

1
2
3 1
4
5
6 2

**Impact of SNR, masker type and noise reduction processing on sentence recognition performance
and listening effort as indicated by the pupil dilation response**

7
8 3
9
10
11 4
12

Barbara Ohlenforst^{1,4}, Dorothea Wendt^{4,5}, Sophia E. Kramer¹, Graham Naylor⁶, Adriana A. Zekveld^{1,2,3},
Thomas Lunner^{2,3,4,5}

13
14 5
15 6
16 7
17

¹Section Ear & Hearing, Dept. of Otolaryngology-Head and Neck Surgery, VU University Medical
Center and Amsterdam Public Health Research Institute, Amsterdam, The Netherlands;

18 8
19 9
20

²Department of Behavioral Sciences and Learning, Linköping University, Sweden;

21 10
22

³Linnaeus Centre HEAD, The Swedish Institute for Disability Research, Linköping and Örebro
Universities, Sweden;

23 11
24 12
25

⁴Eriksholm Research Center, Oticon A/S, Denmark;

26 13
27

⁵Department of Electrical Engineering, Technical University of Denmark, Denmark;

28 14
29

⁶MRC/CSO Institute of Hearing Research, Scottish Section, Glasgow, United Kingdom. Part of the
University of Nottingham.

30 15
31 16
32
33
34 17
35

18 Received, 2017

36 18
37
38

19 This work was supported by grants from the European Commission (FP7-LISTEN607373) and the Oticon
20 Foundation.

39 19
40 20
41
42

21 Corresponding author: Barbara Ohlenforst,

43 21
44

22 Section Ear & Hearing, Dept. of Otolaryngology-Head and Neck Surgery, VU University Medical Center and
23 Amsterdam Public Health Research Institute, Amsterdam, The Netherlands; P.O. Box 7057, 1007 MB, Phone:
24 +31 20 4440952; Fax: +3120 444 2033

45 22
46 23
47 24
48

25 Eriksholm Research Centre, Rørtangvej 20, 3070 Snekkersten, Denmark;

49 25
50

26 E-mail: b.ohlenforst@vumc.nl; Phone: 0045-48298900 Fax: 0045-49223629

51 26
52
53
54
55
56
57
58
59

60
61
62 **Abstract**
63
64

65 28 Recent studies have shown that activating the noise reduction scheme in hearing aids results in a
66
67 29 smaller peak pupil dilation (PPD), indicating reduced listening effort, and 50% and 95% correct
68
69 30 sentence recognition with a 4-talker masker. The objective of this study was to measure the effect of
70
71 31 the noise reduction scheme (on or off) on PPD and sentence recognition across a wide range of
72
73 32 signal-to-noise ratios (SNRs) from +16 dB to -12 dB and two masker types (4-talker and stationary
74
75 33 noise). Relatively low PPDs were observed at very low (-12 dB) and very high (+16 dB to +8 dB) SNRs
76
77 34 presumably due to 'giving up' and 'easy listening', respectively. The maximum PPD was observed
78
79 35 with SNRs at approximately 50% correct sentence recognition. Sentence recognition with both
80
81 36 masker types was significantly improved by the noise reduction scheme, which corresponds to the
82
83 37 shift in performance from SNR function at approximately 5 dB toward a lower SNR. This intelligibility
84
85 38 effect was accompanied by a corresponding effect on the PPD, shifting the peak by approximately 4
86
87 39 dB toward a lower SNR. In addition, with the 4-talker masker, when the noise reduction scheme was
88
89 40 active, the PPD was smaller overall than that when the scheme was inactive. We conclude that with
90
91 41 the 4-talker masker, noise reduction scheme processing provides a listening effort benefit in addition
92
93 42 to any effect associated with improved intelligibility. Thus, the effect of the noise reduction scheme
94
95 43 on listening effort incorporates more than can be explained by intelligibility alone, emphasizing the
96
97 44 potential importance of measuring listening effort in addition to traditional speech reception
98
99 45 measures.
100

101
102
103 46

104
105
106 47

107
108
109 48 **Keywords:** Hearing impairment, speech recognition, noise reduction scheme, hearing aids, pupil
110
111 49 dilation, listening effort, signal-to-noise ratio
112
113
114
115
116
117
118

1. Introduction

Audiological evaluations and research studies investigating hearing aid signal processing have typically focused on changes or benefits in intelligibility but often failed to provide a complete picture of the processes involved in speech recognition (Dillon et al., 1993; Ricketts et al., 2001; Sarampalis et al., 2009). Traditional speech reception measures have been shown to be insensitive to the possible benefits of hearing aid algorithms due to ceiling effects or great variability (Gatehouse et al., 1990). Baer and colleagues (Baer et al., 1993) suggested that the greatest benefit of noise reduction processing in hearing aids may be reduced listening effort rather than enhanced speech intelligibility.

According to the Framework for Understanding Effortful Listening (FUEL) (Pichora-Fuller et al., 2016), listening effort depends on a range of factors, including not only individual factors, such as hearing ability and motivation to continue listening, but also external factors, such as the task demands imposed by the listening situation (Brehm, 1999). Participants may invest less effort in their task performance when the task demands are too high or allocate less cognitive resources under very easy listening conditions (Ohlenforst et al., 2017a). Recently, an increasing number of studies have sought additional methods to gain information about effortful listening as a supplement to traditional audiological measures to assess individual hearing ability (McGarrigle et al., 2014; Ohlenforst et al., 2017b; Pals et al., 2013; Wu et al., 2016). These methods include subjective assessments, such as self-reports and questionnaires (McAuliffe et al., 2012; Panico et al., 2009; Picou et al., 2011); behavioral measures, such as dual-task paradigms or reaction time measures (Fraser et al., 2010; Houben et al., 2013; Tun et al., 2009); and physiological measures, such as the pupil response and functional magnetic resonance imaging (fMRI) or EEG measures (Kuchinsky et al., 2013; Obleser et al., 2012; Petersen et al., 2015). Importantly, the listening conditions may affect listening effort even when speech intelligibility is not affected, such as when speech intelligibility is at a ceiling and hence constitutes an insensitive outcome measure (Koelewijn et al., 2014; Wendt et al.,

178
179
180 75 2017). For example, Wendt et al., (2017) showed that activating the noise reduction scheme at
181
182 76 ceiling performance reduced listening effort, but speech in noise performance was unaffected.
183
184 77 Therefore, simultaneously assessing listening effort and speech performance may uncover challenges
185
186 78 or changes in processing speech that may not be evident with traditional measures.
187
188
189 79 Numerous studies across different research areas have shown that pupil dilation increases as the
190
191 80 processing load imposed by the task demands increases (Beatty, 1982; Engelhardt et al., 2010;
192
193 81 Granholm et al., 1996; Kahneman, 1973; Van Der Meer et al., 2010). Pupillometry has repeatedly
194
195 82 been verified as a valid measure for quantifying the effort required for speech recognition with
196
197 83 background noise (Koelewijn et al., 2012; Koelewijn et al., 2014; Kramer et al., 1997; Ohlenforst et
198
199 84 al., 2017a; Ohlenforst et al., 2017b; Wendt et al., 2017; Zekveld et al., 2011). For instance, the SNR
200
201 85 (ranging from -20 dB to +16 dB) and masker type (stationary and 1-talker masker) have been shown
202
203 86 to affect pupil dilation during listening (Ohlenforst et al., 2017a). Recent studies indicate that effort is
204
205 87 not necessarily monotonically related to the task demands. The changes in effort follow an inverse U-
206
207 88 shaped function, indicating that listeners may exert less effort due to 'giving up' under very difficult
208
209 89 conditions and 'taking it easy' when listening at high SNRs (Ohlenforst et al., 2017a; Wu et al., 2016;
210
211 90 Zekveld et al., 2014). Ohlenforst et al. (Ohlenforst et al., 2017a) investigated the peak pupil dilation
212
213 91 (PPD) across a range of SNRs in hearing-impaired and normal-hearing listeners. These authors
214
215 92 showed that the PPD, which is an indication of the cognitive processing load, was affected by an
216
217 93 interaction between the masker type and hearing status of the individual. In the presence of a
218
219 94 stationary noise masker, the hearing-impaired listeners showed relatively large PPDs across a wide
220
221 95 range of SNRs, while the normal-hearing listeners showed a maximum PPD across a relatively narrow
222
223 96 range of low (challenging) SNRs (Ohlenforst et al., 2017a). With a single-talker masker, the maximum
224
225 97 PPD was in the mid-range of SNRs, while relatively smaller PPDs were observed at low and high SNRs
226
227 98 in both groups of listeners. Interestingly, recent findings across a variety of studies in the field of
228
229 99 listening effort suggest that the allocation of mental resources needed during listening to reach
230
231
232
233
234
235
236

