1    Auditory spatial analysis in reverberant audio-visual multi-talker

2    environments with congruent and incongruent visual room information

3

4    Axel Ahrens

5    Hearing Systems Section, Department of Health Technology, Technical University of Denmark,

6    2800 Kgs. Lyngby, Denmark

7    aahr@dtu.dk

8

9    Kasper Duemose Lund

10    Hearing Systems Section, Department of Health Technology, Technical University of Denmark,

11    2800 Kgs. Lyngby, Denmark

12    kasperduemose@gmail.com

13

14

15

16

17

18

19

20    Date: 13 February 2022

21

22    Running Title: Auditory scene analysis

23    Abstract

24    In multi-talker situation, listeners have the challenge to identify a target speech source out of a

25    mixture of interfering background noises. In the current study it was investigate how listeners

26    analyze audio-visual scenes with varying complexity in terms of number of talkers and

27    reverberation. Furthermore, the visual information of the room was either coherent with the

28    acoustic room or incoherent. The listeners' task was to locate an ongoing speech source in a

29    mixture of other speech sources. The 3D audio-visual scenarios were presented using a

30    loudspeaker array and virtual reality glasses. It was shown that room reverberation as well as the

31    number of talkers in a scene influence the ability to analyze an auditory scene in terms of accuracy

32    and response time. Incongruent visual information of the room did not affect this ability. When

33    few talkers were presented simultaneously, listeners were able to quickly and accurately detect a

34    target talker even in adverse room acoustical conditions. Reverberation started to affect the

35    response time when four or more talkers were presented. The number of talkers became a

36    significant factor for five or more simultaneous talkers.

37

38    Keywords: Speech perception; Virtual Reality; Localization

39

## 40  I.   Introduction

41  The human auditory system has the ability to focus on a speech stream in the presence of

42  interfering speech stimuli. Such a multi-talker scenario has been termed the cocktail-party situation

43  (Bronkhorst, 2000; Cherry, 1953). Many factors are known to reduce the ability to understand

44  speech in such a cocktail-party situation, e.g., the level of the target speech relative to the

45  interferers, the number of talkers, or the type of listening room. These effects are commonly

46  measured by asking the listeners to repeat a word or a sentence or to write down the perceived

47  stimulus. However, in our daily life the task in a cocktail-party situation is usually different, where

48  it is necessary to follow a conversation and to identify a certain topic or continuous speech stream

49  out of an interfering speech mixture. In the current study we investigated the ability of listeners to

50  analyze an acoustic scene with varying complexity in terms of number of interfering talkers, room

51  reverberation and coherency of visual room information.

52  The number of interfering talkers has been shown to influence the intelligibility of a target talker.

53  (S. A. Simpson & Cooke, 2005) showed that the intelligibility decreases when increasing the

54  number of interfering speech sources for up to eight interfering talkers, as the ability to listen into

55  speech gaps is reduced and at the same time the interfering speech remains intelligible and can be

56  confused with the target speech. When further increasing the number of interfering talkers, the

57  intelligibility was shown to improve as the interferers become more noise-like and therefore do

58  not contain understandable speech.

59  Reflections and reverberation are present in nearly all communication scenarios. Room

60  reverberation has been shown to negatively affect speech perception in a number of studies (Best

61  et al., 2015; Bronkhorst & Plomp, 1990; Moncur & Dirks, 1967; Nabelek & Mason, 1981; Nábělek

62  & Pickett, 1974). Particularly, the diffuse reverberation, i.e., the late reverberant tail, has been

63  shown to reduce speech intelligibility, while early reflections do not seem to harm, or might even

64      improve speech perception (Arweiler et al., 2013; Arweiler & Buchholz, 2011; Warzybok et al.,

65      2013).

66      Previous studies have investigated the ability of listeners to identify and locate speech in the

67      presence of other speech sources. (Kopčo et al., 2010) measured the localization accuracy of a

68      digit spoken by a female talker in the presence of words spoken by male interfering talkers.  The

69      target and the interferers were all presented in the frontal area of the listener. They found that the

70      presence of the interferers reduced the localization accuracy. (Buchholz & Best, 2020) measured

71      localization accuracy with a similar target digit as in (Kopčo et al., 2010) but with a more realistic

72      background noise scene. The interfering signals were seven paired conversations (both male and

73      female) at various locations in a simulated cafeteria. Results showed that the localization accuracy

74      was only affected by the noise when the target source was distant but not when it was nearby. This

75      finding suggests an interaction with reverberation, as farther sources have more reverberant energy

76      relative to the direct sound compared to nearby sources.

77      While these studies focused on the ability to locate a speech signal in a speech background,

78      (Hawley et al., 1999) investigated both the localization accuracy of speech as well as the

79      intelligibility. They showed that the inability to correctly locate a source did not limit the ability

80      to correctly understand it. However, the number of interfering sources was limited to three.

81      (Weller et al., 2016) presented a novel method to evaluate the ability to analyze a complex acoustic

82      scene. They asked their listeners to judge the location of all talkers presented in a virtual cocktail-

83      party situation by indicating the gender of the talkers. When varying the number of simultaneously

84      presented talkers, they found that normal-hearing listeners were able to correctly locate and count

85      the number of talkers for up to four sources. When six talkers were presented, the accuracy

86      decreased.

87    Most of the beforementioned studies focused on the ability to localize speech but less to

88    comprehend the speech. However, in a real-world cocktail party, listeners need to perform both

89    tasks to successfully communicate. In the current study, we asked listeners to locate a talker

90    speaking about a certain topic, while presenting a varying number of other simultaneous talkers.

91    Thus, the primary task was to understand the speech and the secondary task to locate the talker.

92    The experiment was conducted in an audio-visual virtual environment using a loudspeaker array

93    and virtual reality glasses. The listeners' task was to indicate a semi-transparent avatar at the

94    location of an acoustic source talking about a topic indicated by an icon. The sources were located

95    at one of fifteen possible locations with 15° horizontal separation. The number of simultaneous

96    speech sources was varied between two and eight. Three virtual rooms were simulated visually

97    and acoustically. Furthermore, a condition with incongruent audio-visual cues was presented by

98    visually showing the anechoic room and acoustically presenting the reverberant room or vice

99    versa.

100

## 101 II.   Methods
102

### 103   A. Participants
104

105    Thirteen Danish native speaking normal-hearing listeners aged 20-26 years participated in the

106    experiment (7 female and 6 male). Participants were paid on an hourly basis and gave consent to

107    an ethics agreement approved by the Science-Ethics Committee for the Capital Region of Denmark

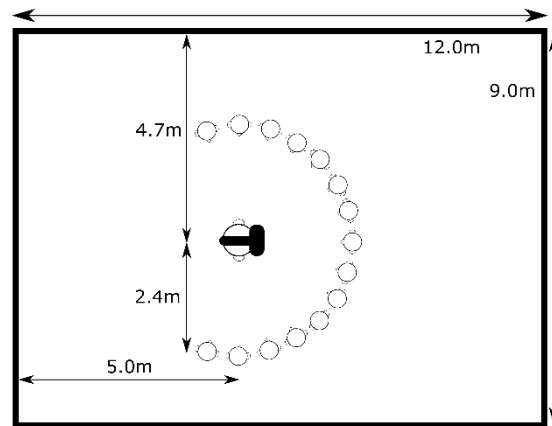108    (reference H-16036391).

