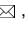1  **A machine learning toolbox for the analysis of sharp-wave ripples reveal**
2  **common features across species**

3  Andrea Navas-Olive[1]*✉, Adrian Rubio[1]*, Saman Abbaspoor[2], Kari L. Hoffman[2,3]
4  and Liset M de la Prida[1,✉,#]

5  [1]Instituto Cajal, CSIC, Madrid 28002. Spain

6  [2]Psychological Sciences, Vanderbilt Brain Institute, Vanderbilt University, USA

7  [3]Biomedical Engineering, Vanderbilt University, USA

8

9

10  * These authors contributed equally

11  ✉Correspondence to: acnavasolive@gmail.com and lmprida@cajal.csic.es

12  [#] Lead author: lmprida@cajal.csic.es

13

14  **Abstract**

15  The study of sharp-wave ripples (SWRs) has advanced our understanding of memory function, and
16  their alteration in neurological conditions such as epilepsy and Alzheimer's disease is considered a
17  biomarker of dysfunction. SWRs exhibit diverse waveforms and properties that cannot be fully
18  characterized by spectral methods alone. Here, we describe a toolbox of machine learning (ML)
19  models for automatic detection and analysis of SWRs. The ML architectures, which resulted from a
20  crowdsourced hackathon, are able to capture a wealth of SWR features recorded in the dorsal
21  hippocampus of mice. When applied to data from the macaque hippocampus, these models were
22  able to generalize detection and revealed shared SWR properties across species. We hereby
23  provide a user-friendly open-source toolbox for model use and extension, which can help to
24  accelerate and standardize SWR research, lowering the threshold for its adoption in biomedical
25  applications.

26

27  **Keywords:** ripples; neural networks; convolutional neural networks; hippocampus; monkey

28

29

30

31

**Introduction**

The study of brain rhythms has bolstered our understanding of the neural basis of cognition. Because these signals emerge from the coordinated activity of multiple neurons, they can be used as biomarkers of the underlying cognitive process[1]. For example, hippocampal sharp-wave ripples (SWRs) represent the most synchronous pattern in the mammalian brain, and are widely considered to contribute to the consolidation of memories[2]. SWRs consist of brief high-frequency oscillations or 'ripples' (100-250Hz), which can be detected around the hippocampal CA1 cell layer during rest or sleep. An avalanche of excitatory inputs from the CA3 region, typically visible as a slower sharp-wave component, triggers ripples locally in CA1[3,4]. Within the ripple event, neural firing patterns that occurred during exploratory behavior are reactivated outside of the experience[5,6], leading the SWR to be used as an index of consolidation-associated reactivation or replay[7–10].

Although SWRs can be detected across an array of recording methods, subfield locations, and species[2,11] their underlying mechanisms and consequent local field potential (LFP) features are understood almost exclusively from measurements in rat and mouse dorsal hippocampal CA1. Even within this region, SWRs exhibit a large diversity of waveforms that presumably reflect the myriad combinations of reactivating ensembles[12–14]. Using spectral methods their characteristics are shown to vary along the long (septotemporal) CA1 axis within animals[15] and most notably with phylogenetic distance across species e.g. when measured in the human versus non-human primates[11,16,17]. Furthermore, in diseases affecting hippocampal function, such as in Temporal Lobe Epilepsy, pathological forms of ripples have been reported [18–21], as well as along aging [22,23]. However, spectral properties alone are suboptimal to separate these events from other types of faster oscillations [24–26].

To address this challenge, many researchers have developed feature-based strategies for detecting LFP oscillations using machine learning (ML) tools[16,27–32]. These novel strategies have accelerated our understanding the underlying mechanisms of SWRs, and the improvement of closed-loop interventions beyond those using spectral features alone[31,33]. Yet these methods have been focused on a single detection method optimized for a single target application, typically either in mouse dorsal CA1 or within lab-specific approaches to detection in brains of humans with epilepsy. As LFP recordings are increasingly common in the clinic, the need to scale analysis from small laboratory animals to the human brain is pressing[10,34–39]. Developing these new tools will provide the community with straightforward methods to identify SWRs from pathological oscillations across the range of recording technologies, sampled regions, and background pathologies. Therefore, there is a broad demand for a consolidated toolbox of ML methods for LFP feature analysis that can be easily applied across species, to aid in understanding of brain function, but also advance biomedical applications.

Here, we develop and analyze a set of ML architectures applied to the problem of SWR identification, and compiled in an open toolbox: https://github.com/PridaLab/rippl-AI. To favor an unbiased screening of potential ML solutions, we ran a hackathon with people from very disparate fields with the mission of detecting SWR using algorithms in a supervised manner. Using community-based solutions in neuroscience is gaining traction due to their ability to foster interdisciplinary and diverse perspectives, and to promote collaboration and data sharing[40–43]. We selected the most promising architectures from the hackathon and standardized them for fair comparisons. We show how the different ML models could bias SWR detection and identify conditions for their optimal performance and stability in the mouose hippocampus (*Mus musculus*). We then extend the analysis to SWRs recorded in the macaque hippocampus (*Macaca mulatta*), to demonstrate the generalizability of SWRs detection methods to the primate order. This proof of principle will foster the development of feature-based detection algorithms for future applications to a range of models and approaches, including the human brain.
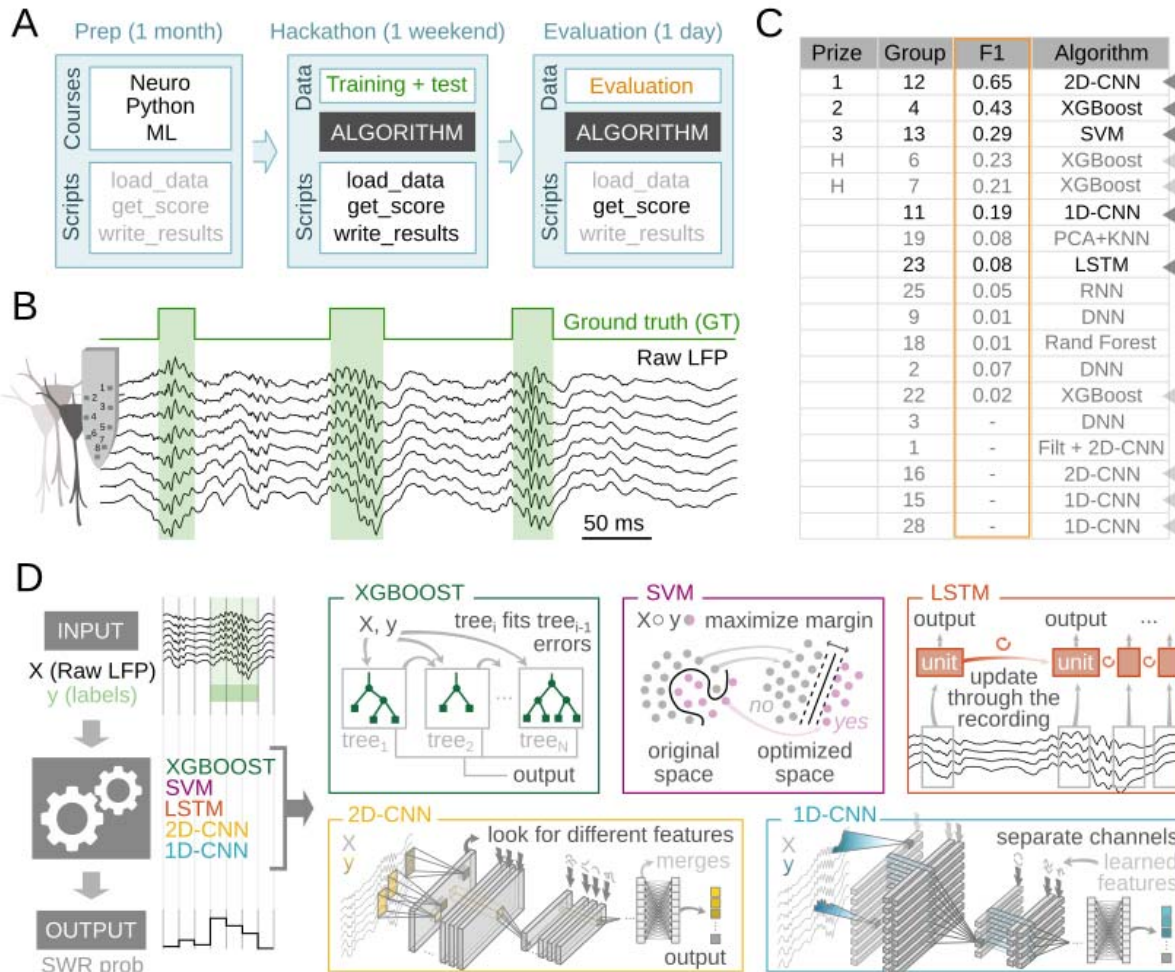

**Results**

**Community-based proposal of ML models of SWR**

To create a diversity of ML supervised models of SWRs, we organized a hackathon that promoted unbiased community-based solutions from scientists unfamiliar with neuroscience research, and SWRs in particular (see Methods). The hackathon challenge was to propose a ML model that

86 successfully identifies SWR in a dataset of high-density LFP recordings from the CA1 dorsal
87 hippocampus of mice, used before for similar purposes[31]. Preparatory courses introduced
88 participants into the main topics required for the challenge (Fig.1A). To standardize the different ML
89 models, they were given access to Python functions for loading the data, to evaluate model
90 performance, and to write results in a common format. Annotated data consisted of raw LFP
91 signals (8-channels) sampled at 30 kHz, and containing SWR events manually tagged by an expert
92 (training set: 1794 events, two sessions from 2 mice; test set: 1275 events; two sessions from 2
93 mice; Fig.1B).

94



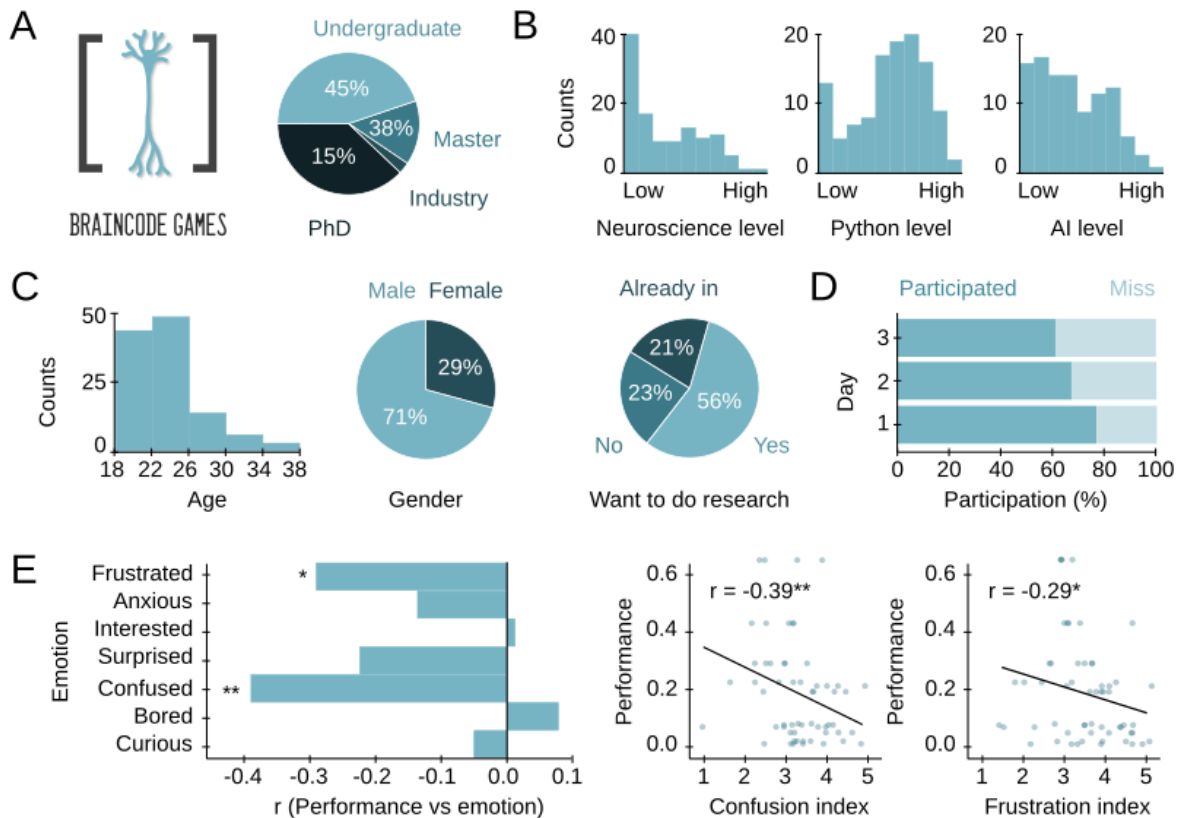| Prize | Group | F1 | Algorithm |
|---|---|---|---|
| 1 | 12 | 0.65 | 2D-CNN |
| 2 | 4 | 0.43 | XGBoost |
| 3 | 13 | 0.29 | SVM |
| H | 6 | 0.23 | XGBoost |
| H | 7 | 0.21 | XGBoost |
|  | 11 | 0.19 | 1D-CNN |
|  | 19 | 0.08 | PCA+KNN |
|  | 23 | 0.08 | LSTM |
|  | 25 | 0.05 | RNN |
|  | 9 | 0.01 | DNN |
|  | 18 | 0.01 | Rand Forest |
|  | 2 | 0.07 | DNN |
|  | 22 | 0.02 | XGBoost |
|  | 3 | - | DNN |
|  | 1 | - | Filt + 2D-CNN |
|  | 16 | - | 2D-CNN |
|  | 15 | - | 1D-CNN |
|  | 28 | - | 1D-CNN |

95
96 **Figure 1: Unbiased community-based proposals of ML models of SWR. A,** Organization of the
97 hackathon. A preparatory phase (Prep) established the basic grounds of the challenge in terms of minimal
98 knowledge about SWR, Python programming and Machine Learning (ML) models. It also looked to
99 standardize scripts and data management. The second phase consisted on the hackathon, which lasted over
100 53h during three days, with participants having access to the annotated training dataset and some Python
101 scripts. During the last evaluation phase, a new test set was released to participants 3 hours before the end
102 of the hackathon. Solutions were ranked using the F1-score (see methods). **B,** Example of the training data
103 consisting on 8 channels of raw LFP (black) sampled at 30 kHz, with the manually tagged ground truth (GT),
104 corresponding to SWR events. **C,** Results from the hackathon. Solutions were ranked by the F1-score. F1
105 represents the harmonic mean between Precision (percentage of good detections) and Recall (percentage of
106 detected GT events). Deep Neural Networks (DNN), Convolutional Neural Networks (CNN), Recurrent
107 Neural Networks (RNN) with/without Long-Short Term Memory (LSTM); Random Forest decision trees (Rand
108 Forest), Extreme Gradient Boosting (XGBoost), Support Vector Machines (SVM), k-Nearest Neighbors (kNN).
109 Chosen solutions are marked with arrowheads. Darker arrows point to the group that got the highest score of
110 each particular architecture; light arrows point repeated architectures. **D**, Schematic representation of the
111 SWR detection strategy and the 5 ML models used in this work.

