if clustering was suboptimal, by probability-weighted averaging
656 of heuristically separated estimate clouds (see *Suppl. Fig. 1B* and `LF_BPExtract` in
657 `C_trackMiceDLC.m`). After that, all markers underwent a temporal and spatial analysis in
658 tandem, both aimed at constructing pieces of unattributed tracks. In the temporal analysis, marker
659 positions that obviously belonged together were first assembled into small spatiotemporal threads
660 with the same, unknown identity (i.e., belonging to the same mouse). The speed and acceleration
661 that these small threads represented were subsequently assessed over time to yield longer
662 threads of marker positions with the same identity (see *Suppl. Fig. 1C* and `LF_tempAttr` in
663 `C_trackMiceDLC.m`). In the spatial analysis, all marker positions were analyzed on a frame-by-
664 frame basis, grouping markers with the same identity based on a logical combination of
665 anatomically permitted inter-marker distances (see *Suppl. Fig. 1D* and `LF_spatAttr` in
666 `C_trackMiceDLC.m`). Next, recognizing the overlap between the unattributed threads and the
667 attributed sets of marker positions, the spatial and temporal information was logically combined
668 in its entirety. This process was mathematically equivalent to a discrete convolution of the sets of
669 integers representing marker identities with the discrete response function describing the thread
670 to which they belong. For example, if a snout and head center clearly belong together in a specific
671 frame, the threads that run through them also belong together (see *Suppl. Fig. 1E* and
672 `LF_attrConv` in `C_trackMiceDLC.m`). Afterwards, the thread ends still loose were connected
673 based on quadratic spatial trajectory estimates for each marker individually. Lastly, all marker
674 trajectories were smoothed with piece-wise shape-preserving cubic interpolation. Reliable
675 tracking was ensured for all recordings by visual inspection and, if required, manual corrections
676 (~10 per trial on average, a major reduction).

677
678 *Manual Visual Animal Tracking:* We manually tracked the spatial locations of all mice during all
679 USVs from the video data to assess the precision of the automatic visual and acoustic tracking.
680 During manual tracking, the observer was presented with a combined display of the vocalization
681 spectrogram and the concurrent video image at the temporal midpoint of each USV (*MultiViewer*,
682 custom-written, MATLAB-based visualization tool). The display included a zoom function for
683 optimal accuracy, as tracking was click-based. Users could also freely scroll in time to ensure
684 consistent animal identities. Only the snout and head center (i.e., midpoint between the ears)
685 needed to be annotated because these points define a vector representing the head location and
686 direction, which was all that was required in subsequent behavioral analyses.

687

688 *Localization of Ultrasonic Vocalizations:* USVs were spatially localized using a hybrid approach
689 that drew on the complementary strengths of the 2 microphone arrays (see *Fig. 2*). For example,
690 the Cam64 array provided excellent localization for USVs with energy below ~90 kHz, due to the
691 increasing noise floor of the MEMS (microelectro-mechanical systems) microphones with sound
692 frequency. Conversely, the 4 high-quality ultrasonic microphones (USM4) have a rather flat noise
693 level as a function of frequency. On the other hand, USM4 will occasionally have glitches in one
694 of the microphones, which can be compensated for in Cam64-based estimates through the
695 number of microphones. As a consequence, the errors of the two methods show an L-shape (see
696 *Fig. 3A*), which highlights the synergy of a hybrid approach.

697 Acoustic localization using the Cam64 recordings was performed on the basis of delay-
698 and-sum beamforming⁸². In beamforming, signals from all microphones are combined to estimate
699 a spatial density that correlates with the probability of a given location being the origin of the
700 sound. Specifically, we computed beamforming estimates for a surface situated 1 cm above and
701 co-centered with the interaction platform, extending to 5 cm beyond all edges of the platform (i.e.,
702 50 x 40 cm in total) at a final resolution of 1 mm in both dimensions. We refer to this density of
703 sound origin as $D_{SO}(x, y)$ where x and y denote spatial coordinates. To prevent noises unrelated
704 to a specific USV from contaminating the location estimate, we limited beamforming to a particular
705 frequency range estimated from the simultaneous data of the USM4 array that enveloped the
706 USV. Spatial density was defined as

$$707 \quad D_{SO}(x, y) = \sum_{f=F_{min}}^{F_{max}} D_{SO}(x, y, f) = \sum_{f=F_{min}}^{F_{max}} \sum_{m=1}^{64} e^{i2\pi f d(m, x, y, z)}$$

708 where $d(m, x, y, z)$ denotes the difference in arrival time at each microphone m for sounds emitted
709 from a location with coordinates (x, y, z) , where z is omitted in $D_{SO}(x, y)$ as it is a fixed distance to
710 the plane of the microphone array. Beamforming was performed in the computational cloud
711 backend provided by the Cam64 manufacturer, the so-called Sorama Portal¹.

