
Research Article: Methods/New Tools | Novel Tools and Methods

SaLSa: a combinatory approach of semi-automatic labeling and long short-term memory to classify behavioral syllables

<https://doi.org/10.1523/ENEURO.0201-23.2023>

Cite as: eNeuro 2023; 10.1523/ENEURO.0201-23.2023

Received: 13 June 2023

Revised: 19 October 2023

Accepted: 9 November 2023

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.eneuro.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

Copyright © 2023 Sakata

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

1 SaLSa: a combinatorial approach of semi-automatic
2 labeling and long short-term memory to classify
3 behavioral syllables
4

5 Abbreviated title: SaLSa for behavioral syllable classification

6
7 Shuzo Sakata

8 Strathclyde Institute of Pharmacy and Biomedical Sciences, University of Strathclyde, 161
9 Cathedral Street, Glasgow G4 0RE, UK

10

11 Author Contributions: SS designed research, performed research and wrote the paper.

12

13 Correspondence should be address to Shuzo Sakata (shuzo.sakata@strath.ac.uk).

14

15 Number of Figures: 6

16 Number of Extended Data Figures: 2

17 Number of Tables: 0

18 Number of words for Abstract: 218

19 Number of words for Significance Statement: 115

20 Number of words for Introduction: 500

21 Number of words for Discussion: 1057

22

23 Acknowledgements: We thank Abigail Hatcher Davies for her technical assistance. This work was
24 supported by MRC (MR/V033964/1) and Horizon2020-ICT (DEEPER, 101016787) to SS.

25

26 Conflict of interest: Author reports no conflict of interest.

27

28 Funding Sources: MRC (MR/V033964/1) and Horizon2020-ICT (DEEPER, 101016787) to SS

29

30 Abstract

31 Accurately and quantitatively describing mouse behavior is an important area. Although advances
32 in machine learning have made it possible to track their behaviors accurately, reliable classification
33 of behavioral sequences or syllables remains a challenge. In this study, we present a novel
34 machine learning approach, called SaLSa (a combination of semi-automatic labeling and long
35 short-term memory-based classification), to classify behavioral syllables of mice exploring an open
36 field. This approach consists of two major steps: first, after tracking multiple body parts, spatial and
37 temporal features of their egocentric coordinates are extracted. A fully automated unsupervised
38 process identifies candidates for behavioral syllables, followed by manual labeling of behavioral
39 syllables using a graphical user interface. Second, a long short-term memory (LSTM) classifier is
40 trained with the labeled data. We found that the classification performance was marked over 97%.
41 It provides a performance equivalent to a state-of-the-art model while classifying some of the
42 syllables. We applied this approach to examine how hyperactivity in a mouse model of Alzheimer's
43 disease (AD) develops with age. When the proportion of each behavioral syllable was compared
44 between genotypes and sexes, we found that the characteristic hyper-locomotion of female AD
45 mice emerges between 4 and 8 months. In contrast, age-related reduction in rearing is common
46 regardless of genotype and sex. Overall, SaLSa enables detailed characterization of mouse
47 behavior.

48

49 Significance Statement

50 Describing complex animal behavior is a challenge. Here, we developed an open-source,
51 combinatory approach to behavioral syllable classification, called SaLSa (a combination of semi-
52 automatic labeling and long short-term memory-based classification). In order to classify
53 behavioral syllables, this approach combines multiple machine learning methods to label video
54 frames semi-automatically and train a deep learning model. To demonstrate SaLSa's versatility, we
55 monitored the exploratory behavior of an Alzheimer's disease mouse model and delineated their
56 complex behaviors. We found that female Alzheimer's mice become hyperactive in the sense that
57 their locomotion behavior, but not other active behaviors, appear more frequently than controls and
58 even male Alzheimer's mice as they age. SaLSa offers a toolkit to analyze complex behaviors.

59 Introduction

60 In modern neuroscience, a goal is to understand the relationship between behavior and neural
61 ensembles. However, accurately and quantitatively describing complex behavior remains
62 challenging. In the past, mouse behaviors have often been assessed using simple, subjective
63 criteria in a series of behavioral tests. However, advances in machine learning have changed the
64 field (Mathis et al., 2020; Pereira et al., 2020; Luxem et al., 2023).

65 To describe mouse behaviors, several steps are required. First, behaviors are video-monitored.
66 The state-of-the-art is to utilize a depth camera (Wiltchko et al., 2020) or multiple cameras (Dunn
67 et al., 2021; Schneider et al., 2022). While some experiments require video recording from the
68 bottom to monitor limb movement (Pereira et al., 2019; Bohoslav et al., 2021; Luxem et al., 2022),
69 most experiments still utilize videos taken from the top. Second, the movement of the mouse's
70 body parts is tracked frame-by-frame to estimate animal posture. Deep learning-based algorithms
71 have been widely adopted for this purpose (Mathis et al., 2018; Pereira et al., 2019). The final step
72 is to identify and classify distinct behavioral sequences or syllables. While a variety of approaches
73 have been developed over the past decade (Kabra et al., 2013; Hsu and Yttri, 2021; Luxem et al.,
74 2023), this final step is still challenging.

75 There are two broad categories of methods used to classify behavioral syllables. The first category
76 is a top-down approach, which involves pre-defining a set of rules to identify behavioral syllables or
77 applying supervised machine learning to classify them (Kabra et al., 2013; Segalin et al., 2021;
78 Harris et al., 2023). The second category is a bottom-up approach, which involves analyzing data
79 patterns using unsupervised classification algorithms (Dunn et al., 2021; Hsu and Yttri, 2021;
80 Gabriel et al., 2022; Jia et al., 2022; Luxem et al., 2022; Weinreb et al., 2023). Either approach has
81 its own advantages and disadvantages. For example, although the top-down approach provides
82 interpretable outcomes by definition, it can be laborious to prepare labeled datasets for model
83 training. Conversely, while the bottom-up approach is less unbiased, setting optimal parameters
84 and providing each syllable to an interpretable behavioral label can be a non-trivial task.