112
113

114    Participants submitted eighteen different solutions (Fig.1C). The most used architecture was the
115    Extreme Gradient Boosting (XGBoost; 4 proposals), a decision tree-based algorithm very popular
116    for its balance between flexibility, accuracy and speed[44] (Fig.1C). Some other popular architectures
117    were one and two-dimensional Convolutional Neural Networks (1D-CNN, 2D-CNN; 3 and 3
118    solutions, respectively), Deep Neural Networks (DNN, 3 solutions)[45], and Recurrent Neural
119    Networks (RNN; 2 solutions)[45] (Fig.1C). RNN were presented in both their standard feed-forward
120    version, and as the Long-Short Term Memory (LSTM) version that includes feedback connections,
121    more suited for processing time series data[46].
122
123    Although all these architectures are neural networks typically used for pattern recognition, the way
124    they process and learn from data is remarkably different. For example, whereas CNNs are based
125    on kernels specialized in spotting particular spatially contiguous features of the input, LSTMs use
126    memory cells that look for time-dependent relationships in the data. Two other algorithms were
127    also submitted: a Support Vector Machine (SVM; 1 solution; Fig.1C) and a clustering-based
128    solution based on dimensionality reduction by Principal Component Analysis (PCA), followed by k-
129    Nearest Neighbors (kNN) clustering (1 solution; Fig.1C). From the 18 solutions submitted, 5 were
130    not functional and could not be scored (Fig.1C, bottom). Analysis of the hackathon experience in
131    relationship to the submitted solutions are summarized in Fig.S1 (see methods for details).



132
133    **Figure S1. Information about the hackathon. A,** A hackathon was organized to seek for community-based
134    solutions to the SWR challenge from people unfamiliar to SWR neurophysiology. Among the 116 participants,
135    there were undergraduate students (45%), Master students (38%), PhD students (15%), and industry
136    workers (3%). **B,** There was a general lack of neuroscience knowledge, although most participants declared
137    a high-level performance in Python. Most groups integrated people with programming abilities and basic ML
138    knowledge. **C,** Participant age (left), gender (middle; 71% male, 29% female participants) and involvement in
139    research (right; 21% already in research; 56% interested in doing basic research; 23% not motivated for
140    basic research activities). **D,** Self-reported participation rate during the three days of the hackathon. **E,**
141    Correlation between the performance metric of the proposed solution and emotional states of participants as
142    quantified from their responses to surveys recorded during the hackathon (Spearman rank-order correlation *,
143    $p<0.05$; **, $p<0.01$). Only performance of functional solutions were used. See Methods for details.
144
145

146  We sought to identify the more promising architectures for a subsequent in depth analysis.
147  Performance of submitted ML models was measured using the F1-score (see next section). The
148  best performances were achieved by the 2D-CNN, one of the XGBoost models, and the SVM
149  algorithm. Since 1D-CNNs and RNNs were submitted by several groups, and given their previously
150  successful application to SWR detection[28,31], we decided to include them as well, resulting in five
151  different machine learning architectures (Fig.1C; dark arrowheads).

152  The goal of the ML models is to identify the presence of a SWR (or part of it) in a given analysis
153  window (Fig.1D, left). The selected ML architectures covered a range of processing strategies
154  (Fig.1D, right). XGBoost is a very popular ML algorithm that uses many decision trees in a parallel
155  fashion, making it one of the fastest algorithms[47]. SVM regression lays within the statistical learning
156  framework, and its objective is to find a new space where samples from different categories (SWRs
157  vs no-SWRs) are maximally separated, making it one of the most robust classification methods[48].
158  LSTMs are especially suited for regression and classification of temporal series like in natural
159  language processing, using a memory-based strategy to extract relationships between non-
160  continuous time points[46]. CNNs represent a very common approach for many detection and
161  classification tasks applied to different data modalities (1D for signals, 2D for images and 3D for
162  video or volumetric reconstructions), and can approach human performance on many tasks[49].
163  While 2D-CNNs process input data by considering adjacency on both dimensions (spatial and
164  temporal, in our case), the 1D-CNN solution treats each channel independently and only considers
165  time adjacency, making them two distinct processing algorithms.
166
167  This community-based ML architecture bank that was produced by participants who were
168  unfamiliar with SWR studies can be used to evaluate the problem of SWR automatic detection in
169  experimental contexts. We next focused on standardizing processing and retraining the different
170  models.
171
172
173  **Standardization and retraining of selected algorithms**

174  After careful examination of the submitted solutions, we noticed that data pre-processing and
175  training strategies were very different between groups. Data characteristics, like the sampling
176  frequency or the number of channels used for detection can influence operation. To standardize
177  analysis, we chose to down sample to 1250 Hz, and normalize input data using z-scores, which
178  account for differences in mean values and standard deviation across experimental sessions.

179  We then retrained the submitted ML architectures using the same training set of the hackathon. We
180  randomly divided the dataset into a training set (70%), and a test set (30%) to evaluate their
181  performance in unseen data prior to a more thorough validation (Fig.2A). We explored a wide
182  range of hyper-parameters for each architecture, which included the number of LFP channels (1, 3
183  or 8), the size of the analysis window (from 6.4 up to 50 ms) and model-specific parameters like
184  "maximum tree depth" for XGBoost, "bidirectionality" for LSTM or "kernel factor" for CNNs (Fig.2A).
185  A trained ML architecture set with a particular combination of its hyper-parameters gives rise to a
186  particular "trained model" (Fig.2A). Because each architecture had different numbers of hyper-
187  parameters, we ended up with different numbers of trained models for each architecture (1944 for
188  XGBoost, 72 for SVM, 2160 for LSTM, 60 for 2D-CNN, and 576 for 1D-CNN). We then used the
189  test set to choose the 50-best models from each architecture, and further tested their performance
190  using a new validation dataset (7586 SWR events; 21 sessions from 8 mice), previously used for
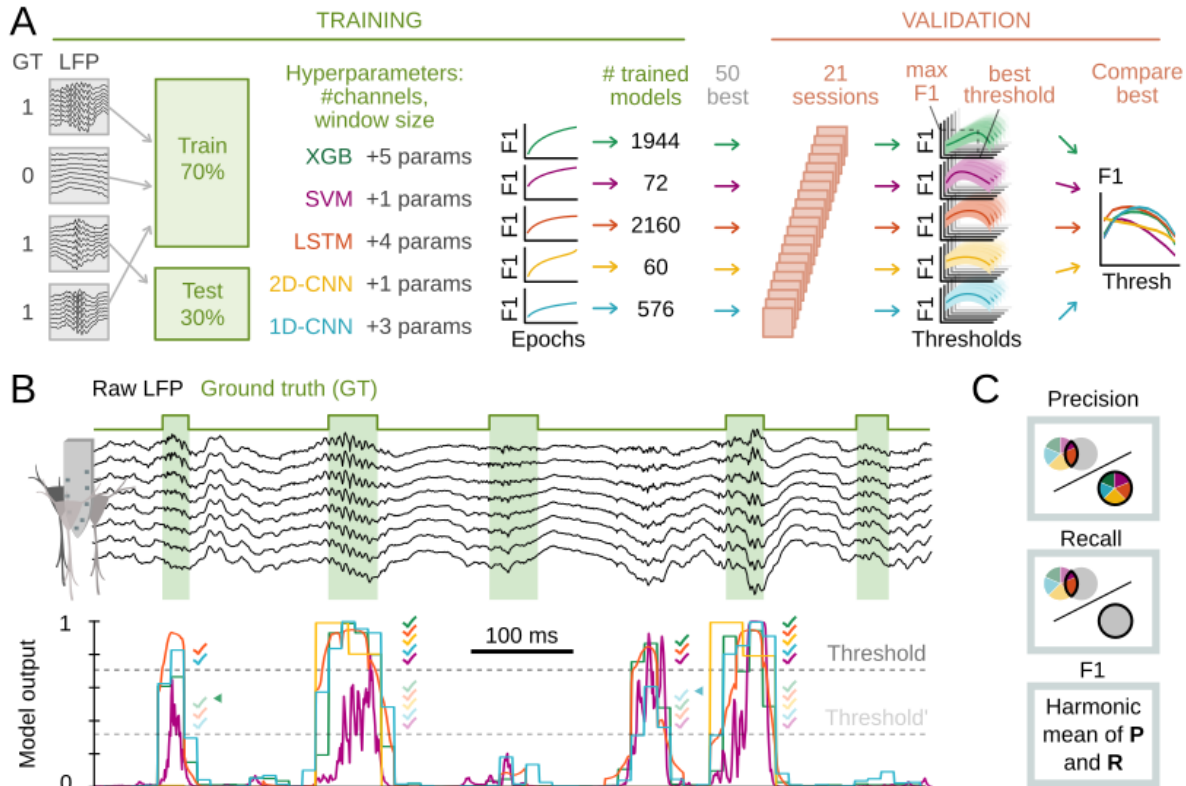191  the 1D-CNN model[31] (Fig.2A, right).
192
193  The goal of training is to make the model output as similar as possible to the ground truth (GT).
194  Because model outputs are continuous numbers between 0 and 1 representing the probability of
195  the presence of the event in the window of analysis, choosing the detection threshold can affect
196  performance (Fig.2B). Lower thresholds would result in more detections (Fig.2B, light-gray
197  discontinuous threshold line), normally implying a larger number of both True and False Positives,
198  while higher thresholds are more conservative at the expenses of False Negatives (Fig.2B, dark-
199  gray threshold line). An ideal model would perform well regardless the threshold, but in practice
200  selecting the threshold that optimizes the True Positive-False Positive trade-off is unavoidable but
201  crucial for experiments. A performance score that takes into account this trade-off is the F1-score,
202  computed as the harmonic mean between Precision (percentage of good detections) and Recall