109

110 ## B. Material

111

112 The speech material for target and interferers was taken from a database of anechoically recorded

113 monologues in Danish (see (Lund et al., 2019) for details[1]). Each monologue was designed with

114 characteristic features in mind, ensuring significant difference of the content. The database consists

115 of ten monologues each spoken by ten native Danish speakers.

116

117 ## C. Audio-visual rooms

118

119 Three different acoustic and visual rooms were used in this study, a high-reverberant room, a mid-

120 reverberant room and an anechoic room. The dimensions of all three rooms remained constant as

121 shown in Figure 1, both acoustically and visually. However, the surface materials differed.
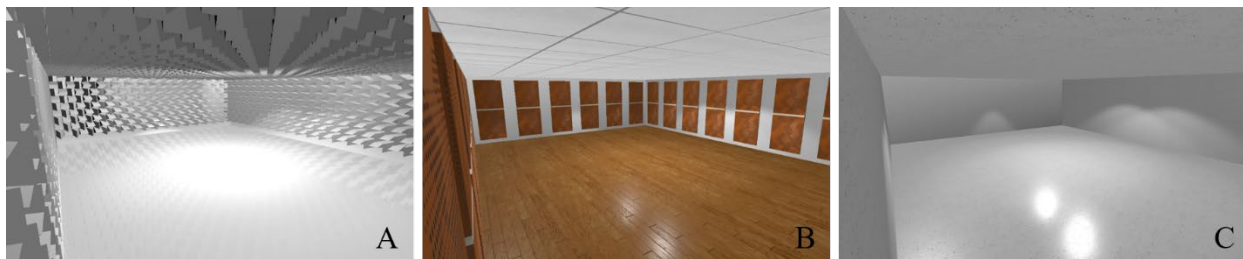
122



124 Figure 1: Top view of the virtual audio-visual room. The listener is wearing VR glasses with a
125 visual simulation of the room including 15 potential talker positions at 2.4m distance in the frontal
126 hemisphere visualized by the head icons. The height of the room is 2.8m.

---

[1] Data available: https://data.dtu.dk/articles/Recordings_of_Danish_Monologues_for_Hearing_Research/9746285

127

128    Figure 2 shows the visual appearances of the three rooms. Figure 2A shows the anechoic room

129    with foam wedges as commonly seen in anechoic chambers. For the acoustic reproduction of this

130    room only the direct sound was reproduced from single loudspeakers. In Figure 2B the mid-

131    reverberant room can be seen. The visual as well as the acoustical properties were similar to a large

132    living room. The highly reverberant room is shown in Figure 2C. It was modelled with bare

133    concrete surfaces to simulate a highly reverberant, yet realistic environment.

134



136    Figure 2: Visual appearance of the three virtual rooms. A: anechoic, B: mid-reverberant, C: high-
137    reverberant. The dimensions in the rooms are identical, while the surface materials differ.

138

139    The rooms were simulated using the room acoustic simulation software Odeon (Odeon A/S, Kgs.

140    Lyngby, Denmark) with the materials and surface absorption coefficients as shown in Table 1. For

141    the anechoic room, only the direct sound was considered. In Figure 3 the reverberation time, clarity

142    and direct-to-reverberant ratio of the three rooms are shown. The reverberation time as well as the

143    clarity were calculated using the ITA-toolbox (Berzborn et al., 2017), the direct-to-reverberant

144    ratio was calculated as the ratio between the direct sound and the reflections. Mind that for the

145    anechoic condition the clarity and direct-to-reverberant ratio are infinite as no reflections are
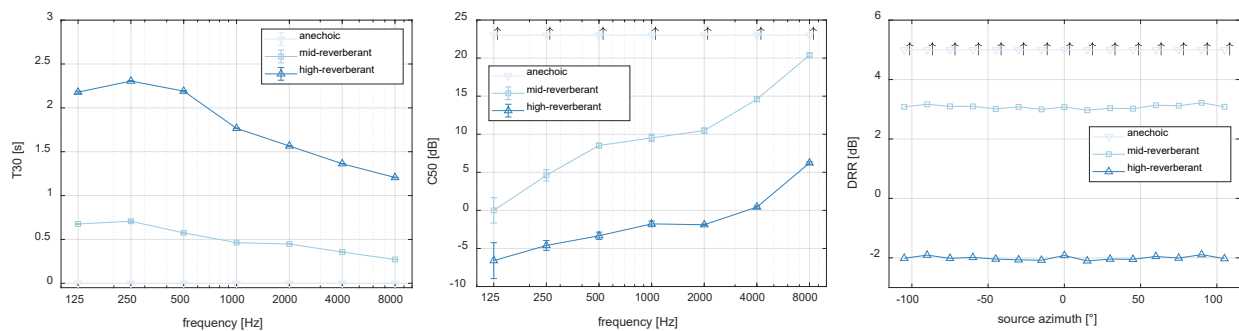
146    present which is indicated with arrows.

147

148     Table 1: Absorption coefficients (α) of the surfaces in the mid-reverberant and high-reverberant room.

| α (mid-rev/high-rev) | 63 Hz | 125 Hz | 250 Hz | 500 Hz | 1 kHz | 2 kHz | 4 kHz | 8 kHz |
|---|---|---|---|---|---|---|---|---|
| Side walls Wooden panels/Brick | 0.2/0.06 | 0.2/0.06 | 0.2/0.06 | 0.3/0.07 | 0.4/0.07 | 0.4/0.07 | 0.5/0.08 | 0.5/0.09 |
| Floor Parquet/Concrete | 0.2/0.05 | 0.2/0.05 | 0.15/0.05 | 0.1/0.05 | 0.1/0.07 | 0.05/0.07 | 0.1/0.07 | 0.1/0.07 |
| Ceiling Gypsumboard/Concrete | 0.3/0.05 | 0.3/0.05 | 0.35/0.05 | 0.4/0.05 | 0.4/0.07 | 0.4/0.07 | 0.5/0.07 | 0.55/0.07 |

149

150



152     Figure 3: Reverberation time (T30), Clarity (C50) and the direct-to-reverberant ratio (DRR) for
153     the three rooms. The T30 and the C50 are shown with respect to octave frequency bands. The DRR
154     is shown with respect to the source azimuth angle. The arrows indicate that the measure is infinite.