85 Here we hypothesize that the long short-term memory (LSTM) model (Hochreiter and Schmidhuber,
86 1997) can be adopted for this purpose due to the following reasons: first, behavioral sequences
87 may contain contextual, long-term dependencies (Musall et al., 2019; Issa et al., 2020). Second,
88 the LSTM model is designed to handle time series data with long-term dependencies and has
89 shown promising results in other fields (Vinyals et al., 2019; Van Houdt et al., 2020). While this
90 model has not been adopted to classify behavioral syllables, it may be a suitable method for
91 addressing the challenges associated with these sequences. Here we introduce SaLSa (a
92 combination of **semi-automatic labeling** and **LSTM-based classification**) (**Figure 1**). This is a
93 combinatorial approach to creating labeled data semi-automatically and training an LSTM classifier.
94 To demonstrate the capability of this approach, we examine how behavioral abnormalities in an
95 Alzheimer's mouse model emerge during aging.

96 Materials and Methods

97 Animals

98 All animal experiments were performed in accordance with the United Kingdom Animals (Scientific
99 Procedures) Act of 1986 Home Office regulations and approved by the University of Strathclyde
100 Animal Welfare and Ethical Review Body and the Home Office (PPL0688994). 5xFAD mice
101 (JAX006554) (Oakley et al., 2006) were bred with wild-type (WT) mice on the C57BL/6J
102 background (>F10). All genotyping was performed by Transnetyx using real-time PCR. 43 mice (13
103 5xFAD male; 11 5xFAD female; 9 WT male; 10 WT female; age range: 1.3 - 9.4 months old) were
104 used. WT mice were littermates. They had ad libitum access to food and water. The animals were
105 housed with their littermates on a 12-h light/dark cycle. All behavioral experiments were performed
106 during the first quarter of the light period (Zeitgeber time 2 - 3).

107

108 Video monitoring of exploratory behaviors

109 The behavioral arena was an acrylic transparent box (40 cm x 40 cm x 40 cm, Displaypro). Paper
110 sheets covered the outside of all side walls and landmark images (e.g., large stripes and crosses)
111 were placed on two walls. Four boxes were placed closely on white tables (Lack Side Table, IKEA).
112 A webcam (C900, NULAXY) was set over the boxes. The pixel resolution was 0.74 pixels/mm at
113 the bottom of the arena. The video was captured at 25 fps by a custom-written program (LabVIEW,
114 National Instruments). For each behavioral session, an animal was placed in the center of the
115 arena and allowed to explore it for 20 min. After the test, the arena was cleaned with 70% ethanol.
116 Only one session was held per day.

117

118 Pose estimation

119 Every video file was cropped into four videos, each containing one box. This was done with
120 custom-written MATLAB code. A DeepLabCut model (Mathis et al., 2018; Lauer et al., 2022) was
121 trained: four videos were chosen from one behavioral session. 50 frames from each video were
122 manually labeled. The labeled body parts consisted of (1) nose, (2) head, (3) left ear, (4) right ear,
123 (5) neck, (6) anterior back (back 1), (7) posterior back (back 2), (8) tail base, (9) mid-tail, and (10)
124 tail tip (**Figure 2A**). The number of training iterations was 410000. All videos were processed with
125 this trained model. Filtered data was used. Because the tail shape did not reflect behavioral
126 syllables, the mid-tail and tail tip were excluded from further analysis in this study. Thus, 8 body
127 parts were analyzed.

128

129 Feature extraction

130 All analyses were implemented in MATLAB (<https://github.com/Sakata-Lab/SaLSa>). All coordinates
131 were converted into egocentric coordinates as the neck was set as the body center and the nose-
132 neck axis was set as the body-center axis. After conversion, comprehensive features were derived.
133 Features can be categorized into spatial and temporal features.

134 The spatial features were as follows: (1) the relative coordinates across body parts, (2) the relative
135 angle of each body part relative to the body-center axis, and (3) the distance between body parts.
136 The temporal features were as follows: (1) the frame-by-frame velocity of each body part, (2) the
137 frame-by-frame velocity of the distance between body parts, (3) the spectrotemporal characteristics
138 of the relative coordinates across body parts, (4) the spectrotemporal characteristics of the relative
139 angle of each body part, (5) the spectrotemporal characteristics of the frame-by-frame velocity of
140 each body part, (6) the spectrotemporal characteristics of the distance between body parts, and (7)
141 the spectrotemporal characteristics of the frame-by-frame velocity of the distance between body
142 parts. To compute spectrotemporal characteristics, we applied a wavelet transformation (*cwt*
143 function in MATLAB). The extracted frequency components were evenly spaced in the wavelet
144 frequency domain: 0.62, 0.88, 1.25, 1.77, 2.50, 3.54, 5.01, 7.09, and 10.0 Hz. Overall, 910 features
145 were extracted.