203  (percentage of detected GT events) (Fig.2C). F1 values of 1.0 would reflect a perfect match
204  between detections and GT, whereas 0.0 reflects a perfect mismatch. Note this was the same
205  score used to rank models in the hackathon.
206
207  After training all architectures by optimizing F1-scores over the test set, we assessed
208  generalization and performance using the validation dataset. We inspected what parametric
209  combinations gave rise to optimal ML models, and found a remarkable variety of distributions
210  (Fig.S2A). All architectures showed a great deal of variability, with almost all available parameter
211  combinations covered. However, some parameters showed biases that depended on the ML
212  architectures, pointing to the necessary requirements for a good performance. For example, all of
213  the 50-best XGBoost models used 8-channels, and in general, more than 1-channel was used
214  across successful architectures (Fig.S2A). Furthermore, different architectures had distinct ranges
215  of parameter values. XGBoost models required longer time windows (25 ms), whereas most SVM
216  models employed shorter windows (<3.2 ms). LSTM, 2D-CNN and 1D-CNNs with variable window
217  sizes all showed very strong performance for >12.8 ms. Finally, LSTM models used both uni- and
218  bi-directional input flow, whereas all of the best models resorted to bidirectionality, suggesting that
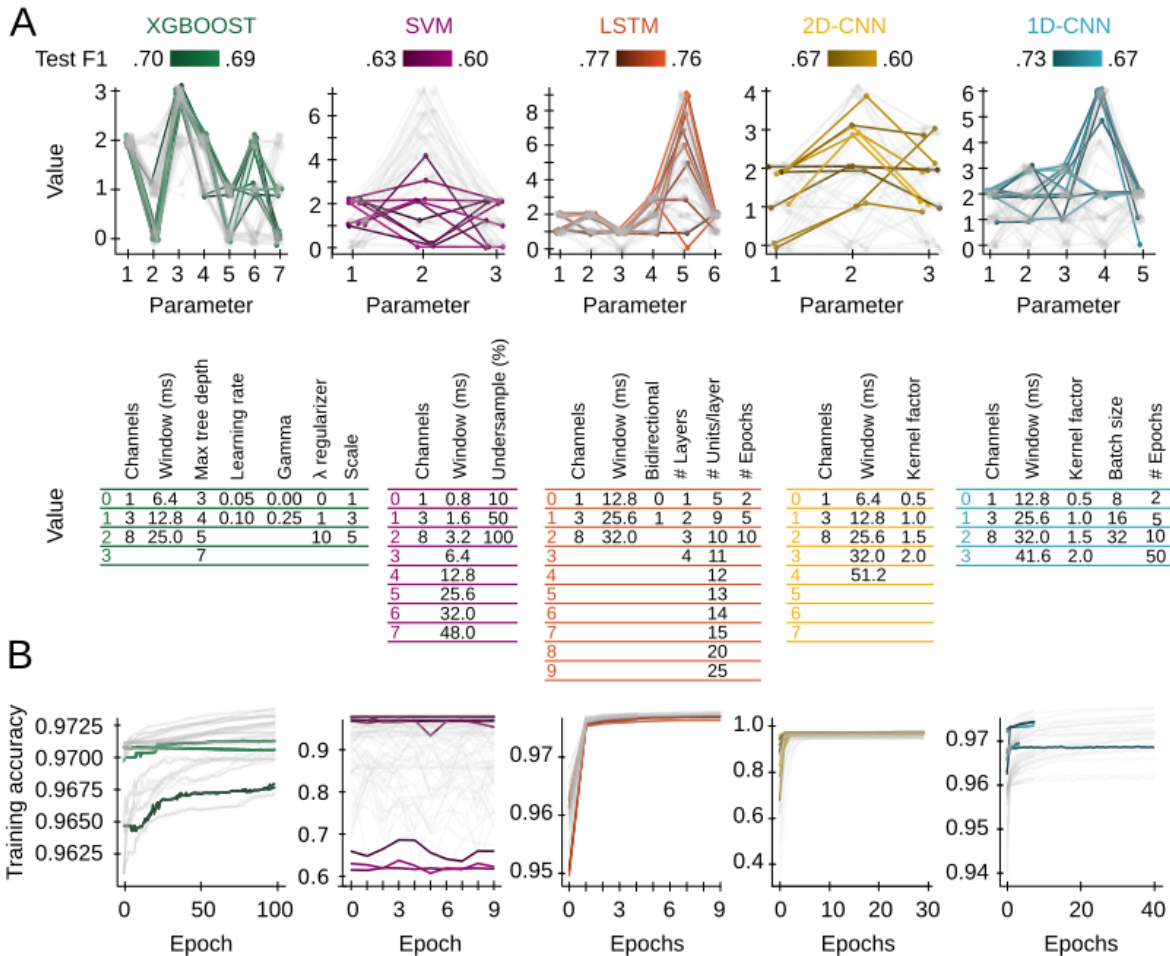219  there should be SWR information also coded in the period preceding an event[50].
220
221  A plug-and-play toolbox to use any of the best 5 models of each architecture for SWR detection is
222  available: https://github.com/PridaLab/rippl-AI.
223



224
225
226  **Figure 2: Training design and performance of ML models. A,** Training and selection criteria scheme. The
227  training dataset used in the hackathon was z-scored and down-sampled to 1250Hz. Training data were
228  shuffled and distributed into train and test subsets (70%-30% respectively). Each architecture was trained to
229  optimize F1 of the test set using several parameters. The 50 best models were tested over a new validation
230  data set (7586 events; 21 sessions from 8 animals), generating an F1 vs threshold curve per model/
231  architecture. Among these 50, the model with highest mean F1 was selected for between-models
232  comparison (right panel). **B,** LFP example of the validation set and the corresponding model outputs per
233  window of analysis. Note different duration of true events. Setting a threshold allows defining the windows
234  containing detected events. Colored ticks represent detections by the different models. Two different
235  thresholds (dark and light gray) can influence what events are detected. Note how detections marked with
236  arrows are dismissed when the threshold increases. Since SWRs constitute about 1-4% of the total

238 detected events are not computed for performance. **C,** Schematic illustration of Precision (percentage of
239 good detections), Recall (percentage of ground truth events that have been detected) and F1-score
240 (harmonic mean between Precision and Recall).

**A**

**XGBOOST** — Test F1 .70 ▬ .69

| | Channels | Window (ms) | Max tree depth | Learning rate | Gamma | λ regularizer | Scale |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 6.4 | 3 | 0.05 | 0.00 | 0 | 1 |
| 1 | 3 | 12.8 | 4 | 0.10 | 0.25 | 1 | 3 |
| 2 | 8 | 25.0 | 5 | | | 10 | 5 |
| 3 | | | 7 | | | | |

**SVM** — .63 ▬ .60

| | Channels | Window (ms) | Undersample (%) |
|---|---|---|---|
| 0 | 1 | 0.8 | 10 |
| 1 | 3 | 1.6 | 50 |
| 2 | 8 | 3.2 | 100 |
| 3 | | 6.4 | |
| 4 | | 12.8 | |
| 5 | | 25.6 | |
| 6 | | 32.0 | |
| 7 | | 48.0 | |

**LSTM** — .77 ▬ .76

| | Channels | Window (ms) | Bidirectional | # Layers | # Units/layer | # Epochs |
|---|---|---|---|---|---|---|
| 0 | 1 | 12.8 | 0 | 1 | 5 | 2 |
| 1 | 3 | 25.6 | 1 | 2 | 9 | 5 |
| 2 | 8 | 32.0 | | 3 | 10 | 10 |
| 3 | | | | 4 | 11 | |
| 4 | | | | | 12 | |
| 5 | | | | | 13 | |
| 6 | | | | | 14 | |
| 7 | | | | | 15 | |
| 8 | | | | | 20 | |
| 9 | | | | | 25 | |

**2D-CNN** — .67 ▬ .60

| | Channels | Window (ms) | Kernel factor |
|---|---|---|---|
| 0 | 1 | 6.4 | 0.5 |
| 1 | 3 | 12.8 | 1.0 |
| 2 | 8 | 25.6 | 1.5 |
| 3 | | 32.0 | 2.0 |
| 4 | | 51.2 | |
| 5 | | | |
| 6 | | | |
| 7 | | | |

**1D-CNN** — .73 ▬ .67

| | Channels | Window (ms) | Kernel factor | Batch size | # Epochs |
|---|---|---|---|---|---|
| 0 | 1 | 12.8 | 0.5 | 8 | 2 |
| 1 | 3 | 25.6 | 1.0 | 16 | 5 |
| 2 | 8 | 32.0 | 1.5 | 32 | 10 |
| 3 | | 41.6 | 2.0 | | 50 |

**B** (Training accuracy vs Epoch plots for each architecture)

**Figure S2: Definition of parameter space in the different ML architectures. A,** Results from the different architectures in the training dataset: XGBOOST, SVM, LSTM, 2D-CNN and 1D-CNN. Tables indicate the different hyper-parameters used to train each architecture. The resulting 10-best models are color-coded by their F1-score in the validation dataset. The remaining 40-best models are shown in light gray. **B,** Evolution of accuracy along training epochs for the ML models shown in A.

### Influence of the temporal and spatial sampling in training performance

252 Next, we sought to evaluate the relationship between model performance, parameters and LFP
253 input characteristics. Given the relevance of the temporal and spatial LFP sampling in the definition
254 of SWRs[31], we started evaluating how the size of the analyzed window and the number of
255 recording channels influenced performance. In order to have as much data as possible, we used
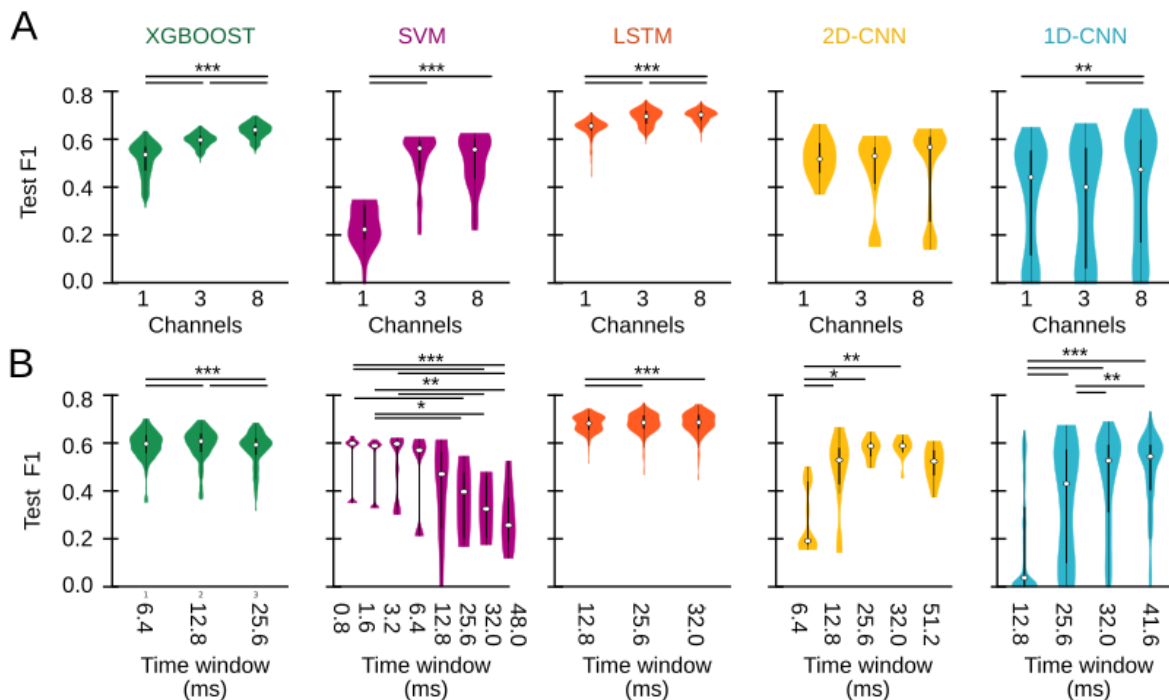256 F1-scores of all the trained models over the test set.

257 We found that XGBoost and LSTM were very stable, with performances changing very little for any
258 combination of window size and the number of channels used, suggesting that these architectures
259 can capture SWR features that are relatively invariant across temporal/spatial windows in the input
260 data (Fig.3A,B). Interestingly, the training parameter that most influenced these two architectures
261 was the number of LFP channels, with 3 and 8 channels providing better performances (Fig.3A).

262 Spatial information was also important for the SVM model, which scored poorly using a single
263 versus several channels (Fig.3A; magenta). As mentioned above, temporal resolution was also
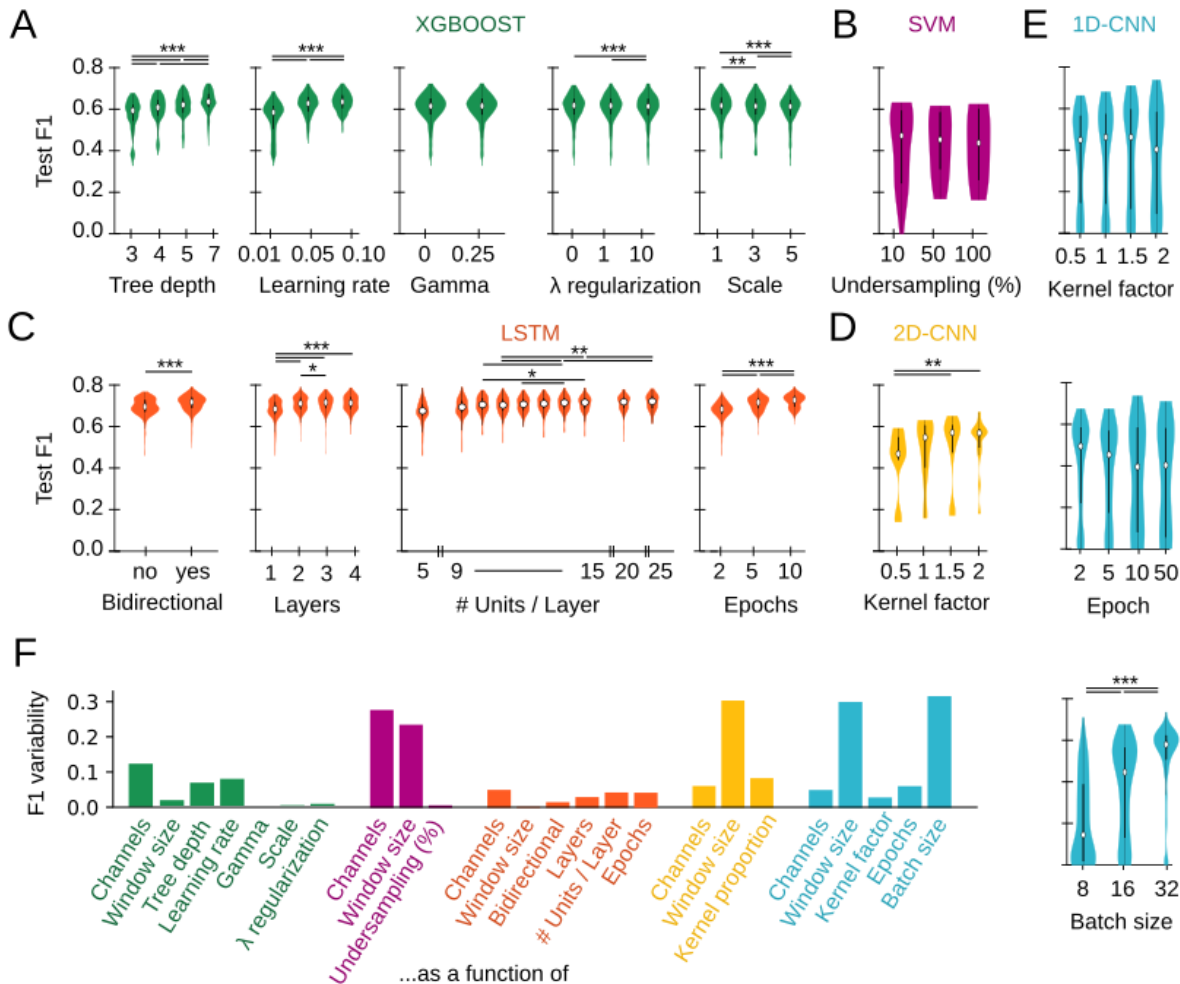264 critical for SVM, which required smaller time windows of <3.2 ms to succeed in detecting SWR

266 performance dropped significantly, indicating that a single SWR cycle and its particular waveform
267 across channels are optimal input information for the SVM architecture to detect events. This effect
268 could be due to the low number of trainable parameters used for SVM (ranging from 1 to 100; see
269 methods), which requires less but more informative data to achieve good performances.