712 The final beamforming estimate was calculated sequentially in 2 steps: first, a coarse
713 estimate with 1 cm resolution was generated over the entire beamforming surface. Second, a
714 fine-grained estimate with 1 mm resolution was generated over a 30 x 30 mm window centered
715 on the peak location of the coarse estimate (see *Fig. 2* for an example). This two-step approach
716 was chosen to optimize performance, as an estimate with 1 mm resolution over the entire
717 beamforming surface would be computationally expensive while failing to produce a better result.

¹ www.sorama.eu/sorama-portal

718 For USVs of sufficient quality (i.e., containing frequency content below ~90 kHz while being
719 sufficiently intense and long), both the coarse and fine estimates of $D_{SO}(x, y)$ contained a peak
720 whose height was typically very large compared to the surrounding values at distances greater
721 than a few cm's. The peak location of the fine-grained estimate was used as the final estimate of
722 the USV's origin. To assess the quality of this location estimate, we computed a signal-to-noise
723 ratio (SNR) per USV as follows:

$$724 \quad SNR_{Cam64}(v) = \frac{\max(D_{SO}(x, y))}{\text{std}(D_{SO}(x, y))}$$

725 where $D_{SO}(x, y)$ is assumed to be calculated for the USV v . The inverse, $1/SNR_{Cam64}$ was used
726 as a proxy for the uncertainty of localization for a given USV.

727 Localization from the USM4 recordings was performed using the SLIM method⁷³. Briefly,
728 SLIM analytically estimates submanifolds (in 2D: surfaces) of a sound's spatial origin for each pair
729 of microphones and combines these into a single estimate by intersecting the manifolds (in 2D:
730 lines). The intersection has an associated uncertainty which scales with the uncertainty of the
731 localization estimate for a given USVs, specifically the uncertainty was defined as the standard
732 deviation of all locations that were >90% times the maximum of the intersection density of all
733 origin curves.

734 Lastly, for each USV where both Cam64 and SLIM location estimates \dot{X}_{Cam64} and
735 \dot{X}_{SLIM} were available, a single estimate \dot{X}_{HyVL} was computed based on the two estimates, spatial
736 uncertainties and their spatial relation to the mice at the current time (see below).

737
738 *USV Assignment:* The final, hybrid location estimate and assignment to a mouse was performed
739 while taking into account the probability of making a false assignment as proposed before⁷⁶,
740 through the calculation of the mouse probability index MPI . While the MPI was previously only
741 used to exclude uncertain assignments (e.g. if two mice are nearly equidistant to the estimated
742 sound location), we also adapted it here to select and combine the location estimates. The MPI_k
743 for each mouse k was computed as

$$744 \quad MPI_k = \frac{P_k}{\sum_{m=1}^n P_m}$$

745 Here, P_k is the probability that the USV in question originated from mouse k computed as
746 $P_k = N(\dot{X}_{Method} - X_{mouth,k}, \sigma_{Method}^2)$, where \dot{X}_{Method} is an estimate of the acoustic origin, $X_{mouth,k}$
747 the position of the mouth of mouse k , and σ_{Method}^2 the uncertainty of the estimate, with \dot{X}_{Method}
748 and σ_{Method}^2 specific to the Method used. $X_{mouth,k}$ was assumed to lie on a line connecting the
749 snout and head-center. For manually tracked recordings, the optimal location on this line was

750 close to the snout (~2% towards the head, where % is relative to the snout-to-head-center tracked
751 distance), while in the automatic tracking it was ahead of the snout tracking point (~15% away
752 from the head). σ_{Method}^2 was computed for each USV as the method's intrinsic per-USV
753 uncertainty estimate. As these uncertainty estimates only correlate with the absolute uncertainty
754 (i.e. in mm's), we scaled them such that their average across all USVs matched the residual error
755 of each method in the Far-condition (all animals >100 mm apart, see Fig. 3C and Oliveira-Stahl
756 et al.⁷³). In this way, the MPI_k for individual USVs took into account the uncertainty of each
757 method: if the uncertainty of one method was higher, probabilities across mice would become
758 more similar and the MPI_k would reduce.