237
238
239 100 speech understanding in daily life listening situations may differ between normal-hearing and
240
241 101 hearing-impaired listeners (Ohlenforst et al., 2017a; Ohlenforst et al., 2017b; Zekveld et al., 2011).
242
243
244 102 Hearing aids are designed to improve the audibility of sounds and facilitate the intelligibility of
245
246 103 speech in both quiet and noisy environments. These improvements may be accompanied by reduced
247
248 104 listening effort. The advanced signal processing in hearing aids includes a digital noise reduction
249
250 105 scheme, which aims to reduce the level of interfering background noise by improving the SNR.
251
252 106 Recent studies indicate that the noise reduction scheme improves the recall of words presented in a
253
254 107 competing multi-talker background (Lunner et al., 2016; Ng et al., 2015; Ng et al., 2013). The
255
256 108 researchers concluded that the noise reduction scheme may reduce the adverse effect of noise on
257
258 109 memory and thereby facilitate the segregation of the target from the multi-talker masker signal. This
259
260
261 110 enhanced memory of the target words was interpreted to represent reduced listening effort (Lunner
262
263 111 et al., 2016; Ng et al., 2015; Ng et al., 2013). Moreover, Wendt et al. (2017) presented speech in a 4-
264
265 112 talker babble masker at two SNRs (SNR50 and SNR95) corresponding to the individual 50% or 95%
266
267 113 sentence recognition level. These authors assessed the effect of the noise reduction scheme by
268
269 114 applying a combination of a digital noise reduction scheme and directional microphones. When the
270
271 115 scheme was activated in the hearing aid, the speech recognition performance at SNR50 was
272
273 116 significantly improved and accompanied by significantly smaller PPDs. Interestingly, activating the
274
275 117 noise reduction scheme did not affect the near-ceiling speech recognition performance at SNR95.
276
277 118 Nevertheless, significantly smaller PPDs were observed, indicating that the noise reduction scheme
278
279 119 had a beneficial effect on listening effort. Thus, measuring listening effort by assessing PPD could
280
281 120 provide a sensitive outcome measure of hearing aid benefit even at high performance level
282
283 121 traditional methods of audiological assessment are not sufficiently sensitive.
284
285
286 122 The studies described above (Ng et al., 2015; Ng et al., 2013; Wendt et al., 2017) indicate that effort
287
288 123 can be reduced with modern hearing aid signal processing. However, knowledge regarding the
289
290 124 benefit of noise reduction processing on listening effort remains very limited as only a few listening
291
292
293
294
295

296
297
298 125 conditions were tested in these studies. In contrast, the effect of noise reduction processing on
299
300 126 intelligibility has been studied by several groups of researchers. In these studies, the inconsistency in
301
302 127 the diverse noise reduction processing schemes studied renders generalization problematic,
303
304 128 especially as processing schemes become increasingly sophisticated over time. Some research
305
306 129 studies have indicated that the application of noise reduction processing may not always be
307
308 130 beneficial for speech intelligibility (Bentler et al., 2008; Nordrum et al., 2006). Such negative effects
309
310 131 suggest that while the background noise may be removed, the target speech might also be degraded.
311
312 132 Stronger or more aggressive signal processing may cause more signal enhancement but could
313
314 133 simultaneously introduce more degradation (Loizou et al., 2011). For example, in a recent study, the
315
316 134 effect of noise reduction processing on sentence recognition was tested in the presence of a
317
318 135 cafeteria background masker (Neher et al., 2013). Simulated hearing aid processing including
319
320 136 coherence-based noise reduction was presented via headphones to hearing-impaired listeners. The
321
322 137 algorithm was designed to suppress the reverberant signal components and diffuse the background
323
324 138 noise at mid to high frequencies but did not include directionality. The results showed that sentence
325
326 139 recognition was unaffected by the moderate noise reduction processing, but the strong noise
327
328 140 reduction processing reduced speech recognition by approximately 5%. The effect was replicated in a
329
330 141 follow up study in which the same acoustic test conditions were used in a group of habitual hearing
331
332 142 aid users (Neher, 2014). Compared to the moderate or no noise reduction processing, the strong
333
334 143 noise reduction processing reduced speech recognition at -4 dB and 0 dB SNR.
335
336
337

338 144 How hearing-impaired listeners invest listening effort across a broader range of listening situations
339
340 145 and how effortful listening relates to performance measures remain unclear. The current study
341
342 146 aimed to examine how a noise reduction scheme influences sentence recognition and listening
343
344 147 effort. The applied noise reduction scheme preserves speech and reduces noise in complex
345
346 148 environments by a fast-acting combination of a beam-former (Kjems et al., 2012) and a single-
347
348 149 channel Wiener post-filter (Jensen et al., 2015) to attenuate interfering sounds. Any effect of the
349
350
351
352
353
354

355
356
357 150 noise reduction processing on intelligibility likely affects the PPD in a corresponding direction as the
358
359 151 intelligibility of speech has a strong and reliable effect on the PPD (Koelewijn et al., 2014; Ohlenforst
360
361 152 et al., 2017a; Zekveld et al., 2014). However, in addition to this intelligibility effect, the noise
362
363 153 reduction processing may have additional effects on the PPD, as suggested by recent studies
364
365 154 investigating listening effort that demonstrated that hearing aid processing has a beneficial effect on
366
367 155 listening effort due to reduced background noise and reduced cognitive effort during speech
368
369 156 processing (Picou et al., 2013; Sarampalis et al., 2009; Wendt et al., 2017). Demonstrating the effect
370
371 157 of noise reduction processing on listening effort combined with simultaneous knowledge regarding
372
373 158 speech in noise performance could further substantiate the value of measuring effort as an extra
374
375 159 dimension in addition to traditional speech reception measures.

376
377
378 160 Recent research found better SRTs in speech recognition in the presence of a single-talker masker
379
380 161 than those in the presence of a stationary noise masker (Koelewijn et al., 2012). The envelope
381
382 162 modulations of the multi-talker masker **might** allow the participants to listen in the energy dips in the
383
384 163 spectral-temporal domain and glimpse parts of the target sentence (Festen et al., 1990; Francart et
385
386 164 al., 2011; Koelewijn et al., 2012; Koelewijn et al., 2014; Rosen et al., 2013). Based on **the**
387
388 165 **characteristics of the masker types and** recent findings, we hypothesize that speech recognition
389
390 166 performance **is better with** the 4-talker masker than that with the stationary noise masker (Koelewijn
391
392 167 et al., 2012; Koelewijn et al., 2014). However, recent studies suggest that the intelligibility of speech
393
394 168 masked by additional interfering speech information may require more mental effort than that with
395
396 169 an energetic mask (Larsby et al., 2008). Informational masking, including lexical interference or the
397
398 170 competition for neural resources, may cause higher listening effort (Beatty, 1982; Koelewijn et al.,
399
400 171 2012; Koelewijn et al., 2014; Scott et al., 2004; Scott et al., 2009). We hypothesized that the better
401
402 172 speech recognition with the 4-taker masker compared to that with the stationary noise masker could
403
404 173 be accompanied by larger PPDs. We hypothesized that sentence recognition is improved and
405
406 174 listening effort is reduced with SNRs corresponding to approximately 50% correct or better
407
408
409
410
411
412
413

414
415
416 175 performance with the active noise reduction compared to the inactive noise reduction scheme. This
417
418 176 hypothesis is motivated by two arguments. First, in a previous study conducted by Wendt and
419
420 177 colleagues (2017), SRT targeting 50% correct performance was significantly improved by the active
421
422 178 noise reduction scheme compared to that with the inactive noise reduction scheme setting. Second,
423
424 179 the segregation between the target and masker signal at very low SNRs might be more difficult for
425
426 180 the algorithm, which might have an impact on the SNR improvement provided by the algorithm.
427
428
429
430
431

181

432 182 **2. Materials and methods**

433 434 183 2.1 Participants

435
436
437 184 Twenty-five experienced hearing aid users were recruited from the Eriksholm Research Centre in
438
439 185 Denmark. On average, the participants had used hearing aids for 7.7 years (SD=3.1 years). The
440
441 186 participants were aged between 46 and 77 years (mean age 64.3 years, SD=9.4) and native Danish
442
443 187 speakers. The audiometric inclusion criterion for the participants was symmetrical, with mild to
444
445 188 moderate sensorineural hearing thresholds. The average pure tone hearing thresholds ranged
446
447 189 between 35 dB and 60 dB HL (see Figure 1), and air-bone gaps less than 10 dB between 500 Hz and
448
449 190 4000 Hz were required in both ears. All the participants had normal or corrected-to-normal vision
450
451 191 and no history of neurological diseases, dyslexia or diabetes mellitus. All the participants provided
452
453 192 written informed consent, and the study was approved by the local regional ethics committee (De
454
455 193 Videnskabetiske Komiteer for Region Hovedstaden).
456
457
458
459
460
461
462
463

