



This is a repository copy of *A machine-learning-based approach to predict early hallmarks of progressive hearing loss*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/229110/>

Version: Accepted Version

Article:

Ceriani, F. orcid.org/0000-0002-5366-341X, Giles, J., Ingham, N.J. et al. (5 more authors) (2025) A machine-learning-based approach to predict early hallmarks of progressive hearing loss. *Hearing Research*, 464. 109328. ISSN 0378-5955

<https://doi.org/10.1016/j.heares.2025.109328>

© 2025 The Authors. Except as otherwise noted, this author-accepted version of a journal article published in *Hearing Research* is made available via the University of Sheffield Research Publications and Copyright Policy under the terms of the Creative Commons Attribution 4.0 International License (CC-BY 4.0), which permits unrestricted use, distribution and reproduction in any medium, provided the original work is properly cited. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>

Reuse

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

1 **A machine-learning-based approach to predict early hallmarks of progressive hearing loss**

2

3 Federico Ceriani^{1,3#}, Joshua Giles^{2,3}, Neil J Ingham^{4,5}, Jing-Yi Jeng¹, Morag A Lewis^{4,5}, Karen P
4 Steel^{4,5}, Mahnaz Arvaneh^{2,3,6}, Walter Marcotti^{1,6}

5

6 ¹*School of Biosciences, University of Sheffield, Sheffield, S10 2TN, UK*

7 ²*Department of Automatic Control and Systems Engineering, University of Sheffield, Sheffield, S1
8 4DT, UK*

9 ³*Centre for Machine Intelligence, University of Sheffield, S10 2TN, UK7*

10 ⁴*Wellcome Sanger Institute, Hinxton, UK*

11 ⁵*Wolfson Sensory, Pain and Regeneration Centre, Guy's Campus, King's College London, London
12 SE1 1UL, UK*

13 ⁶*Neuroscience Institute, University of Sheffield, Sheffield, S10 2TN, UK*

14

15

16

17 **#Corresponding authors:**

18 Federico Ceriani (f.ceriani@sheffield.ac.uk), ORCID identifier: 0000-0002-5366-341X

19

20 **Abstract**

21 Machine learning (ML) techniques are increasingly being used to improve disease diagnosis and
22 treatment. However, the application of these computational approaches to the early diagnosis of age-
23 related hearing loss (ARHL), the most common sensory deficit in adults, remains underexplored.
24 Here, we demonstrate the potential of ML for identifying early signs of ARHL in adult mice. We
25 used auditory brainstem responses (ABRs), which are non-invasive electrophysiological recordings
26 that can be performed in both mice and humans, as a readout of hearing function. We recorded ABRs
27 from C57BL/6N mice (6N), which develop early-onset ARHL due to a hypomorphic allele of
28 *Cadherin23* (*Cdh23^{ahl}*), and from co-isogenic C57BL/6NTac^{*Cdh23⁺*} mice (6N-Repaired), which do not
29 harbour the *Cdh23^{ahl}* allele and maintain good hearing until later in life. We evaluated several ML
30 classifiers across different metrics for their ability to distinguish between the two mouse strains based
31 on ABRs. Remarkably, the models accurately identified mice carrying the *Cdh23^{ahl}* allele even in the
32 absence of obvious signs of hearing loss at 1 month of age, surpassing the classification accuracy of
33 human experts. Feature importance analysis using Shapley values indicated that subtle differences in
34 ABR wave 1 were critical for distinguishing between the two genotypes. This superior performance
35 underscores the potential of ML approaches in detecting subtle phenotypic differences that may elude
36 manual classification. Additionally, we successfully trained regression models capable of predicting
37 ARHL progression rate at older ages from ABRs recorded in younger mice. We propose that ML
38 approaches are suitable for the early diagnosis of ARHL and could potentially improve the success
39 of future treatments in humans by predicting the progression of hearing dysfunction.

40

41

42 **Keywords:** age-related hearing loss, auditory brainstem responses, machine learning, diagnosis.

43

44

45 1. Introduction

46 Progressive hearing loss results in a decrease in hearing sensitivity and ability to understand speech.
47 Among the different forms of progressive hearing loss, age-related hearing loss (ARHL) is the most
48 common sensory deficit in humans, affecting communication and leading to social isolation,
49 depression and diminishing cognitive abilities (Gates & Mills, 2005; Livingston *et al.* 2024).
50 Currently, there are no treatments to prevent or cure ARHL (Wang & Puel, 2020). ARHL is a
51 heterogeneous dysfunction, which results from the cumulative effects of ageing on the auditory
52 system, such as cellular senescence, as well as additional intrinsic (e.g. genetic predisposition, Ingham
53 *et al.* 2019) and extrinsic (e.g. environmental noise) factors. Because of this complex aetiology, the
54 progression of the disease varies between individuals, resulting in different severity and degree of
55 progression of hearing loss. Hearing function in clinical and pre-clinical settings can be examined
56 through a non-invasive electrophysiological test based on the auditory brainstem response (ABR).
57 However, the effects of hearing loss, other than an obvious increase in auditory thresholds, are often
58 difficult to detect using ABR tests. Thus, ARHL is normally diagnosed only after patients start losing
59 key hearing abilities, such as being unable to distinguish words in noisy conditions. This is usually
60 an indication that some severe or irreversible damage has already happened to the sensory cells or
61 neurons that send sound information to the brain. Therefore, as we develop therapies to target ARHL,
62 such as gene-based replacement interventions or small molecules (Lv *et al.* 2024; Schilder *et al.*
63 2024), there is also a pressing need to improve the diagnostic tools to detect and predict the
64 progression of the dysfunction at an early stage. As with any medical condition, treating a disease in
65 its early stages increases the likelihood of successful treatment.

66 Machine learning (ML) techniques are increasingly being explored as tools to improve disease
67 diagnosis and treatment (Goecks *et al.* 2020; Sidney-Gibbons & Sidney-Gibbons, 2019). These
68 techniques leverage advanced algorithms to analyse large datasets, uncovering patterns that may be
69 elusive even to well-trained experts. By identifying complex features in high-dimensional clinical
70 data that correlate strongly with patient phenotypes, ML algorithms can be developed to predict the
71 presence of a disease (Banerjee *et al.* 2023). In the auditory field, significant progress is being made
72 in applying ML to hearing healthcare and research (Chen *et al.* 2021; Shew *et al.* 2019, Cha *et al.*
73 2019, Crowson *et al.* 2023, Chen *et al.* 2024), and there is a growing emphasis on the leveraging of
74 ML-based digital tools to automate hearing assessment (Wasmann *et al.* 2022). However, the
75 potential of these computational techniques to develop diagnostic tools for the early detection of
76 progressive forms of hearing loss remains largely unexplored.

77 Here, we applied ML to ABR data with the goal of detecting early signs of ARHL in mice and
78 forecasting its progression. We recorded ABRs from the commonly used C57BL/6N (6N) mouse

79 strain and from the co-isogenic strain C57BL/6NTac^{Cdh23⁺} (6N-Repaired, [Mianné et al. 2016](#)) at 1, 3,
80 6, 9 and 12 months of age. The 6N mice carry a hypomorphic allele in the Cadherin 23 gene (*Cdh23^{ahl}*,
81 [Johnson et al. 1997](#); [Noben-Trauth et al. 2003](#)), which leads to progressive early-onset hearing loss
82 starting from about 3-6 months of age. Similar to ARHL in humans ([Gates & Mills, 2005](#)), the
83 progression of hearing loss in 6N mice begins at the higher frequencies and worsens over time,
84 resulting in profound hearing loss by 15 months of age ([Jeng et al. 2020a](#); [2020b](#); [Jeng et al. 2021](#)).
85 In contrast, the co-isogenic 6N-Repaired strain, which are corrected for the *Cdh23^{ahl}* mutation using
86 CRISPR/Cas9 ([Mianné et al. 2016](#)), maintains better hearing than 6N mice into old age, especially
87 for tone sensitivity for frequencies of 12 kHz and above ([Mianné et al. 2016](#); [Jeng et al. 2020b](#)). We
88 trained ML models through supervised learning using longitudinal ABR data as input features and
89 genotype (i.e., mouse strain, 6N or 6N-Repaired) as target outputs. We demonstrate that, by
90 recognising anomalies in the ABRs, the ML models were able to detect the mice with the *Cdh23^{ahl}*
91 allele in the very early stages of ARHL. This approach was validated on unseen data of two
92 independently acquired datasets, demonstrating the broad validity and generalisability of our
93 conclusions. Finally, we used ML to forecast the future progression of the hearing capabilities of
94 young adult mice up to 1 year of age. This work highlights the benefit of using ML for the early
95 diagnosis of ARHL, providing a foundation for future studies exploring its applicability to human
96 datasets.

97 2. Methods

98 2.1. Ethical Statement

99 The animal work was licensed by the UK Home Office under the Animals (Scientific Procedures)
100 Act 1986 (Sheffield: PCC8E5E93 and PP1481074; King's College London: P053FFC4C) and was
101 approved by the relevant Ethical Review Committees (University of Sheffield: 180626_Mar). Mice
102 had unlimited access to food and water. For the *in vivo* recording of auditory brainstem responses
103 (ABRs), mice were anaesthetised using intraperitoneal injection of ketamine (100 mg/Kg body
104 weight, Fort Dodge Animal Health, Fort Dodge, USA) and xylazine (10 mg/Kg, Rompun 2%, Bayer
105 HealthCare LLC, NY, USA). At the end of the *in vivo* recordings, mice were either culled by cervical
106 dislocation or recovered from anaesthesia with intraperitoneal injection of atipamezole (1 mg/Kg).
107 Mice under recovery from anaesthesia were returned to their cage, placed on a thermal mat and
108 monitored over the following 2-5 hrs.

109

110 2.2. Auditory brainstem responses

111 Two independent datasets of auditory brainstem responses (ABRs) from different mouse cohorts
112 were used in this study. ABRs from the **primary cohort** were collected at the University of Sheffield
113 from 104 female mice (50 6N and 54 6N-Repaired mice). For all the mice, ABR recordings were
114 performed at 1 month of age, and for some, recordings were also performed at 3, 6, 9 and 12 months
115 of age. These mice were born over a period of 5 months and were housed in the same room within
116 the animal facilities at the University of Sheffield, thus experiencing similar levels of noise exposure
117 throughout the duration of the study. ABRs from the **replication cohort** were collected at King's
118 College London from both males and females at 1 month of age (85 6N and 103 6N-Repaired mice).

119 Following the onset of anaesthesia (see *Ethics statement* above) and the loss of the retraction reflex
120 with a toe pinch, mice were placed onto a heat mat (37°C) in a soundproof chamber (MAC-3 acoustic
121 chamber, IAC Acoustic, UK). Subdermal electrodes were placed under the skin behind the pinna of
122 each ear (reference and ground electrode) and on the vertex of the mouse (active electrode) as
123 previously described (Ingham *et al.* 2019; Ingham *et al.* 2011). Sound stimuli were delivered to the
124 ear by calibrated loudspeakers (MF1-S, Multi Field Speaker, Tucker-Davis Technologies, USA)
125 placed directly in front of the mouse 10 cm (Sheffield) or 20 cm (King's College London) from the
126 nose. Sound pressure was calibrated with a low-noise microphone probe system (ER10B+, Etymotic,
127 USA). Experiments were performed using a customised software (Ingham *et al.* 2011) driving an
128 RZ6 auditory processor (Tucker-Davis Technologies). Auditory thresholds were estimated from the
129 resulting ABR waveform and defined as the lowest sound pressure level (measured in decibel, dB
130 SPL) where any recognisable feature of the waveform was visible. Responses were measured for

131 clicks (which cover a broad range of frequencies, 0.01 ms duration) and pure tones at frequencies of
132 3, 6, 12, 18, 24, 30, 36 and 42 kHz (5 ms duration, 1 ms rise/fall time). Stimulus sound pressure levels
133 were typically 0-95 dB SPL, presented in steps of 5 dB SPL. The brainstem response signal was
134 averaged over 256 repetitions. Tone bursts were 5 ms in duration with a 1 ms on/off ramp time, which
135 was presented at a rate of 42.6/sec. The order of sound stimulus presentation was consistent for all
136 ABR recordings. Responses to click stimuli were recorded first (0 to 95 dB), followed by pure tones
137 from 15 dB to 95 dB at 3 kHz and 6 kHz. Finally, stimuli at varying frequency from high (42 kHz)
138 to low (12 kHz) were presented. This process was repeated in 5 dB increments from 15 dB to 95 dB.

139 Wave 1 amplitudes and latencies were measured using a semiautomatic approach using custom
140 software (doi:10.5281/zenodo.12606227). Automatic identification was manually reviewed and, if
141 required, adjusted to the correct peak. Wave 1 amplitude was calculated as the difference between
142 the amplitude of the first peak and the first trough of the ABR waveform; the latency was calculated
143 as the delay of the Wave 1 peak from the beginning of the recording.

