

Validation of two efficient and robust smartphone-based threshold (GRaBr) and loudness (rACALOS) measures in typical home settings

Chen Xu,¹ Lena Schell-Majoor,¹ and Birger Kollmeier¹

¹ *Medizinische Physik and Cluster of Excellence Hearing4all, Universität Oldenburg, D-26111 Oldenburg, Germany*

Reliable hearing assessment at home can improve accessibility and reduce dependence on in-clinic testing. To be viable, home-based procedures must provide accurate results within short measurement times and remain robust to factors such as ambient noise and variable user attention. This study validated two such procedures—a Graded Response Bracketing method for pure-tone threshold estimation and a reinforced adaptive categorical loudness scaling method for loudness-growth assessment—using remote, smartphone-based testing. Fifteen young adults with normal hearing completed the tasks at home and in the laboratory. Ambient noise levels in home environments were also recorded. Test–retest reliability was assessed by repeating the home measurements on a separate occasion.

Remote measurements closely matched laboratory results, with mean differences below 1 dB for threshold estimation and below 5 dB for loudness scaling. Test–retest differences obtained at home were small, remaining below 2 dB for threshold estimation and below 1 dB for loudness scaling. These findings demonstrate that smartphone-based pure-tone audiometry and loudness-scaling assessments can achieve high accuracy, efficiency, and reliability when using these procedures, provided that basic acoustic-hygiene conditions (e.g., sufficiently low ambient noise) are maintained.

Keywords: remote audiology; ambient noise; validity and reliability; categorical loudness scaling

Corresponding author: Chen Xu

Contact:

chen.xu@uni-oldenburg.de

Department of Medical Physics and Acoustics, Faculty VI

Carl von Ossietzky Universität Oldenburg, 26111, Oldenburg, Germany

ACKNOWLEDGMENTS

This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy – EXC 2177/1 - Project ID 390895286.

DECLARATION OF CONFLICTING INTERESTS

The authors declare that there is no conflict of interest.

1 INTRODUCTION

2 Pure-tone audiometry, a threshold-based testing approach, is widely used to detect hearing
3 loss and guide audiogram-based hearing aid fitting. Its application for mobile and remote testing
4 has therefore received much interest (see Almufarrij et al., 2022 for a review). In contrast, supra-
5 threshold methods such as categorical loudness scaling (CLS, Hellbrück, 1987; Kollmeier, 1997)
6 to characterize supra-threshold deficits and support loudness-based fitting have received much
7 less attention in the literature on mobile or remote applications (Kopun et al., 2022; Xu et al.,
8 2024a). However, this method appears to be particularly suitable for such applications because it
9 is inherently robust to device calibration offsets (Almufarrij et al., 2022; Xu et al., 2025) due to
10 its parameter estimation relying on level differences rather than on absolute levels. Hence, both
11 of these methods are attractive for mobile and remote testing if they are reliable, robust against
12 disturbances, and provide valid estimates of what would be measured in a clinical setting. The
13 current paper aims at validating refined versions of them that fulfill these requirements.

14 Accurate estimation of hearing thresholds remains essential for obtaining reliable reference
15 data in both threshold- and supra-threshold-based tests. Previous work examined the effect of
16 experimenter supervision on smartphone-based pure-tone audiometry and adaptive CLS
17 (ACALOS, Brand and Hohmann, 2002) in a sound-attenuated booth with normal-hearing and
18 hearing-impaired listeners, showing no significant supervision effect (Xu et al., 2024b). To
19 further improve threshold estimation with respect to efficiency and robustness against inattention,
20 we developed the graded adaptive response bracketing (GRaBr) method (Xu et al., 2024a), a
21 model-free adaptive procedure that replaces the conventional yes-no paradigm with a two-
22 interval task in which listeners indicate whether they heard none, one, or two sounds. Compared
23 with the single-interval up-down (SIUD) method, which presents a cue tone in a second interval

24 at a fixed level above the target tone (Lecluyse & Meddis, 2009), GRaBr adaptively adjusts this
25 level difference. Simulation studies, including Xu et al. (2024a), showed that this adaptive
26 strategy yields substantially higher efficiency and robustness, particularly under both long- and
27 short-term inattention. However, a validation of the remote GRaBr listening test with real
28 participants in a home environment—where susceptibility to distractions is higher—against
29 laboratory-based assessments is still lacking, which motivated part of the present study.

30 Furthermore, we consider the adaptive CLS procedure for remote testing which poses
31 challenges in fluctuating ambient noise in home environments that could affect loudness
32 judgments at low sound pressure levels (SPL). Oetting et al. (2014) reported that the mean intra-
33 subject standard deviation of loudness levels close to the threshold was notably high (around 10
34 dB), yielding significant variability in the hearing threshold estimation from loudness judgments
35 near the threshold. In contrast, the mean intra-subject standard deviation at the level
36 corresponding to 50 categorical units on the loudness function (i.e., the estimated uncomfortable
37 loudness level) was below 5 dB, indicating that this measure is more reliable than the threshold
38 estimate. Hence, despite its advantages for assessing supra-threshold hearing problems
39 mentioned above, the conventional ACALOS procedure often fails to accurately predict
40 audiometric thresholds, likely due to differences in stimulus type and psychophysical paradigm
41 (Oetting et al., 2014). To overcome this limitation, we propose the reinforced adaptive
42 categorical loudness scaling (rACALOS) method (Xu, 2025), which incorporates a refined
43 threshold estimation process into ACALOS. This unified approach combines threshold and
44 supra-threshold assessments in a single, time-efficient test, improving estimation accuracy of the
45 hearing threshold and facilitating remote implementation on smartphones.

46 Taken together, the present study aims at validating the performance and test-retest reliability
47 of these novel smartphone-based methods for remote hearing assessment - GRaBr for threshold
48 estimation and rACALOS for categorical loudness scaling - when conducted by listeners in
49 typical home environments in comparison to the reference tests conducted inside the laboratory
50 using the baseline procedures. This study evaluates whether these procedures yield accurate and
51 reliable measures under realistic, uncontrolled acoustic conditions which is a prerequisite for
52 future integration into mobile hearing assessment platforms.

53 **METHODS**

54 **Participants**

55 Fifteen young adults with normal hearing (aged 20 to 35 years; 4 men, 11 women)
56 participated in this study. All participants were members of working groups or students at the
57 University of Oldenburg and were recruited primarily through informal requests within the labs.
58 The three authors did not participate in the study as subjects. All participants self-reported no
59 hearing issues and were classified as being normal hearing (NH) in former studies. Two
60 inclusion criteria were applied: (i) the air-conduction pure-tone average (PTA-4) at 0.5, 1, 2, and
61 4 kHz in the better ear had to be less than or equal to 20 dB HL, and (ii) symmetric hearing,
62 defined as a threshold difference of no more than 20 dB between ears at any test frequency. All
63 15 participants met these criteria. Some listeners (N = 5) received compensation of €12 per hour
64 for their participation, while others took part as part of their work duties. The study was
65 approved by the Research Ethics Committee of the University of Oldenburg (Drs. EK/2023/004).

66 **Reference Laboratory-Based Testing Inside a Booth**

67 We reused the data set from our previous study (Xu et al., 2024b) as the reference laboratory-
68 based listening tests. Full methodological details are provided in Xu et al. (2024b). Following a

69 repeated-measures design, the participants in the present study were the same individuals as
70 those in Xu et al. (2024b). All measurements were conducted by an experimenter inside a sound-
71 attenuated booth. Hearing thresholds were obtained using the SIUD procedure (Lecluyse &
72 Meddis, 2009), and loudness growth functions were assessed using the ACALOS procedure
73 (Brand & Hohmann, 2002). Please note that the SIUD procedure has been validated against the
74 clinical audiograms in Lecluyse and Meddis (2009), hence we assume that thresholds measured
75 by SIUD are comparable to the standard clinical audiogram. Both tests were performed at 0.25, 1,
76 and 4 kHz.