194

195

464 196 2.2 Auditory stimuli

465
466
467
468
469
470
471
472

473
474
475 197 Everyday Danish sentences from the Hearing in Noise Sentence Test (HINT) (Nielsen et al., 2009)
476
477 198 were presented in a spatial setup with five loudspeakers in a sound proof measurement booth as
478
479 199 shown in Figure 2. The target sentences were spoken by a male and presented from a loudspeaker
480
481 200 located at 0 degree azimuth. All the sentences contained five words, 8-9 syllables were included in
482
483 201 each sentence, and the single words did not contain more than four syllables (Nielsen et al., 2009).
484
485 202 The following is an example of a presented sentence: "Filmen er rigtig godt lavet" (translation: "the
486
487 203 movie was well made"). The sentence duration was on average 1.4 seconds. The listeners were
488
489 204 presented with a training list of 20 sentences for each masker type, followed by eight lists of 25
490
491 205 sentences for every SNR. To cover the large number of testing conditions, the sentence material was
492
493 206 re-used across four experimental visits. Recent research assessed the possible learning effects due to
494
495 207 repeated exposure to HINT sentences across three experimental visits with an interval of three
496
497 208 weeks between visits. The results showed that the memory effects of the sentence material are not
498
499 209 significant with limited exposure when the sentences were only presented once during each visit
500
501 210 (Simonsen et al., 2016). The experimental visits in the current study were separated by at least three
502
503 211 weeks, and identical sentence material was not repeated within each visit to prevent learning effects
504
505 212 of the speech material. The speech recognition performance was measured in the presence of a
506
507 213 stationary noise or a 4-talker masker background. The 4-talker masker was made from four single-
508
509 214 talker maskers, including two different male voices and two different female voices. Each separate
510
511 215 talker read a text passage from a newspaper, and one single talker was presented from one
512
513 216 loudspeaker, each positioned at +/- 90 and +/- 150 degree azimuth (Wendt et al., 2017). We
514
515 217 balanced the distribution of the talkers across loudspeakers for each SNR by switching the order of
516
517 218 the talkers. There were never two talkers of the same gender next to each other or on the opposite
518
519 219 position of each loudspeaker. In each trial, the masker started 3 seconds prior to the presentation of
520
521 220 the sentence and ended 3 seconds after the sentence offset. The participants repeated the sentence
522
523 221 aloud once the masker stopped. The same presentation procedure was applied for both masker
524
525
526
527
528
529
530
531

532
533
534 222 types. The long-term average frequency spectrum of both masker types was identical to the
535
536 223 spectrum of the target speech signal, and the masker was always presented at 70 dB SPL. The masker
537
538 224 levels were kept constant to ensure that the noise would not become too loud at low SNRs. Changing
539
540 225 the noise levels might also allow the listeners to estimate the upcoming task difficulty. The same SNR
541
542 226 range was chosen for both masker types. We included a large range of positive SNRs as previous
543
544 227 findings suggested that typical, ecologically sound environments for hearing-impaired listeners occur
545
546 228 at SNRs of approximately +5 dB or better (Festen et al., 1990; Ohlenforst et al., 2017a; Smeds et al.,
547
548 229 2015; Wu et al., 2014; Zekveld et al., 2014). Speech masked with a stationary masker and 4-talker
549
550 230 masker was presented at eight SNRs between -12 dB and +16 dB and distributed in steps of 4 dB. Per
551
552 231 the masker type, 25 sentences were presented for each SNR.
553
554
555
556 232
557
558 233
559
560
561 234 2.3 Noise reduction scheme
562
563

564 235 All the participants wore identical hearing aid models during the sentence recognition test and
565
566 236 examined in the same two different settings. In one setting, the noise reduction scheme was turned
567
568 237 off, but the hearing aid provided audibility based on each individual's hearing threshold via the Voice
569
570 238 Aligned Compression (VAC) rationale (Le Goff, 2015). The VAC amplification rationale is based on a
571
572 239 wide dynamic range compression scheme with compression knee points between 20 and 50 dB SPL
573
574 240 depending on the frequency range and the individuals' hearing thresholds. The hearing aid was set to
575
576 241 mimic the natural acoustic effect of the pinna; thus, the microphone setting was close to
577
578 242 omnidirectional, and no actual noise reduction processing was applied. The other setting involved
579
580 243 activating the noise reduction scheme. In this setting, a fast-acting combination of a minimum
581
582 244 variance distortion-less response (MVDR) beam-former (Kjems et al., 2012) and a single-channel
583
584
585
586
587
588
589
590

591
592
593 245 Wiener post-filter (Jensen et al., 2015) was applied before the VAC. In the algorithm, spatial filtering
594
595 246 and Wiener filtering were applied to attenuate interfering sounds originating behind the listener.
596
597

598 247 The output SNR method suggested by Naylor and Johannesen (2009) was used to directly measure
599
600 248 the SNR effect of the complete noise reduction scheme. The hearing aid was placed in a sound field
601
602 249 and exposed to running speech plus noise mixtures in SNRs ranging from -10 dB SNR to +20 dB in
603
604 250 steps of 5 dB for the two different noise types (speech-weighted unmodulated noise and multi-talker
605
606 251 babble noise). The output SNR method was applied to NR on and off. In the range of -10 dB SNR to
607
608 252 +10 dB SNR, the listeners experience an articulation-index (AI) weighted SNR improvement ranging
609
610 253 from 4.5 dB to 5.2 dB for NR on compared to that for NR off for the speech-weighted noise and an AI
611
612 254 weighted SNR improvement ranging from 4.2 dB to 4.8 dB for the multi-talker babble noise. For SNRs
613
614 255 above +10 dB, the SNR improvement gradually declined to a few dB because the noise estimates in
615
616 256 the noise reduction algorithm decline at high SNRs, and thus, the noise reduction algorithm becomes
617
618 257 less effective.
619
620

621 258
622
623
624 259
625
626

627 260 2.4 Pupillometry

628
629

630 261 During the experiment, the pupil location and pupil size were recorded using an eye tracking system
631
632 262 by SensoMotoric Instruments (SMI, Berlin, Germany, 2D Video-Oculography, version 4), which
633
634 263 applies infrared video tracking to measure the pupil diameter. The eye tracking system had a
635
636 264 sampling frequency of 120 Hz and a spatial resolution of 0.03 mm. The pupil location and pupil size
637
638 265 were recorded by the eye tracker and stored on a connected computer with time stamps
639
640 266 corresponding to the start of each trial, including the masker onset, the sentence onset and the
641
642 267 offset for the post-masker. The experimenter monitored the pupil recordings and applied corrective
643
644 268 actions. In the case that a participant moved his/her head or upper body or the real-time pupil
645
646
647
648
649

650
651
652 269 recordings were missing data regarding the pupil diameter, corrective actions, such as adjusting the
653
654 270 participants' position, the distance to the eye tracker, or light, were applied.
655
656

657 271
658

660 272 2.5 Procedures

661
662

663 273 In total, 17 adults from the Eriksholm pool of participants with recent pure tone audiogram data and
664
665 274 recently made ear impressions (less than 6-month-old) were required to participate in four
666
667 275 experimental visits. We recruited 8 additional participants who required an additional recruitment
668
669 276 visit (total of five visits) to measure the pure tone audiogram and take ear impressions. In total, four
670
671 277 experimental visits, including two visits per masker type, were required for each participant. The
672
673 278 visits were distributed across approximately four months during the fall of 2016 with intervals of at
674
675 279 least three weeks between each visit to avoid learning effects of the sentence material as the
676
677 280 material was repeatedly used (Simonsen et al., 2016). During the four experimental sessions, each
678
679 281 participant sat on a fixed chair in front of the eye tracking system in a sound proof booth. The
680
681 282 experimenter observed the real-time recording of the pupil response from the eye tracking system to
682
683 283 evaluate the pupil recording quality. The height of the chair and the distance to the eye tracker (55
684
685 284 cm +/- 5 cm approximately) were adjusted individually until a stable, continuous pupil response was
686
687 285 measured. The illumination in the measurement booth was fixed during the experiment to an
688
689 286 average of 84.3 lux (SD=3.56 lux). The stationary noise and 4-talker masker were presented at eight
690
691 287 identical SNRs between -12 dB and +16 dB distributed in steps of 4 dB. During each visit, only 1 of the
692
693 288 2 masker types was presented in two blocks of four randomized SNRs. In one block, the noise
694
695 289 reduction scheme was turned on, and in the other block, the noise reduction scheme was turned off.
696
697 290 During each visit, each noise reduction scheme setting (on or off) was tested at four SNR levels. We
698
699 291 balanced the SNR levels for each visit, including two difficult and two easier SNRs (e.g., -12, -4, +4 and
700
701 292 +12 dB SNR or -8, 0, +8 and +16 dB SNR). We balanced the setting of the noise reduction scheme and
702
703
704
705
706
707
708

709
710
711 293 the presented masker types across visits and blocks. Each participant's visit started with a practice
712
713 294 session in which the noise reduction scheme setting was the same as that in the starting block, and
714
715 295 20 sentences at an SNR of +4 dB were tested. The practice session ensured that the participants were
716
717 296 confident with the experimental procedures as it may not be intuitive to inhibit movements and
718
719 297 blinking during the sentence presentation. A sentence was scored as correct if all the words were
720
721 298 correctly repeated.