144 To evaluate the models on a dataset different from the one it was trained on, we found that an
145 alignment procedure was required to maximise model accuracies. This alignment required the shifting
146 of the ABR waveforms from the replication cohort to the left by 0.55 ms. This was likely due to the
147 different distances between mouse and the speaker (~10 cm, accounting for 0.29 ms time difference)
148 and variations in hardware (e.g., differences in electrical delays and timing of sound delivery). The
149 alignment was achieved by removing 54 timepoints at the beginning of the replication cohort trace
150 and at the end of the primary cohort trace (to maintain the same number of features between the two
151 datasets). This procedure, which did not alter the shape and time course of the ABR waveforms, was
152 sufficient to align waveform peaks between the two cohorts. The parameters for feature shifting were
153 determined from ABRs of 6N-Repaired mice from the training datasets from the two cohorts. Note
154 that this transformation, while ensuring correspondence of the features in the two datasets, does not
155 eliminate latency differences between the two strains, as both are shifted by the same amount.

156

157 **2.3. Implementation and evaluation of machine learning models**

158 We developed ML models to address two different tasks: 1) identifying which mice carried the
159 *Cdh23^{ahl}* mutation from their ABRs (classification task) and 2) forecasting the future progression of
160 hearing function in mice (regression task). For the classification task, we used six different classifiers
161 as the basis of our models: Random Convolutional Kernel Transform (ROCKET), Hierarchical Vote
162 Collective of Transformation-based Ensembles V2 (HIVE COTE V2.0), Extreme Gradient Boosting
163 (XGBoost), Random Forest (RF), multilayer perceptron (MLP) and Support Vector Machine (SVM)
164 classifier.

165 The ROCKET algorithm (Dempster *et al.* 2020) classifies time series by applying numerous
166 random filters to the data, extracting key values from the resulting feature maps, and using these
167 values to train a simple, efficient model. This approach is both fast and accurate, transforming
168 complex time series into an easy-to-handle format for effective classification. HIVE COTE V2.0 is a
169 heterogeneous meta ensemble for time series classification (Middlehurst *et al.* 2021) that builds on
170 the ROCKET. This algorithm is more computationally expensive but has shown very high levels of
171 accuracy when utilised to classify other time series data. XGBoost (Extreme Gradient Boosting, Chen
172 *et al.* 2016) and Random Forest are ensemble models that use decision trees as base learners which
173 are widely used due to their accuracy and interpretability. They have proven successful in the
174 classification of electrophysiological data (Edla *et al.* 2018) for a range of applications. Along with
175 these machine learning algorithms the MLP classifier was selected as a simple neural network option
176 to explore. Finally, we selected the SVM classifier as a good base classifier which has been proven
177 to work effectively with minimal computational cost in a wide range of applications. For the
178 regression task, we trained random forest regressor models to predict the outcome of three continuous
179 parameters at different ages (see below). Some of the hyperparameters were tuned using 5-fold grid
180 search cross validation applied to the training set, optimising for F1 score (classification task) or the
181 negative mean squared error (regression task). The hyperparameters of the models are summarised in
182 **S6 Table**.

183 Models were implemented using the *scikit-learn* (Random Forest classifier and regressor, SVM,
184 Multi-layer perceptron, Pedregosa *et al.* 2011), *xgboost* (XGBoost, Chen *et al.* 2016) and *sktime*
185 (HIVE COTE V2.0, ROCKET, Löning *et al.* 2019) python packages. Different ABR waveforms
186 resulting from stimulation with individual intensities/frequencies combinations were concatenated in
187 a single univariate trace and used as input features for the ML models (**Fig 1b, Fig 3a, Fig 5a, Fig**
188 **8d**).

189 For the classification task, ML algorithms were trained through supervised learning using the
190 concatenated ABR waveforms of 1-month old animals as input features and the genotype (6N vs 6N-
191 Repaired) as labels. In all analyses, the “6N” class (i.e., “mice with early-onset ARHL”) was treated
192 as the positive class. All classifiers were preceded by an ANOVA F-test feature selection step, which
193 retained 10% of the features (i.e., ABR timepoints). In this step, the F-statistics scores were calculated
194 between the two classes for every input feature and ranked, and only the 10% top scoring features
195 were preserved as model inputs. By reducing the dimensionality of the dataset, focusing on the most
196 relevant predictors, this method is effective in improving training time and reducing overfitting.

197 We trained and tested the models on data from two laboratories (primary and replication cohorts)
198 either separately or combined. In every instance, the datasets were randomly split into train and test

199 data, with the training data containing 75% of the mice and the test the remaining 25%. For tasks
200 which involved evaluating the models on data of different laboratories, the whole primary and
201 replication cohort datasets were used either for training or testing (see **Results**).

202 We provided two separate evaluations of the models. First, in order to ensure models were working
203 correctly, we performed repeated 5-fold cross validation on the training set collecting 4 different
204 metrics:

- 205 • recall (also called sensitivity, or true positive rate), defined as $\frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$, i.e. the
206 ability of a classifier to correctly identify the positive class (6N);
- 207 • specificity (or true negative rate), defined as $\frac{\text{True negatives}}{\text{True negatives} + \text{False positives}}$, i.e. the ability of a classifier
208 to correctly identify the negative class (6N-Repaired);
- 209 • precision, defined as $\frac{\text{True positives}}{\text{True positives} + \text{False positives}}$, a measure of the accuracy of positive predictions;
- 210 • receiver operating characteristic area under the curve (ROC AUC), i.e. the area under the false
211 positive rate (1-specificity) vs. recall curve, which offers a threshold-independent measure of a
212 model's performance.

213 For all classifiers, the threshold for distinguishing between the two classes was set to a probability
214 of 50%. Models were then trained on the whole training set and confusion matrices were calculated
215 on the test set. We opted to include both cross validation and test set assessment as the relatively
216 small size of the test set limits our ability to draw strong conclusions from its performance alone.
217 However, it allows for comparison with manual classification on the same set of mice.

218 For the regression task, ML algorithms were trained through supervised learning using the
219 concatenated click ABR waveforms recorded at 1 month and 3 months of age as input features and
220 three different parameters measured at 6, 9 or 12 months of age (see **Results**) as targets. For wave 1
221 latency prediction, values for waveforms below the auditory threshold were imputed using the highest
222 latency measured for click stimuli for each mouse. For one mouse for which no detectable ABR signal
223 was present at 12 months of age, missing latency values were imputed using the highest latency for
224 click stimuli observed across all other mice in the dataset. For wave 1 amplitude prediction, values
225 for waveforms below the auditory thresholds were set to zero. Similarly to the classification task,
226 regression models were preceded by a feature selection step based on univariate linear regression
227 tests, which return F-statistics and p-values, and only the 10% top scoring features were preserved as
228 model inputs. The datasets were randomly split into train and test data, with the training data
229 containing 75% of the mice and the test the remaining 25%. Performances of regression models were
230 evaluated as mean absolute error (MAE) averaged across the results of a repeated 5-fold cross

231 validation step on the training set (5 repeats). Moreover, the coefficient of determination (R^2) and
232 MAE were calculated on predictions made on the test set.

233 The Shapley Additive explanations (SHAP) method implemented in the *shap* python module was
234 used for Shapley value estimation (Lundberg & Lee, 2017). The TreeSHAP method was used to
235 estimate Shapley values for tree-based models (Random forest and XGBoost), while KernelSHAP
236 was used for the SVM model. Shapley values were calculated on the test set, using the training set as
237 the background distribution. Feature importances were calculated by averaging the absolute Shapley
238 values computed across all train instances. Feature importances were smoothed with a Savitzky-
239 Golay filter with polynomial order equal to 1 and window size of 0.42 ms (“Global” models) and 0.22
240 ms (“Click” models) for visualisation purposes.

241 Machine learning model implementation, data analysis and figure plotting were conducted using
242 python (version 3.11.8) primarily utilising the scikit-learn (version 1.4.1) and sktime (0.27.0)
243 modules. Computations were performed on a MacBook Pro with M1 processor and 16 GB of RAM
244 (MacOS 15.0, kernel Darwin 24.0.0), and on a workstation equipped with an Intel Xeon Silver 4210R
245 CPU and 256 GB of RAM (Windows 11 Pro for workstations).

246

247 ***2.4. Comparison between ML and manual classification***

248 Three human annotators were asked to blindly label the ABR dataset for comparison with the ML
249 models. Each annotator reviewed ABR data for all samples in both the training and test sets and
250 categorized each instance according to the mouse strain. Annotators were presented with randomised
251 ABR stacks containing responses to click and pure tones (**Fig 3**) or clicks alone (**Fig 5**). Individual
252 ABR waveforms were shown to human experimenters without additional overlays or statistical
253 summaries. Annotators assessed auditory thresholds as part of the classification. Each of the three
254 annotators was an expert in mouse ABR recording and analysis, with specific knowledge of the
255 progressive high frequency hearing loss phenotype of 6N mice compared to 6N-Repaired mice.
256 Predictions on the test set ABRs alone were directly compared to those of the six “Global” and
257 “Click” models (**Fig 3f**, **Fig 5f**). For the “Global” dataset, average predictions (\pm SD) of the three
258 annotators on the whole (train and test) set were as follows: recall 66.7% \pm 13.6%; specificity 82.7%
259 \pm 18.9%; precision 81.9% \pm 14.1%. Average predictions for the “Click” dataset on the whole (train
260 and test) set were: recall 44.7% \pm 8.3%; specificity 67.9% \pm 7.4%; precision 56.4% \pm 1.7%. The
261 results from the three annotators were subsequently evaluated for inter-rater reliability using Fleiss’
262 Kappa, a statistical measure of agreement among multiple raters.

263

264 ***2.5. Statistical analysis***

265 Statistical analysis of experimental data was conducted using either the Aligned Rank Transform
266 (ART) ANOVA, followed by Wilcoxon rank-sum tests with Holm-Bonferroni correction for pairwise
267 post-hoc comparisons, or standard ANOVA with Tukey's Honestly Significant Difference post-hoc
268 test. For model performance comparisons on the primary cohort, cross validation metrics were
269 evaluated using the Friedman test followed by the Nemenyi post-hoc test. To compare model
270 performances across different datasets (primary cohort/replication cohort/combined or
271 “Global”/“Click”) and model types, we used mixed-effects linear models with model type and dataset
272 as fixed effects and the cross-validation splits as random effects to account for non-independence
273 within repeated measures. A significance level of $P < 0.05$ was used to determine statistical
274 significance.

275 3. Results

276 To demonstrate the possibility of ML to detect early signs of a progressive form of deafness, we
277 first acquired ABRs using standardised protocols from a cohort of 1-month-old 6N ($n = 50$) and 6N-
278 Repaired ($n = 54$) mice (**primary cohort, Fig 1a**). To avoid the impact of sex variability in the
279 progression and severity of age-related hearing loss (Nolan, 2020), only female mice were included
280 in the primary cohort. A subset of these mice was aged up to 1 year and ABRs were recorded at
281 regular intervals (3, 6, 9 and 12 months). This approach allowed us to test the efficacy of ML models
282 in learning to distinguish between the two mouse strains based solely on ABRs of 1-month-old mice
283 (classification task). Moreover, we tested whether it was possible to predict the progression and
284 degree of hearing loss at older ages based on ABR data from young adult mice (regression task). To
285 achieve this, we trained several ML algorithms through supervised learning using: 1) ABR data of 1-
286 month old animals as input features and the genotype as labels (classification task) and 2) ABR data
287 recorded at 1 and 3 months of age as input features and ABR characteristics at older ages as target
288 (regression task, **Fig 1b**).

289

290 3.1. Auditory thresholds of 6N and 6N-Repaired mice between 1 and 12 months of age

291 We first determined the progression of hearing loss of our primary cohort of 50 6N and 54 6N-
292 Repaired female mice (**Fig 2**). We found that at 1 month of age both mouse strains showed similar
293 ABR thresholds (**Fig 2a**) and waveforms across most sound stimuli (**S1 Fig**), except for a small
294 increase in the median threshold of 6N mice at the two highest frequencies tested (36 kHz and 42
295 kHz, $P < 0.0001$, pairwise Wilcoxon rank-sum test, ART ANOVA, **Fig 2a**). Difference between the
296 audiograms of the two strains became progressively more evident at older ages (**Fig 2b-e, S2 Fig**).
297 At 3 months, several 6N mice had undetectable ABRs at all intensities for stimuli of 36 and 42 kHz
298 (36 kHz: 19 6N mice out of 50; 42 kHz: 29 mice out of 50) and significantly raised threshold at 30
299 kHz compared with 6N-Repaired mice (**Fig 2b**). Between 6 and 12 months, ABR thresholds were
300 significantly different between the two strains for all stimuli except the lowest frequency tested (3
301 kHz) at 6 and 9 months (**Fig 2c-e**). As previously shown (Jeng *et al.* 2020), we found that the
302 progression of hearing loss was variable in 6N mice (**Fig 2f**), with threshold differences between
303 individual mice of up to 75 dB from 6 months onwards (**Fig 2c-e**).

304

305 3.2. Using ML models to predict the presence of the *Cdh23^{ahl}* allele from ABRs of young mice

306 We then sought to train ML models through supervised learning to classify ABR recordings taken
307 at the earliest timepoint (1 month) based on mouse strain (6N or 6N-Repaired), thereby predicting the
308 presence of the ARHL-linked *Cdh23^{ahl}* allele. To demonstrate the generalisability of this approach,

309 we tested six different ML models: four commonly used classifiers (random forest, XGBoost, support
310 vector machine (SVM) and multi-layer perceptron (MLP)) and two time-series-specific classifiers
311 (HIVE-COTE V2.0 and ROCKET) (see **Methods** for a description of each model and hyperparameter
312 tuning procedure). All classifiers were preceded by an ANOVA F-test feature selection step, which
313 retained 10% of the features (i.e., ABR timepoints, **Fig 3a, Fig 4**, see **Methods**). We randomly split
314 the dataset into a train/validation set and a test set (78 and 26 mice respectively, **Fig 3a**). In all
315 analyses, the “6N” strain (i.e., “mice with early-onset ARHL”) was treated as the positive class, as it
316 represents our primary outcome of interest in evaluating model performance.