155

156     ## D. Task

157

158     The listeners' task was to identify the location of a talker amongst concurrent talker(s) in a virtual

159     audio-visual room according to the story in the monologue. Accuracy and completion time of the

160     task was emphasized by advising the listeners to "find the correct story as fast as possible". The

161     number of concurrent talkers varied between two and eight, thus the number of interfering talkers

162     varied between one and seven. An icon visualizing the target story content was displayed on the

163     backwall in the visual virtual room. The 15 possible talker positions were always represented by

164    semi-transparent humanoid shapes independent of the actual number of concurrent talkers. Figure

165    1 visualizes the possible talker locations between -105° to 105° separated by 15° in the frontal

166    hemisphere at a distance of 2.4 m. The task was performed by pointing at the position where the

167    target talker was perceived. The participants were using a virtual reality controller that included

168    the visual appearance of a laser pointer in the virtual room.

169    For each scene a unique talker, story and position was randomly chosen as the target. Between one

170    and seven masking talkers were included in a similar way. No talker, story or position could occur

171    twice at the same time. For each trial, the acoustic talkers were presented for 120 seconds. The

172    stories were started at a random point in time and were repeated from the beginning after finishing.

173    Thus, no bias towards the beginning of each story was introduced. The listener could indicate the

174    perceived target talker position at any time, even after the audio had stopped. Each individual

175    talker was presented at a sound pressure level of 55 dB SPL.

176    Three congruent audio-visual rooms were used as described above, an anechoic, a mid-reverberant

177    and a high-reverberant room. In addition to the conditions with congruent audio and visual room

178    information, two conditions with incongruent audio-visual cues were considered. These were

179    anechoic acoustics with the appearance of a highly reverberant room and high-reverberant

180    acoustics with the visuals of the anechoic room. Thus, five room conditions were tested. Each of

181    the conditions was repeated three times resulting in 105 trials, five audio-visual conditions and

182    between two and eight concurrent talkers.

183    Prior to the experiment, the listeners performed a familiarization phase, where they were

184    familiarized with the speech material and the story content but not with the task itself. The anechoic

185    version of the ten stories were played back via headphones in a randomized order. Each talker was

186    randomly assigned to one of the stories. Thus, listeners heard each story and each talker once. For

187    the training, listeners were instructed to focus on unique content features or passages of the stories.

188    After completed training listeners were seated in the loudspeaker environment and introduced to

189    the listening task and the interaction method using the VR controller.

190

191    E. Virtual audio-visual setup
192

193    The virtual visual scenes were rendered on the head-mounted display (HMD) of an HTC Vive Pro

194    Eye (HTC Vive system, HTC Corporation, New Taipei City, Taiwan). This system allowed to

195    track the listeners motion and record eye gaze and pupil dilation from inside the HMD with a

196    sampling frequency of up to 120 Hz and an accuracy between 0.5° and 1.1°. The visual virtual

197    scenes were modeled and displayed using Unity (Unity Technologies, San Francisco, California,

198    USA).

199    The acoustic scenes were reproduced on 64-channel spherical loudspeaker array housed in an

200    anechoic chamber (see (Ahrens, Marschall, et al., 2019) for details). The loudspeaker signals were

201    generated using the room acoustic simulation using the LoRA-toolbox (Favrot & Buchholz, 2010).

202    For the loudspeaker playback the nearest loudspeaker mapping was applied, where the direct sound

203    as well as the early reflections are mapped to the nearest loudspeaker. The late reverberant tail is

204    reproduced using $1^{st}$ order ambisonics to achieve a diffuse acoustic field (Favrot & Buchholz,

205    2010).

206    F. Outcome measures and statistical analyses
207

208    To evaluate the listeners' ability to successfully analyze a cocktail-party scenario, two outcome

209    measures were evaluated. First, the ability to correctly identify and locate the target talker. This

210    allows for a binary right/wrong analysis as well as a localization error in degrees. Second, the

211    response time of the listener from audio onset to decision.

212    The outcome measures were analyzed using an analysis of variance of mixed linear models. The

213    computational analyses were done using the statistical computing software R(R Core Team, 2020)

214    and the lmerTest (Kuznetsova et al., 2017) package. Within factor analyses were conducted using

215    marginal means implemented in the emmeans package (Lenth, 2020) with Tukey correction for

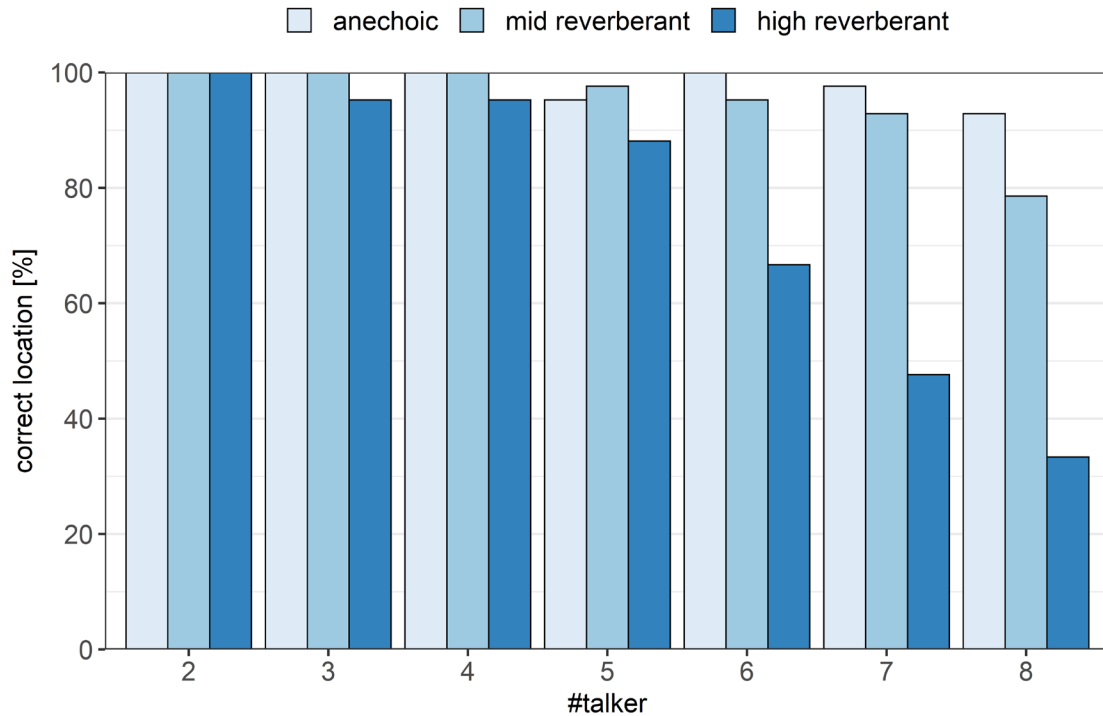216    multiple comparisons.