724 299

727 300

730 301 2.6 Pupil data selection and cleaning

732
733 302 Pupil diameter values more than 2 standard deviations from the mean pupil diameter in a given trial
734
735 303 were defined as blinks. Pupil traces with more than 25% of blinks between the start of baseline (final
736
737 304 second pre-noise before the sentence onset) and the end of the post-masker were excluded from
738
739 305 data analysis. For pupil traces with less than 25% of blinks, the blinks were interpolated linearly
740
741 306 starting with 5 samples before and 7 samples after each blink (Siegle et al., 2008). The pupil response
742
743 307 within each selected and de-blinked trace was smoothed by a 9-point moving average filter. The
744
745 308 reference of the task evoked pupil dilation was the baseline, which corresponded to the average
746
747 309 pupil diameter recorded during the final second of the three second presentation of the masker
748
749 310 before the target speech onset. The PPD was calculated as the maximum pupil dilation between the
750
751 311 onset of the sentence and the offset of the noise relative to the baseline pupil diameter for every
752
753 312 trace (one pupil trace was recorded per sentence). For each participant and each condition, all the
754
755 313 included de-blinked and smoothed traces (≤ 25) were time-aligned and averaged. For each SNR
756
757 314 condition, at least 18 valid pupil traces ($n=25$ traces in total) with less than 25% of blinks were
758
759 315 required per participant to consider the pupil data for the statistical analysis. Eighteen participants
760
761 316 had the required number of valid pupil traces for each of the 32 testing conditions. Six participants
762
763
764
765
766
767

768
769
770 317 had less than 18 valid pupil traces in at least one of the testing conditions, and two participants had
771
772 318 missing data (<18 valid pupil traces) in at 3 test conditions. We calculated the average pupil trace
773
774 319 across all the valid pupil traces per SNR condition and subject. The mean PPD was calculated based
775
776 320 on the averaged pupil trace and thus provided the data for the statistical analysis per SNR and
777
778 321 participant.
779

322

323

780 781 322 782 783 784 323 785 786 787 324 2.7 Statistical analyses

788
789
790 325 Pupil data selection and cleaning were applied to the pupil data from 24 participants (50% female).
791
792 326 One participant was excluded due to unexpected attention problems. We measured 800 pupil traces
793
794 327 during the experimental sessions (excluding the practice traces) per participant, and on average, 38
795
796 328 (SD=12.92) pupil traces were excluded per person. The corresponding sentence recognition scores
797
798 329 for all 800 measured traces were included in the statistical analysis.
799

800
801 330 We applied linear mixed models (LMM) to analyze the data as LMMs tolerate missing values, while
802
803 331 repeated measures ANOVA tests only use complete cases contrary to multilevel analyses. Moreover,
804
805 332 mixed-effects models are more flexible in processing the multilevel structure of the data (i.e., the 8
806
807 333 different SNRs and 2 different hearing aid settings). We averaged over 25 sentences to obtain one
808
809 334 'observation' under each hearing aid setting and listening condition (SNR and masker type), which is
810
811 335 commonly performed in pupillometry research (Koelewijn et al., 2012; Koelewijn et al., 2014;
812
813 336 Ohlenforst et al., 2017a; Zekveld et al., 2011). A linear mixed-effects model was built in R-studio
814
815 337 using the packages lme4 (Bates et al., 2014) and lmerTest (Kuznetsova et al., 2016). The function
816
817 338 lmer was applied to fit the LMM to the data. First, we applied a 3-way LMM ANOVA to statistically
818
819 339 compare the fixed effects of the masker types, SNR and noise reduction setting on the PPD and the
820
821
822
823
824
825
826

827
828
829 340 sentence recognition performance separately to verify the hypothesis that the masker type and SNR
830
831 341 range have an impact on speech recognition performance and the corresponding listening effort. The
832
833 342 probability level of each LMM ANOVA was $p < 0.05$. We did not observe a significant 3-way
834
835 343 interaction effect on the PPD, but we did observe a significant interaction between the SNR and
836
837 344 noise reduction scheme setting. The model was collapsed across masker type, and an additional 2-
838
839 345 way LMM ANOVA was applied to assess the effect of the SNR and noise reduction scheme setting
840
841 346 and the corresponding interaction effect on the PPD.
842
843
844

845 347 The three-way interaction among the masker type, SNR and noise reduction scheme setting on
846
847 348 sentence recognition performance was significant. We created two additional separate LMM
848
849 349 ANOVAs to test the effect of the SNR of each masker type independently (stationary noise and 4-
850
851 350 talker masker) on the percent-correct sentence recognition. In these models, the averaged
852
853 351 percentage of correct sentence recognition scores for each SNR was treated as a dependent
854
855 352 measure, and the participants were treated as a repeated measure, i.e., random effects. The fixed
856
857 353 effects in each separate LMM ANOVA included the categorical variable SNR, the categorical variable
858
859 354 noise reduction scheme setting and the interaction between the SNR and noise reduction scheme
860
861 355 setting. We included the random effect of the SNR and noise reduction scheme as a random slope of
862
863 356 SNR to allow each participant to have their own mean PPD size and effect of SNR or noise reduction
864
865 357 scheme on PPD with both factors nested within participants. The *phia* package, including the
866
867 358 *testInteractions* functions, was used to apply a post hoc interaction analysis. Pairwise comparisons of
868
869 359 the noise reduction scheme setting (on or off) at each SNR level were conducted. The pairwise post-
870
871 360 hocpost hoc analysis was separately applied to both outcome measures (PPD and sentence
872
873 361 recognition performance), and a p-value correction using the Holm method was applied to correct
874
875 362 for the multiple comparisons.
876
877
878

879 363
880
881
882
883
884
885

364 3. Results

365 3.1 Sentence recognition data

366 The results are displayed in Figures 3 and 4. Figure 3 shows the sentence recognition scores across
367 the range of stationary noise masker SNRs with the noise reduction scheme on (solid, gray curve) or
368 off (dashed, gray curve). The sentence recognition scores with the 4-talker masker are shown in
369 Figure 4 with the noise reduction scheme on (solid, gray curve) or off (dashed, gray curve). The error
370 bars represent the standard error of the mean.

371 The 3-way LMM ANOVA revealed significant main effects of SNR ($F_{[7,713]}=1382.5, p<0.001$), noise
372 reduction scheme ($F_{[1,713]}=524.4, p<0.001$), and masker type ($F_{[1,713]}=72.9, p<0.001$), indicating that
373 sentence recognition is affected by differences in the listening conditions (SNR and masker type) and
374 the noise reduction processing algorithm. Furthermore, we found significant interactions between
375 the SNR and noise reduction scheme ($F_{[7,713]}=93.7, p<0.001$), between the SNR and masker type
376 ($F_{[7,713]}=5.73, p<0.001$) and among the SNR, noise reduction scheme and masker type ($F_{[7,713]}=2.82,$
377 $p<0.01$). The interaction between the masker type and noise reduction scheme was not significant.
378 The interaction effects of among the masker type, noise reduction scheme and SNR are larger in the
379 mid-range of SNRs, while at relatively low and high SNRs, floor or ceiling effects of sentence
380 recognition were observed.

381
382 Regarding the stationary noise masker, at relatively high SNRs between +16 dB and +8 dB, the
383 participants achieved 100% sentence recognition independent of the setting of the noise reduction
384 scheme. As the SNR decreased (+8 dB to -8 dB), sentence recognition rapidly decreased until the
385 participants were unable to perform correct sentence recall at -12 dB SNR when the noise reduction
386 scheme was turned off. At -12 dB SNR, the participants could correctly recognize approximately 12%
387 when the noise reduction scheme was turned on. Overall, the sentence recognition curve at the level

945
946
947 388 of 50% correct speech recognition was shifted by approximately 5.5 dB (see Figure 3) toward lower
948
949 389 SNRs when the noise reduction scheme was turned on compared to that when it was turned off. The
950
951 390 LMM ANOVA revealed significant main effects of SNR ($F_{[7,345]}=846.2, p<0.001$) and noise reduction
952
953 391 scheme ($F_{[1,345]}=332.5, p<0.001$) and a significant interaction between the SNR and noise reduction
954
955 392 scheme ($F_{[7,345]}=68.8, p<0.001$). We performed pairwise post hoc comparisons between the two noise
956
957 393 reduction scheme settings (on or off) at each SNR level. Post hoc analysis revealed significant
958
959 394 differences between the noise reduction scheme settings at -12 dB, -8 dB, -4 dB and 0 dB SNR
960
961 395 ($p<0.01$, as indicated by gray diamonds in Figure 3).
962
963
964 396
965
966
967 397 Regarding the 4-talker masker, at SNRs between +16 dB and +8 dB, nearly 100% sentence recognition
968
969 398 was achieved regardless of the noise reduction setting. The overall performance curve was shifted by
970
971 399 approximately 5.1 dB toward the lower SNRs when the noise reduction scheme was turned on
972
973 400 compared to that when it was turned off. By applying an LMM ANOVA, we found significant main
974
975 401 effects of SNR ($F_{[7,345]}=617.3, p<0.001$) and noise reduction scheme ($F_{[1,345]}=223.8, p<0.001$) and a
976
977 402 significant interaction between the SNR and noise reduction scheme ($F_{[7,345]}=36.2, p<0.001$). We
978
979 403 performed pairwise post hoc comparisons between the two noise reduction scheme settings (on or
980
981 404 off) at each SNR level. Significant differences were observed in the sentence recognition performance
982
983 405 between the noise reduction scheme settings at -8 dB, -4 dB, 0 dB and +4 dB SNR ($p<0.01$, as
984
985 406 indicated by gray diamonds in Figure 4).
986
987
988
989 407
990
991
992 408 Arcsine transformation prior to analyzing proportion data, such as the percent of correct responses,
993
994 409 is known to stabilize the variance and normalize proportional data (Studebaker, 1985). We applied
995
996 410 the arcsine transformation to the speech scores and performed the statistical analysis described in
997
998 411 section 2.7 by using LMM ANOVAs of the speech data. The results revealed small differences in the F-
999
1000
1001
1002
1003

1004
1005
1006 412 and p-values compared to those obtained by analyzing the percentage scores. We chose to apply the
1007
1008 413 statistical analysis of the speech data because the prior arcsine transformation did not change the
1009
1010 414 results, and arcsine units are difficult to interpret as they fall into a numeric range that has little
1011
1012 415 intuitive relationship to the proportionate performances.
1013
1014
1015 416