146

147 Unsupervised processes

148 The following steps automatically identified candidates for behavioral syllables based on the
149 extracted features. This step consisted of the following sub-steps: the first step was dimensionality
150 reduction. Principal component analysis (PCA) was performed to reduce the dimension. PCs that
151 explained >85 % variance were used for uniform manifold approximation and projection (UMAP)
152 (<https://www.mathworks.com/matlabcentral/fileexchange/71902>). The parameters, *min_dist*, and
153 *n_neighbors*, were set to 0.3 and 5, respectively. The second step was clustering. In the 2D UMAP
154 space, spectral clustering was performed. The number of clusters was optimized to between 35
155 and 50 clusters using Calinski-Harabasz Criterion (Calinski and Harabasz, 1974). Because the
156 main aim was to extract clusters of certain behavioral syllables with less contamination, we
157 intentionally over-clustered the data. After clustering, each cluster contains a number of sequences
158 of candidate behavioral syllables. The final step is a post-processing. Because short sequences
159 were hard to label manually, fragmented (<0.25 s) sequences were excluded from labeling. This
160 step drastically reduced sequences and clusters. After this processing, a snippet for each cluster
161 was created for manual labeling.

162

163 Manual labeling

164 A custom-written graphical user interface (GUI) was used to label each snippet. According to the
165 initial evaluation, the following six behavioral syllables were often observed: (1) locomotion:
166 walking/running behavior toward one direction, (2) turning: turning behavior at the same position,
167 (3) rearing: rearing behavior toward a wall or at the center of the arena, (4) sniffing: rhythmic
168 sniffing behavior, (5) grooming: grooming behavior, (6) pause: behavior standing still at the same
169 position. Each snippet was labeled as one of these six syllables or miscellaneous where more than
170 two behavioral syllables were contaminated, or initially unrecognized behaviors (such as jumping)
171 were observed. Because this labeling procedure was for the training of an LSTM classifier, we took
172 a conservative approach so that snippets labeled as one of 6 syllables can be less contaminated
173 by other syllables. Out of 43 videos, 29 were manually labeled. 19 videos were used for LSTM
174 model training. 10 additional videos were used for the model performance assessment.

175

176 LSTM training and classification

177 The LSTM classifier was comprised of (1) an input layer, (2) an LSTM layer, (3) a fully connected
178 layer, (4) a Softmax layer, and (5) a classification layer. The input layer consisted of 910 units as
179 910 features were extracted from each video (see above). The LSTM layer contained 200 hidden
180 units. The classifier aimed to classify 6 behavioral syllables. For training, the ADAM optimizer was
181 used with the L_2 norm regularization method. The relationship between three parameters (the
182 gradient threshold, the maximum number of epochs, and the number of hidden units) and
183 evaluation accuracy was systematically assessed as part of optimization. In this study, the gradient
184 threshold was set to 1 and the maximum number of epochs was set to 60. For training, 80% of
185 labeled data was used whereas the remaining labeled data was used for evaluation. The model
186 with the best validation performance was used for further analysis. Once the classifier was trained,
187 all videos were processed to determine behavioral syllables frame-by-frame.

188

189 Classification performance assessment

190 To evaluate the LSTM classifier's performance, additional 10 videos were used. For each
191 behavioral syllable, a receiver operating characteristic (ROC) curve and the area under the curve
192 (AUC) were computed.

193

194 Benchmark testing

195 To compare SaLSa's performance with an existing model, keypoint-MoSeq (Weinreb et al., 2023)
196 was chosen for the following reasons: first, it overperformed other major models, such as B-SOiD
197 (Hsu and Yttri, 2021) and VAME (Luxem et al., 2022). Second, the implementation is

198 straightforward with a limited set of parameters to explore. For training of the keypoint-MoSeq
199 (kpms) model, we took 29 videos used for training and performance assessment for SaLSa since
200 manually labeled data was available. The hyperparameter, κ , was set to $9e4$, and training
201 iterations were set to 250 times. The frequency cutoff for behavioral syllables was set to 0.5%.
202 These sub-threshold syllables were merged as a single syllable for post-hoc analysis. The trained
203 model was used to analyze all these 29 videos to compare the performance of both SaLSa and
204 kpms models with each other. For more direct comparisons between the two models, kpms'
205 syllables were re-assigned based on the comparison data with the manually labeled data: each
206 original syllable was re-assigned as the most presented syllable of six pre-defined syllables (i.e.,
207 locomotion, turning, rearing, sniffing, grooming, and pause).

208 As an additional comparison, a multiclass support vector machine was trained to classify 6
209 syllables with MATLAB's *fitcecoc* function, and performance was assessed. Similar to the LSTM
210 model, 19 labeled videos were used for training whereas the remaining 10 videos were used for
211 performance assessment.

212

213 Quantification of behavioral syllables

214 To quantify the features of each behavioral syllable, three metrics were computed. First, speed
215 (cm/s) was calculated in each episode. This was done by measuring the nose travel distance
216 across frames. Second, the turning angle ($^{\circ}$ /s) was computed for each episode. This was done by
217 calculating the cumulative angle change between the nose and tail base across frames. Finally,
218 "compactness" was defined as the mean distance between two body parts. Smaller compactness
219 results from the squeezed body pose.

220

221 Statistical analysis

222 All statistical analyses were performed in MATLAB (version 2022a). In **Figure 5**, the Shapiro-Wilk
223 test was performed with a Bonferroni correction to check normality. Then, a one-way analysis of
224 variance (ANOVA) with a post-hoc Tukey HSD test was carried out. In **Figure 6**, after performing
225 the Shapiro-Wilk test, an analysis of covariance (ANCOVA) with a post-hoc Tukey HSD test was
226 carried out. A p -value less than 0.05 was considered significant. Otherwise stated, the error bars
227 represent SEM.