270 **Figure 3: Influence of number of channels and analysis window on training performance. A,** Final test
271 F1-score of all trained models depending on the number of input channels: one (pyramidal channel; see
272 methods), three (pyramidal channel and extreme channels), or eight (all channels of the probe). Kruskal-
273 Wallis tests with repeated measures for every architecture: XGBOOST, Chi2(2)=1282.2, p<0.0001; SVM,
274 Chi2(2)=33.1, p<0.0001; LSTM, Chi2(2)=964.4, p<0.0001; 2D-CNN, not significant; 1D-CNN, Chi2(2)=14.6,
275 p=0.0007. Post hoc tests *, p<0.05; **, p<0.01, ***, p<0.001. **B,** Same as panel A, but depending on the time
276 window used for analysis. Kruskal-Wallis tests with repeated measures for every architecture: XGBOOST,
277 Chi2(2)=369.5, p<0.0001; SVM, Chi2(7)=48.8, p<0.0001; LSTM, Chi2(5)=48.0, p<0.0001; 2D-CNN,
278 Chi2(4)=16.5, p=0.0024; 1D-CNN, Chi2(3)=126.5, p<0.0001. Post hoc tests *, p<0.05; **, p<0.01, ***,
279 p<0.001.

280
281

282 Finally, both the 2D- and 1D-CNN models had similar performance for any number of channels,
283 although there was also a trend for higher spatial sampling (Fig.3B, yellow and acqua).
284 Interestingly, both CNN models presented a large F1 dispersion because their performance was
285 very dependent on the window size (Fig.3B). The 2D-CNN model exhibited maximal F1-score for
286 32ms, while most 1D-CNN models best scored for 25 ms (Fig.3B). This may be related to the
287 number of training parameters: the more parameters, the more complex tasks these algorithms
288 can solve, provided the amount of training data is representative enough of the expected variance.
289 This supports accurate detection in longer LFP windows. Examination of the remaining parameters
290 suggested additional differences across architectures (Fig.S3A-E). Interestingly, evaluating their
291 impact on F1-scores confirmed the effect of channels and window size on model behavior
292 (Fig.S3F). For CNN models, the batch size (1D-CNN) and the number of kernels (2D-CNN) were
293 also critical.

**Figure S3: Influence of architecture-specific training parameters on performance. A-E,** F1-scores from the test set for all models of each architecture. All statistical tests were Kruskal Wallis (KW) with repeated measures. **A,** XGBoost training parameters: maximum tree depth (KW: Chi2(3)=1321.6, p<0.0001), learning rate (KW: Chi2(2)=1109.4.6, p<0.0001), gamma (KW not significant), lambda regularization (KW: Chi2(2)=67.8, p<0.0001) and scale (KW: Chi2(2)=111.6, p<0.0001). Post hoc tests *, p<0.05; **, p<0.01, ***, p<0.001. **B,** SVM training parameters: under-sampling percentage (KW not significant). Higher % of undersampling means training the model with higher representativity of GT data. **C,** LSTM training parameters: bidirectionality (KW: Chi2(1)=320.1, p<0.0001), number of layers KW: Chi2(3)=602.4, p<0.0001), number of units per layer (KW: Chi2(9)=543.8, p<0.0001) and training epochs (KW: Chi2(2)=836.1.6, p<0.0001). **D,** 2D-CNN training parameters: number of kernels scaling factor (KW: Chi2(3)=16.0, p=0.0011), number of epochs and batch size (KW not significant). **E,** 1D-CNN training parameters: number of kernels scaling factor (KW not significant), number of training epochs (KW not significant), and batch size (KW: Chi2(2)=196.9, p<0.0001). **F,** F1-score variability as a function of all training parameters. F1 variability was computed as the difference between the maximum and minimum mean F1.
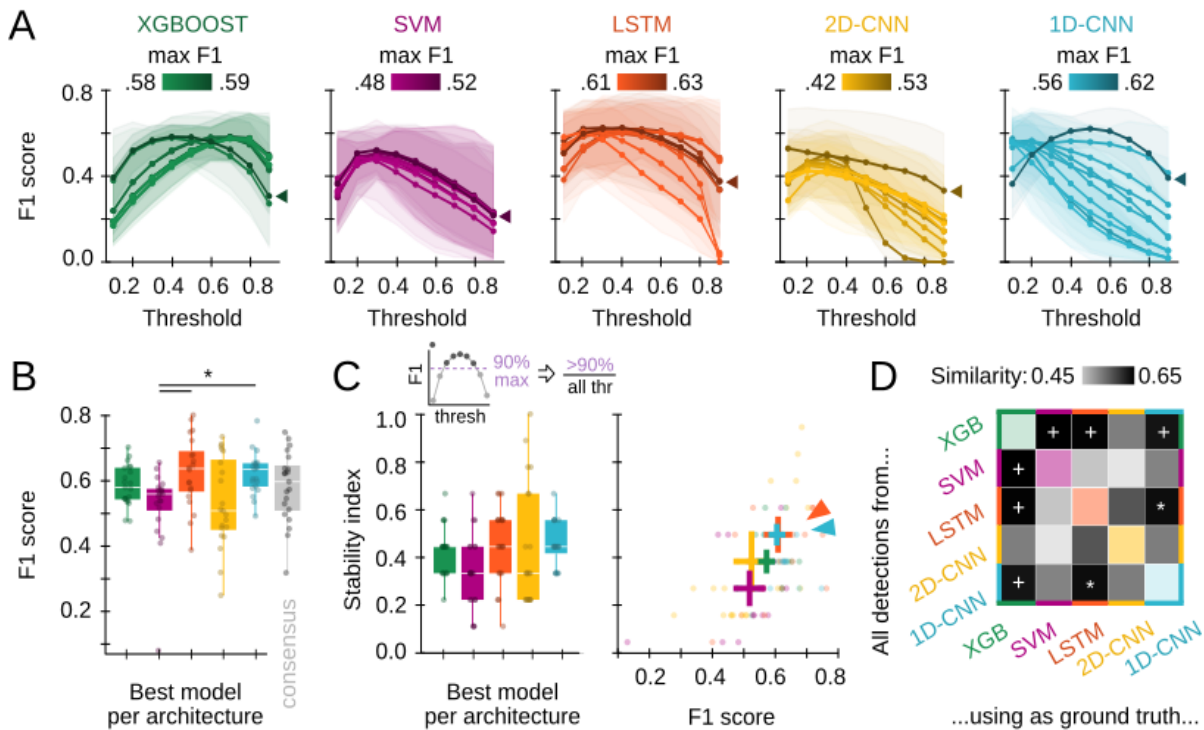
## Comparison between optimized models

The analysis above provided insights on how input characteristics and processing parameters can influence detection performance in different ML models. Understanding how each architecture learns to identify ripple-like events can not only can aid the development of new tools, but unveil what are the key LFP features used for detection. We thus evaluated conditions for their best performance.

For fair comparison between architectures, we selected the 10-best models from the validation set. Remarkably, our previously published 1D-CNN model[31] was among the 10-best 1D-CNN, outperforming other configurations. Plotting F1-scores of all models across a range of thresholds allowed visualization of their performance stability as a function of the probability threshold (Fig.4A).

321 We analyzed their performance along a range of characteristics (performance, robustness, and
322 threshold dependency) to better inform their selection depending on research applications. Five of
323 the 10-best trained models of all architectures are available at https://github.com/PridaLab/rippl-
324 AI/blob/main/optimized_models/

325 The consistency of F1-threshold curves depended on the model architectures (Fig.4A). Most
326 models reached their maximal F1-score at relatively low threshold values of 0.3-0.4 and remained
327 stable until a probability of around 0.5-0.7. Such a behavior indicates robust performance, since
328 even low probability (i.e., relatively uncertain) output predictions overlapped with the ground truth.
329 This property is very useful for online experimental applications, when choosing different
330 thresholds is not manageable, making detection more robust. Interestingly, we found that XGBoost
331 models exhibited good performance at two threshold ranges (0.2-0.4 and 0.6-0.8), depending on
332 how trained models penalized False negative predictions. Similarly, for both CNN architectures, we
333 found several models operating sharply at low thresholds, while others exhibited a relatively stable
334 operation in the 0.4-0.6 range especially for 1D-CNN models. We confirmed the variability of
335 different models within a given architecture by looking at their Precision vs Recall curves for the
336 entire threshold range (Fig.S4A). This variability suggests that even when arising from the same
337 architecture, algorithmic processes and detection strategies by which the different models were
338 detecting SWR events could differ. This may provide a range of models for different applications.



339

340 **Figure 4: Comparison between best performing ML models. A,** F1 against threshold from the 10-best
341 models of each architecture as evaluated in the validation set. Each line represents the performance of one
342 trained model, colored by its maximal F1 (mean from all sessions is plotted in dark color). Data reported as
343 mean±95% confidence interval for validation sessions. Arrows indicate the best model of each architecture.
344 **B,** F1-scores for the best model of panel A. Thresholds used are: 0.4 for XGBoost, 0.5 for SVM, 0.4 for
345 LSTM, 0.1 for CNN2D, 0.5 for CNN1D. Each dot represents a session of the validation set (n=21 sessions; 8
346 mice). In gray, the F1-score for a consensus detector. Kruskal-Wallis, Chi2(5)=26.9, p<0.0001; post hoc tests
347 *, p<0.05; **, p<0.01, ***, p<0.001. **C,** Stability index for the best model of each architecture (left), and the
348 stability index vs the F1 (right). Kruskal-Wallis, Chi2(4)=10.5, p=0.03; post hoc tests. **D,** Similarity between
349 predicted events of different architectures. Models are the same as in panels B-C. To measure the similarity,
350 the mean F1 across validation sessions have been computed, using detected events in the y-axis as
351 detections, and detected events in the x-axis as ground truth. Note the similarity between LSTM and 1D-
352 CNN (white *), and that by XGBoost against SVM, LSTM and 1D-CNN (white +).
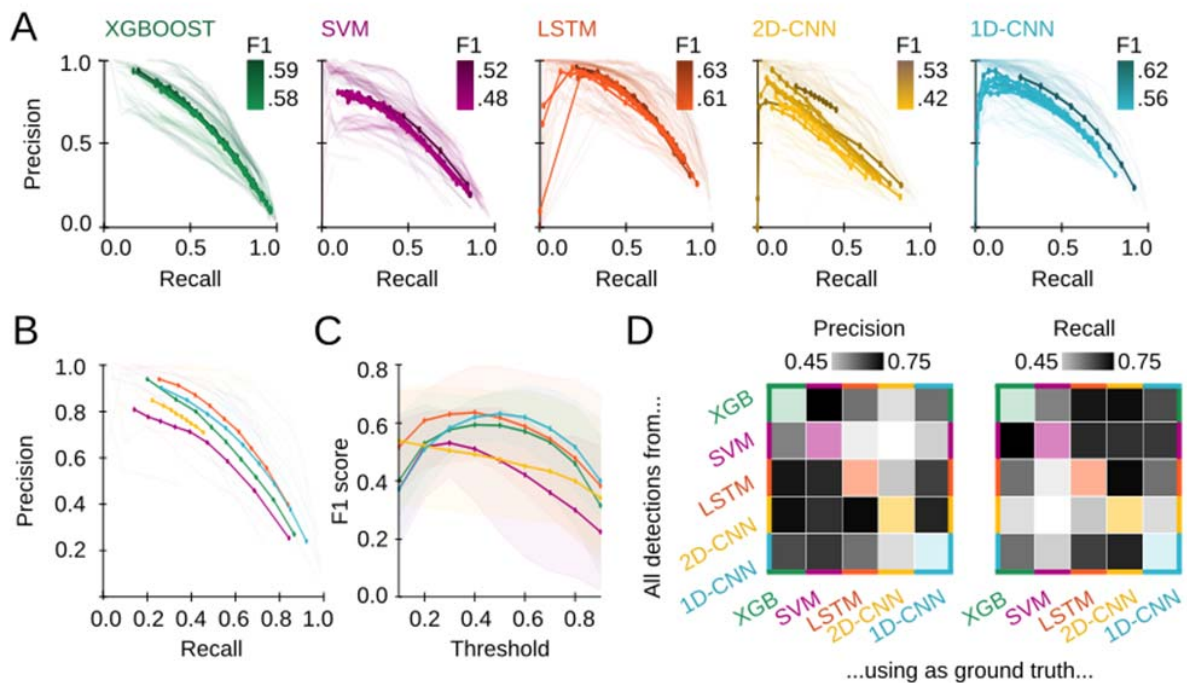
353
354

355    Next, we selected the model that reached the highest F1 value from each architecture (Fig.4A,
356    best models, arrowheads), and compared their scores using all validation sessions (Fig.4B). We
357    found that the LSTM and 1D-CNN best models outperformed other architectures, with mean F1-
358    scores over 0.6 (as a reference, the inter-expert F1-score in our lab is ~0.7[31]. Precision-Recall
359    curves from these two models clearly stood out of the other solutions (Fig.S4B). Importantly, a
360    consensus prediction based on the 5-best models did not perform better than individual
361    architectures alone (Fig.4B; gray).