759 For a given USV, we computed the MPI_k for all mice for both methods. The mouse with
760 the largest MPI_k per method - which coincides with the mouse at the smallest distance to the
761 estimate - was denoted as MPI_{Cam64} and MPI_{SLIM} , respectively. If only one of the two exceeded
762 0.95, this method's estimate was selected. If both exceeded 0.95, then the estimate with the
763 smaller distance to the mouse with the highest MPI_k was chosen. This combination ensured that
764 only reliable assignments were performed, while minimizing the residual error. Similar to
765 Neunuebel et al. 2015, we also excluded estimates that were too far away from any mouse
766 (50 mm). This distance threshold mainly serves to compensate for a deficiency of the MPI : if all
767 mice are far from the estimate, all P_k are extremely small, however, the MPI_k will often exceed
768 0.95. The distance threshold corresponds to setting the individual $P_k = 0$ in the MPI_k , thus
769 excluding candidate mice which are highly unlikely to be the source of the USV. USVs that had
770 no $MPI_k > 0.95$ for either method were excluded from further analysis. The fraction of included
771 USVs is referred to as *selected* in the plots. Maximizing this fraction is essential to perform a
772 complete analysis of vocal communication.

773 We compared the above-described combination strategy to a large number of alternative
774 strategies, including maximum likelihood combination of estimators (Ernst & Banks 2002), or
775 selecting directly based on the largest MPI_k or largest P_k . While all these approaches led to
776 broadly similar results, the described approach achieved the most robust and reliable results (see
777 Discussion for additional details).

778
779 *Audiovisual Alignment:* For both microphone sets, precise measurements of their location in
780 relation to the camera's location were used to position acoustic estimates in the coordinate system
781 of the images provided by the camera. In the final analysis, we noticed for each microphone set
782 small, systematic (0.5-2 mm) shifts in both X and Y. We interpreted these as very small
783 measurement errors in the relative positions of the camera or microphone arrays, and corrected

784 these post-hoc in the setup definition, followed by rerunning all subsequent analysis steps. This
785 reduced all systematic shifts to near 0.

786
787 *Spatial Vocalization Analysis:* To gain insight into the spatial positioning of the interacting mice,
788 we represented the relative animal positions in a polar reference frame centered on the snout of
789 the emitter. In this format, the radial distance corresponded to the snout-snout distance and the
790 radial angle described the relative angle between the gaze direction of the emitter and the snout
791 position of the recipient (i.e., with the line from the head center to the snout of the emitter pointing
792 towards 0°; see also *Fig. 5A*).

793 The position density of the recipient mouse was collected in cumulative fashion, with the
794 polar coordinate system translated appropriately for each USV based on its temporal midpoint.
795 We assumed that the mice had no preference for relative vocalizations to either side of their snout,
796 so all relative spatial positions were agglomerated in the right hemisphere for further analysis. All
797 data points were then binned using a polar, raw-count histogram with bins of 10° and 1 cm.

798
799 *Statistical Analysis*

800 To avoid distributional assumptions, all statistical tests were nonparametric, i.e., Wilcoxon rank
801 sum test for two-group comparisons and Kruskal-Wallis for single factor analysis of variance.
802 Correlations were computed as Spearman's rank-based correlation coefficients. Error bars
803 represent standard errors of the mean (SEM) unless stated otherwise. All statistical analyses
804 were performed in MATLAB v.2018b (The Mathworks, Natick) using functions from the Statistics
805 Toolbox.

806
807 **Acknowledgements**

808
809 We would like to thank Lucas Noldus for suggesting the use of the Sorama Cam64 and Maurice
810 Camp and Toros Senan for technical support relating to the operation and data handling of the
811 Cam64 and the Sorama Portal. We would like to thank Amber van der Stam, Dionne Lenferink,
812 Soha Farboud for assisting with the animal handling and experimental control. BE acknowledges
813 funding from a DCN Internal Grant, funded by Noldus IT b.v. as well as an NWO VIDI grant
814 (016.VIDI.189.052) and a Technology Hotel Grant (ZonMW, 40-43500-98-4141).