228

229 Code availability

230 MATLAB implementations of SaLSa are publicly available (<https://github.com/Sakata-Lab/SaLSa>).

231

232 Results

233 SaLSa (a combination of semi-automatic labeling and LSTM-based classification)

234 The general workflow was as follows (**Figure 1**) (see also Materials and Methods): (1) video
235 recording, (2) body part tracking, (3) feature extraction, (4) classifier training, and (5) classification.
236 The step of classifier training consists of two major components (**Figure 1**): first, labeled data is
237 prepared semi-automatically. This component starts with an unsupervised approach to
238 automatically identify behavioral syllable candidates. This facilitates the subsequent manual
239 labeling step. The second component is the training of an LSTM-based classifier. Using the labeled
240 data, an LSTM classifier is trained to classify sequential data. We call this integrative approach
241 “SaLSa” (a combination of semi-automatic labeling and LSTM-based classification).

242 To deploy SaLSa, we began by collecting 43 videos where mice explored an open arena for 20
243 minutes. As described below, 5xFAD mice and their littermates of both sexes (1.3 – 9.4 months
244 old) were used. To track multiple body parts (**Figure 2A**), a DeepLabCut (DLC) model was trained
245 and all videos were analyzed. The DLC model test error was 1.84 pixels. Because the tail shape
246 did not reflect behavioral syllables, the mid-tail and tail tip were excluded from further analysis in
247 this study. Thus, 8 body parts were analyzed.

248

249 After converting all body part coordinates into egocentric coordinates as the neck was defined as
250 the body center, we extracted 910 features with respect to spatial and temporal features (see
251 Materials and Methods) (**Figure 2B**). **Figure 2C** shows an example sequence of all features. There
252 were notifiable patterns across frames, implying distinct behavioral syllables.

253 To identify syllable candidates, we adopted unsupervised methods, including principal component
254 analysis (PCA), uniform manifold approximation and projection (UMAP), and spectral clustering
255 (**Figures 2D and 2E**). First, 910-dimensional data was reduced into ~15 dimensions by using PCA
256 so that >85% variance could be explained (**Figure 2D**). UMAP was adopted for further reduction to
257 2 dimensions. In this 2-dimensional UMAP space, around 40 clusters were separated using the
258 spectral clustering algorithm. Because we aimed to identify clusters, each containing frames
259 related to a certain behavioral syllable with minimum contamination of other syllables, we
260 intentionally set the number of clusters high (between 35 and 50) (**Figure 2E**).

261 After removing fragmented, short (<0.25 s) sequences, each cluster was manually labeled as one
262 of the following categories: (1) locomotion, (2) turning, (3) rearing, (4) sniffing, (5) grooming, (6)
263 pause, and (7) miscellaneous (**Figure 2E**). Out of 43 videos, 29 videos were manually labeled: 19
264 were used for classifier training whereas 10 were used for classifier evaluation in this study.

265

266 Classification Performance

267 From 19 videos, 38202 frames were labeled as one of six behavioral syllables (**Figure 3A**). While
268 most of the frames were labeled as locomotion, less than 1% of the frames were labeled as sniffing.
269 Despite this uneven distribution of labeled frames, the evaluation performance of the trained LSTM
270 classifier was 97.9%. A range of parameters (i.e., the gradient threshold, the maximum number of
271 epochs, and the number of hidden units) were explored to confirm similar evaluation performance
272 (**Extended Data Figure 3-1**).

273 We further evaluated the classifier's performance by processing additional 10 labeled videos
274 (**Figures 3B and C**). To evaluate evaluation performance, the receiver operating characteristic
275 (ROC) curve (**Figure 3B**) and the area under the curve (AUC) (**Figure 3C**) were calculated for
276 each video and each behavioral syllable. In many cases, the performance was high (>0.95 AUC).
277 As a comparison, a multiclass support vector machine model was also trained. Although the
278 evaluation performance of the trained model was 98.4%, the model failed to generalize the
279 performance to the additional 10 labeled videos (**Extended Data Figure 3-2**). Although the LSTM
280 model's performance was preserved for those videos, the classification performance for sniffing
281 was not good as other syllables (**Figures 3B and C**). Therefore, we excluded sniffing frames for
282 further analysis (**Figures 5 and 6**). Overall, the trained LSTM classifier provided reliable outcomes
283 across most behavioral syllables.

284

285 Comparisons of SaLSa's performance with a state-of-the-art model

286 We compared the performance of SaLSa with that of a state-of-the-art behavioral syllable
287 classification algorithm. Recently, keypoint-MoSeq (kpms) has been introduced and outperforms
288 other models (Weinreb et al., 2023). Therefore, the performance of a kpms model can be a
289 benchmark to evaluate SaLSa. Since we had labelled data from 29 videos, we trained a kpms
290 model with the entire 29 videos. First, we examined how each model classified the labelled data
291 across 6 syllables (**Figures 4A and B**). Consistent with the assessment above, the trained LSTM
292 model provided high classification performance except for sniffing (**Figure 4A**). In the trained kpms
293 model, 15 syllables with one sub-threshold syllable were identified (**Figure 4B**) while all sub-
294 threshold syllables were merged into one syllable. Although some behavioral syllables (locomotion,
295 turning, rearing) were classified into multiple syllables further depending on the direction of the
296 movement, other syllables (sniffing, grooming and pause) tended to be misclassified together. In
297 particular, a significant fraction of grooming behavior was identified as the sub-threshold syllable.
298 We also directly compared the outcomes from two models (**Figure 4C**). The general trend was
299 qualitatively similar to **Figure 4B**. Although locomotion, turning, and rearing were classified into
300 subclasses by the kpms model, other syllables were classified together. To make this trend clearer,
301 syllables from kpms were re-assigned to one of the 6 pre-defined syllables based on the
302 comparison with the labeled data (**Figure 4D**). Then the outcomes of SaLSa and kpms were

303 compared (**Figure 4E**). In this additional analysis, it became clear that kpmis misclassified
304 grooming behavior. Thus, although the state-of-the-art model can classify detailed behavioral
305 syllables in a fully automated fashion, SaLSa can reliably classify major behavioral syllables,
306 including grooming.