362

363    Given the importance of consistent threshold performance for practical applications, we quantified
364    the robustness of F1-threshold curves for the best models using a stability index in the validation
365    dataset (see methods). Models with a stability index of 1.0 provide at least 90% of its maximal
366    performance for any threshold value, a property especially suitable for experimental applications.
367    While the best 2D-CNN model exhibited stability in some validation sessions, the best LSTM and
368    especially the best 1D-CNN best models exhibited more consistent behavior (Fig.4C, left; Fig.S4C).
369    We confirmed this result by plotting the stability index versus F1, where both the best LSTM and
370    1D-CNN best models clearly segregated (Fig.4C, right; arrowhead).



371
372    **Figure S4: Precision-Recall curves of optimized models. A,** Precision (P) vs Recall (R) curves for the 10-
373    best models of each architecture. Each dot represents P-R values for a particular threshold. Each line
374    represents the performance of one trained model, colored by its maximal F1 (mean of all sessions is plotted
375    in dark color; sessions are light colored). **B,** P-R curves for the best model of each architecture (all
376    thresholds). Thick lines represent mean values. Thin lines curves are individual validation sessions. **C,** F1-
377    score as a function of the threshold. Data reported as mean±95% confidence interval for validation sessions.
378    **D,** Similarity between the events predicted by the best model (maximum F1) of each architecture. Models
379    shown are the ones with maximum F1. To measure the similarity, we computed the mean Precision (right)
380    and Recall (left) across validation sessions have been computed, and used detected SWR events of models
381    in the y-axis as detections, and detected events of models in the x-axis as ground truth.

382
383
384    Finally, to evaluate whether the different models were targeting similar or different subsets of SWR
385    events, we compared how similar their detections were. To quantify this similarity, we computed the
386    F1 between both groups of detections, using one of them as the ground truth (Fig.4D). Interestingly,
387    the 1D-CNN and LSTM showed a high level of similarity, in line with their consistent and accurate
388    behavior (Fig.4D, white *). XGBoost scored a high similarity with all other architectures except for
389    the 2D-CNN (Fig.4D, white +). Possibly, this reflects the fact that very few of the XGBoost
390    detections were also predicted also by 2D-CNN, leading to a very low Precision (Fig.S4D). In
391    general, high similarities did not seem to be caused by a particularly high Precision or Recall

392    (model A detects so few events that all coincide with detections of model B), but by a good balance
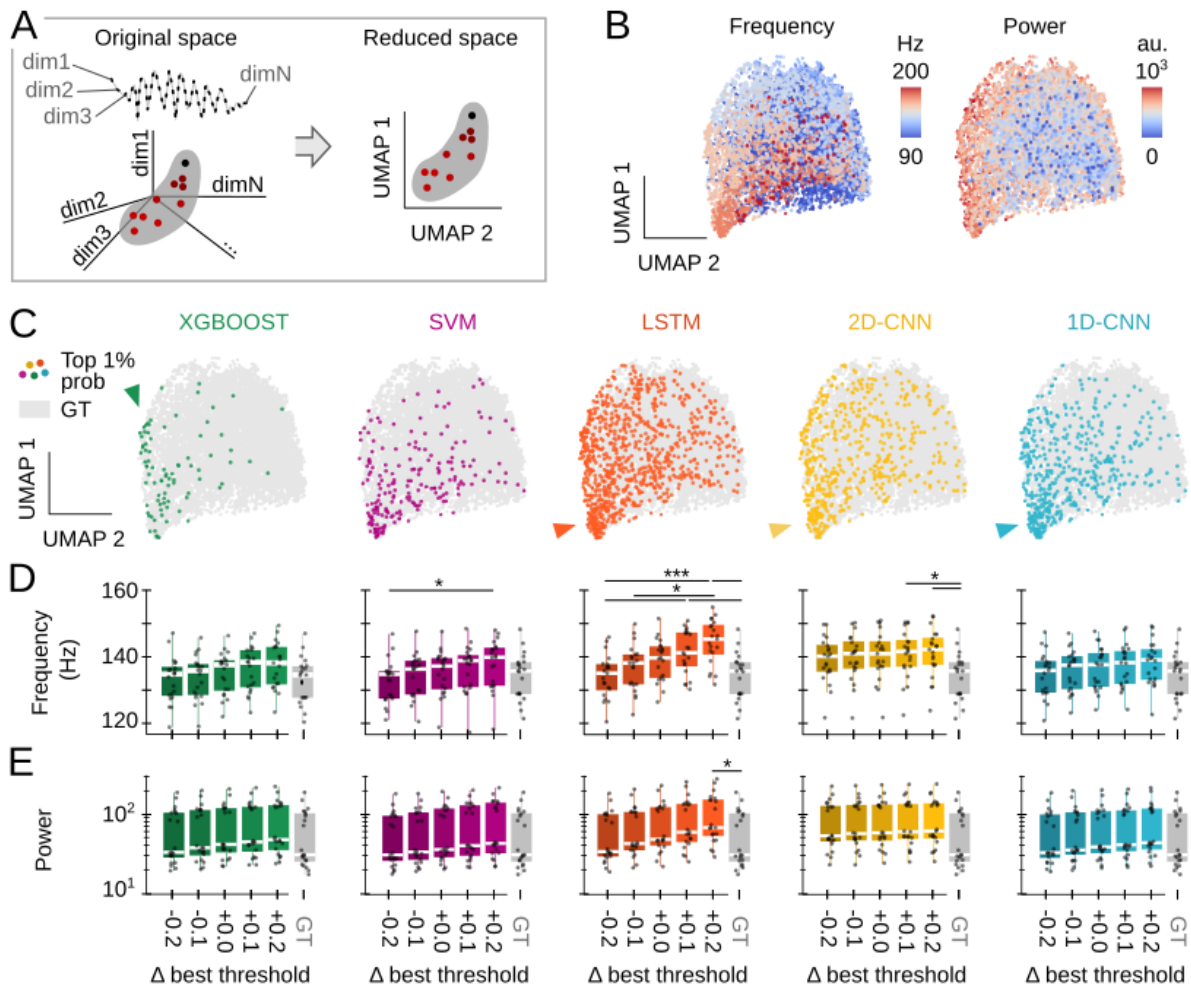393    between both (events of model A and B highly overlap) (Fig.S4D).
394
395
396

### Effect of different ML models on the features of detected SWRs

398    Results above suggest that different models may be relying on different strategies for recognizing
399    SWRs. We thus wondered whether models could be biased towards SWRs with different features
400    (e.g. frequency, amplitude, etc…), and whether these biases could also be reflected over different
401    ranges of output probabilities.

402    In order to evaluate these issues, we resorted to a low-dimensional analysis of SWRs which allows
403    for their unbiased topological characterization[14]. In this strategy, SWR events are considered points
404    in an N-dimensional space, where each dimension X (dimX) represents the LFP value sampled at
405    a given timestamp X (Fig.5A). In our case, as events were GT ripples of 50 ms sampled at 1250Hz
406    (i.e. 63 timestamps), the original space was 63 dimensions. Plotting all SWR events will result in a
407    point cloud, with events sharing similar LFP features lying close to each other, while those of
408    different characteristics distribute separately (Fig.5A). To ease visualization, the SWRs were
409    embedded in a low-dimensional representation using Uniform Manifold Approximation and
410    Projection (UMAP)[14,51].



411

**Figure 5. Effect of ML models and thresholds on the type of detected SWR. A,** Low-dimensional
analysis of SWR features[14]. GT ripples are represented into a high-dimensional space by mapping each
timestamp to a particular dimension. Since the sampling rate is 1250Hz, and windows around SWRs were
cut to 50ms, there are 63 timestamps per event, and so the original space has 63 dimensions. The SWR
cloud is embedded in a low-dimensional space using UMAP. **B,** UMAP embedding projected into the two
first axes. Each dot represents a GT ripple, and its color reflects its frequency (left) and power (right). Note

419    ripples with similar features are close together. **C,** Colored dots superimposed over gray GT data represent
420    the top 1% of detected events for every given architecture, i.e., True Positive events with an output SWR
421    probability above 99% of the maximum probability for that given model. Note that different distributions of
422    events in the cloud reflect biases of ML model used for detection. **D,** Frequency of True Positive SWR
423    detected by each architecture. Each dot represents the mean frequency of detected ripples of one validation
424    session (21 sessions from 8 animals). Kruskal-Wallis tests for every architecture: XGBOOST, not significant;
425    SVM, $Chi2(5)=11.1$, $p=0.049$; LSTM, $Chi2(5)=29.9$, $p<0.0001$; 2D-CNN, $Chi2(5)=13.8$, $p=0.017$; 1D-CNN, not
426    significant. **E**, Spectral power of True Positive events detected by each architecture. Kruskal-Wallis tests for
427    every architecture: XGBOOST, SVM, 2D-CNN and 1D-CNN are not significant; LSTM, $Chi2(5)=14.0$,
428    $p=0.016$. Post hoc tests *, $p<0.05$; **, $p<0.01$, ***, $p<0.001$.
429
430

431    First, we analyzed how ripple frequency and power were distributed in the UMAP embedding by
432    coloring each dot (i.e. each SWR) based on their frequency (Fig.5B, left) and power (Fig.5B, right).
433    As expected from our previous work[14], these features followed different distributions, segregating
434    high-frequencies towards the bottom of the cloud and high-power events radially out (Fig.5B). We
435    then inspected events detected by the best model of each architecture by plotting the top 1%
436    detections, defined as True Positive events for which the model output probability was >99% of its
437    maximum probability (Fig.5C). Interestingly, each model showed different distributions of preferred
438    SWRs. For example, XGBoost was biased towards a subset of high-power and fast SWR events
439    (Fig.5C, green arrowhead), whereas the SVM model exhibited a more heterogenous distrinution. In
440    turn, LSTM and both CNNs assigned higher probabilities to events that had a good frequency-
441    power balance (Fig.5C, orange, yellow and blue arrowheads). Note how these models have more
442    colored events, consistent with their higher stability indices reported above (Fig.4C).

443    To quantify detection biases in each ML model, we analyzed the frequency and power of their True
444    Positive events and compared them against those in the GT. Consistent with the UMAP
445    distributions, SWR frequency was highly dependent on the threshold for SVM, LSTM and 2D-CNN
446    algorithms (Fig.5D). The case of LSTM was particularly striking with differences accumulating for
447    all thresholds. Instead, for the SVM and 2D-CNN biases were significant only when thresholds
448    differed ±0.2 from the optimal value (Fig.5D). As previously reported[31], the 1D-CNN exhibited
449    roughly consistent behavior with SWR features not statistically different from GT events. SWR
450    power exhibited no major dependency on the threshold in any of the models but the LSTM,
451    especially at higher detection thresholds (Fig.5E).

452    Altogether, this analysis suggests that the different ML models can be exploited to detect a wide
453    range of SWRs with different characteristics.

454
455    **Using the toolbox to identifying SWRs in non-human primates**
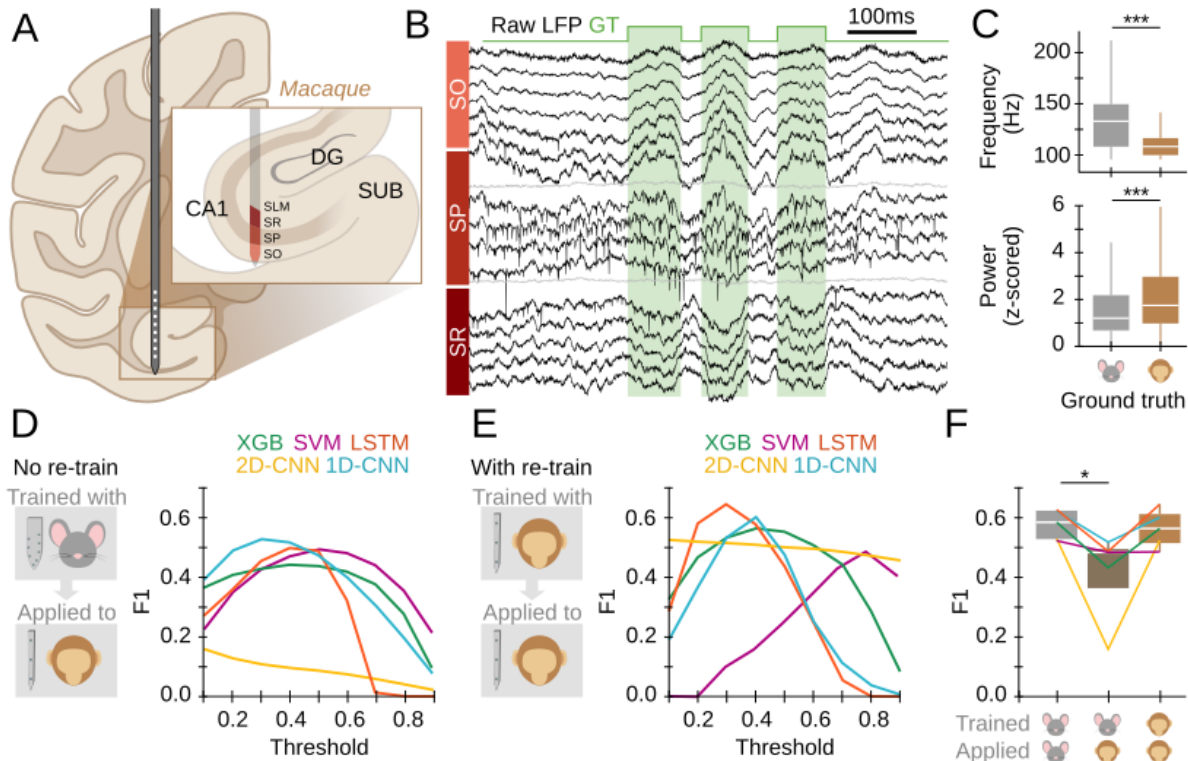456

457    A major motivation of our study is to develop methods which can be generalizable for a wider
458    range of detection contexts, including a greater range of species and biomedical applications. Thus,
459    we applied our ML models to LFP recordings from the hippocampus of the macaque, which shares
460    a high level of genetic, morphological and physiological characteristics with that of its fellow
461    primate, the human, while enabling precise localization of signals roughly comparable to those
462    used for the algorithm development. To accomplish this, we recorded hippocampal LFP signals
463    from a freely moving macaque using a multichannel linear probe[52] (Fig.6A). Unlike the original
464    high-density probes (20 µm), recordings were obtained every 90/60 µm and spanned CA1 layers
465    (Fig.6A). As in mice, SWRs were manually identified (4133 events) to generate the annotated
466    ground truth (Fig.6B). Consistent with the literature[16,17], macaque SWRs had lower frequencies and
467    higher power as compared to mouse ripples (Fig.6C).
468
469    We applied the best model of each architecture trained in head-fixed mice to macaque recordings,
470    and evaluated their performance. For fair comparison, we flipped laminar LFP signals upside down
471    and sampled the channel combination that best matched the characteristic mouse LFP profile (see
472    Methods and layer orientation in Fig.6A). Strikingly, 4/5 models reached a maximum F1 of ~0.5
473    (Fig.6D), close to their maximal performance on mice data (~0.6). SVM, 1D-CNN and LSTM
474    exhibited the best performance, as compared to XGBoost and 2D-CNN (Fig.6D). Importantly, the

475  fact that both LSTM and 1D-CNN have relatively good generalization capability, suggests that they
476  successfully capture shared features of SWRs from mice and macaques.