77 **Equipment, Procedure, and Environment for Remote Testing**

78 Prior to the start of remote testing, a test kit was assembled, which included a smartphone
79 (OnePlus, Android), a USB-C charger, and HD650 circumaural headphones (Sennheiser,
80 Wedemark, Hanover, Germany). The smartphone and headphones were pre-calibrated using a
81 Brüel & Kjær (B&K) artificial ear 4153, a B&K 0.5-inch microphone 4134, a B&K microphone
82 pre-amplifier 2669, and a B&K measuring amplifier 2610, with a target calibration level of 80
83 dB SPL. There was only one test kit in use, and the setup was calibrated once at the beginning of
84 data collection. Upon handing over the test kit, participants received a brief oral explanation of
85 the remote experiments, and consent forms were signed before they began. The oral instructions
86 consisted of three parts: test environment check, internet connection check, and steps for
87 conducting the remote experiments. The details are provided below. Participants could initiate
88 testing at home by connecting to the internet via WLAN and accessing the provided website. For
89 data security, a VPN connection was established using the ‘GlobalProtect’ app when accessing
90 the site. The workflow of the web-based application for remote testing was described in Xu et al.
91 (2024b). A Raspberry Pi 3 Model B (Raspberry Pi Foundation, UK), a Linux-based micro-

92 controller, served as the server hosting the measurement site. All behavioral data were stored on
93 an SD card within the Raspberry Pi, located at the University of Oldenburg.

94 Listeners were asked to complete remote testing at home on two separate days within a
95 week—one for the test measurement and the other for the retest measurement. Each day’s testing
96 lasted approximately half an hour, resulting in a total testing time of one hour per participant.
97 The home environments were primarily located in rural regions of northwestern Germany and
98 cities such as Oldenburg, Cloppenburg, Jever, and Bad Zwischenahn.

99 **Noise Level Measurement for Remote Testing**

100 Ambient noise was measured using the freely available “Decibel X” smartphone app
101 (SkyPaw Co., Ltd.), configured with an A-weighted filter and a 500-ms slow time weighting.
102 The app displayed real-time, average, and maximum noise levels but did not record audio. To
103 calibrate the smartphone’s internal microphone, a Class 3 digital sound level meter (Voltcraft
104 SL-100; ± 2 dB at 1 kHz), previously calibrated in our laboratory, was placed next to the phone
105 while 80-dB SPL narrowband noise was presented via a laptop. The app’s gain was then adjusted
106 to match the digital meter, requiring a correction of 13.7 dB. The same smartphone and
107 headphones were used for all participants to ensure consistent gain.

108 At the beginning of each remote session, participants documented the current ambient noise
109 level. They were instructed to monitor real-time noise using the app and to pause testing
110 whenever levels exceeded 45 dB(A), a threshold selected based on Kopun et al. (2022), who
111 showed that remote loudness-scaling results remain comparable to laboratory measurements
112 below 50 dB(A) ambient noise level. The time and location of each remote session were also
113 recorded.

114 **Remote Listening Tests at Home**

115 A total of 24 sessions were conducted by each subject at home, consisting of 4 listening tests
116 (SIUD, GRaBr, ACALOS, and rACALOS; see details below) across 3 test frequencies (0.25, 1,
117 and 4 kHz) and 2 runs (test and retest), presented in randomized order. Please note that the test
118 and retest measurements are referred to as Run 1 and Run 2, respectively. Participants were
119 instructed to take short breaks between sessions.

120 ***Remote pure-tone audiometry***

121 Remote pure-tone thresholds were measured using two adaptive procedures: the SIUD
122 method and the GRaBr approach (Lecluyse et al., 2009; Xu et al., 2024a). In both tasks, each
123 trial contained a probe tone and a higher-level cue tone, and listeners reported whether they
124 heard zero, one, or two tones. The cue tone was muted on 20% of trials (catch trials). The key
125 difference between the procedures is that SIUD uses a fixed 10-dB level difference between
126 probe and cue tones, whereas GRaBr adapts this difference: it starts at 10 dB, decreases to 5 dB
127 after the first “one-tone” response, and to 2 dB after the second.

128 In SIUD, the step size of the probe tone level began at 10 dB and was reduced to 2 dB after
129 the first “one-tone” response (by setting the level to the midpoint of the preceding two levels). In
130 GRaBr, step sizes were 8, 6, 3, and finally 1 dB after the first, second, and third reversals,
131 respectively. In both procedures, levels were adjusted adaptively—becoming harder after correct
132 responses and easier after incorrect ones. To enable fair comparison, key parameters (e.g.,
133 minimum trials, reversals, starting level) were matched. The initial cue-tone level for both
134 methods was drawn from a discrete uniform distribution of 55–65 dB SPL (1-dB steps). Each
135 track ended after at least 14 reversals and 10 trials, and the first four reversals were discarded.

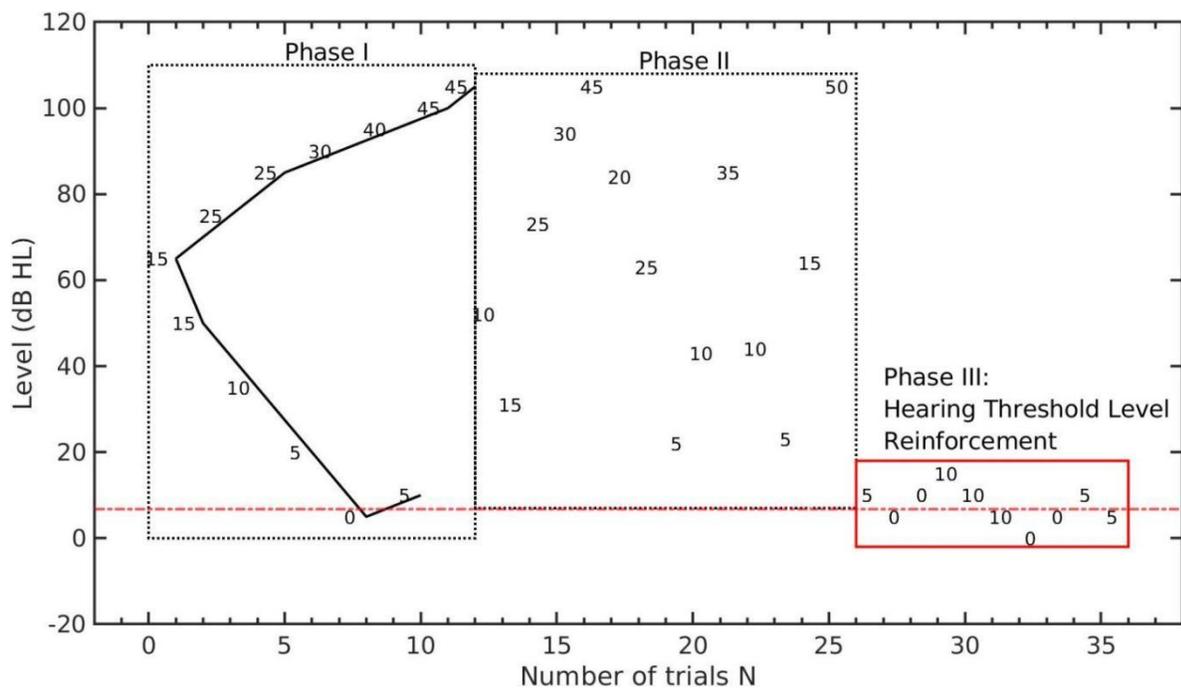
136 Each pure tone lasted 0.2 s, with cosine ramps of 0.02 s and a 0.3 s interval between tones. In
137 SIUD, the correct response rates were fitted to an S-shaped logistic psychometric function, and
138 the level at the 50% correct response point (L_{50}) was estimated as the hearing threshold. For
139 GRaBr, responses from the upper and lower tracks (i.e., the tracks of the cue tones and the probe
140 tones) were fitted to two independent psychometric functions, and the hearing threshold was
141 calculated as the mean level at the 50% correct response point of both functions (i.e.,
142 $0.5*(L_{50,upper} + L_{50,lower})$).

143 *Remote adaptive categorical loudness scaling*

144 The adaptive categorical loudness scaling (ACALOS) method was used to assess the
145 loudness growth function (Brand & Hohmann, 2002; ISO 16832, 2006). In the ACALOS task,
146 participants rated the loudness of stimuli on an 11-point scale with descriptors ranging
147 from 'very soft', "soft", "medium", "loud", and "very loud" with 4 unnamed intermediate
148 categories in between, plus the two limiting categories "not heard", and "too loud". The stimulus
149 levels, ranging from -10 to 105 dB, were presented in a pseudo-random order following an initial
150 estimation of the user-specific dynamic range (Phase I, see Fig. 1), which was updated to obtain
151 a more representative placement of test level in Phase II, encompassing 26 trials. At the end of
152 the procedure, a loudness growth function was modeled by fitting two linear segments and a
153 transition region using a Bezier fit, following the BTUX fitting method (Oetting et al., 2014).
154 Each frequency was assessed in a separate run.