317 Initially, we trained the models using ABRs for the full set of sound stimuli (click and pure tones
318 from 3 to 42 kHz) of 1-month-old mice from the primary cohort (“global” models). As input features,
319 we concatenated the ABR waveforms recorded at various stimulus intensities and frequencies,
320 forming a single univariate time series (**Fig 1b**). We first evaluated the model performances through
321 repeated *k*-fold cross-validation on the 78 mice within the training set (*k*=5 folds and 5 repeats, for a
322 total of 25 splits). Final scores were then calculated by averaging the results of individual splits (see
323 **Methods**). We found that all tested models showed strong overall performances across recall,
324 specificity, precision and receiver operating characteristic area under the curve (ROC AUC) metrics
325 (**Fig 3b-e, S1 Table**). When focusing on recall (i.e., the true positive rate or sensitivity, reflecting the
326 capability of the models to identify “pathological” cases), there was no significant differences among
327 most pairs of models, except between those with the highest score (SVM) and the two tree-based
328 models (**S2 Table**). Moreover, no significant differences were found in the specificity (i.e., true
329 negative rate) and precision scores of the six models ($P=0.0724$ and $P=0.0844$ respectively, Friedman
330 test, **S2 Table**). All models achieved relatively high average ROC AUC scores (**Fig 3e**),
331 demonstrating strong overall discrimination ability between the classes. These results suggest that all
332 models were generally effective in distinguishing between classes, with some models achieving
333 higher discrimination performance on average.

334 The models were then trained on the whole training set (78 mice) and evaluated on the test set (26
335 mice, the same train/test split was kept for all models). We found that the models based on the HIVE
336 COTE V2.0 and ROCKET classifiers showed the best performance and were able to correctly
337 determine the strain of all the 26 mice from their ABR waveforms (100% accuracy: **Fig 3f**). Tree-
338 based models (random forest and XGBoost) were slightly less accurate compared to the other four
339 models and misclassified 4 out of 26 test mice (~15%, **Fig 3f**). In comparison, manual blind dataset
340 labelling by three experimenters demonstrated varied accuracy, with each annotator mislabelling
341 between 5 and 8 mice out of the 26 in the test set, corresponding to an error rate between ~19% and
342 ~31%, with moderate agreement between annotators (Fleiss' Kappa: 0.58).

343 When evaluated on the entire dataset (104 mice), manual classification resulted in a recall of 66.7%
344 \pm 13.6%, specificity of 82.7% \pm 18.9%, and precision of 81.9% \pm 14.1% ($n = 3$ annotators). The lower
345 performance of manual classification can be attributed to its reliance on differences in high-frequency
346 ABRs between the two genotypes. However, at the early age considered, the substantial overlap
347 between the two genotypes reduces the reliability of human-extracted features for accurate genotype
348 differentiation. Notably, at 1 month of age, 25 out of 50 6N mice (50%) had thresholds at 42kHz that
349 were superimposed to Repaired mice (35 to 60 dB SPL, **Fig 2a**). Moreover, there was substantial
350 overlap in the distribution of wave amplitudes and latencies (see **Fig. 6** below). Therefore, human-
351 extracted features may not sufficiently capture the subtle differences between genotypes at this early
352 age since, unlike classification algorithms, manual classification is limited to lower-dimensional
353 representations of the data.

354 To gain an insight into the ML algorithm decision process, we determined the contribution of
355 individual features to the classification task by calculating the mean absolute Shapley values
356 (Lundberg & Lee, 2017). We selected the random forest and XGBoost classifier for this task, as
357 calculating Shapley values was computationally prohibitive for the other four models. We found that
358 features corresponding to higher frequency stimuli (36 and 42 kHz) were the most influential for the
359 classification task (**Fig 4, S3 Fig**). Additionally, features corresponding to click responses had
360 Shapley value elevated across different sound levels, suggesting that subtler differences between the
361 two genotypes may exist in the ABR waveforms associated with these stimuli (**Fig 4**).

362 Overall, these results highlight the potential of ML models to outperform human experts in
363 identifying differences in ABRs due to ARHL, offering an accurate tool for its early detection.

364

365 **3.3. Click ABRs are sufficient to predict the presence of the *Cdh23^{ah1}* allele**

366 To assess the robustness of the models under more challenging conditions, we aimed to restrict the
367 number of input features, simulating a scenario often encountered in clinical settings where higher-
368 frequency tone sensitivity (above 8 kHz, called the extended high frequencies, or EHF) is typically
369 not performed (Hunter *et al.* 2020). Specifically, we asked whether the click ABR alone, which does
370 not display any significant threshold shift until 6 months of age (**Fig 2d**), contained enough
371 information to differentiate between the two mouse strains at 1 month using ML. To test this, we re-
372 trained the six models described above using only the click responses from 1-month-old mice as input
373 features (“Click” models, **Fig 5a**). We found that all the models tested retained good performances
374 across several metrics (**Fig 5b-e, S3 Table, S4 Table**). No significant difference was found in the
375 recall of “Global” and “Click” models ($P=0.2150$, mixed-effects linear model), albeit the former were
376 associated with higher sensitivity and precision ($P < 0.0001$ for both metrics) and ROC-AUC score

377 (P=0.0330, mixed-effects linear models). When evaluated on the test set, models misclassified
378 between 2 to 5 out of 26 mice (error rate between ~8%, MLP and ~19%, random forest) from their
379 click ABR waveforms (**Fig 5f**). These results were consistent with those obtained for the “Global”
380 models that included all frequencies tested (**Fig 3**). In contrast, manual annotation from the three
381 experimenters was much less accurate when the information about high-frequency tones was
382 removed, with each annotator mislabelling between 10 and 12 mice out of the 26 from the test set
383 (error rates between ~38% and ~46%) with very poor agreement between annotators (Fleiss' Kappa:
384 -0.09).

385 As done previously for the “Global” models (**Fig 4**), we sought to interpret the “Click” models by
386 calculating the mean absolute Shapley values for the three algorithms for which the computation was
387 feasible on our hardware (random forest, SVM and XGBoost). We found that the mean absolute
388 Shapley values corresponding to wave 1 and wave 2 of the click ABR waveforms were consistently
389 elevated across the higher sound intensities (**Fig 5g, S3 Fig**), indicating the importance of these
390 features for the identification of early-onset ARHL mice. Wave 1, which reflects the activity of the
391 auditory afferent fibres, was consistently selected by the ANOVA F-test feature selection step across
392 most intensities. This suggests that, even in the absence of an auditory threshold difference,
393 significant variations in the output of the cochlea caused by the *Cdh23^{ahl}* allele may be present at an
394 early age.

395 Overall, these findings indicate that ABR wave 1 features could enable ML models to distinguish
396 between the two mouse strains at an early stage, before threshold differences emerge.

397

398 **3.4. Differences in ABR wave 1 in 1-month-old 6N and 6N-Repaired mice**

399 Next, we tested whether the importance of features in wave 1, which were used by some models to
400 identify 6N and 6N-Repaired mice, were underpinned by differences in the average wave 1 amplitude
401 and latency between the two strains. Using ABR click responses from 1-month-old mice, we found
402 that both wave 1 amplitude and latency differed significantly between the two mouse strains, despite
403 substantial overlap in their distributions (**Fig 6a,d,e**). The difference in average amplitude was
404 maximal at 95 dB SPL (1.6 μ V, 18.4%), while the difference in average latency was maximal at 70
405 dB SPL (56 μ s, 3.9%). Significant differences were also found in latency and amplitude in the
406 individual tone responses (e.g., 18 kHz: **Fig 6b,f,g** and 42 kHz **Fig 6c,h,i**, see also **S4 Fig**). These
407 results indicate that, in a mouse model of early onset progressive hearing loss, ABR waveforms may
408 undergo subtle changes well before a more obvious threshold shift appears. These changes can be
409 detected by ML models, potentially identifying early hallmarks of the dysfunction.

410 To further investigate the role of these features in classification, we compared the performance of a
411 support vector machine (SVM) when trained on only wave 1 latencies and amplitudes or auditory
412 thresholds, compared to when the full ABR waveform was used as input (**S5 Fig**). This analysis was
413 performed on both a “Global” dataset (including wave 1 parameters and thresholds for all click and
414 pure-tone stimuli) and a “Click” dataset (i.e., using only waveforms, thresholds, wave 1 amplitudes
415 and latencies from “Click” ABRs). We found that using the full ABR trace led to significantly higher
416 recall and ROC-AUC scores, reflecting both improved sensitivity in detecting mice carrying the
417 ARHL-linked allele and better discrimination between the two classes across different probability
418 thresholds. Moreover, the full-trace model consistently outperformed wave 1 and threshold-based
419 models on the test dataset, with all models trained using the same train/test split and cross-validation
420 folds (**S5 Fig**). Taken together, these findings indicate that allowing the model to autonomously
421 determine the most relevant features may offer advantages over hypothesis-driven feature selection,
422 leading to improved classification performance.

423

424 **3.5. ML models performances on heterogeneous datasets**

425 To test the ability of the models to generalise to a similar set of ABR data obtained from a different
426 experimental setting, we replicated the previous analysis incorporating into the training/testing data
427 an independently acquired ABR dataset (**replication cohort, Fig 7**). This dataset contained 188 click
428 ABRs of 1-month-old mice (85 6N and 103 6N-Repaired) and, differently from the primary cohort
429 dataset (**Figs 1-6**), was obtained from mice of both sexes.

430 We first tested whether the general approach described above was also applicable to the replication
431 cohort by retraining the ML models using either the primary or the replication cohort datasets, or the
432 two sets combined for training/validation (**Fig 7a-d, S5 Table**). As the replication cohort dataset
433 contained ABRs for stimuli up to 85 dB SPL, we retrained the models from **Fig 5** using only this
434 subset of sound intensities from the primary cohort dataset to allow for comparison. Across the main
435 four metrics (recall, sensitivity, precision, and ROC AUC), no statistically significant differences
436 were found between datasets (primary, replication, or combined cohorts, $P > 0.1$ for all comparisons,
437 mixed-effect linear model, **Fig 7a-d**). Moreover, feature importance analysis using Shapley values
438 highlighted similar features in the primary and replication cohorts (both alone and combined), roughly
439 corresponding to wave 1 and wave 2 (**Fig 7e**) as previously shown (**Fig 5g**). Overall, this result
440 highlights the generalisability of our ML-based approach in capturing relevant patterns in ABR data
441 for the identification of early hallmarks of *Cdh23^{ahl}*-related ARHL.

442 Next, we measured the performance of one of the models (ROCKET) in making predictions on a
443 dataset different from the one it was trained on (**Fig 7f**). We found that the accuracy greatly decreased

444 when making predictions on a different set. A model trained on the primary cohort dataset correctly
445 classified only 101 out of 188 mice (54%) from the replication cohort (recall: 67%, specificity: 43%),
446 while a model trained on 75% of the replication cohort dataset correctly predicted the genotype of 63
447 out of 104 mice (61%) from the primary cohort (recall: 94%, specificity: 30%, **Fig 7f**). In contrast,
448 accuracy on held out test data of either cohort remained higher when the model was trained on a
449 combined dataset (**Fig 7f**). These findings suggest that incorporating data from multiple sources
450 during training is essential to maintain a high prediction accuracy.

451

452 **3.6. ML based prediction of the progression of hearing function in mice**

453 We next examined whether ABR waveforms from young adult mice contained information to
454 predict the future progression of their hearing function. To do so, we used data from the 63 mice (45
455 6N and 18 6N-Repaired mice) from the primary cohort for which ABR measurements were collected
456 at 1, 3, 6, 9 and 12 months of age. ABRs of 6N-Repaired mice were included in the training set to
457 expose the models to data from “good hearing” mice, thus providing a wider range of targets. We
458 sought to train models to predict three parameters: the shift in average thresholds across all stimuli
459 (**Fig 8a**), and wave 1 amplitude (**Fig 8b**) and latency (**Fig 8c**) for click stimuli at three different sound
460 pressure levels (55, 75 and 95 dB SPL). These parameters showed significant age-dependent changes
461 in 6N mice (shift in average threshold: $P < 0.0001$, one-way ANOVA; wave 1 amplitude and latency:
462 $P < 0.0001$ for both, two-way ANOVA). For example, the average threshold increased by 38.0 ± 5.9
463 dB SPL from 1 to 12 months in 6N mice, compared to an increase of only 4.2 ± 4.3 dB SPL in 6N-
464 Repaired mice (**Fig 8a**). Moreover, substantial variability in the progression of these parameters was
465 observed across 6N mice (**Fig 8 a-c**, see also **S6 Fig**).