815 Bibliography

- 816 1. Mahrt, E. J., Perkel, D. J., Tong, L., Rubel, E. W. & Portfors, C. V. Engineered deafness
817 reveals that mouse courtship vocalizations do not require auditory experience. *J. Neurosci.*
818 **33**, 5573–5583 (2013).
- 819 2. Brudzynski, S. M. Biological functions of rat ultrasonic vocalizations, arousal mechanisms,
820 and call initiation. *Brain Sci.* **11**, (2021).
- 821 3. Zaytseva, A. S., Volodin, I. A., Ilchenko, O. G. & Volodina, E. V. Ultrasonic vocalization of
822 pup and adult fat-tailed gerbils (*Pachyuromys duprasi*). *PLoS ONE* **14**, e0219749 (2019).
- 823 4. Volodin, I. A., Dymkaya, M. M., Smorkatcheva, A. V. & Volodina, E. V. Ultrasound from
824 underground: cryptic communication in subterranean wild-living and captive northern mole
825 voles (*Ellobius talpinus*). *Bioacoustics* **31**, 414–434 (2022).
- 826 5. Murrant, M. N. *et al.* Ultrasonic vocalizations emitted by flying squirrels. *PLoS ONE* **8**,
827 e73045 (2013).
- 828 6. Schnitzler, H.-U., Moss, C. F. & Denzinger, A. From spatial orientation to food acquisition in
829 echolocating bats. *Trends Ecol. Evol.* **18**, 386–394 (2003).
- 830 7. Feng, A. S. *et al.* Ultrasonic communication in frogs. *Nature* **440**, 333–336 (2006).
- 831 8. Mourlam, M. J. & Orliac, M. J. Infrasonic and Ultrasonic Hearing Evolved after the
832 Emergence of Modern Whales. *Curr. Biol.* **27**, 1776-1781.e9 (2017).
- 833 9. Bakker, J. & Langermans, J. A. M. Ultrasonic components of vocalizations in marmosets. in
834 *Handbook of Ultrasonic Vocalization - A Window into the Emotional Brain* vol. 25 535–544
835 (Elsevier, 2018).
- 836 10. Ramsier, M. A. *et al.* Primate communication in the pure ultrasound. *Biol. Lett.* **8**, 508–511
837 (2012).
- 838 11. Kikusui, T. *et al.* Cross fostering experiments suggest that mice songs are innate. *PLoS*
839 *ONE* **6**, e17721 (2011).
- 840 12. Litvin, Y., Blanchard, D. C. & Blanchard, R. J. Rat 22kHz ultrasonic vocalizations as alarm

- 841 cries. *Behav. Brain Res.* **182**, 166–172 (2007).
- 842 13. Marconi, M. A., Nicolakis, D., Abbasi, R., Penn, D. J. & Zala, S. M. Ultrasonic courtship
843 vocalizations of male house mice contain distinct individual signatures. *Animal Behaviour*
844 **169**, 169–197 (2020).
- 845 14. Rieger, N. S. & Marler, C. A. The function of ultrasonic vocalizations during territorial
846 defence by pair-bonded male and female California mice. *Animal Behaviour* **135**, 97–108
847 (2018).
- 848 15. Hammerschmidt, K., Radyushkin, K., Ehrenreich, H. & Fischer, J. Female mice respond to
849 male ultrasonic “songs” with approach behaviour. *Biol. Lett.* **5**, 589–592 (2009).
- 850 16. Pultorak, J. D., Matusinec, K. R., Miller, Z. K. & Marler, C. A. Ultrasonic vocalization
851 production and playback predicts intrapair and extrapair social behaviour in a monogamous
852 mouse. *Animal Behaviour* **125**, 13–23 (2017).
- 853 17. Chabout, J., Sarkar, A., Dunson, D. B. & Jarvis, E. D. Male mice song syntax depends on
854 social contexts and influences female preferences. *Front. Behav. Neurosci.* **9**, 76 (2015).
- 855 18. Musolf, K., Meindl, S., Larsen, A. L., Kalcounis-Rueppell, M. C. & Penn, D. J. Ultrasonic
856 Vocalizations of Male Mice Differ among Species and Females Show Assortative
857 Preferences for Male Calls. *PLoS ONE* **10**, e0134123 (2015).
- 858 19. Sugimoto, H. *et al.* A role for strain differences in waveforms of ultrasonic vocalizations
859 during male-female interaction. *PLoS ONE* **6**, e22093 (2011).
- 860 20. Tschida, K. *et al.* A specialized neural circuit gates social vocalizations in the mouse.
861 *Neuron* **103**, 459-472.e4 (2019).
- 862 21. Warren, M. R., Clein, R. S., Spurrier, M. S., Roth, E. D. & Neunuebel, J. P. Ultrashort-
863 range, high-frequency communication by female mice shapes social interactions. *Sci. Rep.*
864 **10**, 2637 (2020).
- 865 22. Nicolakis, D., Marconi, M. A., Zala, S. M. & Penn, D. J. Ultrasonic vocalizations in house
866 mice depend upon genetic relatedness of mating partners and correlate with subsequent