155 An example run of the proposed rACALOS procedure is shown in Fig. 1. The rACALOS
156 followed the same adaptive rules as ACALOS during Phases I and II (see above) but presented
157 additional stimuli near the hearing threshold to better estimate HTL. The starting level of Phase
158 III was set at the minimum level reached in Phases I and II, plus 5 dB. In this phase, a one-up-

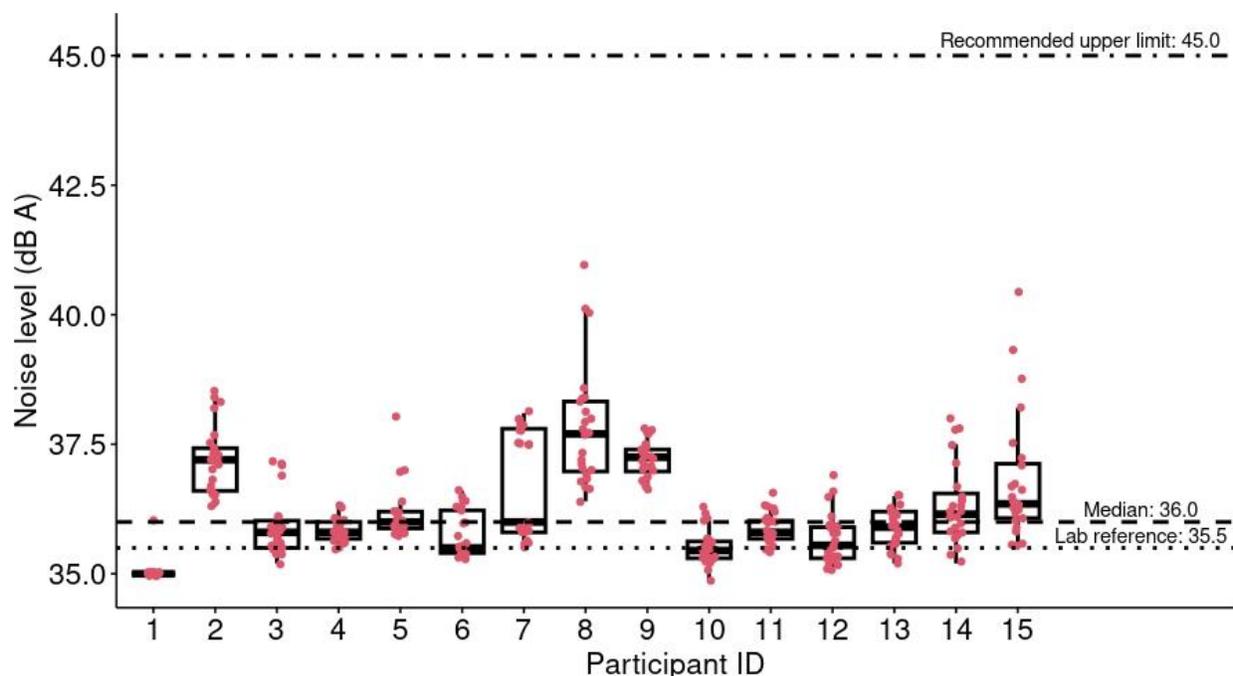
159 one-down adaptive rule was applied: the stimulus level increased by 5 dB if participants
 160 responded with "not heard" and decreased by 5 dB if they selected other loudness categories (e.g.,
 161 "very soft," "medium"). Phase III consisted of 10 trials. The stimuli used were one-third-octave-
 162 band low-noise noises (Kohlrausch et al., 1997). Each noise stimulus had a duration of 1 second
 163 with 0.05-second rise and fall ramps.



164
 165 Fig. 1. An example track of the reinforced adaptive categorical loudness scaling (rACALOS),
 166 where the level (in dB HL) is plotted as a function of the number of trials N. The listener's
 167 response (in categorical units (CU)) is annotated with numbers between 0 ('not heard') and 50
 168 ('too loud'). Left dotted rectangle region: Phase I ('dynamic range estimation'); Middle dotted
 169 rectangle region: Phase II ('presenting and re-estimation'); Right solid red rectangle region:
 170 Phase III ('hearing threshold level reinforcement'); Red dash-dotted line: target threshold. In
 171 Phase III, the step size is set to 5 dB, and the number of trials is set to 10.

172 **RESULTS**

173 **Noise Level Measurements**



174
175 Fig. 2. Ambient noise level (in dB A) measurement across participants (N = 15). Medians, 25th
176 and 75th percentiles, and interquartile ranges (IQR) are visualized in the box-plot while the end
177 of the whiskers denotes the minimum and maximum, indicating the 5th and 95th percentiles
178 respectively. Circles represent individual data points. Dotted line: lab reference (i.e., ambient
179 noise level measured within a booth). Dashed line: median value across subjects. Dot-dashed line:
180 recommended upper limit.

181 Fig. 2 presents a box plot of the ambient noise levels recorded by each participant (N = 15),
182 who documented the noise level a total of 24 times at the start of each measurement session,
183 corresponding to 24 measurement sessions at home within a week. Notably, the noise levels for
184 all participants remained below the recommended upper limit of 45 dB(A). The median noise
185 level across subjects was 36.0 dB, which was approximately 0.5 dB higher than the reference

186 noise level measured inside the sound-attenuated booth. The interquartile range (IQR) across
187 participants was 1.2 dB. Overall, the sound levels in participants' homes were considerably low
188 and comparable to those measured within the booth, indicating a suitable test environment. A
189 few participants (e.g., No. 2 and No. 8) lived near a train station, resulting in slightly elevated
190 noise levels compared to others. Additionally, one participant (No. 1) misinterpreted the task and
191 consistently rounded the recorded noise level to an integer, leading to uniform values across
192 sessions.

193 **Validation of the GRaBr procedure in a home environment**

194 To validate the home-based GRaBr procedure, its thresholds were compared with booth-
195 based thresholds obtained using the SIUD method at 0.25, 1, and 4 kHz. As reported in Table 1,
196 threshold differences were normally distributed (Shapiro–Wilk $p = 0.92$) and showed a
197 negligible mean bias of 0.4 dB. The correlation between home and booth measurements was
198 moderate—likely reflecting the limited variability in this normal-hearing sample—and the
199 RMSE was 7 dB, consistent with previous smartphone-based audiometry studies. Therefore, the
200 GRaBr procedure for home testing is comparable to the reference in-lab threshold measurements,
201 suggesting good validity.

202 Table 1. Validation of the home-based GRaBr procedure for threshold measurements against
203 the in-lab reference measurement using SIUD. Differences in hearing thresholds between the two
204 measurements are shown. Spearman correlation coefficients (R), root-mean-square error
205 (RMSE), and bias are reported for each of the three test frequencies and for their overall average.

| Frequency (Hz) | 250 | 1000 | 4000 | overall |
|----------------|------|------|------|---------|
| R (“spearman”) | 0.29 | 0.61 | 0.29 | 0.47 |
| RMSE (dB) | 8.0 | 4.6 | 7.9 | 7.0 |

| | | | | |
|-----------|-----|-----|------|-----|
| Bias (dB) | 5.5 | 1.0 | -5.3 | 0.4 |
|-----------|-----|-----|------|-----|

206

207 A two-way repeated-measures ANOVA showed no main effect of test environment ($F(1,14)$
208 $= 0.09$, $p = 0.77$). Frequency had a significant effect ($F(2,28) = 20.80$, $p < 0.05$), as did the
209 environment \times frequency interaction ($F(2,28) = 36.01$, $p < 0.05$). Post-hoc tests revealed no
210 significant difference between home and booth thresholds at 1 kHz, while small but significant
211 differences were found at 0.25 and 4 kHz.