466 We trained regression models through supervised learning using click ABR waveforms from 1- and
467 3-month-old mice as input features, while the values of the parameters mentioned above at 6, 9 and
468 12 months were used as targets (**Fig 8d**). We randomly divided the dataset into a train/validation set
469 and a test set (47 and 16 mice respectively; the same training and test mice were kept across all
470 models). As above, we first evaluated model performances through repeated k -fold cross-validation
471 by averaging the mean absolute error (MAE) across splits ($k=5$ folds and 5 repeats, totalling 25 splits,
472 **Fig 8e-g**, see **Methods**). We found that the wave 1 latency model (**Fig 8g**) was the most accurate in
473 predicting the target values (the average MAE was between 3.9% and 7.9% of the mean wave 1
474 latency values at the corresponding age and sound levels, **Fig 8c**). When evaluated on the test set, the
475 models exhibited similar MAE to those calculated in cross-validation, with a coefficient of
476 determination (R^2) ranging from 0.45 to 0.69 (**Fig 8h-j**).

477

478 4. Discussion

479 The application of machine learning (ML) to large biomedical datasets is expected to drive profound
480 changes in clinical diagnosis, delivery of precision medicine and health monitoring (Goecks et al.,
481 2020). In this study, we applied ML to identify early signs of hearing loss and to predict its
482 progression in mice using ABR waveforms as input data. We tested six different ML algorithms on
483 the task of classifying which mice carried *Cdh23^{ahl}* (6N) compared to co-isogenic *Cdh23* repaired
484 mice (6N-Repaired), at a time when hearing thresholds are similar between the two mouse strains
485 (Johnson et al. 1997; Noben-Trauth et al. 2003; Mianné et al. 2016; Jeng et al. 2020a; 2020b). We
486 tested both widely used classifiers (e.g., random forest and SVM) and state-of-the-art algorithms
487 specialised for time series classification (HIVE-COTE V2.0 and ROCKET). The models we
488 implemented retained very good performances even when a restricted set of data (click ABRs) was
489 used instead of the full range of information (pure tone ABRs). This indicates that ML algorithms
490 were able to identify key features associated with hearing loss even in the absence of differences in
491 ABR thresholds, which can more easily be identified by trained experimentalists. None of the six
492 models tested demonstrated a clear advantage in the classification task. Time-series specific
493 classifiers like HIVE COTE V2.0 and ROCKET were the most consistent across different tasks, albeit
494 with longer computational analysis time compared to the other models. In contrast, tree-based models
495 (Random Forest and XGBoost) were less effective, despite offering easier interpretability of their
496 decision processes. Overall, even simpler models, like the one based on the SVM classifier,
497 performed reasonably well in the classification task, suggesting that more complex models might not
498 provide a substantial improvement over simpler ones. However, the consistency of results across
499 different models strengthens the validity of our approach and supports its robustness.

500 Interpretation of the models' decision process for the identification of mice carrying the *Cdh23^{ahl}*
501 allele revealed that the most important features were associated with wave 1, in agreement with the
502 mutation being located in a gene expressed in the cochlear hair cells. Indeed, wave 1 was consistently
503 selected by the ANOVA F-test feature selection step preceding the classifiers, and its amplitude and
504 latency were significantly different between the two genotypes already at one month of age. We
505 validated our ML approach by applying it to a second, more heterogeneous, ABR dataset acquired by
506 another laboratory. We also demonstrated that ML models are well suited to predict the future
507 trajectory of hearing capabilities in mice from early timepoints.

508 509 4.1. Early hallmarks of progressive hearing loss in mice with *Cdh23^{ahl}*

510 In mammals, acoustic information travelling within the cochlear partition is transduced into a
511 receptor potential in the sensory hair cells by the mechanical displacement of the stereociliary bundles

512 projecting from their apical surface (Fettiplace, 2017). Within each hair bundle, individual stereocilia
513 are interconnected by several extracellular linkages (Tilney *et al.* 1992; Goodyear *et al.* 2005). One
514 of these linkages, the tip link, is formed by cadherin 23 and protocadherin 15 (Siemens *et al.* 2004;
515 Ahmed *et al.* 2006; Kazmierczak *et al.* 2007). Tip links transmit force generated during sound-
516 induced displacement of the hair bundles to open mechano-electrical transducer (MET) channels
517 located at the tips of the shorter rows of stereocilia (Beurg *et al.* 2009). MET channel opening leads
518 to the depolarization of hair cells, and fusion of glutamate-filled synaptic vesicles at ribbon synapse
519 active zones, which allow high-rate synaptic transmission onto the auditory afferent fibres (Glowatzki
520 & Fuchs, 2002; Keen & Hudspeth, 2006; Goutman & Glowatzki, 2007).

521 The widely used C57BL/6 mice have a single-nucleotide polymorphism in exon7 of the gene
522 encoding cadherin 23, which affects splicing and leads to skipping of exon 7 (*Cdh23^{ahl}*, Johnson *et al.*
523 *et al.* 2017; Noben-Trauth *et al.* 2003). *Cdh23^{ahl}* has been shown to cause hearing loss starting in the
524 high-frequency cochlear region by about 3 months of age, and then progressing towards the low-
525 frequency, so that C57BL/6 mice become almost completely deaf by 12-18 months of age (Johnson
526 *et al.* 1997; Jeng *et al.* 2020a; Jeng *et al.* 2020b; Jeng *et al.* 2021; Kane *et al.* 2012; Peineau *et al.*
527 2021). Sensitive high frequency ABR thresholds are maintained into old age in co-isogenic repaired
528 C57BL/6 mice (6N-Repaired), in which the *Cdh23^{ahl}* allele was repaired with targeted CRISPR/Cas9
529 gene editing (Mianné *et al.* 2016). The progression of hearing loss of C57BL/6N mice used in this
530 study is consistent with the above previous investigations when using a similar ABR threshold
531 detection approach. However, the larger sample size of our dataset used to train ML algorithms
532 allowed us to identify small differences in the ABR waveform of 6N compared to 6N-Repaired mice
533 already at 1 month of age, which is a time when both mouse strains are considered to have normal
534 hearing (Fig 6). These differences included a significant decrease in the amplitude and increase in the
535 latency of wave 1, which is determined by the synchronous activity of auditory nerve fibres, and it is
536 generally interpreted as a measure of the neural output of the cochlea. Our results are therefore
537 consistent with the idea that a reduction in cochlear output is an early sign of the auditory decline
538 associated with ARHL, preceding threshold elevation (Sergeyenko *et al.* 2013). Moreover, this work
539 indicates that age-related changes to hair cell function in mice harbouring the *Cdh23^{ahl}* allele may
540 occur earlier than previously observed based on changes in synapse count and morphology (Jeng *et al.*
541 2020b; Peineau *et al.* 2021; Stamatakis *et al.* 2006). Further investigations will be required to
542 determine the physiological correlates of this reduction.

543
544
545

546 ***4.2. Wave 1 as an early predictor of progressive hearing loss due to cochlear impairment***

547 The close correspondence between wave 1 and synaptic survival (Kujawa & Liberman 2009)
548 indicates that non-invasive measures of cochlear neural responses, such as the one provided by ABR
549 tests, are a suitable method for early diagnosis of hearing loss even in the presence of normal
550 thresholds, as previously suggested (hidden hearing loss, Sergeyenko *et al.* 2013). Our work
551 demonstrates that ML algorithms can “learn” to single out early differences in ABR data without any
552 a-priori hypothesis of the features that are indicative of hearing loss. The number of features extracted
553 by computational algorithms far exceed those that can be identified by trained
554 experimenters/clinicians (e.g., thresholds, absolute amplitudes and latencies, wave 5/1 ratio and
555 interwave 1-5 latency, Verhulst *et al.* 2016). Moreover, our work demonstrates that feature selection
556 approaches that rely on predefined parameters, such as wave 1 amplitudes and latencies or auditory
557 thresholds, can be less sensitive in identifying early-onset ARHL cases. Conversely, a data-driven
558 approach, in which the model autonomously determines the most relevant features, improves
559 classification performance compared to hypothesis-driven feature selection. Our findings align with
560 prior work showing the potential of machine learning to enhance the objectivity and accuracy of ABR
561 waveform classification (McKearney & MacKinnon 2019). Moreover, ML applications are recently
562 making significant progress in hearing healthcare and research (Lesica *et al.* 2021; Chen *et al.* 2021;
563 Shew *et al.* 2019, Cha *et al.* 2019, Crowson *et al.* 2023, Chen *et al.* 2024).

564 One shortcoming of our approach is that our models demonstrated poor performance when tested on
565 mice from a different cohort, highlighting a key limitation in generalizability and indicating a
566 tendency to overfit to the specific domain used for training. The poor performances are likely linked
567 to variations in equipment and techniques used in ABR recordings, such as electrode placement,
568 mouse position in respect to the speakers, variations in preprocessing pipelines or subtle differences
569 in task execution by the experimenters, leading to a distribution shift between the two datasets.
570 However, we found that incorporating labelled data from the target cohort (Combined dataset) in the
571 training phase was sufficient to maintain classification performance. In the absence of labelled data
572 from the target cohort, transfer learning techniques (e.g. Azab *et al.* 2019; Azab *et al.* 2020, Giles *et*
573 *al.* 2022) could be used to merge knowledge from data collected in different settings to boost model
574 performance across different data sets. For example, domain adaptation techniques such as Domain-
575 Adversarial Neural Networks (DANN, Ganin *et al.* 2016) or CORAL (Sun *et al.* 2017), could help
576 mitigating this issue by aligning feature distributions across datasets, potentially improving model
577 robustness. However, these methods rely on neural networks and are therefore likely to require larger
578 datasets than those used in this study.

579 Furthermore, the models were trained exclusively on mice carrying the *Cdh23^{ahl}* allele and co-
580 isogenic controls, limiting their generalizability to other pathologies affecting ABR waveforms,
581 including other forms of progressive hearing loss. In order to improve the robustness and applicability
582 of the model, future work should focus on expanding the dataset to include a broader cohort of mice,
583 encompassing different ages and genetic backgrounds. Transfer learning techniques could also
584 facilitate the extension of models to the diagnosis of other hearing pathologies associated with
585 changes in ABR waveforms (e.g. [Schaette & McAlpine 2011](#)). A potential approach to further
586 improve model robustness could involve using individual ABR trials rather than averaged waveforms
587 to train/test models, as done in the present work. Resampling individual trials could also provide a
588 more in-depth assessment of the performance of the models under varying conditions.

589 We demonstrated that ML is able not only to identify early signs of hearing loss due to the *Cdh23^{ahl}*
590 allele, but also forecasting the future progression of hearing loss in mice. Interestingly, forecasting
591 the progression of hearing loss was recently applied to a longitudinal study with patients affected with
592 *GJB2*-related sensorineural hearing loss ([Chen et al. 2024](#)). However, translating this approach to
593 humans with unknown mutations linked to hearing loss would face numerous challenges. Variability
594 in ABR waveforms is notably higher in human data than in controlled animal models. Moreover,
595 amplitudes and latencies can vary significantly, both within a single clinic and even more so across
596 different clinics. Electrode montage around the patient head and head size are also known factors
597 influencing ABR measurement in humans ([King & Sininger 1992](#); [Mitchell et al. 1989](#)). These
598 challenges could be overcome by developing a ABR testing pipeline that will allow the acquisition
599 of high-quality, standardised, well curated ABR datasets. Moreover, wave 1 is usually small and more
600 difficult to identify in humans than in mice ([Bramhall 2021](#)). Therefore, other non-invasive
601 measurements of auditory nerve activity, such as the auditory nerve compound action potential
602 recorded with an extra-tympanic electrode ([Eggermont 2017](#)) could be used to develop ML-based
603 diagnostic tools. Finally, multimodal machine learning models that integrate both structured (e.g.,
604 age, gender ([Jerger & Johnson, 1988](#))) and unstructured data (e.g., clinical notes) may be required to
605 achieve reliable predictions.

606 An ML-based approach could also be applied to the identification of new genes involved in
607 progressive auditory dysfunction in large-scale screening studies, which often rely on thresholds as
608 the primary metric ([Bowl et al. 2017](#); [Ingham et al 2019](#)). These, however, are relatively insensitive
609 to primary neuronal degeneration without hair cell loss ([Kujawa & Liberman 2009](#)). By contrast, an
610 ML-based approach could provide a more sensitive tool that does not depend on human-labelled
611 parameters indicative of hearing loss, potentially enabling the discovery of new genes implicated in
612 ARHL.

613 Overall, our findings demonstrate the potential of machine learning applied to ABR data for early
614 detection of hearing loss, providing a framework for developing more sensitive, comprehensive
615 diagnostic tools. While our study focused on a controlled mouse dataset, future work will be necessary
616 to assess the applicability of this approach to human data, where ABR variability across clinics and
617 individuals presents additional challenges.

618 **Declaration of competing interests:** The Authors declare no conflict of interest.

619

620 **CRediT authorship contribution statement**

621 **Federico Ceriani:** conceptualisation, data curation, formal analysis, funding acquisition,
622 methodology, project administration, software, validation, writing – original draft, writing – review
623 & editing. **Joshua Giles:** methodology, funding acquisition, writing – review & editing. **Neil J**
624 **Ingham:** investigation, writing – review & editing. **Jing-Yi Jeng:** funding acquisition, investigation,
625 writing – review & editing. **Morag A Lewis:** investigation, writing – review & editing. **Karen P**
626 **Steel:** funding acquisition, investigation, writing – review & editing. **Mahnaz Arvaneh:**
627 methodology, funding acquisition, supervision, writing – review & editing. **Walter Marcotti:** data
628 curation, funding acquisition, investigation, resources, supervision, writing – original draft, writing –
629 review & editing.