- 867 reproductive success. *Front. Zool.* **17**, 10 (2020).
- 868 23. Rieger, N. S., Monari, P. K., Hartfield, K., Schefelker, J. & Marler, C. A. Pair-bonding leads
869 to convergence in approach behavior to conspecific vocalizations in California mice
870 (*Peromyscus californicus*). *PLoS ONE* **16**, e0255295 (2021).
- 871 24. Kalcounis-Rueppell, M. C. & Petric, R. Female and male adult brush mice (*Peromyscus*
872 *boylii*) use ultrasonic vocalizations in the wild. *Behaviour* **150**, 1747–1766 (2013).
- 873 25. Chen, J. *et al.* Flexible scaling and persistence of social vocal communication. *Nature* **593**,
874 108–113 (2021).
- 875 26. de Chaumont, F., Lemièrre, N., Coqueran, S., Bourgeron, T. & Ey, E. LMT USV Toolbox, a
876 Novel Methodological Approach to Place Mouse Ultrasonic Vocalizations in Their
877 Behavioral Contexts-A Study in Female and Male C57BL/6J Mice and in Shank3 Mutant
878 Females. *Front. Behav. Neurosci.* **15**, 735920 (2021).
- 879 27. Castellucci, G. A., Calbick, D. & McCormick, D. The temporal organization of mouse
880 ultrasonic vocalizations. *PLoS ONE* **13**, e0199929 (2018).
- 881 28. Pultorak, J. D., Alger, S. J., Loria, S. O., Johnson, A. M. & Marler, C. A. Changes in
882 behavior and ultrasonic vocalizations during pair bonding and in response to an infidelity
883 challenge in monogamous california mice. *Front. Ecol. Evol.* **6**, (2018).
- 884 29. Burke, K., Screven, L. A. & Dent, M. L. CBA/CaJ mouse ultrasonic vocalizations depend on
885 prior social experience. *PLoS ONE* **13**, e0197774 (2018).
- 886 30. Zala, S. M., Reitschmidt, D., Noll, A., Balazs, P. & Penn, D. J. Sex-dependent modulation of
887 ultrasonic vocalizations in house mice (*Mus musculus musculus*). *PLoS ONE* **12**, e0188647
888 (2017).
- 889 31. Mun, H.-S., Lipina, T. V. & Roder, J. C. Ultrasonic Vocalizations in Mice During Exploratory
890 Behavior are Context-Dependent. *Front. Behav. Neurosci.* **9**, 316 (2015).
- 891 32. von Merten, S., Hoier, S., Pfeifle, C. & Tautz, D. A role for ultrasonic vocalisation in social
892 communication and divergence of natural populations of the house mouse (*Mus musculus*

- 893 domesticus). *PLoS ONE* **9**, e97244 (2014).
- 894 33. Scattoni, M. L., Crawley, J. & Ricceri, L. Ultrasonic vocalizations: a tool for behavioural
895 phenotyping of mouse models of neurodevelopmental disorders. *Neurosci. Biobehav. Rev.*
896 **33**, 508–515 (2009).
- 897 34. Warren, M. R., Spurrier, M. S., Sangiamo, D. T., Clein, R. S. & Neunuebel, J. P. Mouse
898 vocal emission and acoustic complexity do not scale linearly with the size of a social group.
899 *J. Exp. Biol.* **224**, (2021).
- 900 35. Dou, X., Shirahata, S. & Sugimoto, H. Functional clustering of mouse ultrasonic
901 vocalization data. *PLoS ONE* **13**, e0196834 (2018).
- 902 36. Hoier, S., Pfeifle, C., von Merten, S. & Linnenbrink, M. Communication at the Garden
903 Fence--Context Dependent Vocalization in Female House Mice. *PLoS ONE* **11**, e0152255
904 (2016).
- 905 37. Chabout, J. *et al.* Adult male mice emit context-specific ultrasonic vocalizations that are
906 modulated by prior isolation or group rearing environment. *PLoS ONE* **7**, e29401 (2012).
- 907 38. Hertz, S., Weiner, B., Perets, N. & London, M. Temporal structure of mouse courtship
908 vocalizations facilitates syllable labeling. *Commun. Biol.* **3**, 333 (2020).
- 909 39. Kikusui, T. *et al.* Testosterone increases the emission of ultrasonic vocalizations with
910 different acoustic characteristics in mice. *Front. Psychol.* **12**, 680176 (2021).
- 911 40. Kikusui, T. *et al.* Testosterone regulates the emission of ultrasonic vocalizations and
912 mounting behavior during different developmental periods in mice. *Dev. Psychobiol.* **63**,
913 725–733 (2021).
- 914 41. Timonin, M. E., Kalcounis-Rueppell, M. C. & Marler, C. A. Testosterone pulses at the nest
915 site modify ultrasonic vocalization types in a monogamous and territorial mouse. *Ethology*
916 **124**, 804–815 (2018).
- 917 42. Pultorak, J. D., Fuxjager, M. J., Kalcounis-Rueppell, M. C. & Marler, C. A. Male fidelity
918 expressed through rapid testosterone suppression of ultrasonic vocalizations to novel