212 **Validation of the rACALOS procedure in a home environment**

213 Loudness growth functions measured at home using rACALOS were compared with those
214 obtained in the booth using a standard ACALOS procedure at 0.25, 1, and 4 kHz. Differences
215 followed a normal distribution (Shapiro–Wilk $p = 0.22$), and the overall bias across categorical
216 units (CUs) was small (3.38 dB), indicating good validity of rACALOS relative to standard
217 ACALOS. Agreement metrics varied across loudness levels: correlations for CUs ≥ 35 ranged
218 from 0.57 to 0.62, whereas correlations for softer CUs were lower (< 0.45), accompanied by
219 larger RMSE values. This pattern reflects the higher variability inherent at near-threshold levels
220 rather than a systematic effect of home-environment noise, as deviations occurred in both
221 directions. See Table S1 in the supplementary material for the details of the validation for the
222 home-based rACALOS procedure.

223 A three-way repeated-measures ANOVA assessing test environment, frequency, and CU
224 revealed no significant main effect of environment. Frequency ($F(2,28) = 5.33$, $p < 0.05$) and CU
225 ($F(3,38) = 353.30$, $p < 0.05$) showed significant effects. Post-hoc comparisons indicated that
226 environment-related differences emerged only at a few specific combinations (4 kHz at 5, 25,

227 and 45 CU), while the majority of frequency–CU combinations showed no significant
228 differences.

229 **Test-Retest Reliability for home testing**

230 To assess the reliability of the four home-based listening tests, measurements from the initial
231 run (run 1) were compared with those from the retest (run 2). Details of the statistical analyses in
232 terms of reliability for remote ACALOS and rACALOS testing are provided in the
233 supplementary material (see Table S2).

234 ***SIUD and GRaBr for home testing***

235 As shown in Table 2, the GRaBr procedure showed test-retest intraclass correlation
236 coefficient (ICC) values exceeding 0.75 ($p < 0.05$), indicating good reliability across all three
237 frequencies, whereas the SIUD procedure yielded ICC values ranging from 0.59 to 0.77 ($p <$
238 0.05), reflecting moderate test-retest reliability. This difference was significant ($p < 0.05$), i.e.,
239 GRaBr demonstrated significantly higher test-retest reliability than SIUD based on these metrics.

240 A significant main effect of frequency was observed ($F(2, 28) = 15.46, p < 0.05$). Moreover,
241 pairwise t-tests were performed to assess reliability by comparing the two runs for both adaptive
242 procedures across all three frequencies, showing no significant differences (bias) between runs
243 in most cases, except for GRaBr at 1 kHz ($t(14) = 2.85, p < 0.05$).

244 Table 2. Test-retest reliability between two runs for remote GRaBr and SIUD testing at 0.25, 1,
245 and 4 kHz frequencies. Correlation coefficient (both ‘Pearson’ and ‘Spearman’), RMSE, bias,
246 and intraclass cross-correlation (ICC) between the test and retest measurements are reported.

| | GRaBr (N = 15) | | | SIUD (N = 15) | | |
|----------------|----------------|------|------|---------------|------|------|
| Frequency (Hz) | 250 | 1000 | 4000 | 250 | 1000 | 4000 |
| R (Pearson) | 0.70 | 0.90 | 0.73 | 0.43 | 0.61 | 0.63 |

| | | | | | | |
|--------------------|------|------|------|------|------|------|
| R (Spearman) | 0.47 | 0.78 | 0.65 | 0.40 | 0.57 | 0.65 |
| RMSE | 4.0 | 3.3 | 3.4 | 5.4 | 5.2 | 4.1 |
| Bias | -1.5 | 2.0 | -1.3 | -1.7 | -0.5 | -1.4 |
| ICC2k ^a | 0.77 | 0.88 | 0.83 | 0.59 | 0.77 | 0.74 |

247 ^a ICC2k: Average random raters

248 ^b ns: not significant; * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$; **** $p < 0.0001$

249

250 *ACALOS and rACALOS for home testing*

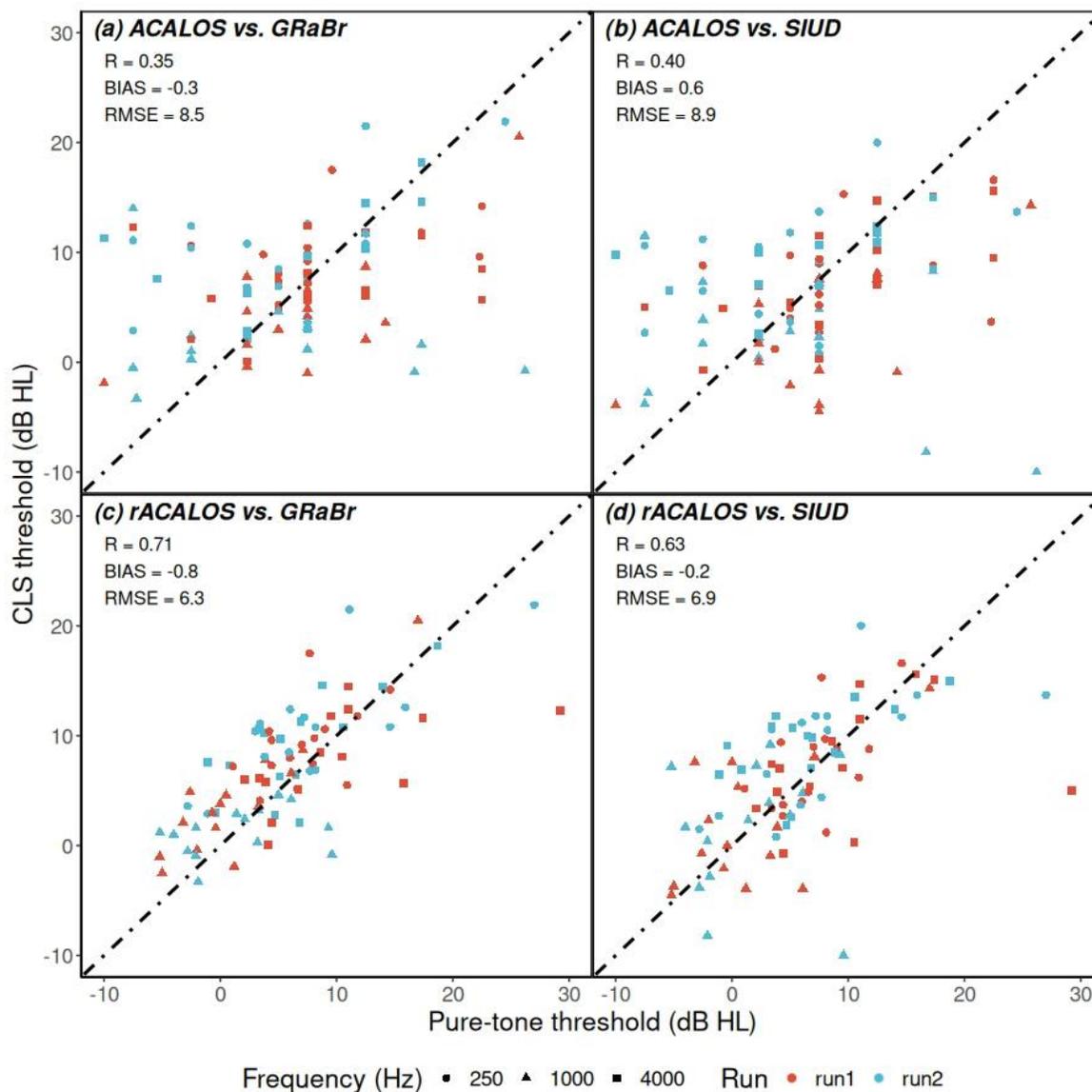
251 The reliability of the ACALOS and rACALOS procedures was assessed using across-run
252 bias (quantified by mean signed difference, MSD) and within-run variability (measured by mean
253 interquartile range, MIQR) (please see Kopun et al. (2022) for the definitions of the MSD and
254 MIQR). Both adaptive procedures demonstrated an MSD of less than 5 dB at all frequencies,
255 indicating a small across-run bias. Most MIQR values did not exceed 10 dB for either procedure
256 at the three frequencies, although they were typically larger than 10 dB at 5, 10, and 15 CU,
257 reflecting a consistent within-run variability. Overall, these metrics suggested that both
258 ACALOS and rACALOS exhibited strong reliability.