307

308 Quantification of behavioral syllables

309 Using the trained LSTM classifier, we processed all frames across all 43 videos. To examine if
310 each behavioral syllable has characteristic features, we quantified three simple metrics: speed
311 (**Figure 5A**), turning angle (**Figure 5B**), and compactness (**Figure 5C**). Speed was defined as the
312 nose speed in each behavioral syllable episode. The turning angle was calculated as the
313 cumulative turning angle of the nose-tail base axis in each episode. Compactness was defined as
314 the average pair-wise distance between body parts. We computed the median value of each metric
315 in each video and compared their averages across behavioral syllables (**Figures 5A-C**).

316 As expected, locomotion was the fastest syllable ($F_{4,210} = 730$, $p < 0.0001$, one-way ANOVA with
317 post-hoc Tukey HSD test) whereas grooming and pause were the slowest syllables even
318 compared with rearing ($p < 0.0005$) (**Figure 5A**). The metric of turning angle also provided
319 expected outcomes (**Figure 5B**): turning exhibited the largest turning angle ($F_{4,210} = 456$, $p < 0.01$,
320 one-way ANOVA with post-hoc Tukey HSD test) whereas grooming and pause showed a smaller
321 angle than rearing ($p < 0.0001$). Grooming and rearing were the most compact syllables ($F_{4,210} =$
322 75.2 , $p < 0.0001$, one-way ANOVA with post-hoc Tukey HSD test) (**Figure 5C**). Because our video
323 monitoring was top-down, rearing typically squeezes their pose in 2-D images. On the other hand,
324 locomotion was the most stretched pose compared to others ($p < 0.001$).

325 To assess the general structure of animals' exploratory behaviors in this particular experimental
326 setting, we computed the fraction and average episode duration of each syllable (**Figures 5D and**
327 **E**). Animals typically spend more time on locomotion and rearing ($F_{4,210} = 129$, $p < 0.0001$, one-
328 way ANOVA with post-hoc Tukey HSD test) and less on pause ($p < 0.0001$) (**Figure 5D**). However,
329 while the effect of behavioral syllables on the duration was significant ($F_{4,210} = 66.1$, $p < 0.0001$,
330 one-way ANOVA), the duration of grooming was comparable to pause ($p = 0.18$, post-hoc Tukey
331 HSD test) (**Figure 5E**). As expected, turning was the shortest ($p < 0.0001$). Overall, the quantities
332 of each behavioral syllable are consistent with our intuition from behavioral syllables.

333

334 Age-related and sex-specific changes in behavioral syllables of 5xFAD mice

335 To apply our approach, we examined how abnormalities in 5xFAD mice's exploratory behavior
336 emerge as they age (**Figure 6**). To this end, we simply compared the fraction of each behavioral
337 syllable as a function of age. The interaction effect between animal groups and age on locomotion

338 was significant ($F_3 = 5.7$, $p < 0.005$, ANCOVA): female 5xFAD mice exhibited significantly higher
339 hyper-locomotion compared to female controls ($p < 0.05$, post-hoc Tukey HSD test) (**Figure 6A**).
340 Consistent with this, the interaction effect on pause was also significant ($F_3 = 3.8$, $p < 0.05$,
341 ANCOVA) (**Figure 6E**). On the other hand, although the interaction effect on rearing was not
342 significant ($F_3 = 1.85$, $p = 0.155$, ANCOVA), the effect of age was significant ($F_1 = 15.4$, $p < 0.0005$),
343 meaning that the fraction of rearing decreased with age regardless animal groups (**Figure 6C**). We
344 did not see any significant interaction effects on turning and grooming ($F_3 = 2.87$, $p = 0.050$ for
345 turning; $F_3 = 1.09$, $p = 0.36$ for grooming, ANCOVA) (**Figures 6B and D**). Although hyperactivity of
346 female 5xFAD mice was well described, our approach could dissect detailed behavioral
347 abnormalities.

348 Discussion

349 In the present study, we developed SaLSa, a combination of semi-automatic labeling and LSTM-
350 based classification. The semi-automatic process facilitates the preparation of labeled data for
351 LSTM training whereas LSTM-based classification provides accurate and generalizable behavioral
352 syllable classification. Applying this approach, we found that hyper-locomotion in female 5xFAD
353 mice emerges between 4 and 8 months old whereas other active behaviors, such as rearing, are
354 not affected by genotype or sex. Given its versatility, SaLSa can classify behavioral syllables
355 without an expensive experimental setup.

356

357 Comparisons to other approaches

358 Over the last decade, a range of approaches has been developed to classify behavioral syllables
359 (Kabra et al., 2013; Pereira et al., 2020; Wiltshcko et al., 2020; Dunn et al., 2021; Hsu and Yttri,
360 2021; Segalin et al., 2021; Jia et al., 2022; Luxem et al., 2022; Harris et al., 2023; Luxem et al.,
361 2023; Weinreb et al., 2023). Since these approaches including SaLSa are applied after body parts
362 detection, they can be applied to videos taken in relatively dark environments as long as body
363 parts are detected reliably. SaLSa is the first approach to utilize LSTM for this purpose, to the best
364 of our knowledge. Deep learning-based approaches have increasingly been adopted for this
365 purpose (Marks et al., 2022; Harris et al., 2023). Although convolutional deep learning models are
366 powerful to classify and segment images, they lack an intrinsic mechanism to hold contextual
367 information. Recurrent neural networks are suitable to process time series data, including
368 behavioral tracking data (Luxem et al., 2022). One advantage of LSTM over conventional recurrent
369 networks is that it can learn long-term dependencies on the data (Hochreiter and Schmidhuber,
370 1997). Because the brain can deal with several orders of magnitudes of time depending on its
371 computational goals (Issa et al., 2020), adopting LSTMs is an extension of ongoing efforts to
372 characterize natural behaviors comprehensively. Another advantage of LSTM is its generalizability
373 like many other deep learning approaches. Once it has been trained, newly acquired videos can be
374 processed by simply extracting features as we demonstrated.