630

631 **Acknowledgements**

632 This work was supported by an Innovation Seed Fund from the RNID (F115) to FC and JJ and by
633 a grant from the Sheffield Neuroscience Institute to FC, WM, MA and JG.

634 FC was supported by a BBSRC grant (BB/V006681/1) to WM and FC.

635 JJ was supported by the RNID-Dunhill Medical Trust Fellowship (PA28).

636 KPS, MAL and NJI were supported by Wellcome (098051; 100669; 089622; 221769/Z/20/Z).

637 The authors thank Steven M. Barnes and Matthew A Loczki at the University of Sheffield for their
638 assistance with the mouse husbandry.

639

640 **Data availability**

641 The data that support the findings is available upon request.

642 The code is available on GitHub at <https://github.com/fedeceri85/abr-ml-analysis-paper>.

643 **References**

- 644 Azab AM, Mihaylova L, Ahmadi H, Arvaneh M. Robust common spatial patterns estimation using
645 dynamic time warping to improve BCI systems. In: ICASSP 2019–2019 IEEE International
646 Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE; 2019. p. 3897–901.
- 647 Azab AM, Ahmadi H, Mihaylova L, Arvaneh M. Dynamic time warping-based transfer learning
648 for improving common spatial patterns in brain–computer interface. *J Neural Eng.* 2020;17:016061.
- 649 Ahmed ZM, Goodyear R, Riazuddin S, Lagziel A, Legan PK, Behra M, Burgess SM, Lilley KS,
650 Wilcox ER, Riazuddin S, Griffith AJ. The tip-link antigen, a protein associated with the transduction
651 complex of sensory hair cells, is protocadherin-15. *J Neurosci.* 2006;26:7022–34.
- 652 Banerjee J, Taroni JN, Allaway RJ, Prasad DV, Guinney J, Greene C. Machine learning in rare
653 disease. *Nat Methods.* 2023;20:803–14.
- 654 Beurg M, Fettiplace R, Nam JH, Ricci AJ. Localization of inner hair cell mechanotransducer
655 channels using high-speed calcium imaging. *Nat Neurosci.* 2009;12:553–8.
- 656 Bramhall NF. Use of the auditory brainstem response for assessment of cochlear synaptopathy in
657 humans. *J Acoust Soc Am.* 2021;150:4440–51.
- 658 Bowl MR, Simon MM, Ingham NJ, Greenaway S, Santos L, Cater H, et al. A large scale hearing
659 loss screen reveals an extensive unexplored genetic landscape for auditory dysfunction. *Nat Commun.*
660 2017;8:886.
- 661 Cha D, Pae C, Seong SB, Choi JY, Park HJ. Automated diagnosis of ear disease using ensemble
662 deep learning with a big otoendoscopy image database. *EBioMedicine.* 2019;45:606–14.
- 663 Chen C, Zhan L, Pan X, Wang Z, Guo X, Qin H, Xiong F, Shi W, Shi M, Ji F, Wang Q, Yu N, Xiao
664 R. Automatic recognition of auditory brainstem response characteristic waveform based on
665 bidirectional long short-term memory. *Front Med.* 2021;7:613708.
- 666 Chen PY, Yang TW, Tseng YS, Tsai CY, Yeh CS, Lee YH, Lin PH, Lin TC, Wu YJ, Yang TH,
667 Chiang YT, Hsu JS, Hsu CJ, Chen PL, Chou CF, Wu CC. Machine learning-based longitudinal
668 prediction for GJB2-related sensorineural hearing loss. *Comput Biol Med.* 2024;176:108597.
- 669 Chen T, Guestrin C. Xgboost: A scalable tree boosting system. In: Proceedings of the 22nd ACM
670 SIGKDD International Conference on Knowledge Discovery and Data Mining. 2016. p. 785–94.
- 671 Crowson MG, Bates DW, Suresh K, Cohen MS, Hartnick CJ. “Human vs Machine” Validation of
672 a Deep Learning Algorithm for Pediatric Middle Ear Infection Diagnosis. *Otolaryngol Head Neck*
673 *Surg.* 2023;169:41–6.
- 674 Dempster A, Petitjean F, Webb GI. ROCKET: exceptionally fast and accurate time series
675 classification using random convolutional kernels. *Data Min Knowl Discov.* 2020;34:1454–95.
- 676 Edla DR, Mangalorekar K, Dhavalikar G, Dodia S. Classification of EEG data for human mental
677 state analysis using Random Forest Classifier. *Procedia Comput Sci.* 2018;132:1523–32.
- 678 Eggermont JJ. Ups and downs in 75 years of electrocochleography. *Front Syst Neurosci.* 2017;11:2.
- 679 Fettiplace R. Hair Cell Transduction, Tuning, and Synaptic Transmission in the Mammalian
680 Cochlea. *Comprehensive Physiology.* 2017;7:1197-227.
- 681 Gates GA, Mills JH. Presbycusis. *Lancet.* 2005;366:1111–20.
- 682 Ganin Y, Ustinova E, Ajakan H, Germain P, Larochelle H, Laviolette F, March M, Lempitsky V.
683 Domain-adversarial training of neural networks. *Journal of machine learning research.* 2016;17:1-35.
- 684 Giles J, Ang KK, Phua KS, Arvaneh M. A transfer learning algorithm to reduce brain-computer
685 interface calibration time for long-term users. *Front Neuroergonomics.* 2022;3:837307.

686 Glowatzki E, Fuchs PA. Transmitter release at the hair cell ribbon synapse. *Nat Neurosci.*
687 2002;5:147–54.

688 Goecks J, Jalili V, Heiser LM, Gray JW. How machine learning will transform biomedicine. *Cell.*
689 2020;181:92–101.

690 Goodyear RJ, Marcotti W, Kros CJ, Richardson GP. Development and properties of stereociliary
691 link types in hair cells of the mouse cochlea. *J Comp Neurol.* 2005;485:75–85.

692 Goutman JD, Glowatzki E. Time course and calcium dependence of transmitter release at a single
693 ribbon synapse. *Proc Natl Acad Sci U S A.* 2007;104:16341–6.

694 Hunter LL, Monson BB, Moore DR, Dhar S, Wright BA, Munro KJ, Zadeh LM, Blankenship CM,
695 Stiepan SM, Siegel JH. Extended high frequency hearing and speech perception implications in adults
696 and children. *Hear Res.* 2020;397:107922.

697 Ingham NJ, Pearson SA, Vancollie VE, Rook V, Lewis MA, Chen J, Buniello A, Martelletti E,
698 Preite L, Lam CC, Weiss FD, Powis Z, Suwannarat P, Lelliott CJ, Dawson SJ, White JK, Steel KP.
699 Mouse screen reveals multiple new genes underlying mouse and human hearing loss. *PLoS Biol.*
700 2019;17:e3000194.

701 Ingham NJ, Pearson S, Steel KP. Using the auditory brainstem response (ABR) to determine
702 sensitivity of hearing in mutant mice. *Curr Protoc Mouse Biol.* 2011;1:279–87.

703 Jeng JY, Johnson SL, Carlton AJ, De Tomasi L, Goodyear RJ, De Faveri F, Furness DN, Wells S,
704 Brown SD, Holley MC, Richardson GP, Mustapha M, Bowl MR, Marcotti W. Age-related changes
705 in the biophysical and morphological characteristics of mouse cochlear outer hair cells. *J Physiol.*
706 2020a;598:3891–910.

707 Jeng JY, Ceriani F, Olt J, Brown SD, Holley MC, Bowl MR, Johnson SL, Marcotti W.
708 Pathophysiological changes in inner hair cell ribbon synapses in the ageing mammalian cochlea. *J*
709 *Physiol.* 2020b;598:4339–55.

710 Jeng JY, Carlton AJ, Johnson SL, Brown SD, Holley MC, Bowl MR, Marcotti W. Biophysical and
711 morphological changes in inner hair cells and their efferent innervation in the ageing mouse cochlea.
712 *J Physiol.* 2021;599:269–87.

713 Jerger J, Johnson K. Interactions of age, gender, and sensorineural hearing loss on ABR latency.
714 *Ear Hear.* 1988;9:168–76.

715 Johnson KR, Erway LC, Cook SA, Willott JF, Zheng QY. A major gene affecting age-related
716 hearing loss in C57BL/6J mice. *Hear Res.* 1997;114:83–92.

717 Kane KL, Longo-Guess CM, Gagnon LH, Ding D, Salvi RJ, Johnson KR. Genetic background
718 effects on age-related hearing loss associated with Cdh23 variants in mice. *Hear Res.* 2012;283:80–
719 8.

720 Kazmierczak P, Sakaguchi H, Tokita J, Wilson-Kubalek EM, Milligan RA, Müller U, Kachar B.
721 Cadherin 23 and protocadherin 15 interact to form tip-link filaments in sensory hair cells. *Nature.*
722 2007;449:87–91.

723 Keen EC, Hudspeth A. Transfer characteristics of the hair cell’s afferent synapse. *Proc Natl Acad*
724 *Sci U S A.* 2006;103:5537–42.

725 King AJ, Sininger YS. Electrode configuration for auditory brainstem response audiometry. *Am J*
726 *Audiol.* 1992;1:63–7.

727 Kujawa SG, Liberman MC. Adding insult to injury: cochlear nerve degeneration after “temporary”
728 noise-induced hearing loss. *J Neurosci.* 2009;29:14077–85.

729 Livingston G, Huntley J, Liu KY, Costafreda SG, Selbæk G, Alladi S, et al. Dementia prevention,
730 intervention, and care: 2024 report of the Lancet standing Commission. *Lancet*. 2024;404:572–628.

731 Löning M, Bagnall A, Ganesh S, Kazakov V, Lines J, Király FJ. sktime: A unified interface for
732 machine learning with time series. *arXiv preprint arXiv:1909.07872*. 2019.

733 Lv J, Wang H, Cheng X, Chen Y, Wang D, Zhang L, et al. AAV1-hOTOF gene therapy for
734 autosomal recessive deafness 9: a single-arm trial. *Lancet*. 2024; S0140–6736(23)02874–X.

735 Lesica NA, Mehta N, Manjaly JG, Deng L, Wilson BS, Zeng FG. Harnessing the power of artificial
736 intelligence to transform hearing healthcare and research. *Nat Mach Intell*. 2021;3:840–9.

737 Lundberg SM, Lee SI. A unified approach to interpreting model predictions. *Adv Neural Inf Process*
738 *Syst*. 2017;30.

739 McKearney RM, MacKinnon RC. Objective auditory brainstem response classification using
740 machine learning. *International journal of audiology*. 2019;58:224-30.

741 Mianné J, Chessum L, Kumar S, Aguilar C, Codner G, Hutchison M, et al. Correction of the
742 auditory phenotype in C57BL/6N mice via CRISPR/Cas9-mediated homology directed repair.
743 *Genome Med*. 2016;8:16.

744 Middlehurst M, Large J, Flynn M, Lines J, Bostrom A, Bagnall A. HIVE-COTE 2.0: a new meta
745 ensemble for time series classification. *Mach Learn*. 2021;110:3211–43.

746 Mitchell C, Phillips DS, Trune DR. Variables affecting the auditory brainstem response: audiogram,
747 age, gender and head size. *Hear Res*. 1989;40:75–85.

748 Noben-Trauth K, Zheng QY, Johnson KR. Association of cadherin 23 with polygenic inheritance
749 and genetic modification of sensorineural hearing loss. *Nat Genet*. 2003;35:21–3.

750 Nolan LS. Age-related hearing loss: Why we need to think about sex as a biological variable. *J*
751 *Neurosci Res*. 2020;98:1705–20.

752 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine
753 learning in Python. *J Mach Learn Res*. 2011;12:2825–30.

754 Peineau T, Belleudy S, Pietropaolo S, Bouleau Y, Dulon D. Synaptic release potentiation at aging
755 auditory ribbon synapses. *Front Aging Neurosci*. 2021;13:756449.

756 Schaeffe R, McAlpine D. Tinnitus with a normal audiogram: physiological evidence for hidden
757 hearing loss and computational model. *J Neurosci*. 2011;31:13452–7.

758 Schilder AGM, Wolpert S, Saeed S, Middelink LM, Edge ASB, Blackshaw H, et al. A phase I/IIa
759 safety and efficacy trial of intratympanic gamma-secretase inhibitor as a regenerative drug treatment
760 for sensorineural hearing loss. *Nat Commun*. 2024;15:1896.

761 Sergeyenko Y, Lall K, Liberman MC, Kujawa SG. Age-related cochlear synaptopathy: an early-
762 onset contributor to auditory functional decline. *J Neurosci*. 2013;33:13686–94.

763 Shew M, New J, Wichova H, Koestler DC, Staecker H. Using machine learning to predict
764 sensorineural hearing loss based on perilymph micro RNA expression profile. *Sci Rep*. 2019;9:3393.

765 Sidey-Gibbons JA, Sidey-Gibbons CJ. Machine learning in medicine: a practical introduction.
766 *BMC Med Res Methodol*. 2019;19:1–18.

767 Siemens J, Lillo C, Dumont RA, Reynolds A, Williams DS, Gillespie PG, et al. Cadherin 23 is a
768 component of the tip link in hair-cell stereocilia. *Nature*. 2004;428:950–5.

769 Stamatakis S, Francis HW, Lehar M, May BJ. Synaptic alterations at inner hair cells precede spiral
770 ganglion cell loss in aging C57BL/6J mice. *Hear Res*. 2006;221:104–18.