259 A repeated measures ANOVA revealed a significant main effect of the procedure, indicating
260 a statistically significant difference between ACALOS and rACALOS ($p < 0.05$). Compared to
261 the ACALOS procedure, the rACALOS procedure yielded lower mean sound levels (by 1.9 dB)
262 across participants and CU. Since the rACALOS and ACALOS procedures are identical in
263 Phases I and II, this difference is likely attributable to the additional trials included in Phase III
264 of the rACALOS procedure (see Fig. 1).

265 No significant effect was found for frequency ($F(2, 28) = 2.51, p = 0.10$), and as expected,
266 the two runs (test and retest measurements) did not differ ($F(1, 14) = 1.97, p = 0.18$). A
267 subsequent post-hoc t-test compared median levels of the ACALOS and rACALOS procedures

268 between runs 1 and 2 across three frequencies and 11 categories, indicating that median levels
269 for run 1 did not significantly differ from those for run 2 in most cases (31 out of 33 groups of
270 comparison = 3 levels of frequency * 11 levels of CU), except for two groups (measurements at
271 0.25 kHz for 25CU ($t(29) = 2.08$) and 40 CU ($t(29) = 2.21$)).

272 Accuracy of HTL Estimation for the rACALOS procedure



273
274 Fig. 3. Scatter plot for comparison between pure-tone (abscissa) and CLS (ordinate) thresholds
275 (defined as the level at 2.5 CU of the loudness growth functions) in dB HL of $N = 15$ individual

276 listeners. Frequency is labeled with different shapes while the run is denoted with different
277 colors (run1: red, run2: blue). A set of statistical metrics (R, Bias, and RMSE) are reported in the
278 top-left corner. For rACALOS, 10 additional trials with a step size of 5 dB were used.

279 Pure-tone audiometric thresholds are plotted against CLS thresholds (i.e., levels at 2.5 CU of
280 the loudness growth functions) for two runs and three frequencies in Fig. 3. Compared to
281 ACALOS, the majority of rACALOS points were consistently and closely clustered around the
282 diagonal line, indicating that thresholds estimated by the rACALOS method aligned more
283 closely with pure-tone thresholds than those from baseline ACALOS and, hence, provide
284 improved accuracy in threshold estimation. Quantitatively, R values increased by 36% for
285 GRaBr and 23% for SIUD when ACALOS was reinforced near the hearing threshold level.
286 Additionally, RMSE values for the rACALOS method decreased by approximately 2 dB
287 compared to the baseline, while biases remained unchanged. Overall, the reinforcement of
288 baseline ACALOS positively influenced cross-correlation and reduced error.

289 The highest correlation coefficient and lowest RMSE were observed between GRaBr
290 thresholds and rACALOS, followed by SIUD and rACALOS. In contrast, the unmodified
291 ACALOS procedure showed lower correlation coefficients and higher RMSEs for both threshold
292 estimation methods, indicating the superior performance of rACALOS, as confirmed by t-tests (p
293 < 0.05).

294 **DISCUSSION**

295 This study validated the performance of two novel smartphone-based procedures—GRaBr
296 for pure-tone audiometry and rACALOS for categorical loudness scaling—conducted by
297 listeners independently in typical home environments. The results demonstrate that accurate and

298 reliable measurements can be obtained without professional supervision or sound-treated booths,
299 confirming the feasibility of remote hearing assessments using standard smartphones.

300 **Ambient Sound Levels in Home Settings**

301 The median background level measured across participants' homes was 36.0 dB(A),
302 comparable to quiet residential environments and well below limits defined by ANSI S3.1–1999
303 (R2018) for covered-ear conditions. These values fall within the updated permissible ranges for
304 circumaural earphones proposed by Margolis et al. (2022), ensuring that ambient sound did not
305 measurably affect the results. While noise levels were not experimentally manipulated,
306 participants were advised to continuously monitor them to confirm that real-world domestic
307 conditions remained suitable for valid measurements. Consequently, the obtained home results
308 were closely aligned with those from the booth.

309 The recorded noise levels were lower than those in earlier remote audiometry studies (Storey
310 et al., 2014; Brennan-Jones et al., 2016; Swanepoel et al., 2015) and comparable to the quietest
311 non-sound-treated settings reported by Serpanos et al. (2022) and Bean et al. (2022). It is likely
312 that our participants conducted the smartphone-based listening tests at home in rural areas during
313 the morning or evening, whereas other studies typically test in clinical settings located in urban
314 areas during the daytime, which tend to be noisier. Overall, the measured noise conditions
315 confirm that typical homes can provide sufficiently quiet environments for reliable mobile
316 hearing assessments.

317 **Pure-Tone Audiometry (GRaBr)**

318 Smartphone-based audiometry with GRaBr showed good agreement with in-booth reference
319 measurements, with a mean bias of 0.4 dB. While both SIUD ($0.59 < ICC < 0.77$) and GRaBr
320 performed reliably, GRaBr achieved higher test–retest consistency ($ICC > 0.75$, $p < 0.05$) and

321 lower within-subject variability, confirming its improved stability predicted by simulation (Xu et
322 al., 2024a). These findings extend prior validations of boothless audiometry (Maclennan-Smith
323 et al., 2013; Swanepoel et al., 2015; Serpanos et al., 2022) to an unsupervised, mobile context.

324 The correlation between at-home and in-booth thresholds ($R = 0.47$) was lower than in
325 studies involving hearing-impaired listeners ($R > 0.9$; Maclennan-Smith et al., 2013), likely
326 because our normal-hearing sample exhibited limited threshold variability. Future refinements -
327 such as automated quality control or remote calibration - may further optimize the GRaBr for
328 remote threshold testing.

329 **(reinforced) Adaptive Categorical Loudness Scaling**

330 Both ACALOS and its reinforced version, rACALOS, yielded valid and reliable loudness
331 functions in typical home conditions. The bias between in-booth and at-home measurements (3.4
332 dB) was smaller than that reported by Kopun et al. (2022), indicating strong agreement across
333 environments. The enhanced rACALOS design produced smaller within-run (IQR) and across-
334 run (MSD) variability than both baseline ACALOS and previously published CLS methods
335 (Rasetshwane et al., 2015; Fultz et al., 2020; Kopun et al., 2022). At 5 CU, the mean IQR values
336 for rACALOS (7.7 dB at 1 kHz and 6.7 dB at 4 kHz) were notably lower than those in earlier
337 studies (≈ 10 – 15 dB), reflecting improved stability near the hearing threshold.

338 Across-run bias was also smaller than in prior work, likely due to the lower average ambient
339 noise level and the reinforced threshold estimation sequence. The additional trials near threshold
340 in rACALOS substantially reduced estimation uncertainty, resulting in the smallest overall
341 variability among tested methods. Compared with alternative adaptive CLS frameworks (Fultz et
342 al., 2020), rACALOS offers higher precision and efficiency, validating its use as a robust supra-
343 threshold measurement for mobile and unsupervised hearing assessments.

344 **Accuracy of HTL Estimation**

345 Table 3 presents a comparison between our current study and several state-of-the-art works
346 (Fultz et al., 2020; Trevino et al., 2016; Sanchez-Lopez et al., 2021) by evaluating the cross-
347 correlation between CLS and pure-tone thresholds. Multiple CLS methods, including Fixed-level
348 (FL), Maximum expected information-Median (MEL-Med), Maximum expected information-
349 Maximum likelihood (MEL-ML), Slope-adaptive (SA), ACALOS, and rACALOS, were used to
350 estimate thresholds, which were then compared with pure-tone thresholds measured using
351 various audiometric methods such as a clinical audiometer, SIUD, and GRaBr. In the studies by
352 Fultz et al. (2020) and Trevino et al. (2016), correlation coefficients R values (according to
353 “spearman”) ranged from 0.21 to 0.26 for the threshold estimated from all four CLS methods ,
354 indicating a relatively weak cross-correlation. Additionally, the RMSEs and biases in these
355 studies were notably large, suggesting that CLS thresholds did not align well with pure-tone
356 thresholds. In contrast, Sanchez-Lopez et al. (2021) applied a baseline ACALOS method using
357 the same audiometric procedure as Fultz et al. (2020), and while the R-value did not significantly
358 improve, both RMSE and bias were notably reduced. In our study, we employed SIUD and
359 GRaBr to measure pure-tone thresholds, yielding a stronger cross-correlation and smaller bias,
360 although the RMSE was slightly larger or comparable to that reported by Sanchez-Lopez et al.
361 (2021).