- 919 females in the monogamous California mouse. *Horm. Behav.* **70**, 47–56 (2015).
- 920 43. Guoynes, C. D. & Marler, C. A. An acute dose of intranasal oxytocin rapidly increases
921 maternal communication and maintains maternal care in primiparous postpartum California
922 mice. *PLoS ONE* **16**, e0244033 (2021).
- 923 44. Tsuji, T. *et al.* Oxytocin administration modulates the complex type of ultrasonic
924 vocalisation of mice pups prenatally exposed to valproic acid. *Neurosci. Lett.* **758**, 135985
925 (2021).
- 926 45. Tsuji, C., Fujisaku, T. & Tsuji, T. Oxytocin ameliorates maternal separation-induced
927 ultrasonic vocalisation calls in mouse pups prenatally exposed to valproic acid. *J.*
928 *Neuroendocrinol.* **32**, e12850 (2020).
- 929 46. Michael, V. *et al.* Circuit and synaptic organization of forebrain-to-midbrain pathways that
930 promote and suppress vocalization. *eLife* **9**, (2020).
- 931 47. Gao, S.-C., Wei, Y.-C., Wang, S.-R. & Xu, X.-H. Medial preoptic area modulates courtship
932 ultrasonic vocalization in adult male mice. *Neurosci. Bull.* **35**, 697–708 (2019).
- 933 48. Tasaka, G.-I. *et al.* Genetic tagging of active neurons in auditory cortex reveals maternal
934 plasticity of coding ultrasonic vocalizations. *Nat. Commun.* **9**, 871 (2018).
- 935 49. Fröhlich, H., Rafiullah, R., Schmitt, N., Abele, S. & Rappold, G. A. Foxp1 expression is
936 essential for sex-specific murine neonatal ultrasonic vocalization. *Hum. Mol. Genet.* **26**,
937 1511–1521 (2017).
- 938 50. Shepard, K. N., Chong, K. K. & Liu, R. C. Contrast Enhancement without Transient Map
939 Expansion for Species-Specific Vocalizations in Core Auditory Cortex during Learning.
940 *eNeuro* **3**, (2016).
- 941 51. Arriaga, G. & Jarvis, E. D. Mouse vocal communication system: are ultrasounds learned or
942 innate? *Brain Lang.* **124**, 96–116 (2013).
- 943 52. Fujita, E., Tanabe, Y., Imhof, B. A., Momoi, M. Y. & Momoi, T. Cadm1-expressing synapses
944 on Purkinje cell dendrites are involved in mouse ultrasonic vocalization activity. *PLoS ONE*

- 945 7, e30151 (2012).
- 946 53. Wang, H., Liang, S., Burgdorf, J., Wess, J. & Yeomans, J. Ultrasonic vocalizations induced
947 by sex and amphetamine in M2, M4, M5 muscarinic and D2 dopamine receptor knockout
948 mice. *PLoS ONE* **3**, e1893 (2008).
- 949 54. Yang, X., Guo, D., Li, K. & Shi, L. Altered postnatal developmental patterns of ultrasonic
950 vocalizations in Dock4 knockout mice. *Behav. Brain Res.* **406**, 113232 (2021).
- 951 55. Binder, M. S., Shi, H. D. & Bordey, A. CD-1 Outbred Mice Produce Less Variable Ultrasonic
952 Vocalizations Than FVB Inbred Mice, While Displaying a Similar Developmental Trajectory.
953 *Front. Psychiatry* **12**, 687060 (2021).
- 954 56. Hepbasli, D. *et al.* Genotype- and Age-Dependent Differences in Ultrasound Vocalizations
955 of SPRED2 Mutant Mice Revealed by Machine Deep Learning. *Brain Sci.* **11**, (2021).
- 956 57. Agarwalla, S. *et al.* Male-specific alterations in structure of isolation call sequences of
957 mouse pups with 16p11.2 deletion. *Genes Brain Behav.* **19**, e12681 (2020).
- 958 58. Tsai, P. T. *et al.* Autistic-like behaviour and cerebellar dysfunction in Purkinje cell Tsc1
959 mutant mice. *Nature* **488**, 647–651 (2012).
- 960 59. Hodges, S. L., Nolan, S. O., Reynolds, C. D. & Lugo, J. N. Spectral and temporal properties
961 of calls reveal deficits in ultrasonic vocalizations of adult Fmr1 knockout mice. *Behav. Brain*
962 *Res.* **332**, 50–58 (2017).
- 963 60. Silverman, J. L., Yang, M., Lord, C. & Crawley, J. N. Behavioural phenotyping assays for
964 mouse models of autism. *Nat. Rev. Neurosci.* **11**, 490–502 (2010).
- 965 61. Ciucci, M. R. *et al.* Reduction of dopamine synaptic activity: degradation of 50-kHz
966 ultrasonic vocalization in rats. *Behav. Neurosci.* **123**, 328–336 (2009).
- 967 62. Palmateer, J. *et al.* Ultrasonic vocalization in murine experimental stroke: A mechanistic
968 model of aphasia. *Restor. Neurol. Neurosci.* **34**, 287–295 (2016).
- 969 63. Erata, E. *et al.* Cnksr2 Loss in Mice Leads to Increased Neural Activity and Behavioral
970 Phenotypes of Epilepsy-Aphasia Syndrome. *J. Neurosci.* **41**, 9633–9649 (2021).