## 217 III.    Results
218

### 219    A. Coherent audio-visual room information
220

221    Figure 4 shows the percentage of correctly located stories. Each bar contains 42 datapoints across

222    the 14 participants and three repetitions. When few talkers are in a scene, the participants were

223    able to accurately locate the correct story in all reverberation conditions. In scenes with more than

224    five talkers, the accuracy in the high-reverberant condition (dark blue) decreases. In the mid-

225    reverberant condition such a decrease can only be observed when eight talkers are in a scene. In

226    the anechoic condition, the participants were able to accurately locate the target story for all
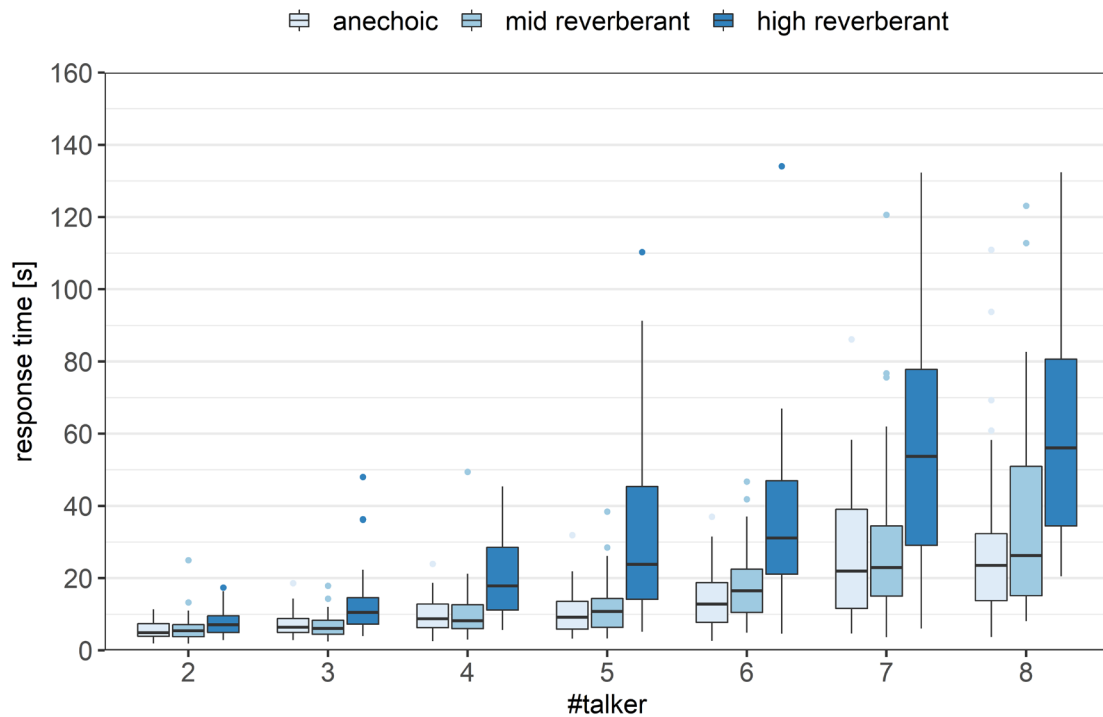
227    numbers of talkers.

228

Figure 4: The percentage of correct response locations. Each bar contains 42 datapoints across subjects and repetitions. The three colors indicate the room conditions.

Figure 5 shows the response time of the correct responses when two to eight talkers were presented simultaneously. The response time is displayed for the audio-visually coherent room conditions with varying reverberation times indicated with the different colors. With an increasing number of simultaneous talkers, the time needed to identify the target talker increased [$F_{(6,755.2)}=73.1$, $p<0.0001$]. The response time was also found to be dependent on the reverberation time [$F_{(2,755.6)}=83.1$, $p<0.0001$]. Furthermore, the interaction term between the number of talkers and the reverberation time was found significant [$F_{(12,754.8)}=5.4$, $p<0.0001$]. Specifically, the high-reverberant condition was found to lead to a higher response time when four or more talkers were presented [$p<0.05$] but not with less than four talkers [$p>0.5$]. The differences between the high-reverberant condition and the anechoic/mid-reverberant condition increases with larger numbers

243    of talkers. No significant differences between the anechoic and the mid reverberant condition was

244    found [p>0.1] across all number of talkers.
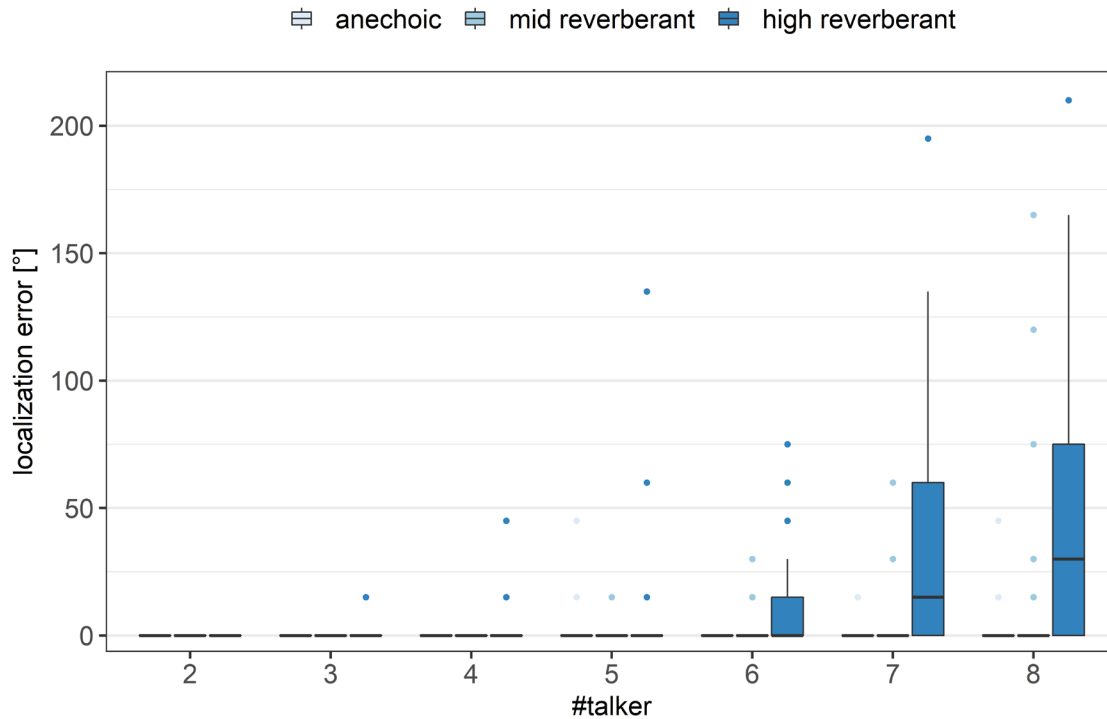
245



246

247    Figure 5: Response time with respect to the number of talkers in a scene of all correct responses.
248    The colors indicate the room reverberation conditions. The boxes cover the range between the 25th
249    and the 75th percentile. The horizontal line in the boxes indicates the median. The whiskers extend
250    to 1.5 times the inter-quartile range. Outliers are indicated as dots.

251

252    In Figure 6 the localization error is shown. In the high-reverberation condition an increasing mean

253    localization error was found for six and more talkers, with the eight-talker setting resulting in a

254    median error of 30°, i.e., two potential positions error from the target location. In the anechoic and

255    mid-reverberant conditions only few errors were found, indicated as outliers in Figure 6.