375 On the other hand, LSTM requires a large amount of labeled data for training. To mitigate this
376 issue, we adopted semi-automatic labeling. By extracting features unbiasedly, behavioral syllable
377 candidates were automatically identified. Compared to an approach that creates snippets randomly
378 (such as MuViLab), our approach reduces manual curation time. In the present study, a 20-minute
379 video took several minutes. This allows labeling a number of videos easily to prepare training data.
380 Thus, our approach is unique in the sense that it combines both unsupervised methods and a
381 supervised LSTM model to classify behavioral syllables.

382 In a broader context, SaLSa takes a top-down approach where pre-defined behavioral syllables
383 are identified semi-automatically and classified by a deep learning model. In the future, it may be
384 interesting to integrate a fully automated bottom-up model with an LSTM classifier.

385 Behavioral abnormalities in 5xFAD

386 As an application of SaLSa, we investigated age- and sex-specific changes in the exploratory
387 behaviors of 5xFAD mice. It has been well documented that female 5xFAD mice exhibit
388 hyperactivity (Oblak et al., 2021). Our approach demonstrates that hyperactivity consists of hyper-
389 locomotion whereas other active behaviors, such as turning and rearing, are similar across animal
390 groups. In particular, rearing decreases with age regardless of genotype and sex, which has not
391 been documented in this mouse model before.

392 The underlying mechanisms of sex-specific hyper-locomotion in 5xFAD mice are unknown. In this
393 mouse model, amyloid plaques can be seen in the hippocampus (primarily the subiculum) as early
394 as 2 months old and the pathology appears across brain regions as they age (Oakley et al., 2006;
395 Oblak et al., 2021). Sex differences in amyloid pathology can also be apparent in the cortical
396 subplate even at 3 months old and seen across multiple brain regions at 4 months old (Oblak et al.,
397 2021). Consistent with this sex-specific pathological progression, a transcriptomic analysis also
398 revealed sex differences in a wide range of molecular pathways (Oblak et al., 2021). In the future, it
399 would be crucial to determine how these sex-specific pathological features link to dysfunctions in
400 neural circuit activity, which lead to age-related, sex-specific hyper-locomotion. Despite this
401 challenge, given the simplicity of our experimental setup, similar approaches can be applied to
402 other animal models.

403

404 Limitations of the study

405 Our study has at least four limitations: first, our pre-set behavioral syllables were limited. We are
406 also aware that some animals jumped or exhibited complex behaviors, such as mixed turning and
407 rearing behaviors. This will require more labeled data or additional post-hoc analysis. For example,
408 based on the predicted score from a classifier, such complex behavioral syllables may be defined.
409 Additionally, sniffing behaviors could not be identified and classified accurately. This could be
410 partly because of the limited resolution of our camera and the accuracy of body-part tracking.
411 Increasing the resolution and adding extra body parts to track may improve this aspect.

412 Second, the model must be re-trained when videos are taken at different frame rates or from a
413 different experimental setup. Because our temporal features assume a certain frame rate (25 fps),
414 a change in frame rate leads to a change in the number of features. To deal with this issue, re-
415 sampling can be considered.

416 Third, while the present study replicated the results (i.e., sex-specific hyperactivity) of the recent
417 comprehensive analysis in 5xFAD mice with the detailed classification of behavioral syllables, the

418 estrous cycle might be a potential confounding factor even though a recent study could not find this
419 to be the case (Levy et al., 2023). Increasing the sampling size by monitoring the estrous cycle will
420 address this issue in the future.

421 Finally, as widely appreciated, deep learning models are challenging to interpret. In this study, it is
422 not straightforward to determine what spatiotemporal features contribute to classification. Several
423 approaches may be considered, such as Local Interpretable Model-Agnostic Explanations (LIME)
424 and visualization.

425

426 **Conclusions**

427 Behavioral syllable classification is important. In the present study, we developed a combinatory
428 approach to semi-automatic labeling and LSTM-based classification, called SaLSa. Our approach
429 can assist manual curation to prepare labeled data semi-automatically. The LSTM classifier reliably
430 classified behavioral syllables in new datasets, which were not used for training. It also provides
431 comparable performance with the state-of-the-art model. Thus, our approach adds a versatile tool
432 for behavioral syllable classification. Combining other advanced technologies, SaLSa facilitates the
433 effort to better understand the neural basis of complex behavior.