362 Considering all the studies, the rACALOS method consistently produces thresholds closest to
363 pure-tone thresholds, outperforming other CLS methods. This finding is consistent with
364 computer simulations that examined threshold variability across different parameter
365 combinations of rACALOS compared to the original ACALOS procedure. However, it is
366 important to note that rACALOS requires more measurement time due to the increased number

367 of trials focused on converging near the HTL. Additionally, while pure-tone thresholds obtained
368 using clinical audiometers are still widely regarded as the ‘gold standard’, more precise (see Xu
369 et al., 2024a; Lecluyse & Meddis., 2009) and criterion-free methods such as SIUD and GRaBr
370 may produce stronger correlations with CLS thresholds. It is also crucial to recognize that this
371 comparison is based on a small sample of young NH listeners, and the conclusions may differ if
372 HI listeners are included or if a larger participant pool is studied. This consideration is
373 particularly relevant for potential discrepancies between the narrowband noise thresholds
374 estimated by the CLS methods used here and the pulsed pure-tone thresholds assessed via
375 audiograms. While threshold differences in our study sample of young NH listeners were
376 minimal, variations in stimulus characteristics—such as spectral extent and modulation
377 spectrum—may yield threshold differences in naïve listeners with hearing impairments.
378 Nonetheless, we also expect these differences to be minimal, as the low-noise, third-octave-band
379 noise utilized here shares key perceptual characteristics with a frequency-modulated sinusoid
380 with minor envelope fluctuations and an instantaneous frequency confined well within a critical
381 band (cf. Zwicker & Fastl, 2013).

382 Table 3. Comparison of threshold estimates obtained from different pure-tone audiometry
383 methods and categorical loudness scaling (CLS) methods across several state-of-the-art studies
384 and the present study. Spearman correlation coefficients (R), root-mean-square error (RMSE),
385 bias, and sample size (N) are reported for each method combination. The highest R value and the
386 lowest RMSE and bias within each comparison set are highlighted in bold. (FL = Fixed-level;
387 MEL-Med = Maximum-expected-information, median; MEL-ML = Maximum-expected-
388 information, maximum likelihood; SA = Slope-adaptive, as defined in Fultz et al. (2020))

| | Audiometric method | CLS method | N | R (Spearman) | RMSE (dB) | Bias (dB) |
|--|--------------------|------------|----|--------------|------------|-------------|
| Fultz et al. 2020; Trevino et al. 2016 | Audiometer | FL | 17 | 0.21 | 12.2 | -6.9 |
| | | MEL-Med | | 0.26 | 25.3 | -18.0 |
| | | MEL-ML | | 0.26 | 15.5 | -10.6 |
| | | SA | | 0.21 | 15.7 | -8.4 |
| Sanchez-Lopez et al. 2021 | Audiometer | ACALOS | 11 | 0.24 | 7.1 | -2.3 |
| current | SIUD | ACALOS | 15 | 0.44 | 9.4 | 1.5 |
| | GRaBr | | | 0.38 | 9.0 | 1.0 |
| | SIUD | rACALOS | | 0.59 | 7.8 | 0.5 |
| | GRaBr | | | 0.71 | 6.9 | 0.04 |

389

390 **Advantages of rACALOS**

391 **Increased time efficiency:** The rACALOS procedure combines two listening tests—pure-
392 tone audiometry and ACALOS—into a single, integrated protocol. This approach significantly
393 reduces the measurement time required for participants by eliminating the need for separate tests
394 with separate instructions and training procedures.

395 **Improved HTL accuracy:** Compared to the original ACALOS, rACALOS includes
396 additional trials near the hearing threshold level (HTL), enhancing the precision of HTL
397 estimation (see Table 3 for details). These modifications enable the seamless integration of
398 audiometric measurement into the ACALOS framework.

399 **Consistent user interface and minimal training requirements:** The rACALOS procedure
400 employs the same interface and measurement paradigm as ACALOS. Consequently, participants
401 already familiar with ACALOS require no additional training, while new users need to become
402 familiar with only a single, consistent interface for both listening tests.

403 **Limitations and Outlook**

404 In this study, we conducted smartphone-based listening tests outside of a sound booth,
405 preceded by ambient noise level measurements. Given the recommendation for conducting tests
406 in rather quiet acoustical conditions, the testing environments generally exhibited a low
407 background noise level. However, many individuals live in urban regions with significant vehicle
408 or industrial noise, where real-world environments are typically much noisier. Testing in such
409 noisy conditions warrants further investigation. Potential solutions, such as circumaural muffs or
410 noise-canceling earphones (NCE), could prove effective. For instance, Saliba et al. (2017)
411 evaluated mobile-based audiometry under 50 dB A background noise, using passive and active
412 noise cancellation by placing circumaural muffs over insert headphones, successfully reducing
413 noise. Similarly, Clark et al. (2017) tested NCE (BoseQuietComfort 15) in a patient consultation
414 room and found that NCE sufficiently attenuated ambient noise below the ANSI standards.

415 Another key concern for out-of-booth audiometric tests is distraction. As noted by Margolis
416 et al. (2022), background noise not only causes direct masking but also acts as a source of
417 distraction. Their study demonstrated that increasing background noise levels led to elevated
418 hearing thresholds and higher subjective ratings of distraction. Xu et al. (2024a) further
419 supported these findings, characterizing distraction from internal noise (e.g., background noise)
420 as long-term inattention. They also proposed and simulated short-term inattention—where
421 listeners are distracted by external events—during mobile hearing tests, though this has yet to be
422 validated with human participants.

423 Another limitation of this study is the use of an integrated microphone for noise measurement.
424 Studies like Kopun et al. (2022) recommend using an external microphone, such as the MicW

425 iBoundary, which provides higher accuracy in capturing frequency characteristics and calibration
426 precision compared to the internal microphone used here.

427 Enhanced calibration of smartphone microphones can, in principle, be achieved using an
428 external reference sound, such as a whistle tone produced by a standard empty beer bottle
429 (Scharf et al., 2024). However, accurate calibration of the playback level remains a challenge and
430 is essential for precise auditory measurements. While pragmatic approaches such as biological
431 calibration—using a normal-hearing reference listener to define 0 dB HL (Honeth, 2010;
432 Masalski et al., 2014)—may offer a practical alternative, this aspect was beyond the scope of the
433 present proof-of-concept study with its limited sample size. Future work could focus on data-
434 driven estimation of calibration coefficients to enable fully unsupervised field applications (Xu et
435 al., 2025).

436 In the future, we plan to expand our study by increasing the number of normal-hearing (NH)
437 participants and incorporating hearing-impaired (HI) participants. Compared to NH listeners, we
438 expect the validity of pure-tone audiometry and ACALOS tests in HI listeners to be comparable,
439 as supported by previous studies (e.g., Hazan et al., 2022; Bean et al., 2022; Xu et al., 2024b),
440 suggesting that hearing loss does not significantly affect test validity. Regarding reliability, HI
441 listeners are expected to show similar or sometimes even higher test-retest reliability in
442 audiometry than NH listeners (Hazan et al., 2022), likely due to their elevated hearing thresholds,
443 which reduce the impact of ambient noise. Similarly, in the ACALOS procedure, HI listeners
444 should exhibit similar or even greater reliability at lower levels of the loudness growth function,
445 as they are less affected by background noise. This may have a limited impact on the accuracy of
446 the loudness growth slope fitted to the data, which is generally increased in hearing-impaired
447 listeners exhibiting recruitment. However, since the slope estimate is derived from loudness

448 judgments across multiple supra-threshold levels, it is considered a reliable and robust measure,
449 even when obtained through self-administered smartphone assessments.

450 Finally, Shen et al. (2018) and Kursun et al. (2023) introduced a quick categorical loudness
451 scaling (qCLS) procedure based on a Bayesian adaptive method, which can estimate equal
452 loudness contours within just 5 minutes. Given its efficiency and accuracy, incorporating qCLS
453 into future smartphone-based loudness tests is worth considering. However, it remains uncertain
454 whether qCLS can estimate hearing thresholds as precisely as the rACALOS developed in this
455 study, highlighting the need for further research to evaluate its threshold accuracy in comparison.