256

Figure 6: Localization error with respect to the number of talkers. The three colors indicate the room conditions. The boxes cover the range between the 25th and the 75th percentile. The horizontal line in the boxes indicates the median. The whiskers extend to 1.5 times the inter-quartile range. Outliers are indicated as dots.
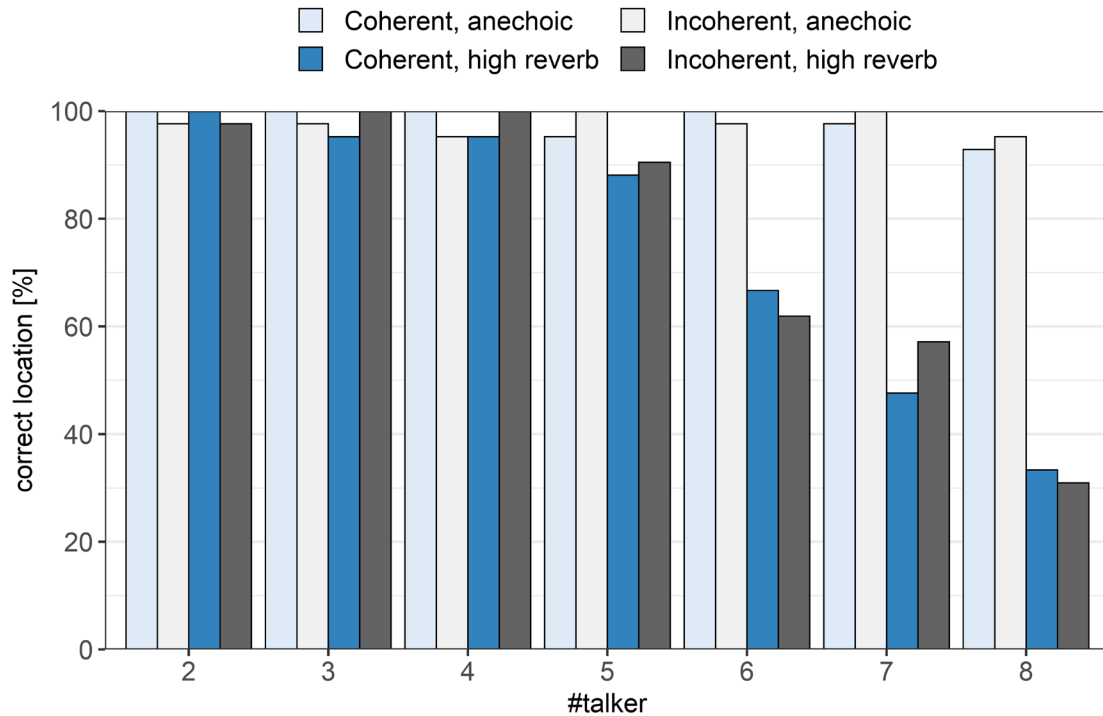
## B. Incoherent audio-visual room information

Figure 7 shows the percentage of correctly identified stories, comparing the coherent and the

incoherent audio-visual conditions with and without reverberation. The light blue/grey bars

indicate the acoustically anechoic conditions and the dark bars the acoustically reverberant

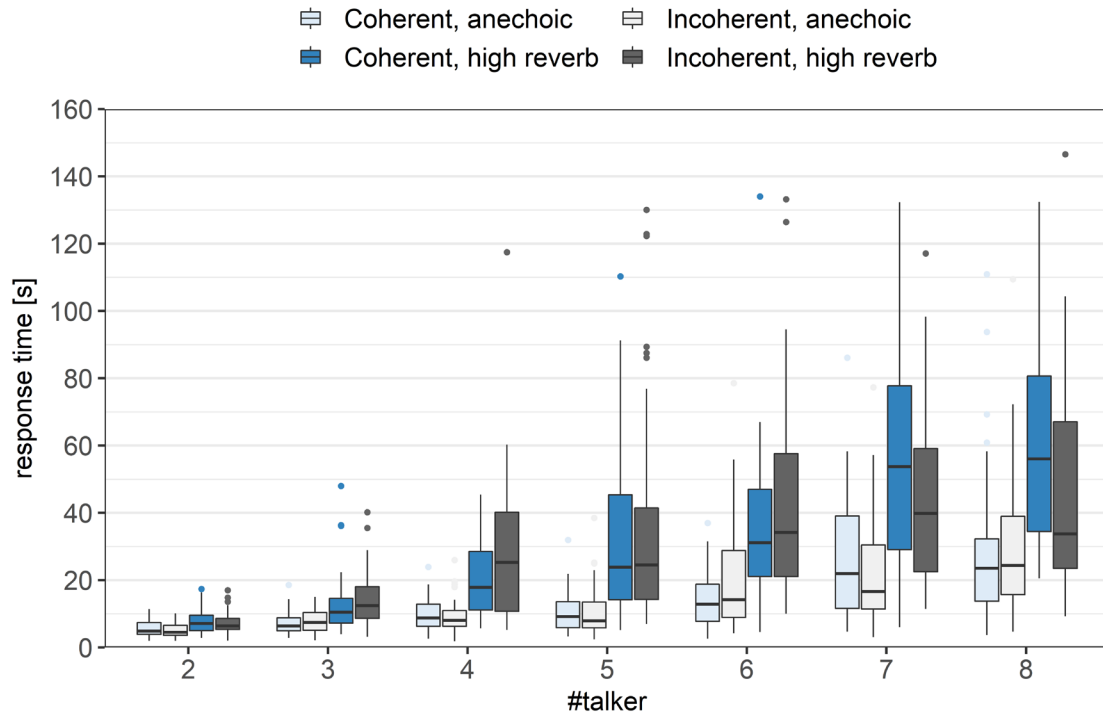conditions. No differences arise from the audio-visual incongruency.

Figure 7: The percentage of correct response locations comparing the coherent and incoherent audio-visual conditions. Each bar contains 42 datapoints across subjects and repetitions. The three colors indicate the room conditions.

Figure 8 shows the response times for the incongruent audio-visual conditions (grey boxes), i.e., the conditions with anechoic acoustic stimuli and the visuals of the reverberant room (light grey) and with high acoustic reverberation and the visuals of the anechoic room (dark grey). Additionally, the response times from the coherent anechoic and reverberant conditions are shown (blue boxes, as in Figure 5). No significant difference was found between the congruent and the incongruent condition [$p > 0.12$].

Figure 8: Response time with respect to the number of talkers in a scene. The light-blue and light-grey boxes indicate the anechoic room acoustic condition with coherent and incoherent visual information, respectively. The dark-blue and dark-grey boxes indicate the high reverberant room acoustic condition. The boxes cover the range between the 25th and the 75th percentile. The horizontal line in the boxes indicates the median. The whiskers extend to 1.5 times the inter-quartile range. Outliers are indicated as dots.