434 **References**

- 435 Bohoslav JP, Wimalasena NK, Clausing KJ, Dai YY, Yarmolinsky DA, Cruz T, Kashlan AD, Chiappe ME, Orefice
436 LL, Woolf CJ, Harvey CD (2021) DeepEthogram, a machine learning pipeline for
437 supervised behavior classification from raw pixels. *Elife* 10.
- 438 Calinski T, Harabasz J (1974) A dendrite method for cluster analysis. *Communications in Statistics*
439 3:1-27.
- 440 Dunn TW, Marshall JD, Severson KS, Aldarondo DE, Hildebrand DGC, Chettih SN, Wang WL,
441 Gellis AJ, Carlson DE, Aronov D, Freiwald WA, Wang F, Olveczky BP (2021) Geometric
442 deep learning enables 3D kinematic profiling across species and environments. *Nat*
443 *Methods* 18:564-573.
- 444 Gabriel CJ, Zeidler Z, Jin B, Guo C, Goodpaster CM, Kashay AQ, Wu A, Delaney M, Cheung J,
445 DiFazio LE, Sharpe MJ, Aharoni D, Wilke SA, DeNardo LA (2022) BehaviorDEPOT is a
446 simple, flexible tool for automated behavioral detection based on markerless pose tracking.
447 *Elife* 11.
- 448 Harris C, Finn KR, Kieseler ML, Maechler MR, Tse PU (2023) DeepAction: a MATLAB toolbox for
449 automated classification of animal behavior in video. *Sci Rep* 13:2688.
- 450 Hochreiter S, Schmidhuber J (1997) Long short-term memory. *Neural Comput* 9:1735-1780.
- 451 Hsu AI, Yttri EA (2021) B-SOiD, an open-source unsupervised algorithm for identification and fast
452 prediction of behaviors. *Nat Commun* 12:5188.
- 453 Issa JB, Tocker G, Hasselmo ME, Heys JG, Dombeck DA (2020) Navigating Through Time: A
454 Spatial Navigation Perspective on How the Brain May Encode Time. *Annu Rev Neurosci*
455 43:73-93.
- 456 Jia Y, Li S, Guo X, Lei B, Hu J, Xu XH, Zhang W (2022) Selfee, self-supervised features extraction
457 of animal behaviors. *Elife* 11.
- 458 Kabra M, Robie AA, Rivera-Alba M, Branson S, Branson K (2013) JAABA: interactive machine
459 learning for automatic annotation of animal behavior. *Nat Methods* 10:64-67.
- 460 Lauer J, Zhou M, Ye S, Menegas W, Schneider S, Nath T, Rahman MM, Di Santo V, Soberanes D,
461 Feng G, Murthy VN, Lauder G, Dulac C, Mathis MW, Mathis A (2022) Multi-animal pose
462 estimation, identification and tracking with DeepLabCut. *Nat Methods* 19:496-504.
- 463 Levy DR, Hunter N, Lin S, Robinson EM, Gillis W, Conlin EB, Anyoha R, Shansky RM, Datta SR
464 (2023) Mouse spontaneous behavior reflects individual variation rather than estrous state.
465 *Curr Biol* 33:1358-1364 e1354.
- 466 Luxem K, Mocellin P, Fuhrmann F, Kursch J, Miller SR, Palop JJ, Remy S, Bauer P (2022)
467 Identifying behavioral structure from deep variational embeddings of animal motion.
468 *Commun Biol* 5:1267.
- 469 Luxem K, Sun JJ, Bradley SP, Krishnan K, Yttri E, Zimmermann J, Pereira TD, Laubach M (2023)
470 Open-source tools for behavioral video analysis: Setup, methods, and best practices. *eLife*
471 12:e79305.
- 472 Marks M, Qiuhan J, Sturman O, von Ziegler L, Kollmorgen S, von der Behrens W, Mante V,
473 Bohacek J, Yanik MF (2022) Deep-learning based identification, tracking, pose estimation,
474 and behavior classification of interacting primates and mice in complex environments. *Nat*
475 *Mach Intell* 4:331-340.
- 476 Mathis A, Schneider S, Lauer J, Mathis MW (2020) A Primer on Motion Capture with Deep
477 Learning: Principles, Pitfalls, and Perspectives. *Neuron* 108:44-65.
- 478 Mathis A, Mamidanna P, Cury KM, Abe T, Murthy VN, Mathis MW, Bethge M (2018) DeepLabCut:
479 markerless pose estimation of user-defined body parts with deep learning. *Nat Neurosci*
480 21:1281-1289.
- 481 Musall S, Kaufman MT, Juavinett AL, Gluf S, Churchland AK (2019) Single-trial neural dynamics
482 are dominated by richly varied movements. *Nat Neurosci* 22:1677-1686.
- 483 Oakley H, Cole SL, Logan S, Maus E, Shao P, Craft J, Guillozet-Bongaarts A, Ohno M, Disterhoft J,
484 Van Eldik L, Berry R, Vassar R (2006) Intraneuronal beta-amyloid aggregates,
485 neurodegeneration, and neuron loss in transgenic mice with five familial Alzheimer's
486 disease mutations: potential factors in amyloid plaque formation. *J Neurosci* 26:10129-
487 10140.
- 488 Oblak AL et al. (2021) Comprehensive Evaluation of the 5XFAD Mouse Model for Preclinical
489 Testing Applications: A MODEL-AD Study. *Front Aging Neurosci* 13:713726.

- 490 Pereira TD, Shaevitz JW, Murthy M (2020) Quantifying behavior to understand the brain. *Nat*
491 *Neurosci* 23:1537-1549.
- 492 Pereira TD, Aldarondo DE, Willmore L, Kislin M, Wang SS, Murthy M, Shaevitz JW (2019) Fast
493 animal pose estimation using deep neural networks. *Nat Methods* 16:117-125.
- 494 Schneider A, Zimmermann C, Alyahyay M, Steenbergen F, Brox T, Diester I (2022) 3D pose
495 estimation enables virtual head fixation in freely moving rats. *Neuron* 110:2080-2093 e2010.
- 496 Segalin C, Williams J, Karigo T, Hui M, Zelikowsky M, Sun JJ, Perona P, Anderson DJ, Kennedy A
497 (2021) The Mouse Action Recognition System (MARS) software pipeline for automated
498 analysis of social behaviors in mice. *Elife* 10.
- 499 Van Houdt G, Mosquera C, Nápoles G (2020) A review on the long short-term memory model.
500 *Artificial Intelligence Review* 53:5929-5955.
- 501 Vinyals O et al. (2019) Grandmaster level in StarCraft II using multi-agent reinforcement learning.
502 *Nature* 575:350-354.
- 503 Weinreb C, Osman MAM, Zhang L, Lin S, Pearl J, Annapragada S, Conlin E, Gillis WF, Jay M,
504 Shaokai Y, Mathis A, Mathis MW, Pereira T, Linderman SW, Datta SR (2023) Keypoint-
505 MoSeq: parsing behavior by linking point tracking to pose dynamics.
506 [bioRxiv:2023.2003.2016.532307](https://doi.org/10.1101/2023.2003.2016.532307).
- 507 Wiltschko AB, Tsukahara T, Zeine A, Anyoha R, Gillis WF, Markowitz JE, Peterson RE, Katon J,
508 Johnson MJ, Datta SR (2020) Revealing the structure of pharmacobehavioral space
509 through motion sequencing. *Nat Neurosci* 23:1433-1443.