1018 417 3.2 Pupil data

1021 418 Figure 3 shows the PPD data under the stationary noise masker conditions, and Figure 4 shows the
1022
1023 419 PPD data under the 4-talker masker conditions across SNRs. The 3-way LMM ANOVA revealed the
1024
1025 420 significant main effects of the SNR ($F_{[7,699.1]}=26.82, p<0.001$), noise reduction scheme ($F_{[1,699.1]}=25.34,$
1026
1027 421 $p<0.001$), and masker type ($F_{[1,699.1]}=21.37, p<0.01$), and a significant interaction was observed
1028
1029 422 between the SNR and noise reduction scheme ($F_{[7,699.1]}=9.97, p<0.01$). No significant interaction was
1030
1031 423 observed between the masker type and SNR or masker type and noise reduction scheme. The
1032
1033 424 interaction effect between the SNR and noise reduction scheme suggests that the SNR-dependency
1034
1035 425 of the PPD differs when the noise reduction scheme is on from that when the scheme is off. We did
1036
1037 426 not test two separate models for each masker type per sentence recognition performance. In an
1038
1039 427 additional 2-way LMM ANOVA that collapsed across the level of masker type, the noise reduction
1040
1041 428 scheme setting and masker type were not significant, which is similar to the interaction with the SNR.
1042
1043 429 The 2-way LMM ANOVA revealed a significant main effect of noise reduction scheme setting
1044
1045 430 ($F_{[1,715.05]}=25.08, p<0.001$), a significant main effect of SNR ($F_{[7,715.07]}=25.94, p<0.001$) and a significant
1046
1047 431 interaction effect between the noise reduction scheme setting and SNR ($F_{[7,715.05]}=9.72, p<0.001$) on
1048
1049 432 the PPD. Pairwise post hoc comparisons of the two noise reduction scheme settings (on or off) were
1050
1051 433 applied at each SNR level. Significant differences were observed between the noise reduction
1052
1053 434 scheme settings in the PPD measured at -8 dB, -4 dB, 0 dB and +4 dB SNR.
1054
1055
1056
1057 435

1063
1064
1065 436 Figure 3 shows the averaged PPD data across SNRs for the stationary noise masker when the noise
1066
1067 437 reduction scheme was active (black, solid line) and when the noise reduction scheme was inactive
1068
1069 438 (black, dashed line). The PPD plateaued with relatively high SNRs between +16 and +8 dB where high
1070
1071 439 performance was reached independently of the noise reduction scheme setting. When the noise
1072
1073 440 reduction scheme was turned off, as the SNR further decreased, a steady increase in PPD was
1074
1075 441 observed until a maximum PPD was reached at -4 dB SNR. The corresponding sentence recognition
1076
1077 442 was approximately 38% correct. The maximum PPD was shifted by 4 dB toward lower SNRs when the
1078
1079 443 noise reduction scheme was turned on, and this maximum corresponded to an approximately 52%
1080
1081 444 correct sentence recognition. At the lowest SNR of -12 dB, relatively lower PPDs were observed
1082
1083 445 under both noise reduction scheme settings.

1084
1085
1086 446 Figure 4 shows the PPD data across SNRs with the noise reduction scheme on (black, solid curve) or
1087
1088 447 off (black, dashed curve) under the 4-talker masker condition. The PPD measured with high SNRs
1089
1090 448 between +16 dB and +8 dB was overall consistent but larger when the noise reduction scheme was
1091
1092 449 off compared than that when it was on. Further decreases in the SNRs resulted in continuous
1093
1094 450 increases in the PPD until the maximum PPD was reached between -4 dB and 0 dB SNR when the
1095
1096 451 noise reduction scheme was off and between -8 and -4 dB SNR when the noise reduction scheme
1097
1098 452 was on. The range of the maximum PPD was shifted by approximately 4 dB toward the lower SNRs
1099
1100 453 when the noise reduction scheme was turned on compared to that when it was turned off.
1101
1102
1103
1104 454

1105 1106 1107 455 3.3 Summary of the results

1108
1109
1110 456 The preceding statistical analyses support the following summary of the results: The effect of the
1111
1112 457 noise reduction scheme applied in this study on sentence recognition was to shift the performance
1113
1114 458 function across SNRs by approximately 5.5 dB for the stationary masker and approximately 5.1 dB for
1115
1116 459 the 4-talker masker toward the lower SNRs. For both masker types, the effect of the noise reduction
1117
1118
1119
1120
1121

1122
1123
1124 460 scheme on listening effort (as measured by the PPD) was to shift the peak of the PPD function across
1125
1126 461 SNRs by approximately 4 dB toward the lower SNR. In addition, in the case of the 4-talker masker,
1127
1128 462 the noise reduction scheme lowered the average PPD by approximately 35% compared to the
1129
1130 463 inactive noise reduction scheme.
1131
1132

1133 464

1136 465

1139 466 4. Discussion

1140
1141
1142 467 In the present study, the effect of a noise reduction scheme on sentence recognition and PPD was
1143
1144 468 examined across a range of SNRs with two masker types. For both masker types, the noise reduction
1145
1146 469 scheme had a large beneficial effect on sentence recognition, which was accompanied by a
1147
1148 470 corresponding effect on listening effort, as indicated by the PPD.
1149
1150

1151 471

1154 472 4.1 Relationship among noise reduction scheme, SNR and speech recognition

1155
1156 473 For the stationary and 4-talker maskers, the sentence recognition performance was significantly
1157
1158 474 improved when the noise reduction scheme was active compared to that when it was inactive. The
1159
1160 475 results showed improved sentence recognition not only at performance levels of approximately 50%
1161
1162 476 and higher but also at lower sentence recognition performances. Notably, sentence recognition was
1163
1164 477 mainly improved across a large range of negative SNRs between 0 dB and -12 dB. The findings of the
1165
1166 478 present study confirm and extend the previously shown benefits of a noise reduction scheme on
1167
1168 479 sentence recognition with an approximately 50% successful performance rate (Wendt et al., 2017) at
1169
1170 480 higher and lower performance levels. Additionally, the present study confirmed that the currently
1171
1172 481 tested noise reduction scheme can significantly improve speech intelligibility in very challenging
1173
1174
1175
1176
1177
1178
1179
1180

1181
1182
1183 482 sound environments. Hence, this finding might allow hearing-impaired listeners to participate in
1184
1185 483 communication situations that might otherwise be impossibly challenging.
1186
1187

484

1191 485 4.2 Relationship among noise reduction scheme, SNR and PPD

1192
1193
1194 486 In line with recent research (Ohlenforst et al., 2017a; Zekveld et al., 2014), the present results
1195
1196 487 confirm that the changes in speech recognition are accompanied by changes in PPD. We found the
1197
1198 488 maximum PPD with SNRs producing approximately 50% correct sentence recognition and relatively
1199
1200 489 smaller PPDs at very low and very high SNRs. The indication that listening effort follows an inverted
1201
1202 490 U-shape across a range of SNRs also supports the findings reported in a recent study (Wu et al., 2016)
1203
1204 491 in which dual-task paradigms were applied to assess listening effort across a wide range of SNRs. Wu
1205
1206 492 et al. (2016) found that second-task performance (reaction time) was the worst (i.e., longest) at SNRs
1207
1208 493 for 30-50% speech recognition and better at both lower and higher SNRs. The change in the PPD
1209
1210 494 function at positive SNRs when the percent-correct sentence recognition is saturated might be
1211
1212 495 affected by the type of speech material used in the sentence recognition test. The transfer function
1213
1214 496 of the speech intelligibility index is modifiable depending on the tested sentence material, and more
1215
1216 497 difficult speech material can change the transfer function. Thus, the transfer function at positive
1217
1218 498 SNRs might already be saturated for speech intelligibility index values that are not at the level of
1219
1220 499 saturation. However, we designed this experiment to intentionally reach a ceiling in performance,
1221
1222 500 although with very positive SNRs, a ceiling effect is achieved regardless of the presented speech
1223
1224 501 material.
1225
1226

1227
1228 502 The statistical analysis revealed that the level of the SNR and the noise reduction scheme setting
1229
1230 503 significantly affected the PPD. The impact of the masker type on the PPD was rather small, which
1231
1232 504 might contrast with previous studies reporting that listening effort required for speech recognition is
1233
1234 505 altered by the type of background masker (e.g., Koelewijn et al., 2012; Koelewijn et al., 2014).
1235
1236
1237
1238
1239

1240
1241
1242 506 Koelewijn and colleagues reported significantly larger pupil dilation responses for masker types
1243
1244 507 containing speech information, and the increase in effort was mainly explained by the semantic
1245
1246 508 inference with the target. However, Koelewijn and colleagues examined the impact of masker types
1247
1248 509 on the PPD at similar intelligibility levels corresponding to 50% correct speech recognition. Therefore,
1249
1250 510 comparisons between the PPDs of the different masker types were drawn at varying SNRs. Our data
1251
1252 511 indicate that the PPDs are strongly affected by the SNRs, which is in line with the results of previous
1253
1254 512 studies (Zekveld and Kramer, 2014; Ohlenforst et al., 2017). Hence, the differentiation between the
1255
1256 513 effect of the SNR and masker type is not possible based on these aforementioned studies by
1257
1258 514 Koelewijn and colleagues. Our results suggest that when examining the PPD across a range of
1259
1260 515 intelligibility varying between 0 to 100%, a non-linear change in the PPD, with maximum PPDs
1261
1262 516 occurring at approximately 50% recognition, could be observed independently of the masker type.
1263
1264 517 Furthermore, the impact of the masker type might be less pronounced when testing fixed SNRs,
1265
1266 518 which is in line with the results of previous work (see Wendt et al., 2018 in press).
1267
1268
1269