296     IV.    Discussion

297

298     In the current study we investigated the ability of normal-hearing listeners to identify and locate a

299     story in the presence of other stories. The task of the listeners was to locate a target story in the

300     presence of a varying number of simultaneous interfering talkers. Furthermore, the effect of audio-

301     visual room information was investigated, by testing different audio-visually coherent and

302     incoherent reverberant environments. The data showed that the localization accuracy and the

303     response time are affected by the number of simultaneous talkers as well as by reverberation. With

304     an increase of number of interfering talkers and an increase of reverberation time the performance

305     of the listeners decreased. Presenting incoherent audio-visual room information did not affect the

306     outcome measures.

307     A. Effect of number of talkers

308     Several factors are likely to affect the increase in response time with increasing number of talkers.

309     In the present study the speech level of each talker was kept constant independent of the number

310     of talkers, and therefore, the signal-to-noise ratio (SNR) decreases. Thus, the intelligibility is

311     expected to drop with the number of simultaneous talkers. However, the effective SNR is

312     constantly changing with head-motion and fluctuations in the signals (Grange & Culling, 2016).

313     The head-motion introduces a variation of the target and interferer angles relative to the head and

314     thus head-shadow and interaural time differences vary. Both head-shadow and interaural time

315     differences have been shown to be utilized to separate target and interfering speech sources

316     (Bronkhorst, 2000; Culling et al., 2004). Fluctuations in the speech signals allow for dip-listening

317     which can significantly improve the SNR in some time-frequency bins. Such glimpses can help to

318     better understand speech (Glyde et al., 2013; Miller & Licklider, 1950). When many speech

319     sources are presented, such glimpses are usually reduced (Cooke, 2006; Freyman et al., 2004).

320  Another effect that likely influences the response time is the amount of informational masking,

321  i.e., confusions between the target and the interferers (Carhart et al., 1969; Durlach et al., 2003;

322  Kidd et al., 2008; Watson, 2005). Previous studies have argued that the amount of informational

323  masking decreases with increasing number of simultaneous talkers (Carhart et al., 1975; Freyman

324  et al., 2004; S. A. Simpson & Cooke, 2005). However, in the current study the target speaker needs

325  to be identified by understanding the speech and to do so, listeners also need to understand the

326  content of the interferers. Thus, the listener needs to employ a strategy to search through the

327  auditory scene and while performing the search an interfering talker becomes a temporary target

328  talker. Therefore, the definition of informational masking that was already controversial in classic

329  speech perception tasks (Durlach et al., 2003; Kidd et al., 2008; Watson, 2005) becomes even more

330  complex. How the listeners perform this task and which search strategies they employ, remains an

331  open question and is out of the scope of the current study.

332  B. Effect of Reverberation

333  Reverberation was found to affect the response time only between the mid-reverberation and the

334  high-reverberation conditions, and when there were four or more talker in a scene. In literature, it

335  is reported that reverberation affects speech intelligibility more with few interfering talkers

336  because potential speech gaps and pauses get 'filled' with the reverberant energy (Bolt &

337  MacDonald, 1949; Xia et al., 2018). Such gaps generally do not exist with many overlapping

338  speech sources (Cooke, 2006; Freyman et al., 2004). A potential explanation for the disagreement

339  is that the task remains fairly easy with additional reverberation when few talkers are in a scene

340  and thus, the effect of reverberation is masked.

341  No difference in response time was observed between the anechoic and the mid-reverberant

342  conditions. The inexistent difference between the anechoic and the mid-reverberant condition

343  contradicts results from previous studies where differences in speech perception between mildly

344    reverberant conditions and anechoic conditions were found (Ahrens, Marschall, et al., 2019;

345    Duquesnoy & Plomp, 1980; Plomp, 1976). The reason for this discrepancy could be that the test

346    paradigm might not be as sensitive to capture small differences in reverberation time, as traditional

347    speech tests. However, (Kopčo et al., 2010) discussed a similar finding that mild reverberation

348    does not affect the speech localization in background speech by comparing their study with data

349    from (B. D. Simpson et al., 2006). This raises the question if there is an effect of mild reverberation

350    on speech intelligibility in everyday situations or if this effect can only be observed in artificial

351    listening scenarios in the laboratory.

352

### 353    C. Experimental paradigm

354    The spatial scene analysis method employed in this study was similar to (Weller et al., 2016). The

355    most significant difference between the approaches is that in the current study the target speech

356    stimulus needed to be understood while the task in (Weller et al., 2016) was to judge the gender

357    of all talkers presented in a scene. Consequently, they used the total number of perceived talkers

358    as their main outcome measure, while we used the response time. Furthermore, in their study the

359    participants needed to translate the spatial percept from an egocentric auditory perception onto a

360    top-down view interface. This translation was not needed in the current study as virtual reality was

361    employed as a user interface.

362    While the use of virtual reality can allow for a more user-friendly interface, virtual reality could

363    also introduce issues to an experiment. For example, the auditory percept might be affected by the

364    physical presence of the headset which has been shown to be negligible for setups with far spaced

365    sources (Ahrens, Lund, et al., 2019; Gupta et al., 2018). Furthermore, virtual reality glasses might

366    alter the participant's behavior due to their physical appearance but also because the visual world

367   is not an exact copy of the real world. However, the influence is likely negligible in this

368   experimental setup.

369   Contrary to classical speech perception studies where a %-correct or a reception threshold is

370   determined, in the present study the response time was used as the main outcome measure.

371   (Drullman & Bronkhorst, 2000) used a similar speech localization/identification paradigm with

372   sentences and words instead of ongoing speech. They showed that the trend of change in

373   intelligibility with increasing number of talkers was similar to the trend of the response times, i.e.,

374   with more interfering talkers the intelligibility decreases, and the response time increases. While

375   the material and the task were not fully comparable between these studies, one can expect a

376   correlation between speech intelligibility and response time.

377

378   ## D. Effect of incoherent AV
379   Visual information is known to affect speech perception (McGurk & MacDonald, 1976). However,

380   the effect of visual room information on auditory perception remains unclear. Previous studies

381   showed that visual information of the room can improve auditory distance perception (Calcagno

382   et al., 2012) and incongruent audio-visual cues can disrupt distance or externalization percepts

383   (Gil-Carvajal et al., 2016). However, visual information has been shown to not affect the percept

384   of reverberation (Schutte et al., 2019), which is in line with the results from the current study.

385

386   ## E.  Limitations
387   The speech material (10 stories spoken by 10 talkers) was recorded specifically for this study with

388   the aim to have distinctly different content that can be visualized with an icon. Furthermore, we

389   aimed for natural speech as opposed to highly controlled recordings with professional speakers.

390   This approach also comes with disadvantage; for example some stories or talkers might be easier

391 to understand than others. However, as stories and talkers were chosen randomly, their influence

392 is likely to be little over the sufficiently large number of iterations.