- 971 64. Menuet, C. *et al.* Age-related impairment of ultrasonic vocalization in Tau.P301L mice:
972 possible implication for progressive language disorders. *PLoS ONE* **6**, e25770 (2011).
- 973 65. Palazzo, E., Fu, Y., Ji, G., Maione, S. & Neugebauer, V. Group III mGluR7 and mGluR8 in
974 the amygdala differentially modulate nocifensive and affective pain behaviors.
975 *Neuropharmacology* **55**, 537–545 (2008).
- 976 66. Moskal, J. R. & Burgdorf, J. Ultrasonic vocalizations in rats as a measure of emotional
977 responses to stress: models of anxiety and depression. in *Handbook of Ultrasonic*
978 *Vocalization - A Window into the Emotional Brain* vol. 25 413–421 (Elsevier, 2018).
- 979 67. Coffey, K. R., Marx, R. G. & Neumaier, J. F. DeepSqueak: a deep learning-based system
980 for detection and analysis of ultrasonic vocalizations. *Neuropsychopharmacology* **44**, 859–
981 868 (2019).
- 982 68. Fonseca, A. H., Santana, G. M., Bosque Ortiz, G. M., Bampi, S. & Dietrich, M. O. Analysis
983 of ultrasonic vocalizations from mice using computer vision and machine learning. *eLife* **10**,
984 (2021).
- 985 69. Zala, S. M., Reitschmidt, D., Noll, A., Balazs, P. & Penn, D. J. Automatic mouse ultrasound
986 detector (A-MUD): A new tool for processing rodent vocalizations. *PLoS ONE* **12**,
987 e0181200 (2017).
- 988 70. Van Segbroeck, M., Knoll, A. T., Levitt, P. & Narayanan, S. MUPET-Mouse Ultrasonic
989 Profile EXtraction: A Signal Processing Tool for Rapid and Unsupervised Analysis of
990 Ultrasonic Vocalizations. *Neuron* **94**, 465-485.e5 (2017).
- 991 71. Chabout, J., Jones-Macopson, J. & Jarvis, E. D. Eliciting and analyzing male mouse
992 ultrasonic vocalization (USV) songs. *J. Vis. Exp.* (2017) doi:10.3791/54137.
- 993 72. Ivanenko, A., Watkins, P., van Gerven, M. A. J., Hammerschmidt, K. & Englitz, B.
994 Classifying sex and strain from mouse ultrasonic vocalizations using deep learning. *PLoS*
995 *Comput. Biol.* **16**, e1007918 (2020).
- 996 73. Oliveira-Stahl, G. *et al.* High-precision spatial analysis of mouse courtship vocalization