477
478
479



480
481  **Figure 6. Extending sharp-wave ripple detection to non-human primates. A,** Linear multichannel probes
482  were used to obtain LFP recordings from the anterior hippocampus of a freely moving monkey. **B,** SWR
483  events were manually tagged (4133 events) as in mouse data. **C,** Significant differences between SWR
484  recorded in mice and monkey. Kruskal-Wallis Chi2=1649, p<0.0001 for frequency; Kruskal-Wallis Chi2=407,
485  p<0.0001 for power. Posthoc, ***, p<0.001. Data from the GT in both cases. **D**. The best model of each
486  architecture trained in mouse data was applied to detect SWRs on the macaque data. Input data consisted
487  of 5 LFP channels of SO, SP and SR, and 3 interpolated channels (see methods for details). We evaluated
488  all models by computing F1-score against the ground truth (GT). Note relatively good results from non-
489  retrained ML models. **E,** Results of model re-training using macaque data. Data were split into a training and
490  test dataset (50% and 20% respectively), used to train the models; and a validation set (30%), used to
491  compute the F1 (left panel). **F,** F1-scores for the maximal performance of each model before and after re-
492  training. Kruskal-Wallis test, Chi2(2)=8.06, p=0.018. Post hoc tests *, p<0.05.

493
494
495
496  We next chose to re-train the 5 models with the macaque dataset, using 50% for training and 20%
497  for testing. The remaining 30% was used for validation to compute the final F1. For re-training, we
498  reset all trainable parameters (internal weights) but kept all architectural hyper-parameters fixed
499  (number of input number of channels, input window length, number of layers, etc…). Performances
500  improved after retraining for 4/5 models, reaching a F1 increase of +0.3 for 2D-CNN (Fig.6E). The
501  best model was LSTM, followed by 1D-CNN and XGBoost. SVM was the only model that did not
502  improve after retraining, but exhibited a shift towards larger thresholds. Furthermore, performance
503  of macaque SWR detection after re-training reached the mouse level (Fig.6E), suggesting that
504  these models identified similar key features in both species, and could readily be trained to similar
505  levels of accuracy across mice and monkeys. A user-friendly open python notebook to re-train any
506  of the 5 models and use it for event detection is available at https://github.com/PridaLab/rippl-
507  AI/blob/main/examples_retraining.ipynb

508
509
510

**Discussion**

Here, we provide a pool of models for automatic SWR detection based on different ML architectures. These include some of the most used ML solutions, such as XGBoost, SVM, 1D- and 2D-CNN and LSTM. The models, which resulted from unbiased community-based proposals, are able to capture a wealth of SWR features recorded in the dorsal hippocampus of head-fixed mice. When applied to LFP recordings from a freely moving macaque, these models were able to generalize detection.

The need for detecting and classifying high-frequency oscillations such as SWR has accelerated over recent years for advanced biomedical applications[28,33,35,41,53]. Identification of these events can help to delineate normal from pathological epileptogenic territories[18,54,55], and to develop closed-loop intervention strategies for boosting memory function[33,35]. However, spectral-based methods have revealed suboptimal and the community is actively seeking for novel feature-based strategies. Recently, solutions based in ML methods have started to emerge[25,28,31,54]. Using these tools will drive advances not only in online detection of SWRs, but also their unbiased categorization for better mechanistic understanding[11,13,31,56,57], including their functional ties to visuospatial and episodic memory[10,11,16,34,38,39].

Amongst the 5 ML models examined here we found the LSTM and 1D-CNN to provide the best performance and reliability using rodent data. The other models exhibited roughly similar behavior depending on the input parameter selection (recording channels and analysis windows). While in general, we found that all of them performed better with high-density multi-channel recordings (8 channels), some of them (e.g. 2D-CNN) exhibited similar results while operating over data sampled with 1 to 3 channels. This suggests they may be able to identify characteristic features with reduced spatial information, which could facilitate applications to human recordings[19,37].

Detection of SWR candidates with ML models is based on using a probability threshold. We found that the different models exhibited a degree of sensitivity to threshold selection, with LSTM, XGBoost and 1D-CNN providing a wider range of operational stability. This suggests there is a larger range of thresholds in these models which provide relatively similar performance. Instead, SVM and the 2-CNN better operate in a very narrow threshold range. This is very important for online applications, when threshold selection can affect experimental results in real time[25].

The different ML models are biased towards SWRs with slightly different properties, probably reflecting their internal representations of SWR characteristic LFP features[31]. During training, each model learns to identify what specific LFP features made ripples distinguishable from background LFP signals, so that during SWR detection, the presence of those features raises their output probability. The fact that the properties of detected SWR depend on the probability threshold for SVM and 2D-CNN suggests that frequency and relative power are some of the LFP features these models identified during training. On the contrary, XGBoost, LSTM and 1D-CNN models, which showed less bias, may be capturing other LFP features such as the spatial profile. This is consistent with results from the analysis of the influence of spatial sampling in training performance in these ML models.

When applied to data from the macaque anterior hippocampus, we found that models trained with LFP signals from the dorsal hippocampus of mice can perform relatively well, especially considering established differences in frequency and in LFP shape in monkey and human[10,16,17]. After re-training, their operation improved significantly, reaching the inter-experts' performance levels at 0.7[31]. This demonstrates the strong capability of the ML models to generalize and suggests the existence of shared features across species. This is of particular importance, because many human applications may not have the exact spatial localization or the same electrode types, in some cases even within studies, and so any effective ML applications will need a high degree of generalizability. It also demonstrates the proof of principle for applying to a wider range of measurements, including other animal models and ripple-adjacent pathologies such as MTL seizures[54].

More testing along these lines will identify the extent of generalizability across different permutations of species, location, electrode sampling and type, to find the limits of these ML models. To enable such developments, we made several of the 10-best trained models and our

568 coding strategies for detection and retraining openly available to the research community at
569 https://github.com/PridaLab/rippl-AI. They can be tested through open-source notebooks that are
570 ready to use, with enough examples to illustrate their operation capability. Although the notebooks
571 provide easily readable code, they may not be optimal for further code development. That is why
572 the core functions are written as separate Python modules. Users can test these models for SWR
573 detection by loading their own data and defining the channels. The ripple_AI repository has a wide
574 variety of SWR detection tools that include optional supervised detection curation, and a graphical
575 user interface for a quick visual exploration of detected events depending on the threshold chosen,
576 as well as the option of retraining a model with the user's own data.
577
578 This collection of resources joins to the many other community-based approaches for model
579 benchmarking[30,41,o53,58]. Crowdsourced solutions are becoming a tool to advance solutions of
580 particularly difficult problems which require knowledge integration[40,43]. This provides the field with
581 a set of platforms for detecting events from diverse datasets using traditional and state-of-the-art
582 algorithms (e.g., our own ripple-AI toolbox, and https://www.sharpwaveripples.org/). Our toolbox
583 goes beyond SWR detection, easing development of personalized ML models to detect other
584 electrophysiological events of interest[32]. This may be critical in experimental and/or clinical cases,
585 where other detection criteria, i.e. F-values, than those maximizing performance may be more
586 important. For instance, different experiments may call for avoiding either type I or type II errors,
587 and hence the balance between Precision and Recall. Such a versatility of our toolbox may be
588 further exploited to accelerate our understanding of hippocampal function and to support the
589 development of biomedical applications.
590

**References**

592
593  1. da Silva, F. L. EEG and MEG: Relevance to Neuroscience. *Neuron* **80**, 1112–1128 (2013).
594  2. Buzsáki, G. Hippocampal sharp wave-ripple: A cognitive biomarker for episodic memory and
595     planning. *Hippocampus* **25**, 1073–188 (2015).
596  3. Csicsvari, J., Hirase, H., Mamiya, A. & Buzsáki, G. Ensemble patterns of hippocampal CA3-
597     CA1 neurons during sharp wave-associated population events. *Neuron* **28**, (2000).
598  4. Stark, E. *et al.* Pyramidal cell-interneuron interactions underlie hippocampal ripple
599     oscillations. *Neuron* **83**, 467–480 (2014).
600  5. Buzsáki, G., Lai-Wo S., L. & Vanderwolf, C. H. Cellular bases of hippocampal EEG in the
601     behaving rat. *Brain Research Reviews* **6**, 139–171 (1983).
602  6. Kudrimoti, H. S., Barnes, C. A. & McNaughton, B. L. Reactivation of hippocampal cell
603     assemblies: effects of behavioral state, experience, and EEG dynamics. *J. Neurosci.* **19**,
604     4090–4101 (1999).
605  7. Genzel, L. *et al.* A consensus statement: defining terms for reactivation analysis. *Philos.*
606     *Trans. R. Soc. Lond. B. Biol. Sci.* **375**, (2020).
607  8. Joo, H. R. & Frank, L. M. The hippocampal sharp wave–ripple in memory retrieval for
608     immediate use and consolidation. *Nature Reviews Neuroscience* **19**, 744–757 (2018).
609  9. Pfeiffer, B. E. The content of hippocampal 'replay'. *Hippocampus* (2017).
610     doi:10.1002/hipo.22824
611  10. Mil, A. *et al.* Replay of cortical spiking sequences during human memory retrieval. *Science*
612      **367**, 1128–1130 (2020).
613  11. Liu, A. A. *et al.* A consensus statement on detection of hippocampal sharp wave ripples and
614      differentiation from other fast oscillations. *Nat. Commun. 2022 131* **13**, 1–14 (2022).
615  12. Reichinnek, S., Künsting, T., Draguhn, A. & Both, M. Field potential signature of distinct
616      multicellular activity patterns in the mouse hippocampus. *J. Neurosci.* **30**, 15441–9 (2010).
617  13. Ramirez-Villegas, J. F., Logothetis, N. K. & Besserve, M. Diversity of sharp-wave-ripple LFP
618      signatures reveals differentiated brain-wide dynamical events. *Proc. Natl. Acad. Sci. U. S. A.*
619      **112**, E6379-87 (2015).
620  14. Sebastian, E. R. *et al.* Topological analysis reveals input mechanisms behind feature
621      variations of sharp-wave ripples. *Under Revis.* (2023).
622  15. Patel, J., Schomburg, E. W., Berényi, A., Fujisawa, S. & Buzsáki, G. Local generation and
623      propagation of ripples along the septotemporal axis of the hippocampus. *J. Neurosci.* **33**,
624      17029–41 (2013).
625  16. Leonard, T. K. *et al.* Sharp Wave Ripples during Visual Exploration in the Primate

626          Hippocampus. *J. Neurosci.* **35**, 14771–14782 (2015).

627   17.   Skaggs, W. E. *et al.* EEG sharp waves and sparse ensemble unit activity in the macaque
628          hippocampus. *J. Neurophysiol.* **98**, 898–910 (2007).

629   18.   Bragin, A., Engel, J., Wilson, C. L., Fried, I. & Mathern, G. W. Hippocampal and entorhinal
630          cortex high-frequency oscillations (100-500 Hz) in human epileptic brain and in kainic acid-
631          treated rats with chronic seizures. *Epilepsia* **40**, 127–137 (1999).

632   19.   Worrell, G. A. *et al.* High-frequency oscillations in human temporal lobe: Simultaneous
633          microwire and clinical macroelectrode recordings. *Brain* **131**, 928–937 (2008).

634   20.   Alvarado-Rojas, C. *et al.* Different mechanisms of ripple-like oscillations in the human
635          epileptic subiculum. *Ann. Neurol.* **77**, 281–290 (2015).

636   21.   Valero, M. *et al.* Mechanisms for Selective Single-Cell Reactivation during Offline Sharp-
637          Wave Ripples and Their Distortion by Fast Ripples. *Neuron* **94**, (2017).

638   22.   Cowen, S. L., Gray, D. T., Wiegand, J. P. L., Schimanski, L. A. & Barnes, C. A. Age-
639          associated changes in waking hippocampal sharp-wave ripples. *Hippocampus* **30**, 28–38
640          (2020).

641   23.   Born, H. A. *et al.* Genetic suppression of transgenic APP rescues Hypersynchronous
642          network activity in a mouse model of Alzheimer's disease. *J. Neurosci.* **34**, 3826–3840
643          (2014).