510

511 **Legends**

512 **Figure 1. SaLSa.** After recording videos and processing them with DeepLabCut (“pose
513 estimation”), spatial and temporal features are extracted from the egocentric coordinates of tracked
514 body parts (feature extraction). Based on a set of videos, an LSTM classifier is trained (model
515 training). This component consists of two parts: first, through fully automated unsupervised
516 processes including dimensionality reduction and clustering, behavioral syllable candidates are
517 identified. Using a graphical user interface, the identified candidates are manually labeled (semi-
518 automatic labeling). For the training and evaluation of an LSTM classifier, labeled data will be used.
519 Once the classifier is trained, the entire sequence of extracted features is processed to classify
520 behavioral syllables.

521

522 **Figure 2. Semi-automatic labeling.** (A) A frame with tracked body parts. Although 10 body parts
523 were tracked, the mid-tail and tail tip were excluded for quantitative analyses. (B) The categories of
524 extracted features. (C) A chunk of normalized feature values. Each feature was Z-scored for a
525 visualization purpose. Dotted lines separate different feature categories. The feature category
526 number corresponds to the number indicated in (B). (D) Cumulative distribution of explained
527 variance across principal components (PCs). The threshold for including PCs was set at 85%
528 variance explained. (E) UMAP representation of the entire video sequences and example frame
529 sequences of labeled behavioral syllables. Data reduced by principal component analysis (PCA)
530 was further reduced to 2 dimensions by UMAP. A spectral clustering algorithm was applied. By
531 removing fragmented (<0.25 s) sequences, each cluster was manually annotated. The example
532 sequences indicate a down-sampled sequence of each labeled behavioral syllable. Please note
533 that although the mid-tail and tail tip were shown in sample frames, because the tail shape did not
534 reflect behavioral syllables, the mid-tail and tail tip were excluded from all quantitative analyses
535 including the UMAP analysis.

536

537 **Figure 3. LSTM training data and performance.** (A) The proportion of each behavioral syllable
538 for LSTM training. (B) Receiver operating characteristic curves for each behavioral syllable based
539 on independently labeled 10 videos. FA, false alarm. (C) The area under the curve (AUC) values
540 across behavioral syllables. In **Extended Data Figure 3-1**, a range of parameters in the LSTM
541 model were systematically assessed. The evaluation performance of a multiclass support vector
542 machine is shown in **Extended Data Figure 3-2**.

543

544 **Extended Data Figure 3-1. Systematic comparison of three parameters for evaluation**
545 **accuracy and training duration.** (A and B) The evaluation accuracy of LSTM models with

546 variable maximum numbers of epochs and hidden units. The gradient threshold was set at 1 in (A)
547 and 2 in (B). (C and D) Training duration across conditions.

548

549 **Extended Data Figure 3-2. Performance of a multiclass support vector machine.** (A) Receiver
550 operating characteristic curves for each behavioral syllable based on independently labeled 10
551 videos. FA, false alarm. (B) The area under the curve (AUC) values across behavioral syllables.

552

553 **Figure 4. Benchmark testing of SaLSa.** (A) Confusion matrix between manually labelled data (y-
554 axis) and SaLSa outputs (x-axis). The values indicate what fraction of manually labelled data was
555 classified as each syllable by SaLSa. (B) Confusion matrix between manually labelled data and
556 keypoint-MoSeq (kpms) outputs. Syllable 16 is a sub-threshold (0.5% cutoff) class. (C) Confusion
557 matrix between outputs from SaLSa and kpms. The values indicate what fraction of frames
558 classified by SaLSa were classified by kpms. (D) Confusion matrix between manually labelled data
559 and re-assigned kpms syllables. The original kpms syllables were re-assigned to one of 6 syllables
560 based on comparison to manually labelled data (B). (E) Confusion matrix between SaLSa output
561 and the re-assigned kpms syllables.

562

563 **Figure 5. Quantification of behavioral syllables.** (A) The median speed per episode of each
564 video across behavioral syllables. (B) The median turning angle per episode of each video across
565 behavioral syllables. (C) The median compactness per episode of each video across behavioral
566 syllables. (D) The fraction of each behavioral syllable. (E) The average episode duration of each
567 video across behavioral syllables. *inset*, *p*-values of post-hoc pair-wise comparisons ($n = 43$, two-
568 way ANOVA with post-hoc Tukey HSD test). L, locomotion; T, turning; R, rearing; G, grooming; P,
569 pause.

570

571 **Figure 6. Age-dependent effects of genotype and sex on behavioral contents in 5xFAD mice.**
572 The fraction of each behavioral syllable as a function of age. Behavioral syllables are (A)
573 locomotion, (B) turning, (C) rearing, (D) grooming, and (E) pause. *P*-values of ANCOVA are shown.

574