- 997 behavior reveals sex and strain differences. *BioRxiv* (2021)
998 doi:10.1101/2021.10.22.464496.
- 999 74. Heckman, J. J. *et al.* High-precision spatial localization of mouse vocalizations during social
1000 interaction. *Sci. Rep.* **7**, 3017 (2017).
- 1001 75. Warren, M. R., Sangiamo, D. T. & Neunuebel, J. P. High channel count microphone array
1002 accurately and precisely localizes ultrasonic signals from freely-moving mice. *J. Neurosci.*
1003 *Methods* **297**, 44–60 (2018).
- 1004 76. Neunuebel, J. P., Taylor, A. L., Arthur, B. J. & Egnor, S. E. R. Female mice ultrasonically
1005 interact with males during courtship displays. *eLife* **4**, (2015).
- 1006 77. Mahrt, E., Agarwal, A., Perkel, D., Portfors, C. & Elemans, C. P. H. Mice produce ultrasonic
1007 vocalizations by intra-laryngeal planar impinging jets. *Curr. Biol.* **26**, R880–R881 (2016).
- 1008 78. Liu, R. C., Miller, K. D., Merzenich, M. M. & Schreiner, C. E. Acoustic variability and
1009 distinguishability among mouse ultrasound vocalizations. *J. Acoust. Soc. Am.* **114**, 3412–
1010 3422 (2003).
- 1011 79. Holy, T. E. & Guo, Z. Ultrasonic songs of male mice. *PLoS Biol.* **3**, e386 (2005).
- 1012 80. Barnes, T. D., Rieger, M. A., Dougherty, J. D. & Holy, T. E. Group and individual variability
1013 in mouse pup isolation calls recorded on the same day show stability. *Front. Behav.*
1014 *Neurosci.* **11**, 243 (2017).
- 1015 81. Musolf, K., Hoffmann, F. & Penn, D. J. Ultrasonic courtship vocalizations in wild house
1016 mice, *Mus musculus musculus*. *Animal Behaviour* **79**, 757–764 (2010).
- 1017 82. Van Veen, B. D. & Buckley, K. M. Beamforming: a versatile approach to spatial filtering.
1018 *IEEE ASSP Mag.* **5**, 4–24 (1988).
- 1019 83. Mathis, A. *et al.* DeepLabCut: markerless pose estimation of user-defined body parts with
1020 deep learning. *Nat. Neurosci.* **21**, 1281–1289 (2018).
- 1021 84. Rotschafer, S. E., Trujillo, M. S., Dansie, L. E., Ethell, I. M. & Razak, K. A. Minocycline
1022 treatment reverses ultrasonic vocalization production deficit in a mouse model of Fragile X

- 1023 Syndrome. *Brain Res.* **1439**, 7–14 (2012).
- 1024 85. Choi, H., Park, S. & Kim, D. Two genetic loci control syllable sequences of ultrasonic
1025 courtship vocalizations in inbred mice. *BMC Neurosci.* **12**, 104 (2011).
- 1026 86. Pomerantz, S. M. & Clemens, L. G. Ultrasonic vocalizations in male deer mice
1027 (*Peromyscus maniculatus bairdi*): their role in male sexual behavior. *Physiol. Behav.* **27**,
1028 869–872 (1981).
- 1029 87. Nunez, A. A., Nyby, J. & Whitney, G. The effects of testosterone, estradiol, and
1030 dihydrotestosterone on male mouse (*Mus musculus*) ultrasonic vocalizations. *Horm. Behav.*
1031 **11**, 264–272 (1978).
- 1032 88. Sangiamo, D. T., Warren, M. R. & Neunuebel, J. P. Ultrasonic signals associated with
1033 different types of social behavior of mice. *Nat. Neurosci.* **23**, 411–422 (2020).
- 1034 89. Cha Zhang, Florencio, D., Ba, D. E. & Zhengyou Zhang. Maximum likelihood sound source
1035 localization and beamforming for directional microphone arrays in distributed meetings.
1036 *IEEE Trans. Multimedia* **10**, 538–548 (2008).
- 1037 90. Shemesh, Y. *et al.* High-order social interactions in groups of mice. *eLife* **2**, e00759 (2013).
- 1038 91. Lee, N. S. & Beery, A. K. Neural circuits underlying rodent sociality: A comparative
1039 approach. *Curr. Top. Behav. Neurosci.* **43**, 211–238 (2019).
- 1040 92. Kappel, S., Hawkins, P. & Mendl, M. T. To group or not to group? good practice for housing
1041 male laboratory mice. *Animals (Basel)* **7**, (2017).
- 1042 93. Kondrakiewicz, K., Kostecky, M., Szadzińska, W. & Knapska, E. Ecological validity of social
1043 interaction tests in rats and mice. *Genes Brain Behav.* **18**, e12525 (2019).
- 1044 94. Keesom, S. M., Finton, C. J., Sell, G. L. & Hurley, L. M. Early-Life Social Isolation
1045 Influences Mouse Ultrasonic Vocalizations during Male-Male Social Encounters. *PLoS ONE*
1046 **12**, e0169705 (2017).
- 1047 95. Portfors, C. V. Types and functions of ultrasonic vocalizations in laboratory rats and mice.
1048 *J. Am. Assoc. Lab. Anim. Sci.* **46**, 28–34 (2007).

- 1049 96. Ernst, M. O. & Banks, M. S. Humans integrate visual and haptic information in a statistically
1050 optimal fashion. *Nature* **415**, 429–433 (2002).
- 1051 97. Urbanus, B. H. A., Peter, S., Fisher, S. E. & De Zeeuw, C. I. Region-specific Foxp2
1052 deletions in cortex, striatum or cerebellum cannot explain vocalization deficits observed in
1053 spontaneous global knockouts. *Sci. Rep.* **10**, 21631 (2020).