1270 519
1271
1272
1273 520 One strength of the present study is the replication of previous findings, demonstrating the beneficial
1274
1275 521 effect of a noise reduction scheme in hearing aids on sentence recognition and the PPD (Wendt et
1276
1277 522 al., 2017). There were several factors that were kept constant between the setup of the recent study
1278
1279 523 by Wendt and colleagues (2017) and the current study. In both studies, the same noise reduction
1280
1281 524 scheme was tested during a sentence recognition task with identical stimulus material (HINT
1282
1283 525 sentences in a 4-talker masker). Additionally, the large number of listeners (n=17) that participated in
1284
1285 526 the study by Wendt et al., (2017) were used in the present study. Both studies contribute to the field
1286
1287 527 of hearing research and listening effort by providing new valuable knowledge showing the possible
1288
1289 528 benefits of a noise reduction scheme for hearing-impaired listeners wearing hearing aids.
1290
1291

1292 529
1293
1294
1295
1296
1297
1298

5. Conclusion

The present study demonstrates that a noise reduction scheme in commercial hearing aids can reduce the effort required during speech recognition in stationary noise and a 4-talker masker. With both maskers, the noise reduction processing resulted in a shift in the performance (sentence recognition) function toward lower (more challenging) SNRs, and a corresponding shift in the PPD function was observed. For the 4-talker masker, in addition to the speech recognition-related reduction in the PPD, a main effect of noise reduction processing on the PPD was observed, indicating that the cognitive processing load and some aspects of listening effort may be reduced independent of the SNR. These results also confirm previous findings by showing that for hearing-impaired listeners using hearing aids during speech recognition, listening effort changes in a non-monotonic way as a function of the SNR. This knowledge is essential for future research in the field of listening effort and the hearing aid industry for improving the development of better hearing aid algorithms.

6. Acknowledgments

The authors would like to thank Per Bruun Brockhoff from the Institute for Mathematics and Computer Science at the Technical University of Denmark (DTU Compute) and Birgit Lissenberg-Witte from the department of Epidemiology and Biostatistics at the VU in Amsterdam for their support and advice with the statistical analyses. We would also like to thank Renskje Hietkamp for her support with the participant recruitment and data collection and Nicolas Le Goff and Jesper Jensen for the fruitful discussions about the tested hearing aid technology. We would like to thank Jacob Aderhold for technical support and advice with the hearing aids used in this study and Yang Wang for fruitful teamwork and discussions throughout the study. Finally, we wish to thank all the participants, the

1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416

554 European Commission (grant FP7-LISTEN607373) and the Oticon Foundation for supporting this
555 study. Co-author GN was supported by the UK Medical Research Council (grant U135097131) and a
556 grant from the Chief Scientist Office.

557

558 **Appendix**

559

Table 1: Beta estimates of the sentence recognition performance scores and PPD at each SNR level show the mean differences between the inactive and active noise reduction scheme setting. The SNR levels are compared to the lowest SNR at -12 dB.

SNRs [dB] compared to the reference SNR of -12 dB	-8	-4	0	+4	+8	+12	+16
Beta estimates of performance with the stationary noise masker	-33.68	-40.57	-5.68	7.92	11.03	11.67	11.40
Beta estimates of performance with the 4-talker masker	-26.57	-46.05	-20.73	-4.07	4.57	4.40	5.97
Beta estimates of the PPD collapsed across stationary noise masker and 4-talker masker	-0.03	0.04	0.08	0.05	0.03	0.01	0.03

560

561

1476
1477
1478 **562 References**
1479

- 1480 563 Baer, T., Moore, B.C.J., Gatehouse, S. 1993. Spectral contrast enhancement of speech in noise for
1481 564 listeners with sensorineural hearing impairment: Effects on intelligibility, quality, and
1482 565 response times. . *Journal of Rehabilitation Research and Development* 30, 49-72.
1483 566 Bates, D., Mächler, M., Bolker, B., Walker, S. 2014. Fitting linear mixed-effects models using lme4.
1484 567 arXiv preprint arXiv:1406.5823.
1485 568 Beatty, J. 1982. Task-evoked pupillary responses, processing load, and the structure of processing
1487 569 resources. *Psychological bulletin* 91, 276.
1488 570 Bentler, R., Wu, Y.-H., Kettel, J., Hurtig, R. 2008. Digital noise reduction: Outcomes from laboratory
1490 571 and field studies. . *International journal of audiology* 47, 447-460.
1491 572 Borch Petersen, E., Wöstmann, M., Obleser, J., Lunner, T. 2017. Neural tracking of attended versus
1492 573 ignored speech is differentially affected by hearing loss. . *Journal of Neurophysiology* 117,
1493 574 18-27.
1494 575 Brehm, J.W. 1999. The intensity of emotion. *Personality and Social Psychology Review* 3, 2-22.
1495 576 Dillon, H., Lovegrove, R. 1993. Single microphone noise reduction systems for hearing aids: A review
1497 577 and an evaluation. Boston: Allyn and Bacon.
1498 578 Engelhardt, P.E., Ferreira, F., Patsenko, E.G. 2010. Pupillometry reveals processing load during spoken
1499 579 language comprehension. *The Quarterly Journal of Experimental Psychology* 63, 639-645.
1500 580 Festen, J.M., Plomp, R. 1990. Effects of fluctuating noise and interfering speech on the
1501 581 speech-reception threshold for impaired and normal hearing. *The Journal of the Acoustical*
1502 582 *Society of America* 88, 1725-1736.
1503 583 Francart, T., Van Wieringen, A., Wouters, J. 2011. Comparison of fluctuating maskers for speech
1505 584 recognition tests. *International journal of audiology* 50, 2-13.
1506 585 Fraser, S., Gagné, J.-P., Alepins, M., Dubois, P. 2010. Evaluating the effort expended to understand
1507 586 speech in noise using a dual-task paradigm: The effects of providing visual speech cues.
1508 587 *Journal of Speech, Language, and Hearing Research* 53, 18-33.
1510 588 Gatehouse, S., Gordon, J. 1990. Response times to speech stimuli as measures of benefit from
1511 589 amplification. . *British Journal of Audiology*. 24, 63-68.
1512 590 Granholm, E., Asarnow, R.F., Sarkin, A.J., Dykes, K.L. 1996. Pupillary responses index cognitive
1513 591 resource limitations. *Psychophysiology* 33, 457-461.
1514 592 Houben, R., v. Doorn-Bierman, M., Dreschler, W.A. 2013. Using response time to speech as a
1515 593 measure for listening effort. . *International journal of audiology* 52, 735-761.
1517 594 Jensen, J., Pedersen, M.S. 2015. Analysis of beamformer directed single-channel noise reduction
1518 595 system for hearing aid applications, *Acoustics, Speech and Signal Processing (ICASSP), 2015*
1519 596 *IEEE International Conference on*. IEEE. pp. 5728-5732.
1520 597 Kahneman, D. 1973. *Attention and effort* Prentice-Hall Englewood Cliffs, NJ.
1522 598 Kjemis, U., Jensen, J. 2012. Maximum likelihood based noise covariance matrix estimation for multi-
1523 599 microphone speech enhancement, *Signal Processing Conference (EUSIPCO), 2012*
1524 600 *Proceedings of the 20th European*. IEEE. pp. 295-299.
1525 601 Koelewijn, T., Zekveld, A.A., Festen, J.M., Kramer, S.E. 2012. Pupil dilation uncovers extra listening
1526 602 effort in the presence of a single-talker masker. *Ear and Hearing* 33, 291-300.
1527 603 Koelewijn, T., Shinn-Cunningham, B.G., Zekveld, A.A., Kramer, S.E. 2014. The pupil response is
1528 604 sensitive to divided attention during speech processing. *Hearing research* 312, 114-120.
1529
1530
1531
1532
1533
1534