456 **CONCLUSION**

457 This study validates two efficient and robust smartphone-based methods for hearing
458 assessment—GRaBr for pure-tone threshold estimation and rACALOS for categorical loudness
459 scaling—under typical home conditions but with calibrated hardware. The findings demonstrate
460 that both measures provide reliable and valid outcomes comparable to those obtained in
461 controlled laboratory settings.

462 In the validation experiment, pure-tone audiometry using GRaBr and categorical loudness
463 scaling using rACALOS yielded equivalent results between home and sound-attenuated
464 environments at 0.25, 1, and 4 kHz, confirming the validity of remote testing. Test–retest results
465 further indicate moderate-to-good reliability across sessions, supporting the consistency of
466 home-based measurements.

467 GRaBr outperformed the SIUD method in reliability and accuracy across all tested
468 frequencies, making it a preferred approach for mobile threshold estimation. Similarly,
469 rACALOS produced threshold estimates that were closer to audiometric thresholds measured by
470 SIUD and GRaBr than those from the baseline ACALOS procedure, demonstrating its advantage

471 in improving hearing threshold estimation and test efficiency by combining threshold and
472 loudness assessments.

473 Overall, the results support the feasibility of using GRaBr and rACALOS for reliable and
474 efficient hearing assessments in real-world home environments, marking a step forward toward
475 validated smartphone-based audiometry for large-scale or remote applications.

476 GLOSSARY

| Abbreviation | Meaning |
|--------------|--|
| ACALOS | adaptive categorical loudness scaling |
| ANOVA | analysis of variance |
| B&K | Brüel&Kjaer |
| BTUX | fitting method for loudness function in ACALOS |
| CLS | categorical loudness scaling |
| CU | categorical units |
| FL | fixed-level procedure |
| GRaBr | graded response bracketing |
| HI | hearing impaired |
| HTL | hearing threshold level (at 2.5 CU on the loudness function) |
| ICC | intraclass cross-correlation |
| IQR | interquartile ranges |
| MEL-Med | maximum expected information-median |
| MEL-ML | maximum expected information-maximum likelihood |
| MIQR | mean interquartile range |
| MPANLs | maximum permissible ambient noise levels |
| MSD | mean signed difference |
| NCE | noise reduction earphones |
| NH | normal hearing |
| PTA | pure-tone average |

| | |
|---------|--|
| qCLS | quick categorical loudness scaling |
| rACALOS | reinforced adaptive categorical loudness scaling |
| RMSE | root mean squared error |
| SA | slope-adaptive procedure |
| SIUD | single interval up and down |
| SPL | sound pressure level |

477

478 **REFERENCES (BIBLIOGRAPHIC)**

- 479 Akeroyd, M. A., Arlinger, S., Bentler, R. A., Boothroyd, A., Dillier, N., Dreschler, W. A., ... &
480 Kollmeier, B. (2015). International Collegium of Rehabilitative Audiology (ICRA)
481 recommendations for the construction of multilingual speech tests: ICRA Working Group
482 on Multilingual Speech Tests. *International journal of audiology*, 54(sup2), 17-22.
- 483 Almufarrij, I., Dillon, H., Dawes, P., Moore, D. R., Yeung, W., Charalambous, A. P., ... &
484 Munro, K. J. (2022). Web-and app-based tools for remote hearing assessment: a scoping
485 review. *International Journal of Audiology*, 1-14.
- 486 American National Standards Institute. Maximum Permissible Ambient Noise Levels for
487 Audiometric Test Rooms. (ANSI S3.1–R2018). New York, NY: American National
488 Standards Institute; 2018
- 489 Bean, B. N., Roberts, R. A., Picou, E. M., Angley, G. P., & Edwards, A. J. (2022). Automated
490 audiometry in quiet and simulated exam room noise for listeners with normal hearing and
491 impaired hearing. *Journal of the American Academy of Audiology*, 33(01), 006-013.
- 492 Behar, A. (2021). Audiometric tests without booths. *International Journal of Environmental
493 Research and Public Health*, 18(6), 3073.

- 494 Bianco, R., Mills, G., de Kerangal, M., Rosen, S., & Chait, M. (2021). Reward enhances online
495 participants' engagement with a demanding auditory task. *Trends in Hearing*, 25,
496 23312165211025941.
- 497 Bland, J. M., & Altman, D. G. (1999). Measuring agreement in method comparison studies.
498 *Statistical methods in medical research*, 8(2), 135-160.
- 499 Bland, J. M., & Altman, D. G. (2007). Agreement between methods of measurement with
500 multiple observations per individual. *Journal of biopharmaceutical statistics*, 17(4), 571-
501 582.
- 502 Brand, T., & Hohmann, V. (2002). An adaptive procedure for categorical loudness scaling. *The*
503 *Journal of the Acoustical Society of America*, 112(4), 1597-1604.
- 504 Brand, T., 2000. Analysis and Optimization of Psychophysical Procedures in Audi-ology.
505 Universität Oldenburg, Germany. PhD thesis.
- 506 Brennan-Jones, C. G., Eikelboom, R. H., Swanepoel, D. W., Friedland, P. L., & Atlas, M. D.
507 (2016). Clinical validation of automated audiometry with continuous noise-monitoring in
508 a clinically heterogeneous population outside a sound-treated environment. *International*
509 *journal of audiology*, 55(9), 507-513.
- 510 Buhl, M., Akin, G., Saak, S., Eysholdt, U., Radeloff, A., Kollmeier, B., & Hildebrandt, A. (2022).
511 Expert validation of prediction models for a clinical decision-support system in audiology.
512 *Frontiers in Neurology*, 13, 960012.
- 513 Clark, J. G., Brady, M., Earl, B. R., Scheifele, P. M., Snyder, L., & Clark, S. D. (2017). Use of
514 noise cancellation earphones in out-of-booth audiometric evaluations. *International*
515 *Journal of Audiology*, 56(12), 989-996.

516 Fultz, S. E., Neely, S. T., Kopun, J. G., & Rasetshwane, D. M. (2020). Maximum expected
517 information approach for improving efficiency of categorical loudness scaling. *Frontiers*
518 *in Psychology*, 11, 578352.

519 Giavarina, D. (2015). Understanding bland altman analysis. *Biochemia medica*, 25(2), 141-151.

520 Hazan, A., Luberadzka, J., Rivilla, J., Snik, A., Albers, B., Méndez, N., ... & Kinsbergen, J.
521 (2022). Home-Based Audiometry With a Smartphone App: Reliable Results?. *American*
522 *Journal of Audiology*, 31(3S), 914-922.

523 Hellbrück, J. (1987). How to measure loudness under natural conditions?. *The Japanese Journal*
524 *of Ergonomics*, 23(5), 307-310.

525 Honeth, L., Bexelius, C., Eriksson, M., Sandin, S., Litton, J. E., Rosenhall, U., ... & Bagger-
526 Sjöbäck, D. (2010). An internet-based hearing test for simple audiometry in nonclinical
527 settings: preliminary validation and proof of principle. *Otology & Neurotology*, 31(5),
528 708-714.

529 ISO 16832, 2006. Acoustics Loudness Scaling by Means of Categories. Standard of the
530 International Organization for Standardization, Geneva, Switzerland.

531 Kohlrausch, A., Fassel, R., Van Der Heijden, M., Kortekaas, R., Van De Par, S., Oxenham, A. J.,
532 & Püschel, D. (1997). Detection of tones in low-noise noise: Further evidence for the role
533 of envelope fluctuations. *Acta Acustica united with Acustica*, 83(4), 659-669.

534 Kollmeier, B. (Ed.). (1997). Hörflächenskalierung: Grundlagen und Anwendung der kategorialen
535 Lautheitsskalierung für Hördiagnostik und Hörgeräte-Versorgung. Median-Verlag von
536 Killisch-Horn.