393 One aim of this study was to develop a test paradigm that is more like real-life listening than most

394 current speech intelligibility tests. While the task of understanding and locating a speech stream

395 out of interfering speech is more similar to traditional speech tests, it is by no means a replications

396 of a realistic cocktail-party situation. Firstly, all talkers are located at the same distance and with

397 the same speech level and face the listener. This decision was made to not give any level,

398 directional or direct-to-reverberant energy cues other than the information from the room

399 reflections and the talkers themselves. Secondly, the visual avatars are highly conceptualized

400 human bodies. Technology does not yet allow to visualize highly realistic human avatars with

401 conventional computational power and effort. When using avatars that share similarities with real

402 humans but evidently are not, viewers might get distracted (compare uncanny valley, (Diel et al.,

403 2022)). Thirdly, lip-movements have not been included in this study. This choice was made

404 because lip-movement simulations are not, as to the knowledge of the authors, evaluated for

405 hearing research purposes. Additionally, the aim of the avatars was more to be a 'response-box'

406 than an actual simulation of a human talker.

407

408 V. Conclusions

409 In the present study we investigated the ability of listeners to analyze a spatial scene with multiple

410 talkers. A varying number of simultaneously spoken stories was presented in different reverberant

411 environments and listeners were asked to locate a target story. Results showed that the number of

412 simultaneous talkers affected the correct identification as well as the response time. Reverberation

413    only affected the outcome measures when the reverberation time was high but not with moderate

414    reverberation.

415

416    Acknowledgement

420

421    References

422

423    Ahrens, A., Lund, K. D. K. D., Marschall, M., & Dau, T. (2019). Sound source localization with
424        varying amount of visual information in virtual reality. *PLOS ONE*, *14*(3), e0214603.
425        https://doi.org/10.1371/journal.pone.0214603

426    Ahrens, A., Marschall, M., & Dau, T. (2019). Measuring and modeling speech intelligibility in
427        real and loudspeaker-based virtual sound environments. *Hearing Research*, *377*, 307–317.
428        https://doi.org/10.1016/j.heares.2019.02.003

429    Arweiler, I., & Buchholz, J. M. (2011). The influence of spectral characteristics of early reflections
430        on speech intelligibility. *The Journal of the Acoustical Society of America*, *130*(2), 996–1005.
431        https://doi.org/10.1121/1.3609258

432    Arweiler, I., Buchholz, J. M., & Dau, T. (2013). The influence of masker type on early reflection
433        processing and speech intelligibility (L). *The Journal of the Acoustical Society of America*,
434        *133*(1), 13–16. https://doi.org/10.1121/1.4770249

435    Berzborn, M., Bomhardt, R., Klein, J., Richter, J. G., & Vorländer, M. (2017). The ITA-Toolbox :
436        An Open Source MATLAB Toolbox for Acoustic Measurements and Signal Processing.
437        *Fortschritte Der Akustik*, 222–225. http://www.ita-toolbox.org/publications/ITA-
438        Toolbox_paper2017.pdf

439    Best, V., Keidser, G., Buchholz, J. M., & Freeston, K. (2015). An examination of speech reception
440        thresholds measured in a simulated reverberant cafeteria environment. *International Journal
441        of Audiology*, *54*(10), 682–690. https://doi.org/10.3109/14992027.2015.1028656

442    Bolt, R. H., & MacDonald, A. D. (1949). Theory of Speech Masking by Reverberation. *The
443        Journal of the Acoustical Society of America*, *21*(6), 577–580.
444        https://doi.org/10.1121/1.1906551

445   Bronkhorst, A. W. (2000). The Cocktail Party Phenomenon: A Review of Research on Speech
446       Intelligibility in Multiple-Talker Conditions. *Acta Acustica United with Acustica*, *86*(1), 117–
447       128.

448   Bronkhorst, A. W., & Plomp, R. (1990). A Clinical Test for the Assessment of Binaural Speech
449       Perception in Noise. *International Journal of Audiology*, *29*(5), 275–285.
450       https://doi.org/10.3109/00206099009072858

451   Buchholz, J. M., & Best, V. (2020). Speech detection and localization in a reverberant multitalker
452       environment by normal-hearing and hearing-impaired listeners. *The Journal of the Acoustical
453       Society of America*, *147*(3), 1469–1477. https://doi.org/10.1121/10.0000844

454   Calcagno, E. R., Abregú, E. L., Eguía, M. C., & Vergara, R. (2012). The role of vision in auditory
455       distance perception. *Perception*, *41*(2), 175–192. https://doi.org/10.1068/p7153

456   Carhart, R., Johnson, C., & Goodman, J. (1975). Perceptual masking of spondees by combinations
457       of talkers. *The Journal of the Acoustical Society of America*, *58*(S1), S35–S35.
458       https://doi.org/10.1121/1.2002082

459   Carhart, R., Tillman, T. W., & Greetis, E. S. (1969). Perceptual Masking in Multiple Sound
460       Backgrounds. *The Journal of the Acoustical Society of America*, *45*(3), 694–703.
461       https://doi.org/10.1121/1.1911445

462   Cherry, E. C. (1953). Some Experiments on the Recognition of Speech, with One and with Two
463       Ears. *The Journal of the Acoustical Society of America*, *25*(5), 975–979.
464       https://doi.org/10.1121/1.1907229

465   Cooke, M. (2006). A glimpsing model of speech perception in noise. *The Journal of the Acoustical
466       Society of America*, *119*(3), 1562–1573. https://doi.org/10.1121/1.2166600

467   Culling, J. F., Hawley, M. L., & Litovsky, R. Y. (2004). The role of head-induced interaural time
468       and level differences in the speech reception threshold for multiple interfering sound sources.
469       *The Journal of the Acoustical Society of America*, *116*(2), 1057–1065.
470       https://doi.org/10.1121/1.1772396

471   Diel, A., Weigelt, S., & Macdorman, K. F. (2022). A meta-analysis of the uncanny valley's
472       independent and dependent variables. *ACM Transactions on Human–Robot Interaction*,
473       *11*(1). https://doi.org/10.1145/3470742

474   Drullman, R., & Bronkhorst, A. W. (2000). Multichannel speech intelligibility and talker
475       recognition using monaural, binaural, and three-dimensional auditory presentation. *The
476       Journal of the Acoustical Society of America*, *107*(4), 2224–2235.
477       https://doi.org/10.1121/1.428503

478   Duquesnoy, A. J., & Plomp, R. (1980). Effect of reverberation and noise on the intelligibility of
479       sentences in cases of presbyacusis. *The Journal of the Acoustical Society of America*, *68*(2),
480       537–544. https://doi.org/10.1121/1.384767

481    Durlach, N. I., Mason, C. R., Kidd, G., Arbogast, T. L., Colburn, H. S., & Shinn-Cunningham, B.
482        G. (2003). Note on informational masking (L). *The Journal of the Acoustical Society of*
483        *America*, *113*(6), 2984. https://doi.org/10.1121/1.1570435

484    Favrot, S., & Buchholz, J. M. (2010). LoRA: A loudspeaker-based room auralization system. *Acta*
485        *Acustica United with Acustica*, *96*(2), 364–375. https://doi.org/10.3813/AAA.918285