- 1535
1536
1537 605 Kramer, S.E., Kapteyn, T.S., Festen, J.M., Kuik, D.J. 1997. Assessing aspects of auditory handicap by
1538 606 means of pupil dilatation. *Audiology* 36, 155-164.
- 1539 607 Kuchinsky, S.E., Ahlstrom, J.B., Vaden Jr., K.I., Cute, S.L., Humes, L.E., Dubno, J.R., Eckert, M.A. 2013.
1540 608 Pupil size varies with word listening and response selection difficulty in older adults with
1541 609 hearing loss. . *Psychophysiology* 50, 23-34.
- 1542 610 Kuznetsova, A., Bruun Brockhoff, P., H., B.C.R. 2016. lmerTest: Tests in Linear Mixed Effects Models. R
1543 611 package version 2.0-33., <https://CRAN.R-project.org/package=lmerTest>.
1544
1545 612 Larsby, B., Hällgren, M., Lyxell, B. 2008. The interference of different background noises on speech
1546 613 processing in elderly hearing impaired subjects. . *International journal of audiology* 47(sup2),
1547 614 S83-S90.
- 1550 615 Le Goff, N. 2015. Amplifying Soft Sounds - A Personal Matter. Oticon Whitepaper. Whitepaper,
1551 616 Oticon A/S.
- 1552 617 Loizou, P.C., Kim, G. 2011. Reasons why current speech-enhancement algorithms do not improve
1553 618 speech intelligibility and suggested solutions. . *IEEE Transactions on Audio, Speech, and*
1554 619 *Language Processing* 19, 47-56.
- 1556 620 Lunner, T., Rudner, M., Rosenbom, T., Ågren, J., Ng, E.H.N. 2016. Using Speech Recall in Hearing Aid
1557 621 Fitting and Outcome Evaluation Under Ecological Test Conditions. . *Ear and Hearing* 37, 145S-
1558 622 154S.
- 1559 623 McAuliffe, M.J., Wilding, P.J., Rickard, N.A., O'Beirne, G.A. 2012. Effect of Speaker Age and Speech
1560 624 Recognition and Perceived Listening Effort in Older Adults With Hearing Loss. . *Journal of*
1561 625 *Speech, Language, and Hearing Research* 55, 838-847.
- 1563 626 McGarrigle, R., Munro, K.J., Stewart, A.J., Dawes, P. 2014. Listening effort and fatigue: are we talking
1564 627 about the same thing? *International journal of audiology*.
- 1565 628 **Naylor G & Johannesson RB (2009) Long-Term Signal-to-Noise Ratio at the Input and Output of**
1566 629 **Amplitude-Compression SystemsJ Am Acad Audiol 20:161-171 (2009) DOI:**
1567 630 **10.3766/jaaa.20.3.2**
- 1569 631 Neher, T. 2014. Relating hearing loss and executive functions to hearing aid users' preference for, and
1570 632 speech recognition with, different combinations of binaural noise reduction and microphone
1571 633 directionality. . *Frontiers in Neuroscience* 8.
- 1572 634 Neher, T., Grimm, G., Hohmann, V., Kollmeier, B. 2013. Do Hearing Loss and Cognitive Function
1573 635 Modulate Benefit From Different Binaural Noise-Reduction Settings? . *Ear and Hearing* 35,
1574 636 52-62.
- 1576 637 Ng, E.H.N., Rudner, M., Lunner, T., Rönnerberg, J. 2015. Noise reduction improves memory for target
1577 638 language speech in competing native but not foreign language speech. *Ear and hearing* 36,
1578 639 82-91.
- 1580 640 Ng, E.H.N., Rudner, M., Lunner, T., Pedersen, M.S., Rönnerberg, J. 2013. Effects of noise and working
1581 641 memory capacity on memory processing of speech for hearing-aid users. *International*
1582 642 *Journal of Audiology* 52, 433-441.
- 1583 643 Nielsen, J.B., Dau, T. 2009. Development of a Danish speech intelligibility test. *International journal of*
1584 644 *audiology* 48, 729-741.
- 1585 645 Nordrum, S., Erler, S., Garstecki, D., D., S. 2006. Comparison of performance on the hearing in noise
1586 646 test using directional microphones and digital noise reduction algorithms. . *American Journal*
1587 647 *of Audiology* 15, 81-91.
- 1589
1590
1591
1592
1593

1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652

- 648 Obleser, J., Wöstmann, M., Hellbernd, N., Wilsch, A., Maess, B. 2012. Adverse Listening Conditions
649 and Memory Load Drive a Common Alpha Oscillatory Network. . The Journal of Neuroscience
650 32, 12376-12383.
- 651 Ohlenforst, B., Zekveld, A.A., Lunner, T., Wendt, D., Naylor, G., Wang, Y., Versfeld, N.J., Kramer, S.E.
652 2017a. Impact of stimulus-related factors and hearing impairment on listening effort as
653 indicated by pupil dilation. Hearing Research.
- 654 Ohlenforst, B., Zekveld, A.A., Jansma, E.P., Wang, Y., Naylor, G., Lorens, A., Lunner, T., Kramer, S.E.
655 2017b. Effects of Hearing Impairment and Hearing Aid Amplification on Listening Effort: A
656 Systematic Review. Ear and hearing 38, 267.
- 657 Pals, C., Sarampalis, A., Başkent, D. 2013. Listening effort with cochlear implant simulations. Journal
658 of Speech, Language, and Hearing Research 56, 1075-1084.
- 659 Panico, J., C., H.E. 2009. Influence of Text Type, Topic Familiarity, and Stuttering Frequency on
660 Listener Recall, Comprehension, and Mental Effort. Journal of Speech, Language, and
661 Hearing Research 52, 534-546.
- 662 Petersen, E.B., Wöstmann, M., Obleser, J., Stenfelt, S., Lunner, T. 2015. Hearing loss impacts neural
663 alpha oscillations under adverse listening conditions. . Frontiers in Psychology 6.
- 664 Pichora-Fuller, M.K., Kramer, S.E., Eckert, M.A., Edwards, B., Hornsby, B.W., Humes, L.E., Lemke, U.,
665 Lunner, T., Matthen, M., Mackersie, C.L. 2016. Hearing impairment and cognitive energy: The
666 framework for understanding effortful listening (FUEL). Ear and hearing 37, 5S-27S.
- 667 Picou, E.M., Ricketts, T.A., Hornsby, B.W.Y. 2011. Visual Cues and Listening Effort: Individual
668 Variability. Journal of Speech, Language, and Hearing Research 54, 1416-1430.
- 669 Picou, E.M., Ricketts, T.A., Hornsby, B.W. 2013. How hearing aids, background noise, and visual cues
670 influence objective listening effort. . Ear and Hearing 34, e52-e64.
- 671 Ricketts, T., Lindley, G., Henry, P. 2001. Impact of compression and hearing aid style on directional
672 hearing aid benefit and performance. . Ear and Hearing 22, 348-361.
- 673 Rosen, S., Souza, P., Ekelund, C., Majeed, A.A. 2013. Listening to speech in a background of other
674 talkers: Effects of talker number and noise vocoding. The Journal of the Acoustical Society of
675 America 133, 2431-2443.
- 676 Sarampalis, A., Kalluri, S., Edwards, B., Hafter, E. 2009. Objective Measures of Listening Effort: Effects
677 of Background Noise and Noise Reduction. Journal of Speech, Language, and Hearing
678 Research 52, 1230-1240.
- 679 Scott, S.K., Rosen, S., Wickham, L., Wise, R.J. 2004. A positron emission tomography study of the
680 neural basis of informational and energetic masking effects in speech perception. The Journal
681 of the Acoustical Society of America 115, 813-821.
- 682 Scott, S.K., Rosen, S., Beaman, C.P., Davis, J.P., Wise, R.J. 2009. The neural processing of masked
683 speech: evidence for different mechanisms in the left and right temporal lobes. The Journal
684 of the Acoustical Society of America 125, 1737-1743.
- 685 Shinn-Cunningham, B.G., Best, V. 2008. Selective attention in normal and impaired hearing. . Trends
686 in amplification 12, 283:299.
- 687 Siegle, G.J., Ichikawa, N., Steinhauer, S. 2008. Blink before and after you think: blinks occur prior to
688 and following cognitive load indexed by pupillary responses. Psychophysiology 45, 679-687.
- 689 Simonsen, L.B., Hietkamp, R.K., Bramsløw, L. 2016. Learning effects of repeated exposure to Hearing
690 In Noise Test, Annual Conference of the British Society of Audiology, Coventry, UK.

1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711

- 691 Smeds, K., Wolters, F., Rung, M. 2015. Estimation of signal-to-noise ratios in realistic sound scenarios.
692 *Journal of the American Academy of Audiology* 26, 183-196.
- 693 **Studebaker, G. A. (1985). A rationalized arcsine transform. *Journal of Speech, Language, and Hearing***
694 ***Research*, 28(3), 455-462.**
- 695 Tun, P.A., McCoy, S., Wingfield, A. 2009. Aging, Hearing Acuity, and the Attentional Costs of Effortful
696 Listening. . *Psychology and Aging* 24, 761-766.
- 697 Van Der Meer, E., Beyer, R., Horn, J., Foth, M., Bornemann, B., Ries, J., Kramer, J., Warmuth, E.,
698 Heekeren, H.R., Wartenburger, I. 2010. Resource allocation and fluid intelligence: Insights
699 from pupillometry. *Psychophysiology* 47, 158-169.
- 700 Wendt, D., Hietkamp, R.K., Lunner, T. 2017. Impact of noise and noise reduction on processing effort:
701 A pupillometry study. *The Journal of the Acoustical Society of America* 141, 4040-4040.
- 702 **Wendt, D., Koelewijn, T., Książek, P., Kramer, S.E, Lunner, T. 2017. Toward a more comprehensive**
703 **understanding of the impact of masker type and signal-to-noise ratio on the pupillary**
704 **response while performing a speech-in-noise test. *Accepted for Hearing Research***
- 705 Wu, Y.-H., Stangl, E., Zhang, X., Perkins, J., Eilers, E. 2016. Psychometric functions of dual-task
706 paradigms for measuring listening effort. *Ear and hearing* 37, 660-670.
- 707 Wu, Y.-H., Aksan, N., Rizzo, M., Stangl, E., Zhang, X., Bentler, R. 2014. Measuring listening effort:
708 Driving simulator vs. simple dual-task paradigm. *Ear and hearing* 35, 623.
- 709 Zekveld, A.A., Kramer, S.E. 2014. Cognitive processing load across a wide range of listening
710 conditions: Insights from pupillometry. *Psychophysiology* 51, 277-284.
- 711 Zekveld, A.A., Kramer, S.E., Festen, J.M. 2011. Cognitive load during speech perception in noise: the
712 influence of age, hearing loss, and cognition on the pupil response. *Ear and hearing* 32, 498-
713 510.