771 Sun B, Feng J, Saenko K. Correlation alignment for unsupervised domain adaptation. Domain
772 adaptation in computer vision applications. arXiv:1612.01939.

773 Tilney LG, Tilney MS, DeRosier DJ. Actin filaments, stereocilia, and hair cells: how cells count
774 and measure. *Annu Rev Cell Biol.* 1992;8:257–74.

775 Verhulst S, Jagadeesh A, Mauermann M, Ernst F. Individual differences in auditory brainstem
776 response wave characteristics: relations to different aspects of peripheral hearing loss. *Trends Hear.*
777 2016;20:2331216516672186.

778 Wang J, Puel JL. Presbycusis: an update on cochlear mechanisms and therapies. *J Clin Med.* 2020;
779 9:218.

780 Wasmann JW, Pragt L, Eikelboom R, Swanepoel DW. Digital approaches to automated and
781 machine learning assessments of hearing: scoping review. *Journal of Medical Internet Research.*
782 2022;24:e32581.

783

784 **Figure legends**

785

786 **Fig 1. Schematic representation of the work.**

787 (a) Auditory brainstem responses (ABRs) were recorded from a cohort of 1-month-old 6N (50 mice)
788 and 6N-Repaired mice (54 mice). Some mice were also tested at 3, 6, 9 and 12 months of age.
789 Anaesthetised animals were presented with auditory stimuli comprising clicks and pure tones of
790 frequencies ranging from 3 to 42 kHz and stimulus intensities ranging from 15 dB SPL to 95 dB SPL
791 in 5 dB SPL increments. Bottom panel shows an example of an ABR recording from a 6N mouse.
792 Each of the traces in the matrix lasts 10 ms and represents a response to one combination of stimulus
793 intensity and click/tone frequency. (b) Responses to individual intensities/frequencies combinations
794 were concatenated in a single trace and used as input features for the machine learning (ML) models.
795

796 **Fig 2. Age-dependent change in ABR thresholds in 6N and 6N-Repaired mice.**

797 (a-e) ABR thresholds for click and pure tone stimuli at different ages. Solid points represent median
798 \pm median absolute deviation, while individual lines represent audiograms of individual mice. The
799 number of mice for each genotype is indicated in parenthesis. The same mice cohort was repeatedly
800 tested at different ages. Note that not all mice were investigated at all age timepoints. Significant
801 differences between the two genotypes are indicated next to the data points (*: $P < 0.0001$, pairwise
802 Wilcoxon rank-sum test, ART ANOVA). At 1 month of age, auditory thresholds were already
803 significantly different between the two strains ($P < 0.0001$, ART ANOVA). (f) Rasterplots showing
804 auditory thresholds as a function of age for four different stimuli (click, 18, 30, 42 kHz) in the 63
805 mice (45 6N mice and 18 6N-Repaired mice) that were evaluated at all ages tested. Note the
806 progressive increase in auditory thresholds of 6N mice (red band) compared to 6N-Repaired mice
807 (blue band).

808

809 **Fig 3. ML models can accurately predict the presence of the *Cdh23^{ahl}* allele from ABR**
810 **waveforms early on.**

811 (a) Data flow of the ML training/ testing process. The ABR dataset from 1-month old mice in the
812 primary cohort was randomly split into a train/validation set (75% of mice) and a test (hold-out) set
813 (25% of mice). The same train/test split was consistently used across all the models. The models
814 consisted of an ANOVA feature selection step, where the 10% top scoring features (timepoints) were
815 selected (see **Methods**), followed by one of six classifiers. Models were initially evaluated using
816 repeated stratified k -fold cross-validation ($k = 5$ splits and 5 repeats). The 25 scores produced by this
817 step were averaged to provide a measure of the overall performances of the models. Finally, the

818 models were trained on the entire training/validation set from the primary cohort and final scores
819 were obtained by testing predictions on the held-out test set. **(b-e)** Average metrics of the six models
820 trained on the whole ABR (click and 8 tones) as estimated in the 5×5 cross-validation step. Solid dots
821 represent the mean ± SD. Smaller dots indicate the scores from individual folds in the cross-validation
822 step (25 scores per model, see also **S1 Table** and **S2 Table** for statistical comparisons). **(f)** Confusion
823 matrices highlighting the performances of the models trained on the whole ABR (click and 8 tones)
824 on the final test set. HC: HiveCoteV2.0; MLP: multilayer perceptron. RF: random forest; Rckt:
825 ROCKET; SVM: support vector classifier; XGB: XGBoost. Rep: 6N-Repaired. NPV: negative
826 predictive value.

827

828 **Fig 4. Most important features for *Cdh23^{ahl}* prediction highlighted by the models.**

829 Matrix of average ABR waveforms (104 mice in the primary cohort at 1-month-old from both the 6N
830 and 6N-Repaired strains). Note that the y-axis scale is adjusted independently to the minimum and
831 maximum for each trace. The part of the traces in blue indicates parts of the ABR selected by the
832 ANOVA F-test feature selection step preceding the classifiers. Each trace is superimposed to a colour-
833 coded raster plot representing the normalised mean absolute Shapley values. Higher Shapley values
834 indicate the most influential features for model prediction. The displayed raster plots were calculated
835 as averages of the normalised Shapley values for the random forest (RF) and XGBoost (XGB) models
836 (see **S3 Fig** for the Shapley values of the two individual models). Shapley values were calculated on
837 the test set, using the training set as background distribution. This analysis indicated that responses
838 to higher frequencies (36 and 42 kHz) at high sound levels are the most important features for these
839 two models, followed by features associated to click stimuli.

840

841 **Fig 5. Click responses alone are sufficient for predicting the presence of *Cdh23^{ahl}* from ABRs**
842 **(a)** Input features for “Click” models. The greyed-out trace (i.e., tone ABRs) indicates features not
843 used for training/testing the models (compared to “Global” models, **Fig 1b**, **Fig 3a**). **(b-e)** Average
844 performances of the six models trained on the click ABR alone as estimated in the cross-validation
845 step. Solid dots represent the mean ± SD. Smaller dots indicate the score of individual folds in the
846 cross-validation step (25 scores per model, see also **S3 Table** and **S4 Table** for statistical
847 comparisons). **(f)**, Confusion matrices highlighting the performances of the models trained on the
848 click ABR on the test set. HC: HiveCoteV2.0; MLP: multilayer perceptron. RF: random forest; Rckt:
849 ROCKET; SVM: support vector classifier; XGB: XGBoost. Rep: 6N-Repaired. NPV: negative
850 predictive value. **(g)** Average click ABR waveform (104 mice from 1-month-old of both 6N and 6N-
851 Repaired strains) superimposed to a colour-coded raster plot representing the normalised mean

852 absolute Shapley values. Shapley values from three different models (random forest, SVM and
853 XGBoost) were normalised and averaged (see **S3 Fig** for the Shapley values of the individual
854 models). The part of the traces in blue indicates parts of the ABR selected by the ANOVA F-test
855 feature selection step preceding the classifiers. This analysis highlighted parts of wave 1 and wave 2
856 at sound intensities above 50 dB SPL as the features with the highest importance for model
857 predictions.

858

859 **Fig 6. Differences in ABR Wave 1 in 1-month-old 6N and 6N-Repaired mice.**

860 (a-c) Comparisons of average ABR waveforms from 50 6N and 54 6N-Repaired mice for click stimuli
861 (a), 18 kHz (b) and 42 kHz (c) tones at 95 dB SPL. The traces on the right provide a magnified view
862 of the dashed in the traces on the left, highlighting the subtle differences in wave 1 between the
863 average waveforms of the two genotypes (arrows in panel a). (d-i), Average wave 1 amplitude (d,f,h)
864 and latency (e,g,i) as a function of sound level for 6N and 6N-Repaired mice for click, 18 kHz and
865 42 kHz sound stimuli. Significant differences between the two mouse strains: $P < 0.0001$ (for the
866 three stimuli for both amplitude and latency, two-way ANOVA, panels d-i). Solid lines represent the
867 mean \pm SD, while lighter traces show individual mice.

868

869 **Fig 7. ML model performances decrease when predicting an external dataset.**

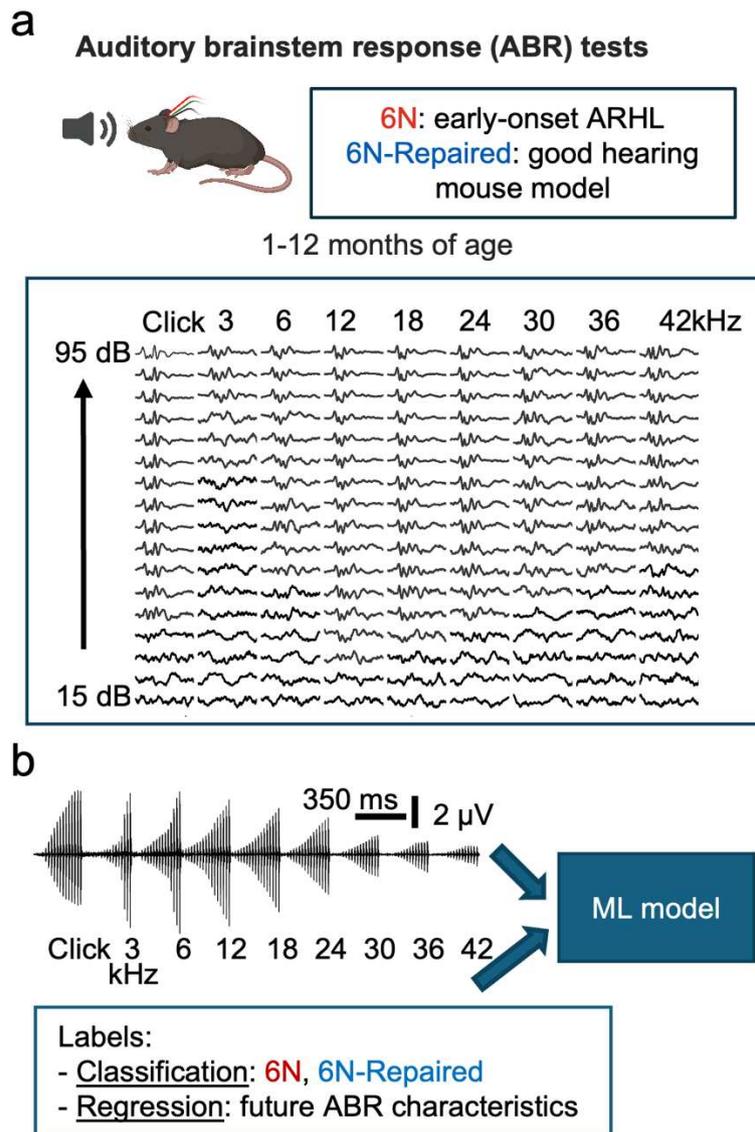
870 (a-d) Average metrics of the six models trained on click ABR as estimated in the cross-validation
871 step. Models were trained and validated on each of two independently acquired datasets (primary and
872 replication cohort) or on a combined dataset from one-month-old mice. To align with the acquisition
873 protocols of the two datasets, only sound intensities from 15 dB SPL to 85 dB SPL were used. Solid
874 dots represent the mean \pm SD while the smaller dots indicate the scores from individual folds in the
875 cross-validation step (25 scores per model per dataset, see also **S5 Table**). (e) Average click ABR
876 waveforms from three datasets (primary cohort, replication cohort, combined) superimposed to a
877 colour-coded raster plot representing the normalised mean absolute Shapley values. The values of
878 three different models (random forest, SVM and XGBoost) were normalised and averaged. Models
879 were trained and tested on either the primary cohort dataset (left), replication cohort dataset (center)
880 or both datasets combined (right). This analysis consistently identified parts of wave 1/ wave 2 as the
881 features with the highest importance for model predictions across datasets. Primary cohort: 50 6N
882 and 54 6N-Repaired (Rep.) mice; replication cohort: 85 6N and 103 6N-Repaired mice. (f) Array of
883 confusion matrices displaying performances of ROCKET models with different combination of the
884 three datasets used for training and/or testing. For tasks which involved evaluating the models on data
885 from different laboratories (i.e. training on replication cohort/testing on primary cohort or vice versa),
886 the whole primary and replication cohort datasets were used either for training or testing. Note that

887 results for training/testing on primary cohort data (top left confusion matrix) are slightly different
888 from the same model in **Fig 5f** due to the difference in input features between the two (15-95 dB SPL
889 click ABRs vs 15-85 dB SPL click ABRs). HC: HiveCoteV2.0; MLP: multilayer perceptron. RF:
890 random forest; Rckt: ROCKET; SVM: support vector classifier; XGB: XGBoost. Rep: 6N-Repaired.
891 NPV: negative predictive value.