486    Freyman, R. L., Balakrishnan, U., & Helfer, K. S. (2004). Effect of number of masking talkers and
487        auditory priming on informational masking in speech recognition. *The Journal of the*
488        *Acoustical Society of America*, *115*(5), 2246–2256. https://doi.org/10.1121/1.1689343

489    Gil-Carvajal, J. C., Cubick, J., Santurette, S., & Dau, T. (2016). Spatial Hearing with Incongruent
490        Visual or Auditory Room Cues. *Scientific Reports*, *6*. https://doi.org/10.1038/srep37342

491    Glyde, H., Buchholz, J., Dillon, H., Best, V., Hickson, L., & Cameron, S. (2013). The effect of
492        better-ear glimpsing on spatial release from masking. *The Journal of the Acoustical Society*
493        *of America*, *134*(4), 2937–2945. https://doi.org/10.1121/1.4817930

494    Grange, J. A., & Culling, J. F. (2016). The benefit of head orientation to speech intelligibility in
495        noise. *The Journal of the Acoustical Society of America*, *139*(2), 703–712.
496        https://doi.org/10.1121/1.4941655

497    Gupta, R., Ranjan, R., He, J., & Gan, W.-S. (2018). Investigation of effect of VR/AR headgear on
498        Head related transfer functions for natural listening. *AES International Conference on Audio*
499        *for Virtual and Augmented Reality*. http://www.aes.org/e-lib/browse.cfm?elib=19697

500    Hawley, M. L., Litovsky, R. Y., & Colburn, H. S. (1999). Speech intelligibility and localization in
501        a multi-source environment. *The Journal of the Acoustical Society of America*, *105*(6), 3436–
502        3448. https://doi.org/10.1121/1.424670

503    Kidd, G., Mason, C. R., Richards, V. M., Gallun, F. J., & Durlach, N. I. (2008). Informational
504        Masking. In W. A. Yost, A. N. Popper, & R. R. Fay (Eds.), *Auditory Perception of Sound*
505        *Sources. Springer Handbook of Auditory Research* (Vol. 29, pp. 143–189).
506        https://doi.org/10.1007/978-0-387-71305-2_6

507    Kopčo, N., Best, V., & Carlile, S. (2010). Speech localization in a multitalker mixture. *The Journal*
508        *of the Acoustical Society of America*, *127*(3), 1450–1457. https://doi.org/10.1121/1.3290996

509    Kuznetsova, A., Brockhoff, P. B., & Christensen, R. H. B. (2017). lmerTest Package: Tests in
510        Linear Mixed Effects Models . *Journal of Statistical Software*, *82*(13).
511        https://doi.org/10.18637/jss.v082.i13

512    Lenth, R. (2020). *emmeans: Estimated Marginal Means, aka Least-Squares Means*. https://cran.r-
513        project.org/package=emmeans

514    Lund, K. D., Ahrens, A., & Dau, T. (2019). A method for evaluating audio-visual scene analysis
515        in multi-talker environments. *Proceedings of the International Symposium on Auditory and*
516        *Audiological Research, Vol. 7: Auditory Learning in Biological and Artificial Systems*.

517    McGurk, H., & MacDonald, J. (1976). Hearing lips and seeing voices. *Nature*, *264*(5588), 746–
518        748. https://doi.org/10.1038/264746a0

519    Miller, G. A., & Licklider, J. C. R. (1950). The Intelligibility of Interrupted Speech. *Journal of the*
520        *Acoustical Society of America*, *22*(2), 167–173. https://doi.org/10.1121/1.1906584

521    Moncur, J. P., & Dirks, D. (1967). Binaural and Monaural Speech Intelligibility in Reverberation.
522        *Journal   of   Speech   Language   and   Hearing   Research*,   *10*(2),   186.
523        https://doi.org/10.1044/jshr.1002.186

524    Nabelek, A. K., & Mason, D. (1981). Effect of Noise and Reverberation on Binaural and Monaural
525        Word Identification by Subjects with Various Audiograms. *Journal of Speech, Language,*
526        *and Hearing Research*, *24*(3), 375–383. https://doi.org/10.1044/jshr.2403.375

527    Nábělek, A. K., & Pickett, J. M. (1974). Reception of consonants in a classroom as affected by
528        monaural and binaural listening, noise, reverberation, and hearing aids. *The Journal of the*
529        *Acoustical Society of America*, *56*(2), 628–639. https://doi.org/10.1121/1.1903301

530    Plomp, R. (1976). Binaural and Monaural Speech Intelligibility of Connected Discourse in
531        Reverberation as a Function of Azimuth of a Single Competing Sound Source (Speech or
532        Noise).   *Acta   Acustica   United   with   Acustica*,   *34*(4),   200–211.
533        http://www.ingentaconnect.com/content/dav/aaua/1976/00000034/00000004/art00004

534    R Core Team. (2020). *R: A Language and Environment for Statistical Computing*. https://www.r-
535        project.org/

536    Schutte, M., Ewert, S. D., & Wiegrebe, L. (2019). The percept of reverberation is not affected by
537        visual room impression in virtual environments. *The Journal of the Acoustical Society of*
538        *America*, *145*(3). https://doi.org/10.1121/1.5093642

539    Simpson, B. D., Brungart, D. S., Iyer, N., Gilkey, R. H., & Hamil, J. T. (2006). DETECTION
540        AND LOCALIZATION OF SPEECH IN THE PRESENCE OF COMPETING SPEECH
541        SIGNALS. *Proceedings of the 12th International Conference on Auditory Display*.

542    Simpson, S. A., & Cooke, M. (2005). Consonant identification in N-talker babble is a
543        nonmonotonic function of N. *The Journal of the Acoustical Society of America*, *118*(5), 2775–
544        2778. https://doi.org/10.1121/1.2062650

545    Warzybok, A., Rennies, J., Brand, T., Doclo, S., & Kollmeier, B. (2013). Effects of spatial and
546        temporal integration of a single early reflection on speech intelligibility. *The Journal of the*
547        *Acoustical Society of America*, *133*(1), 269–282. https://doi.org/10.1121/1.4768880

548    Watson, C. S. (2005). Some Comments on Informational Masking. *Acta Acustica United with*
549        *Acustica*, *91*(2005), 502–512.

550    Weller, T., Best, V., Buchholz, J. M., & Young, T. (2016). A Method for Assessing Auditory
551        Spatial Analysis in Reverberant Multitalker Environments. *Journal of the American Academy*
552        *of Audiology*, *27*(7), 601–611. https://doi.org/10.3766/jaaa.15109

553    Xia, J., Xu, B., Pentony, S., Xu, J., & Swaminathan, J. (2018). Effects of reverberation and noise

554        on speech intelligibility in normal-hearing and aided hearing-impaired listeners. *The Journal*

555        *of the Acoustical Society of America*, *143*(3), 1523–1533. https://doi.org/10.1121/1.5026788

556