- 537 Kollmeier, B., Gilkey, R. H., & Sieben, U. K. (1988). Adaptive staircase techniques in
538 psychoacoustics: A comparison of human data and a mathematical model. *The Journal of*
539 *the Acoustical Society of America*, 83(5), 1852-1862.
- 540 Koo, T. K., & Li, M. Y. (2016). A guideline of selecting and reporting intraclass correlation
541 coefficients for reliability research. *Journal of chiropractic medicine*, 15(2), 155-163.
- 542 Kopun, J. G., Turner, M., Harris, S. E., Kamerer, A. M., Neely, S. T., & Rasetshwane, D. M.
543 (2022). Evaluation of Remote Categorical Loudness Scaling. *American journal of*
544 *audiology*, 31(1), 45-56.
- 545 Kursun, Bertan & Petersen, Erik & Shen, Yi. (2023). Exploring Self-directed Hearing-aid Fitting
546 with No Booth And No Audiogram. 10.13140/RG.2.2.19575.19360.
- 547 Lecluyse, W., & Meddis, R. (2009). A simple single-interval adaptive procedure for estimating
548 thresholds in normal and impaired listeners. *The Journal of the Acoustical Society of*
549 *America*, 126(5), 2570-2579.
- 550 MacLennan-Smith, F., Swanepoel, D. W., & Hall III, J. W. (2013). Validity of diagnostic pure-
551 tone audiometry without a sound-treated environment in older adults. *International*
552 *journal of audiology*, 52(2), 66-73.
- 553 Margolis, R. H., Saly, G. L., & Wilson, R. H. (2022). Ambient Noise Monitoring during Pure-
554 Tone Audiometry. *Journal of the American Academy of Audiology*, 33(01), 045-056.
- 555 Masalski, M., Grysiński, T., & Kręcicki, T. (2014). Biological calibration for web-based hearing
556 tests: evaluation of the methods. *Journal of medical Internet research*, 16(1), e2798.
- 557 Meinke, D. K., & Martin, W. H. (2023). Boothless audiometry: Ambient noise considerations.
558 *The Journal of the Acoustical Society of America*, 153(1), 26-39.

- 559 Min, S. H., & Zhou, J. (2021). SmploT: An R package for easy and elegant data visualization.
560 *Frontiers in Genetics*, 12, 2582.
- 561 Oetting, D., Brand, T., & Ewert, S. D. (2014). Optimized loudness-function estimation for
562 categorical loudness scaling data. *Hearing Research*, 316, 16-27.
- 563 Ooster, J., Krueger, M., Bach, J. H., Wagener, K. C., Kollmeier, B., & Meyer, B. T. (2020).
564 Speech audiometry at home: automated listening tests via smart speakers with normal-
565 hearing and hearing-impaired listeners. *Trends in Hearing*, 24, 2331216520970011.
- 566 Peng, Z. E., Waz, S., Buss, E., Shen, Y., Richards, V., Bharadwaj, H., ... & Venezia, J. H. (2022).
567 Remote testing for psychological and physiological acoustics. *The Journal of the*
568 *Acoustical Society of America*, 151(5), 3116-3128.
- 569 Rasetshwane, D. M., Trevino, A. C., Gombert, J. N., Liebig-Trehearn, L., Kopun, J. G., Jesteadt,
570 W., ... & Gorga, M. P. (2015). Categorical loudness scaling and equal-loudness contours
571 in listeners with normal hearing and hearing loss. *The Journal of the Acoustical Society*
572 *of America*, 137(4), 1899-1913.
- 573 Revelle, W. (2018). *psych: Procedures for psychological, psychometric, and personality research*.
- 574 Robler, S. K., Coco, L., & Krumm, M. (2022). Telehealth solutions for assessing auditory
575 outcomes related to noise and ototoxic exposures in clinic and research. *The Journal of*
576 *the Acoustical Society of America*, 152(3), 1737-1754.
- 577 Saliba, J., Al-Reefi, M., Carriere, J. S., Verma, N., Provencal, C., & Rappaport, J. M. (2017).
578 Accuracy of mobile-based audiometry in the evaluation of hearing loss in quiet and noisy
579 environments. *Otolaryngology–Head and Neck Surgery*, 156(4), 706-711.
- 580 Sanchez-Lopez, R., Nielsen, S. G., El-Haj-Ali, M., Bianchi, F., Fereczkowski, M., Cañete, O.
581 M., ... & Santurette, S. (2021). Auditory tests for characterizing hearing deficits in

582 listeners with various hearing abilities: The BEAR test battery. *Frontiers in neuroscience*,
583 15, 724007.

584 Scharf, M. K., Huber, R., Schulte, M., & Kollmeier, B. (2024). Microphone calibration
585 estimation for mobile audiological tests with resonating bottles. *International Journal of*
586 *Audiology*, 1-7. DOI: 10.1080/14992027.2024.2395416

587 Serpanos, Y. C., Hobbs, M., Nunez, K., Gambino, L., & Butler, J. (2022). Adapting audiology
588 procedures during the pandemic: Validity and efficacy of testing outside a sound booth.
589 *American Journal of Audiology*, 31(1), 91-100.

590 Shen, Y., Zhang, C., & Zhang, Z. (2018). Feasibility of interleaved Bayesian adaptive procedures
591 in estimating the equal-loudness contour. *The Journal of the Acoustical Society of*
592 *America*, 144(4), 2363-2374.

593 Storey, K. K., Muñoz, K., Nelson, L., Larsen, J., & White, K. (2014). Ambient noise impact on
594 accuracy of automated hearing assessment. *International Journal of Audiology*, 53(10),
595 730-736.

596 Swanepoel, D. W., Matthysen, C., Eikelboom, R. H., Clark, J. L., & Hall III, J. W. (2015). Pure-
597 tone audiometry outside a sound booth using earphone attenuation, integrated noise
598 monitoring, and automation. *International Journal of Audiology*, 54(11), 777-785.

599 Swanepoel, D. W., Mngemane, S., Molemong, S., Mkwanzazi, H., & Tutshini, S. (2010). Hearing
600 assessment—reliability, accuracy, and efficiency of automated audiometry. *Telemedicine*
601 *and e-Health*, 16(5), 557-563.

602 Thai-Van, H., Joly, C. A., Idriss, S., Melki, J. B., Desmettre, M., Bonneuil, M., ... & Reynard, P.
603 (2023). Online digital audiometry vs. conventional audiometry: a multi-centre
604 comparative clinical study. *International Journal of Audiology*, 62(4), 362-367.

- 605 Trevino, A. C., Jesteadt, W., & Neely, S. T. (2016). Development of a multi-category
606 psychometric function to model categorical loudness measurements. *The Journal of the*
607 *Acoustical Society of America*, 140(4), 2571-2583.
- 608 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D. A., François, R., ... & Yutani,
609 H. (2019). Welcome to the Tidyverse. *Journal of open source software*, 4(43), 1686.
- 610 Wiseman, K., Slotkin, J., Spratford, M., Haggerty, A., Heusinkvelt, M., Weintraub, S., ... &
611 McCreery, R. (2023). Validation of a tablet-based assessment of auditory sensitivity for
612 researchers. *Behavior research methods*, 55(6), 2838-2852.
- 613 Xu, C. (2025). *Crucial Elements of a Virtual Hearing Clinic on Mobile Devices - Psychophysics,*
614 *Diagnostic Parameter Estimation and Validation* (Doctoral dissertation, University of
615 Oldenburg Press (UOLP)), ISBN 978-3-8142-2424-4. doi:
616 10.13140/RG.2.2.14023.00169
- 617 Xu, C., & Kollmeier, B. (2025). Calibration offset estimation in mobile hearing tests via
618 categorical loudness scaling. *arXiv preprint arXiv:2508.14824*.
- 619 Xu, C., Hülsmeyer, D., Buhl, M., & Kollmeier, B. (2024a). How Does Inattention Influence the
620 Robustness and Efficiency of Adaptive Procedures in the Context of Psychoacoustic
621 Assessments via Smartphone?. *Trends in Hearing*, 28, 23312165241288051.
- 622 Xu, C., Schell-Majoor, L., & Kollmeier, B. (2024b). Development and verification of non-
623 supervised smartphone-based methods for assessing pure-tone thresholds and loudness
624 perception. *International Journal of Audiology*, 1-11.
- 625 Zhao, S., Brown, C. A., Holt, L. L., & Dick, F. (2022). Robust and Efficient Online Auditory
626 Psychophysics. *Trends in hearing*, 26, 23312165221118792.

627 Zwicker, E., & Fastl, H. (2013). Psychoacoustics: Facts and models (Vol. 22). Springer Science
628 & Business Media.