644   24.   Engel, J., Bragin, A., Staba, R. & Mody, I. High-frequency oscillations: What is normal and
645          what is not? *Epilepsia* **50**, 598–604 (2009).

646   25.   Sethi, A. & Kemere, C. Real time algorithms for sharp wave ripple detection. *Annu. Int. Conf.*
647          *IEEE Eng. Med. Biol. Soc. IEEE Eng. Med. Biol. Soc. Annu. Int. Conf.* **2014**, 2637–2640
648          (2014).

649   26.   Liu, A. A. *et al.* A consensus statement on detection of hippocampal sharp wave ripples and
650          differentiation from other fast oscillations. *Nat. Commun. 2022 131* **13**, 1–14 (2022).

651   27.   Kulkarni, P. M. *et al.* A deep learning approach for real-time detection of sleep spindles. *J.*
652          *Neural Eng.* **16**, 36004 (2019).

653   28.   Hagen, E. *et al.* RippleNet: a Recurrent Neural Network for Sharp Wave Ripple (SPW-R)
654          Detection. *Neuroinformatics* **19**, (2021).

655   29.   Nadalin, J. K. *et al.* Application of a convolutional neural network for fully-automated
656          detection of spike ripples in the scalp electroencephalogram. *J. Neurosci. Methods* **360**,
657          109239 (2021).

658   30.   Valenchon, N. *et al.* The Portiloop: A deep learning-based open science tool for closed-loop
659          brain stimulation. *PLoS One* **17**, e0270696 (2022).

660   31.   Navas-Olive, A., Amaducci, R., Jurado-Parras, M.-T., Sebastian, E. R. & de la Prida, L. M.
661          Deep learning based feature extraction for prediction and interpretation of sharp-wave
662          ripples in the rodent hippocampus. *Elife* **11**, (2022).

663   32.   Frey, M. *et al.* Interpreting wide-band neural activity using convolutional neural networks.
664          *Elife* **10**, (2021).

665   33.   Talakoub, O., Gomez Palacio Schjetnan, A., Valiante, T. A., Popovic, M. R. & Hoffman, K. L.
666          Closed-Loop Interruption of Hippocampal Ripples through Fornix Stimulation in the Non-
667          Human Primate. *Brain Stimul.* **9**, 911–918 (2016).

668   34.   Norman, Y. *et al.* Hippocampal sharp-wave ripples linked to visual episodic recollection in
669          humans. *Science (80-. ).* **365**, (2019).

670   35.   Geva-Sagiv, M. *et al.* Augmenting hippocampal-prefrontal neuronal synchrony during sleep
671          enhances memory consolidation in humans. *Nat. Neurosci.* **26**, 1100–1110 (2023).

672   36.   Tong, A. P. S., Vaz, A. P., Wittig, J. H., Inati, S. K. & Zaghloul, K. A. Ripples reflect a
673          spectrum of synchronous spiking activity in human anterior temporal lobe. *Elife* **10**, (2021).

674   37.   Curot, J. *et al.* Local neuronal excitation and global inhibition during epileptic fast ripples in
675          humans. *Brain* **146**, 561–575 (2023).

676   38.   Leonard, T. K. & Hoffman, K. L. Sharp-Wave Ripples in Primates Are Enhanced near
677          Remembered Visual Objects. *Curr. Biol.* **27**, 257–262 (2017).

678   39.   Hussin, A. T., Leonard, T. K. & Hoffman, K. L. Sharp-wave ripple features in macaques
679          depend on behavioral state and cell-type specific firing. *Hippocampus* **30**, 50–59 (2020).

680   40.   Berens, P. *et al.* Community-based benchmarking improves spike rate inference from two-
681          photon calcium imaging data. *PLoS Comput. Biol.* **14**, (2018).

682   41.   Kuhlmann, L. *et al.* Epilepsyecosystem.org: crowd-sourcing reproducible seizure prediction
683          with long-term human intracranial EEG. *Brain* **141**, 2619–2630 (2018).

684   42.   Wheeler, D. W. *et al.* Hippocampome.org: A knowledge base of neuron types in the rodent
685         hippocampus. *Elife* **4**, (2015).
686   43.   de la Prida, L. M. & Ascoli, G. A. Explorers of the cells: Toward cross-platform knowledge
687         integration to evaluate neuronal function. *Neuron* **109**, 3535–3537 (2021).
688   44.   Chen, T. & Guestrin, C. XGBoost: A Scalable Tree Boosting System. *Proc. ACM SIGKDD*
689         *Int. Conf. Knowl. Discov. Data Min.* **13-17-Augu**, 785–794 (2016).
690   45.   Schmidhuber, J. Deep learning in neural networks: an overview. *Neural Netw.* **61**, 85–117
691         (2015).
692   46.   Hochreiter, S. & Schmidhuber, J. Long Short-Term Memory. *Neural Comput.* **9**, 1735–1780
693         (1997).
694   47.   H, F. J. Greedy Function Approximation: A Gradient Boosting Machine on JSTOR. *Ann.*
695         *Stat.* **29**, 1189–1232 (2001).
696   48.   Cortes, C., Vapnik, V. & Saitta, L. Support-vector networks. *Mach. Learn. 1995 203* **20**,
697         273–297 (1995).
698   49.   Cun, L. *et al.* Handwritten Digit Recognition with a Back-Propagation Network. *Adv. Neural*
699         *Inf. Process. Syst.* **2**, 396--404 (1990).
700   50.   de la Prida, L. M. *et al.* Threshold behavior in the initiation of hippocampal population bursts.
701         *Neuron* **49**, 131–142 (2006).
702   51.   McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection
703         for Dimension Reduction. (2018).
704   52.   Abbaspoor, S., Hussin, A. T. & Hoffman, K. L. Theta- and gamma-band oscillatory
705         uncoupling in the macaque hippocampus. *Elife* **12**, (2023).
706   53.   Dutta, S., Ackermann, E. & Kemere, C. Analysis of an open source, closed-loop, realtime
707         system for hippocampal sharp-wave ripple disruption. *J. Neural Eng.* **16**, (2019).
708   54.   Blanco, J. A. *et al.* Unsupervised classification of high-frequency oscillations in human
709         neocortical epilepsy and control patients. *J. Neurophysiol.* **104**, 2900–2912 (2010).
710   55.   Kucewicz, M. T. *et al.* High frequency oscillations are associated with cognitive processing
711         in human recognition memory. *Brain* **137**, 2231–2244 (2014).
712   56.   Liu, X. *et al.* E-Cannula reveals anatomical diversity in sharp-wave ripples as a driver for the
713         recruitment of distinct hippocampal assemblies. *Cell Rep.* **41**, (2022).
714   57.   Valero, M. *et al.* Mechanisms for Selective Single-Cell Reactivation during Offline Sharp-
715         Wave Ripples and Their Distortion by Fast Ripples. *Neuron* **94**, 1234-1247.e7 (2017).
716   58.   Modi, B. *et al.* Benchmarking algorithms that automatically detect sharp wave ripples. *Under*
717         *Rev.* (2023).
718

**Methods**

**The hackathon**

In order to explore a wide variety of ML solutions to the problem of SWR detection, we organized a hackathon (https://thebraincodegames.github.io/index_en.html). We specifically targeted people unfamiliar with SWR studies, who could provide unbiased solutions to the challenge. A secondary goal of the hackathon was to promote their interest and engagement at the interface between Neuroscience and Artificial Intelligence especially for future young scientists. The event was held in Madrid in October 2021, using remote web-platforms. Some of us (ANO) coordinated the event. Consent to participate and to share relevant personal data was obtained prior to the event. All participants were informed on the goal of the hackathon and agreed that their solutions were subject to subsequent investigation and modification.

The hackathon comprised 36 teams of 2 to 5 people (71% males, 29% female), for 116 participants in total. They represent 45% of Undergraduate students, 38% Master students, 15% PhD students and 3% non-academic workers (Fig.S1A). On average, they were young in their professional career with 77% of participants being research-oriented (Fig.S1A). Previous to the hackathon, we monitored the participants' self-declared knowledge level on Neuroscience, Python programming and ML in general using a survey (Fig.S1B). To provide a homogenous floor to address the challenge, we organized three online seminars to cover each of the three topics one month before the activity. Seminars were recorded and made available for review along the experience.

The hackathon was held during one weekend (Friday to Sunday), during which groups had to design and train a ML algorithm to detect SWRs. To standardize the different algorithms for future comparison, they were given Python functions to load the data, compute a performance score, and write results in a common format. Data sets were available from a public research-oriented repository at Figshare (https://figshare.com/authors/Liset_M_de_la_Prida/402282). Participants were given a training set to train their algorithms, and a test set to run validation tests. Data consisted on raw 8-channel LFP signals from the hippocampal CA1 region, recorded with high-density probes, which was used before for similar purposes (Navas-Olive et al. 2022). SWR were manually tagged to be used as ground truth (training set: 1794 events, two sessions from two mice; test set: 1275 events; two sessions from two animals). Since participants had two days to design and train solutions, groups were allowed to interact with us to ask for technical questions and clarification.

We monitored participant's engagement throughout the hackathon using short questionnaires. This allowed us to check their motivation and other emotional states (i.e., frustration, interest, etc…). Some people dropped out along the days of the hackathon (Fig.S1D). We found many participants felt confused and frustrated with the challenge, and this correlated with their performance, as a posterior analysis suggested (Fig.S1E).

**Datasets and ground truth**

Participants of the hackathon were provided with an annotated dataset consisting of raw LFP signals (8-channels) sampled at 30,000 Hz. SWR events were manually tagged by an expert who for each event identified their start and end. The start of the SWR was defined near the first ripple of the sharp-wave onset. The end of the event was defined at the latest ripple or when sharp-wave resumed. The training set consisted of two recording sessions from 2 mice (Navas-Olive et al., 2022). They contained 1794 manually tagged SWRs. The test set consisted of two recording sessions from another 2 mice and contained 1275 SWR events.

For posterior analysis of the results of the hackathon, we used a validation dataset consisting on the 2 test sessions mentioned before plus another 19 sessions for a total of 21 sessions from 8 different mice. They all contained a total of 7423 manually tagged SWR.

The ground truth, i.e. the analysis windows containing a SWR event, was generated for all sessions with the help of a Matlab 2019b tool that considered the window size.

**ML models specifications**

Five architectures were selected out of the 18 solutions submitted to the hackathon: XGBoost, SVM, LSTM, 2D-CNN and 1D-CNN. For the purpose of fair comparisons, they were retrained and tested using homogenized pre-processing steps and data management strategies (see below).

772  We used Python 3.9.13 with libraries Numpy 1.19.5, Pickle 4.0 and H5Py 3.1.0. To build the
773  different neural networks, we used the Tensorflow 2.5.3 library, with Keras 2.5.0 as the application-
774  programming interface. XGBoost 1.6.1 was used to train and test the boosted decision trees
775  classifiers. Scikit-learn 1.1.2 and Imbalanced-learn 0.9.1 were used to train support vector machine
776  classifiers. Analysis and training of the models were conducted on a personal computer (i7-11800H
777  Intel processor with 16 GB RAM and Windows 10).

778  **Data preparation**

779  For subsequent training and analysis of the architectures selected from the hackathon, all data was
780  pre-processed. From each recording session two matrices were extracted, X, with the raw LFP
781  data, shaped (# of timestamps, # of channels) and Y, the ground truth generated from the expert
782  tagging (# of timestamps). A timestamp of Y is 1 if a SWR event is present.

783  Values for matrix X were subsampled at 1250 Hz, taking into consideration that SWRs are events
784  that have frequencies in the range of 150 to 250 Hz. Before retraining the algorithms, data was z-
785  scored with the mean and standard deviation of the whole session.

786  **Training, validation, and test split**

787  For retraining the architectures, the same training dataset provided in the hackathon was used (2
788  sessions from 2 mice; 1794 SWR events). For initial testing, these two sessions were split
789  according to a70/30 train/validation design. To evaluate the generalization capabilities of the
790  models when presented with unseen data, we used several validation sessions, which provide the
791  necessary animal-to-animal, as well as within animal (sessions) variability. Validation sessions
792  included the 2 test dataset provided in the hackathon and 19 additional sessions (21 sessions from
793  8 mice, 7423 SWR events)

794  For retraining, the two training sessions were concatenated and divided into 60 seconds epochs.
795  Each epoch was assigned randomly to the train or validation set, following the desired split
796  proportion. The data was reshaped to be compatible with the required input dimensionality of each
797  architecture (see below). In order to evaluate model performance, two different datasets were used:
798  the test set described above (used for an initial screening of the 50-best models for each
799  architecture), and the validation set (used for generalization purposes).

800  Identification of SWR events in the data was implemented using analysis windows of different sizes.
801  To identify SWR events detected by the ML models, we set a probability threshold to identify
802  windows with positive and negative predictions. GT was annotated in the different analysis
803  windows of each session. Accordingly, predictions were classified in four categories: True Positive
804  (TP), when the prediction was positive and the GT window did contain a SWR event; False
805  Positive (FP), when the prediction was positive in a window that did not contain any SWR; False
806  Negative (FN), when the prediction was negative in a window with a SWR; and True Negative (TN)
807  when the prediction was negative and the window did not contain any SWR event.

808  If a positive prediction had a match with any window containing a SWR it was considered a TP, or it
809  was classified as FP otherwise. All true events that did not have any matching positive prediction
810  were considered FN. Negative predictions with no matching true events windows were TN.