892

893 **Fig 8. ML approach to predict the progression of hearing function in mice.**

894 (a-c) Change in ABR properties over time for 6N (top panels) and 6N-Repaired (bottom panels)
895 mice (45 6N and 18 6N-Repaired mice). Panel (a) displays the change in average threshold, calculated
896 as the difference between the average thresholds for click and eight tones relative to the value at 1
897 month of age. Wave 1 amplitude and latency for click stimuli at three different sound levels (55 dB,
898 blue; 75 dB, orange and 95 dB, green) are shown in panels (b) and (c), respectively. (d) Scheme of
899 the regression model. The models included an ANOVA feature selection step, selecting 10% of
900 timepoints as features (see **Methods**), followed by a random forest regression model. The input of
901 the models consisted of concatenated click ABRs at 1 and 3 months of age, while the targets were the
902 values of the three ABR parameters described in panels (a-c) at 6, 9 and 12 months of age.
903 Training/cross validation set: 47 mice; test (hold-out) set: 16 mice. The same
904 training/validation/testing split was kept for all the models. (e-g) Mean absolute error (MAE) of the
905 regression models for changes in average thresholds (e), wave 1 amplitude (f) and latency (g) of click
906 stimuli, as estimated in the cross-validation step. The target age is indicated on the x axis. Filled dots
907 represent the mean \pm SD. Smaller dots show the score of individual folds in the cross-validation step
908 (5 folds, 5 repeats, 25 scores per age and sound level). (h-j) scatter plots displaying predicted versus
909 real values for the change in average threshold (h) and wave 1 amplitudes (i) and latencies (j) at the
910 indicated sound levels for the test set (16 mice). Each triplet of connected symbols represents an
911 individual mouse in the test set. The coefficient of determination (R^2) and the mean absolute error
912 (MAE) of each model are indicated at the top. The dashed grey lines represent ideal predictions.



914
915
916

Figure 1.

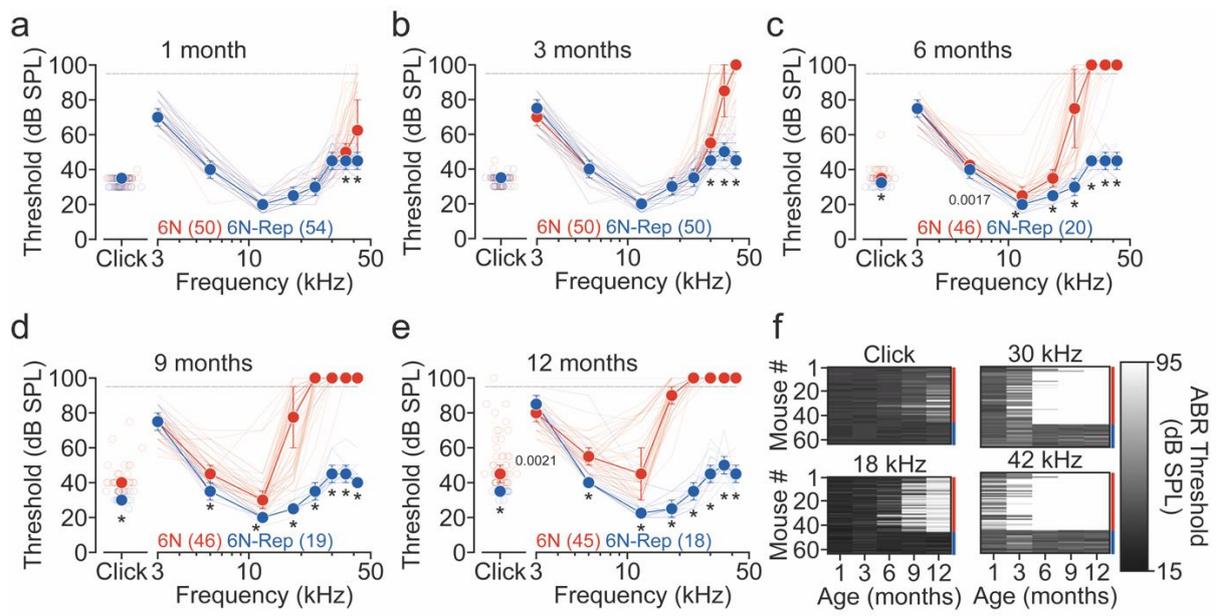


Figure 2.

917

918

919

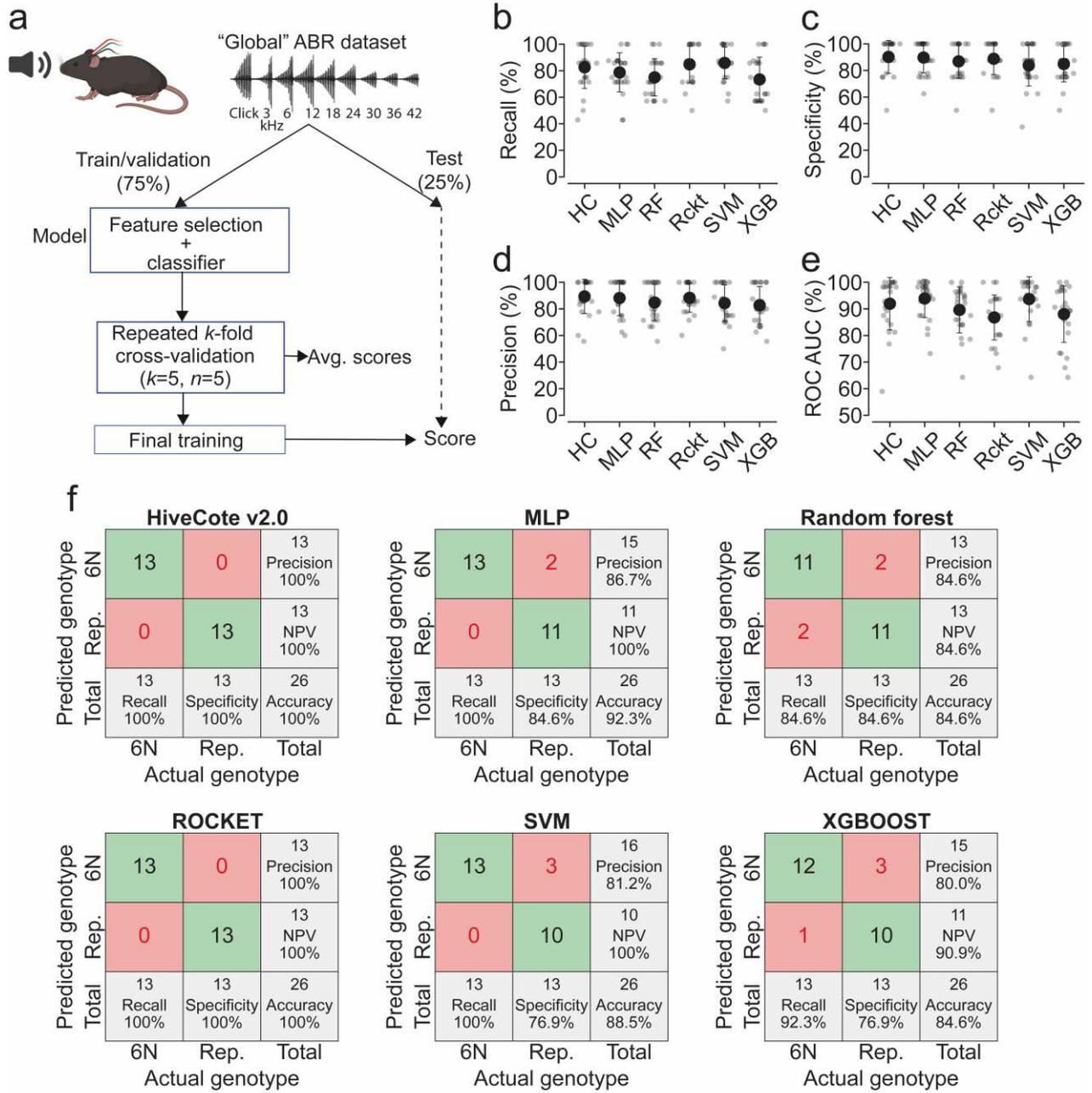
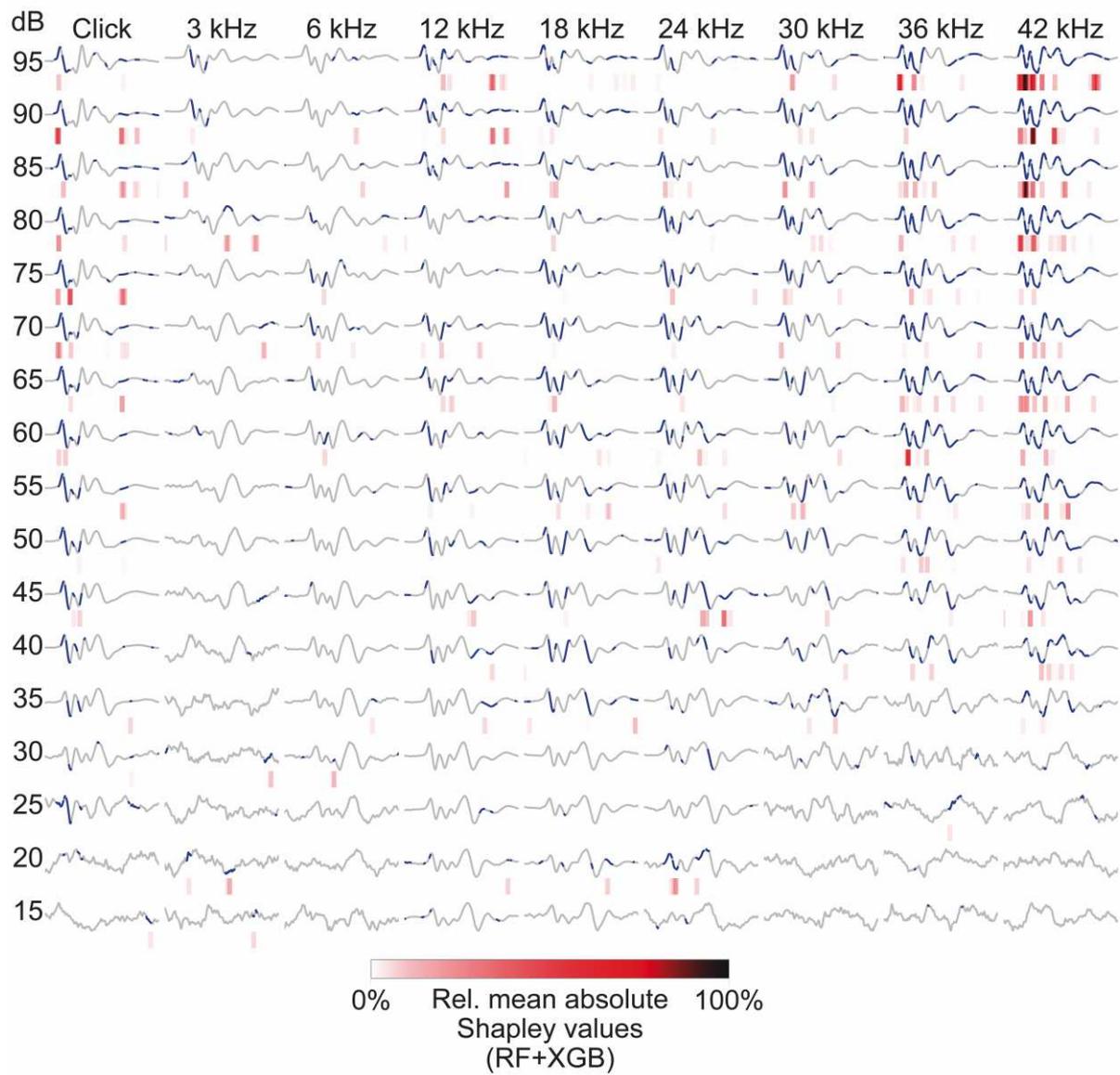
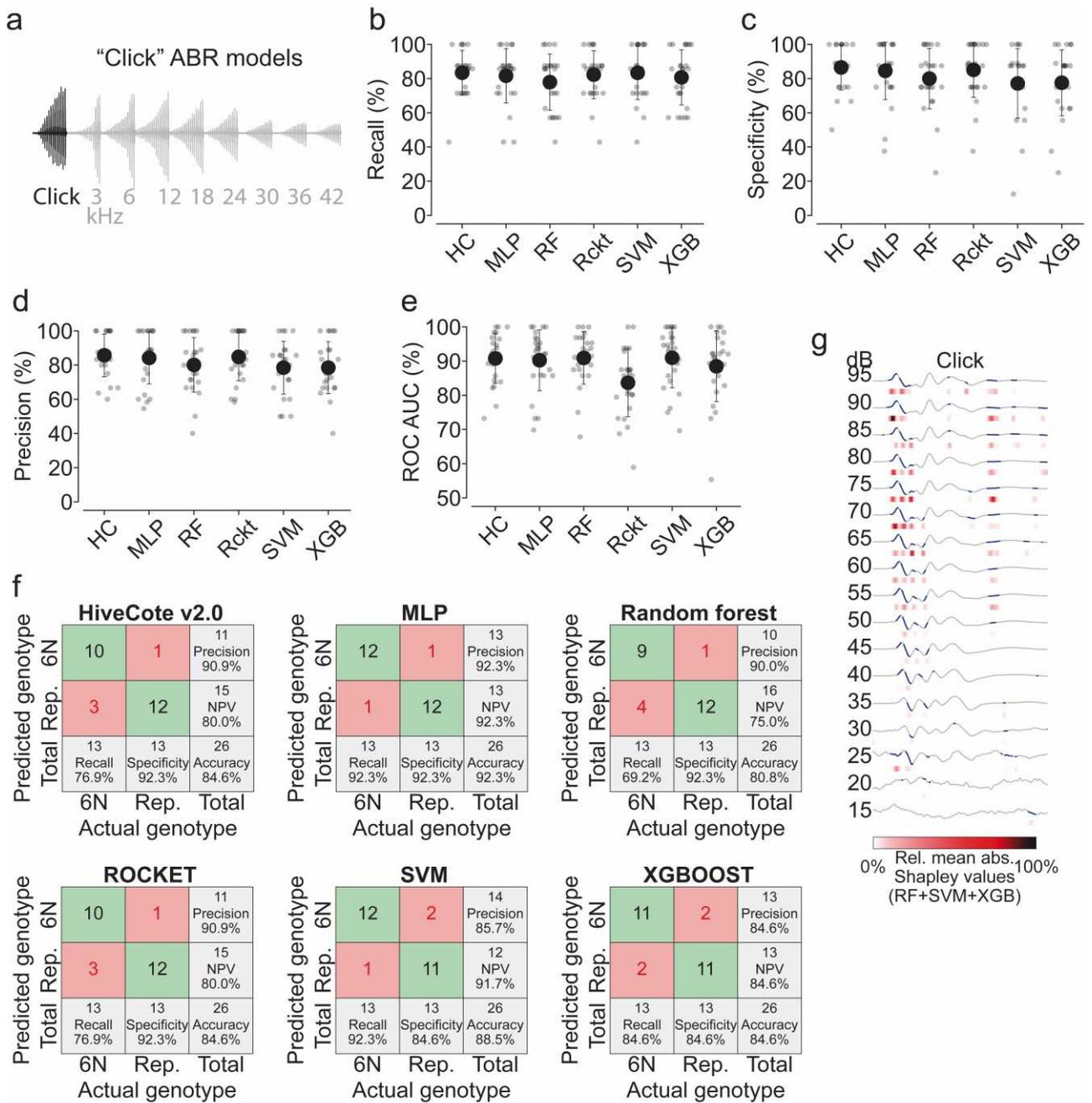


Figure 3.



923
924
925

Figure 4.



926

927

Figure 5.

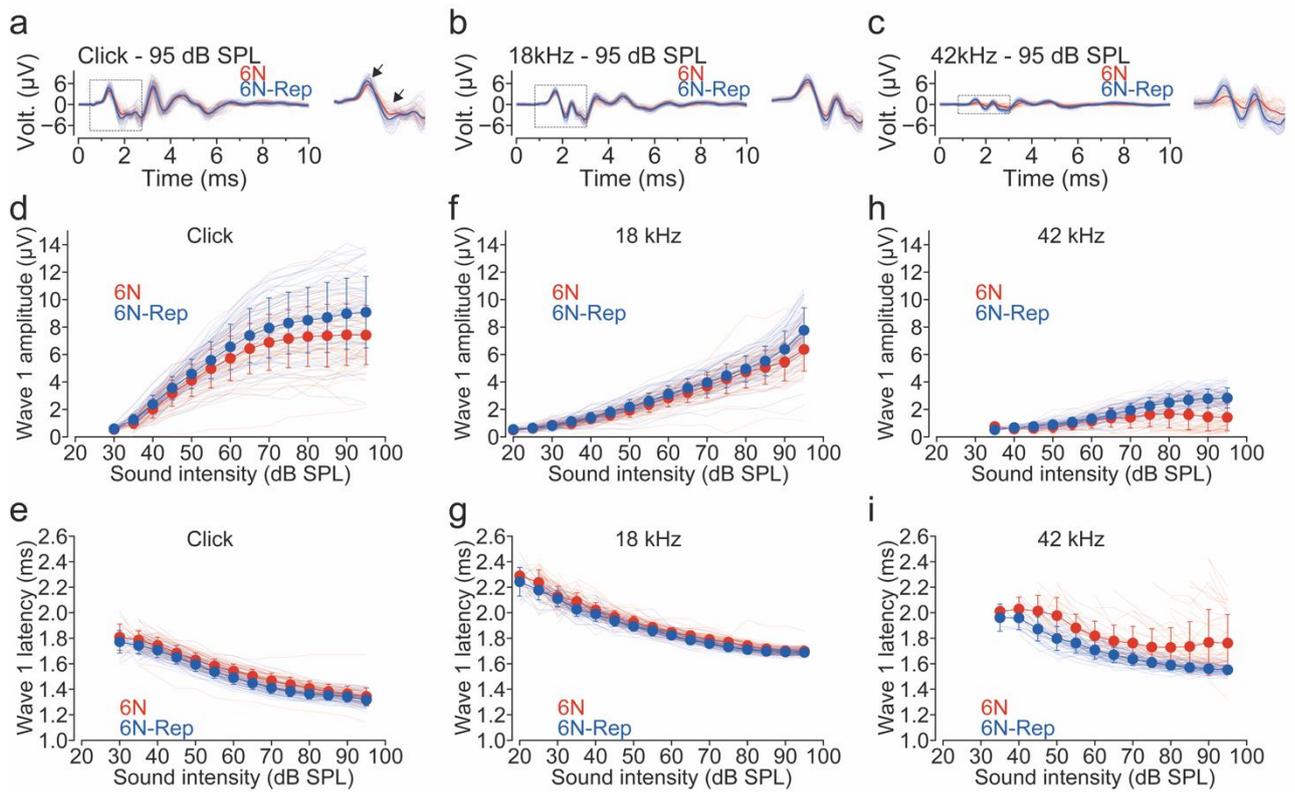
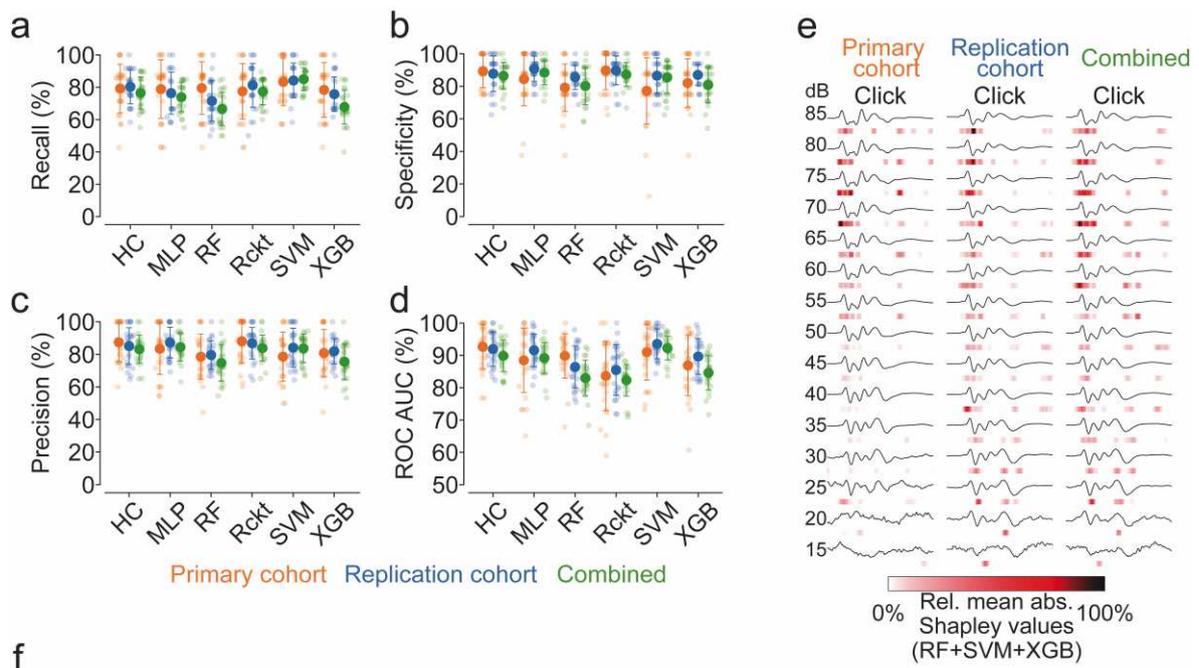


Figure 6.

928

929

930



f

Model: ROCKET

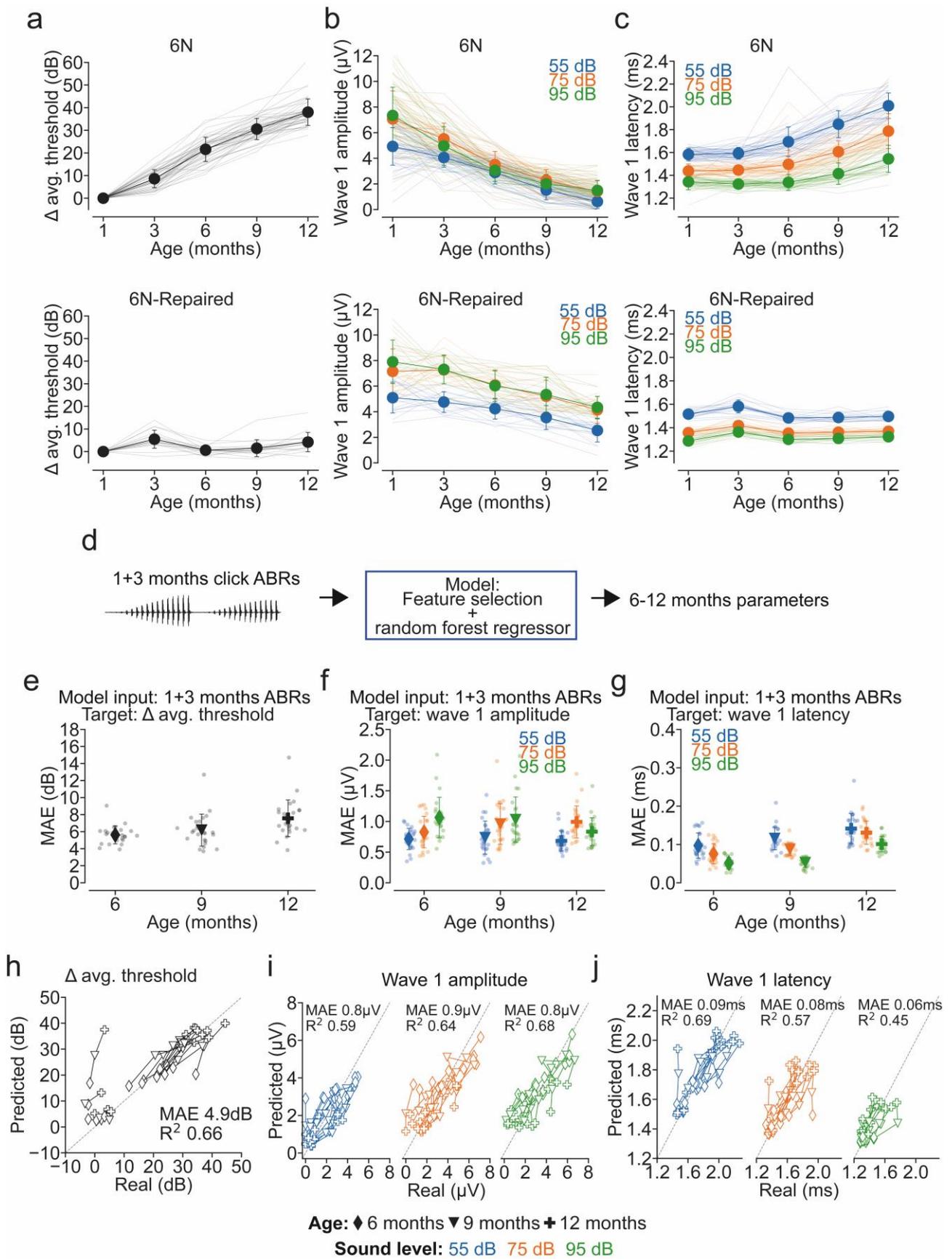
Testing set	Primary cohort			Replication cohort			Combined			
	Predicted genotype	6N	Rep.	Predicted genotype	6N	Rep.	Predicted genotype	6N	Rep.	
Primary cohort	6N	11	2	6N	47	38	6N	11	2	
	Rep.	2	11	Rep.	3	16	Rep.	2	11	
	Total	13	13	Total	50	54	Total	13	13	
		Recall	Specificity	Accuracy	Recall	Specificity	Accuracy	Recall	Specificity	Accuracy
		84.6%	84.6%	84.6%	94.0%	29.6%	60.6%	84.6%	84.6%	84.6%
		6N	Rep.	Total	6N	Rep.	Total	6N	Rep.	Total
		Actual genotype			Actual genotype			Actual genotype		
Replication cohort	6N	57	59	6N	16	1	6N	15	2	
	Rep.	28	44	Rep.	8	22	Rep.	9	21	
	Total	85	103	Total	24	23	Total	24	23	
		Recall	Specificity	Accuracy	Recall	Specificity	Accuracy	Recall	Specificity	Accuracy
		67.1%	42.7%	53.7%	66.7%	95.7%	80.9%	62.5%	91.3%	76.6%
		6N	Rep.	Total	6N	Rep.	Total	6N	Rep.	Total
		Actual genotype			Actual genotype			Actual genotype		
		116	72	188	17	30	47	17	30	47
		Precision	NPV	Accuracy	Precision	NPV	Accuracy	Precision	NPV	Accuracy
		84.6%	84.6%	84.6%	55.3%	84.2%	60.6%	84.6%	84.6%	84.6%
		6N	Rep.	Total	6N	Rep.	Total	6N	Rep.	Total
		Actual genotype			Actual genotype			Actual genotype		

Figure 7.

931

932

933



934

935

Figure 8.