811  With predicted and true events classified into those four categories, there are three measures than
812  can be used to evaluate the performance of the model. Precision (P), which was computed as the
813  total number of TPs divided by TPs and FPs, represents the percentage of predictions that were
814  correct.

$$Precision = \frac{TP}{TP + FP}$$

815  Recall (R), which was calculated as TPs divided by TPs and FNs, represents that percentage of
816  true events that were correctly predicted.

$$Recall = \frac{TP}{TP + FN}$$

817  Finally, the F1-score, calculated as the harmonic mean of Precision and Recall, represents the
818  network performance, penalizing imbalanced models.

$$F1 = \frac{2 + (Precision * Recall)}{Precision + Recall}$$

819 To ease subsequent evaluation of ML models for SWR analysis we provide open-access to codes
820 for retraining strategies: https://github.com/PridaLab/rippl-AI

**Parameter fitting**

822 Different combinations of parameters and hyper-parameters were tested for each architecture
823 during the training phase (1944 for XGBoost, 72 for SVM, 2160 for LSTM, 60 for 2D-CNN, 576 for
824 1D-CNN).

825 Two parameters were shared across all architectures: the number of channels and the number of
826 timestamps in the analysis window (referred as the window size). These parameters define the
827 dimensionality of the input data (# timesteps x # channels), i.e. the number of input features.

828 The number of channels to be used was set at 1, 3 or 8. When 1 channel was chosen, it was that
829 corresponding to the CA1 pyramidal layer channel, defined as the channel with most power in the
830 ripple bandwidth (150-250 Hz). The superficial, pyramidal, and deep channels were used as 3
831 channels. All the channels in the shank were used for the 8-channels input configuration.

832 The number of timestamps defines the window size. The tested values depended on each
833 architecture, and ranged between windows of 0.8 to 51.2 milliseconds. The rest of the parameters
834 were specific for each architecture (see below).

835 The F1-score metric for the training and test set was calculated to compare the performance of the
836 models, with the test F1 serving as a priori metric of the generalization of the models, allowing for a
837 selection of models without performing a complete validation.

838 For each model, a test-F1 array was calculated with different thresholds (generally, from 0.1 to 0.9
839 with 0.1 increments), and the highest value for each model was used for comparison among
840 models of the same architecture. As a result, the 50-best performing models were selected after
841 the initial retrained test.

**Validation process**

843 The aim of validation is to find the model that generalizes best to unseen data for each architecture.
844 With that in mind, defining a metric that takes this into account is not a straight-forward task.

845 To weigh each validation session (21) independently, a F1 array was calculated for each individual
846 session, resulting in matrix of 21 per number of threshold-values (#th). The mean of sessions gives
847 us a #th array that quantifies the performance/generalization of the model as a function of the
848 chosen threshold. The maximum value of this array will represent the best performance that could
849 be achieved with this model if the threshold is correctly selected. This single value is what will be
850 compared. Using this strategy, we narrowed down available models to the 10-best of each
851 architecture, before selecting the best model.

**XGBoost**

853 Based in the Gradient Boosting Decision Trees algorithm, this architecture trains a tree with a
854 subset of samples, and then calculates its output[44]. The misclassified samples are used to train a
855 new tree. The process is repeated until a predefined number of classifiers are trained. The final
856 model output is the weighted combination of individual outputs.

857 In the training process, we worked with quantitative features (LFP values per channels) and a
858 threshold value for a specific feature was considered in each training step. If this division correctly
859 classifies some samples of the subset, two new nodes were generated in the next tree level, where
860 the operation was repeated until the maximum tree depth was achieved, and a new tree with the
861 misclassified samples is generated. The input is one dimensional (# of channels x # of timesteps)
862 and produces a single output.

863 Specific parameters of XGBOOST are Maximum depth, the maximum levels for each tree, may
864 lead to overfitting. Learning rate, which controls the influence of each individual model in the
865 ensemble of trees. Gamma is the minimum loss reduction required to make a further partition on a
866 leaf node, with larger values leading to conservative models. Parameter λ contributes to the

21

867 regularization, with larger values preventing overfitting. Scale is used in imbalanced problems, the
868 larger the more penalized false negatives are during training.

869 Trained models had a number of trainable parameters ranging from 1500 to 17900.

**SVM**

871 Support Vector Machine is a classical classifier that searches for a hyperplane in the input
872 dimensionality that maximizes the separation between different classes. This is only possible in
873 lineal separation problems, so some misclassifications are allowed in real tasks. Usually, SVM
874 performs a transformation on the original data using a kernel (linear or otherwise) that increases
875 the data dimensionality but facilitates classification.

876 During training, the parameters that define the separation hyperplane are updated until the
877 maximum number of iterations is achieved, or the rate of change in the parameters go below a
878 threshold. The input is one-dimensional (# of channels x # of timesteps) and produces a single
879 output.

880 Specific parameters of SVM are the kernel type. Using nonlinear kernels resulted in an explosive
881 growth in training and predicting times, due to the enormous number of training data points. Only
882 the linear kernel produced manageable times. The under-sample proportion rules out negative
883 samples (windows without ripple) until the desired balance is achieved: 1 indicates the same
884 number of positives and negatives.

885 Trained models had a number of trainable parameters ranging from 1 to 480.

**LSTM**

887 Recurrent Neural Networks (RNNs) are a subtype of NNs especially suited to work with temporal
888 series of data, extracting the "hidden" relations and tendencies between non-contiguous instants.
889 Long Short-Term Memory (LSTMs) are RNNs with modifications that prevent some associated
890 problems[46].

891 During training three sets of weights and biases are updated in each LSTM unit, associated with
892 different "gates" (Forget, input and output). To prevent overfitting, two layers of dropout (DP) and
893 batch normalization (batchNorm) were inserted between LSTM layers. DP randomly prevents
894 some outputs to propagate to the next layer. BatchNorm normalizes the output of the previous
895 layer. The final layer is a dense layer that outputs the event probability.  The input is two-
896 dimensional (# of timesteps, # of channels) and produces a probability for each timestep. After
897 each window the internal weights are reset.

898 Specific LSTM parameters: bidirectional if the model process the windows forwards and backwards
899 simultaneously; # of layers is the number of LSTM layers; # of units the number of LSTM units in
900 each layer, and # of epochs, which is the number of times the training data is used to perform
901 training.

902 Trained models had a number of trainable parameters ranging from 156 to 52851.

**2D-CNN**

904 Convolutional Neural Networks use convolutional layers consisting of kernels (spatial filters) to
905 extract the relevant features of an image[49]. Successive layers use this as inputs to compute
906 general features of the image. This 2D-CNN moves the kernels along the two axes, temporal
907 (timesteps) and spatial (channels). The first half of the architecture includes MaxPooling layers that
908 reduce the dimensionality and prevent overfitting. A batchNorm layer follows every convolutional
909 layer. Finally, a dense layer produces the event probability of the window.

910 During training, the weights and biases of every kernel are updated to minimize the loss function,
911 with was taken as the binary cross entropy:

$$H_p(q) = \frac{-1}{N} \sum_{i=1}^{N} y_i \cdot log\big(p(y_i)\big) + (1 - y_i) \cdot log\big(1 - p(y_i)\big)$$

912 N is the number of windows in the training set, $y_i$ is the label of the *i* window and $p(y_i)$ is the
913 probability of ripple that the model predicts. The input is # of timesteps and # of channels; and
914 produces a single probability for each window.

915 The 2D-CNN was tested with a fixed number of layers and kernel dimension. The kernel factor
916 parameter determined the number of kernels in this structure: 32xkf (2x2), 16xkf (2x2), 8xkf (3x2),
917 16xkf (4x1), 16xkf (6x1) and 8xkf (8x1). In parenthesis the size of the kernels in each layer.

918 Trained models had a number of trainable parameters ranging from 1713 to 24513.

919 **1D-CNN**

920 This model is also a convolutional neural network, but the kernels only move along the temporal
921 axis while processing spatial information. The number of layers and the kernel size was fixed. The
922 tested models had 7 sets of 1D convolutional layer, batchNorm and LeakyRelu layer, followed up
923 by a dense sigmoid activation unit. This model is similar to our previous CNN solution[31].

924 During training, the weights and biases of the layers were also updated with the objective of
925 minimizing the binary cross entropy. The input is # of timesteps and # of channels, and produces a
926 single probability for each window.

927 The specific parameters for 1D-CNN included the kernel factor, which defined the number of
928 kernels in each conv layer. The size and stride for each layer was equal and fixed. The size of the
929 kernels in the first layer was defined as the length of the input window divided by 8. Structure: 4x*kf*
930 (# timesteps//8 x # timesteps//8), 2x*kf* (1x1), 8x*kf* (2x2), 4x1 (1x1), 16x*kf* (2x2), 8x*kf* (1x1) and
931 32x*kf* (2x2). Parameters also include # of epochs, as the number of times the training data is used
932 to perform training, and # of training batch samples, which is the number of windows that are
933 processed before parameter updating.

934 Trained models had a number of trainable parameters ranging from 342 to 4253.

935 **Characterization of SWR features**

936 SWR properties (ripple frequency and power) were computed using a 100 ms window around the
937 center of the event, measured at the pyramidal channel of the raw LFP. Preferred frequency was
938 computed first by calculating the power spectrum of the 100 ms interval using the enlarged
939 bandpass filter 70 and 400 Hz, and then looking for the frequency of the maximum power. In order
940 to account for the exponential power decay in higher frequencies, we subtracted a fitted
941 exponential curve ('fitnlm' from MATLAB toolbox) before looking for the preferred frequency. To
942 estimate the ripple power, the spectral contribution was computed as the sum of the power values
943 for all frequencies lower than 100 Hz normalized by the sum of all power values for all frequencies
944 (of note, no subtraction was applied to this power spectrum).

945 **Dimensionality reduction using UMAP**

946 To classify SWR, we used topological approaches[14]. The UMAP version 0.5.1 (https://umap-
947 learn.readthedocs.io/en/latest/) in Python 3.8.10 Anaconda was used, which is known to properly
948 preserve local and global distances while embedding data in a lower dimensional space. In all cas-
949 es, we used default values for reconstruction parameters. Algorithms were initialized randomly.
950 UMAP provided robust results independent on initialization.

951 **Prediction and retraining of non-human primate data-set**

952 To study the generalization capabilities of the different architectures, we used data from a freely
953 moving macaque targeting similar CA1, as completed in our mouse data (methods are described in
954 reference[52]). Recordings were obtained with a 64-ch linear polymer probe (custom 'deep array
955 probe', Diagnostic Biochips) that recorded across the CA1 layers of the anterior hippocampus
956 (Fig.6A) where layers were identifiable relative to the main pyramidal layer, which contains the
957 greatest unit activity and SWP power. LFP signals were sampled at 30 kHz using a Freelynx
958 wireless acquisition system (Neuralynx, Inc). Data corresponds to periods of immobility for a
959 duration of almost 2 hours and 40 minutes, predominantly comprised of sleep in overnight housing.
960 LFP intervals presenting a high level of noise across all channels was not considered for analysis.

961 Similar to the procedures used in mice, SWR beginning and ending times were manually tagged
962 (ground truth). First, the best model of each architecture, already trained with the mouse data, was

963 used to predict the output of the primate data with no retraining. For this purpose, we used
964 recordings of different channels around the CA1 pyramidal channel, and matched to meet the
965 laminar organization of the dorsal mouse hippocampus. Specifically, we used a CA1 radiatum
966 channel, 720 µm from the pyramidal layer, three channels in the pyramidal layer, at +90µm, +0µm
967 and -90µm from the pyramidal channel, and a stratum oriens channel 720µm from the pyramidal
968 channel. The pyramidal channel was defined at the site with the maximal ripple power. We
969 complemented these 5 recordings with 3 more interpolated signals, making a total of 8 input
970 channels [oriens, interpolated, pyramidal, pyramidal, pyramidal, interpolated, interpolated, radiatum]
971 using a linear interpolation script available at Github: https://github.com/PridaLab/rippl-
972 AI/blob/main/aux_fcn.py. The applied pre-processing was the same as with the mice data:
973 subsampling to 1250Hz and z-score normalization.

974 With the aim of studying the effect of retraining with completely different data, we retrained the
975 models. Data was split in three sets (50% training, 20% test, 30% validation), and used to retrain
976 and validate the models. For re-training, we reset all trainable parameters (internal weights) but
977 kept all architectural hyper-parameters fixed (input number of channels, input window length,
978 number of layers, etc…) as with the mouse data, making the re-training process much faster than
979 the original training that required a deep hyper-parametric search (per model re-train: 2min for
980 XGBoost, 10-30min for SVM, 3-20min for LSTM, 1-10 min for 2D-CNN and 1-15 min for 1D-CNN).

981 **Code and data availability**

982 Data and codes used in this study are available. The training and test set data are available at
983 https://figshare.com/authors/Liset_M_de_la_Prida/402282 and listed independently as follows:

984 M de la Prida, Liset (2021): Amigo2_2019-07-11_11-57-07. figshare. Dataset.
985 https://doi.org/10.6084/m9.figshare.16847521.v2
986 M de la Prida, Liset (2021): Som2_2019-07-24_12-01-49. figshare. Dataset.
987 https://doi.org/10.6084/m9.figshare.16856137.v2
988 M de la Prida, Liset (2021): Dlx1_2021-02-12_12-46-54. figshare. Dataset.
989 https://doi.org/10.6084/m9.figshare.14959449.v4
990 M de la Prida, Liset (2021): Thy7_2020-11-11_16-05-00. figshare. Dataset.
991 https://doi.org/10.6084/m9.figshare.14960085.v1
992
993 Codes for some of the best trained models of all architectures are available in an open-source
994 repository https://github.com/PridaLab/rippl-AI and documented in open-source notebooks for
995 model retraining https://github.com/PridaLab/rippl-AI/blob/main/examples_retraining.ipynb and for
996 SWR detection https://github.com/PridaLab/rippl-AI/blob/main/examples_detection.ipynb
997
998 **Acknowledgements**

1007