714
715

716

Figure legends

718 Fig. 1: Averaged pure tone hearing thresholds of the left and right ears across frequencies
719 (125 Hz to 8 kHz) among the twenty-four hearing-impaired participants. Error bars show the standard
720 deviations of the mean.

721 Fig. 2: Spatial loudspeaker setup as used in Wendt et al., 2017. Target speech was presented
722 from the front. Masker signals were presented at 90, 150, 210 and 270 degree azimuth. The
723 stationary noise masker was presented as four individual point sources. For the four-talker masker,
724 one single talker was presented from one loudspeaker each.

725 Fig. 3: Peak pupil dilation (PPD) (black color) and percentage-correct sentence recognition
726 scores (gray color) are shown on the right y-axis across the signal-to-noise ratios (SNRs) with the
727 stationary masker and the noise reduction scheme turned on or off. Error bars represent the
728 standard error of the mean. Dark gray diamonds at -12, -8, -4, 0 and +4 dB SNR represent significant
729 differences in sentence recognition performance between the active and inactive noise reduction in
730 the pairwise comparison at each SNR level ($p < 0.01$).

731 Fig. 4: Peak pupil dilation (PPD) (black color) and the percentage of correct sentence
732 recognition scores (gray color) are shown on the right y-axis across the signal-to-noise ratios (SNRs)
733 with the 4-talker masker and noise reduction scheme on or off. Error bars represent the standard
734 error of the mean. Dark gray diamonds at -8, -4, 0 and +4 dB SNR represent significant differences in
735 sentence recognition performance between the active and inactive noise reduction in the pairwise
736 comparison at each SNR level ($p < 0.01$).

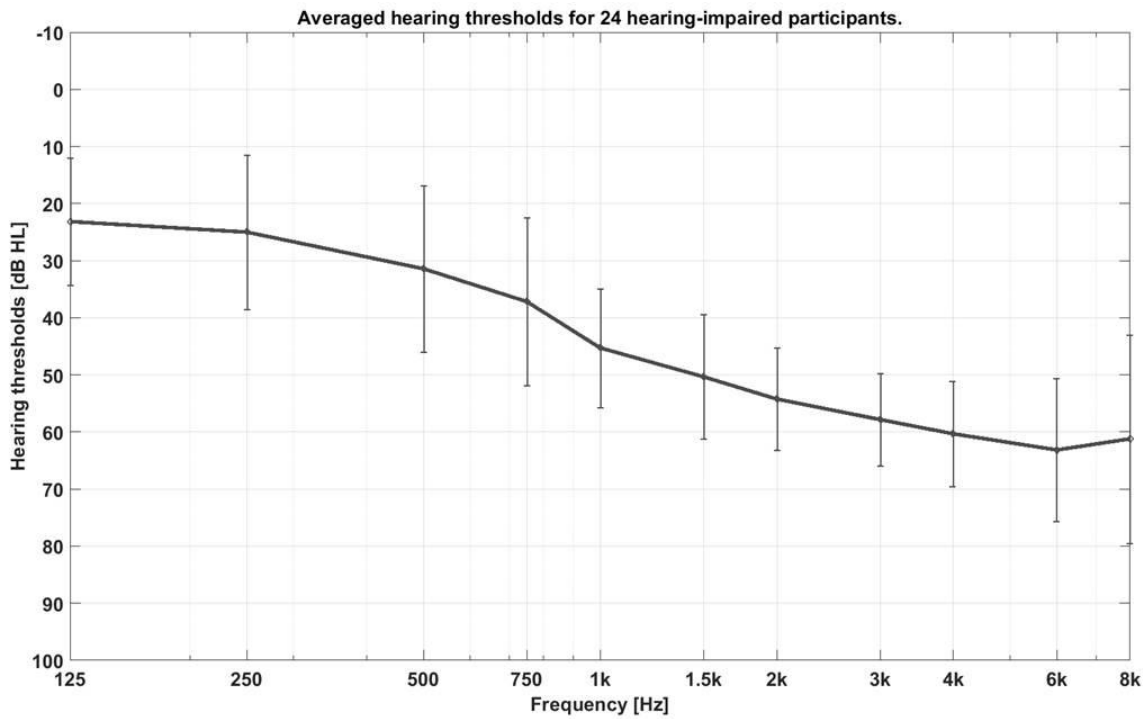


Fig. 1: Averaged pure tone hearing thresholds of the left and right ears across frequencies (125 Hz to 8 kHz) among the twenty-four hearing-impaired participants. Error bars show the standard deviations of the mean.

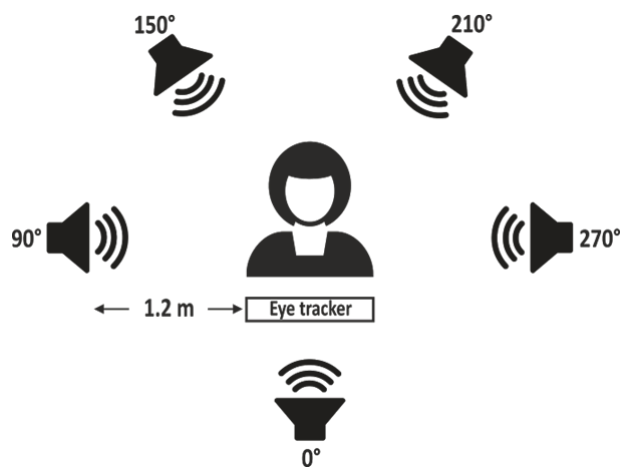


Fig. 2: Spatial loudspeaker setup as used in Wendt et al., 2017. Target speech was presented from the front. Masker signals were presented at 90, 150, 210 and 270 degree azimuth. The stationary noise masker was presented as four individual point sources. For the four-talker masker, one single talker was presented from one loudspeaker each.

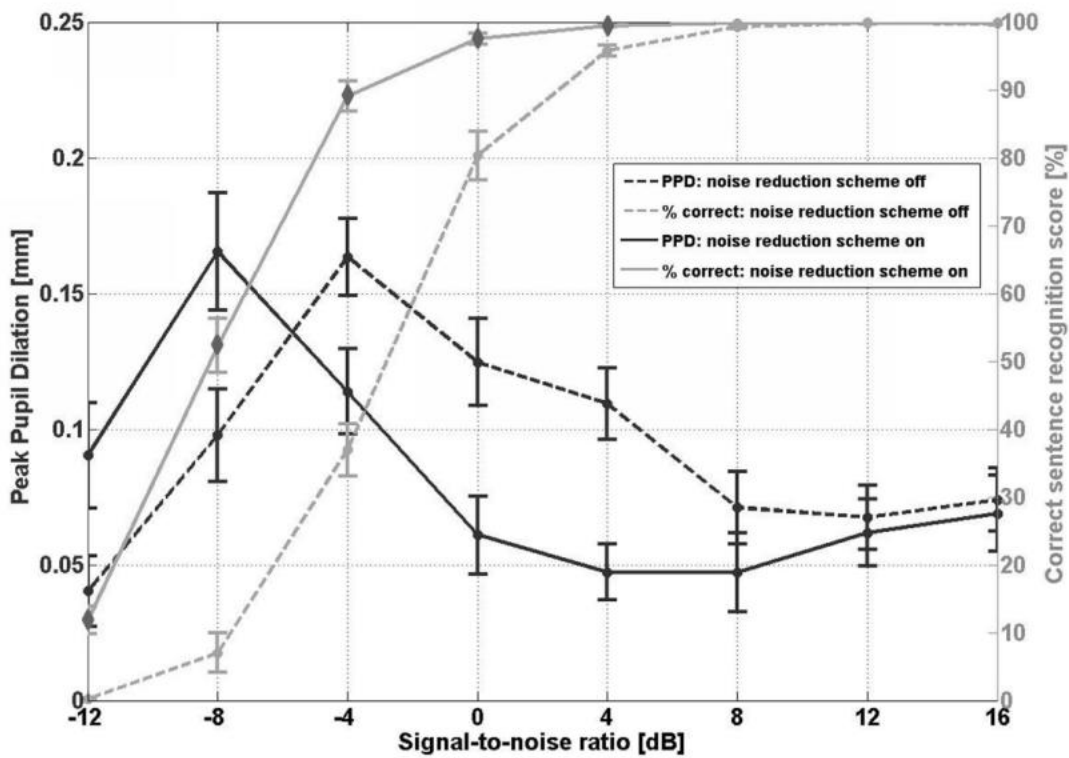


Fig. 3: Peak pupil dilatation (PPD) (black color) and percentage-correct sentence recognition scores (gray color) are shown on the right y-axis across the signal-to-noise ratios (SNRs) with the stationary masker and the noise reduction scheme turned on or off. Error bars represent the standard error of the mean. Dark gray diamonds at -12, -8, -4, 0 and +4 dB SNR represent significant differences in sentence recognition performance between the active and inactive noise reduction in the pairwise comparison at each SNR level ($p < 0.01$).

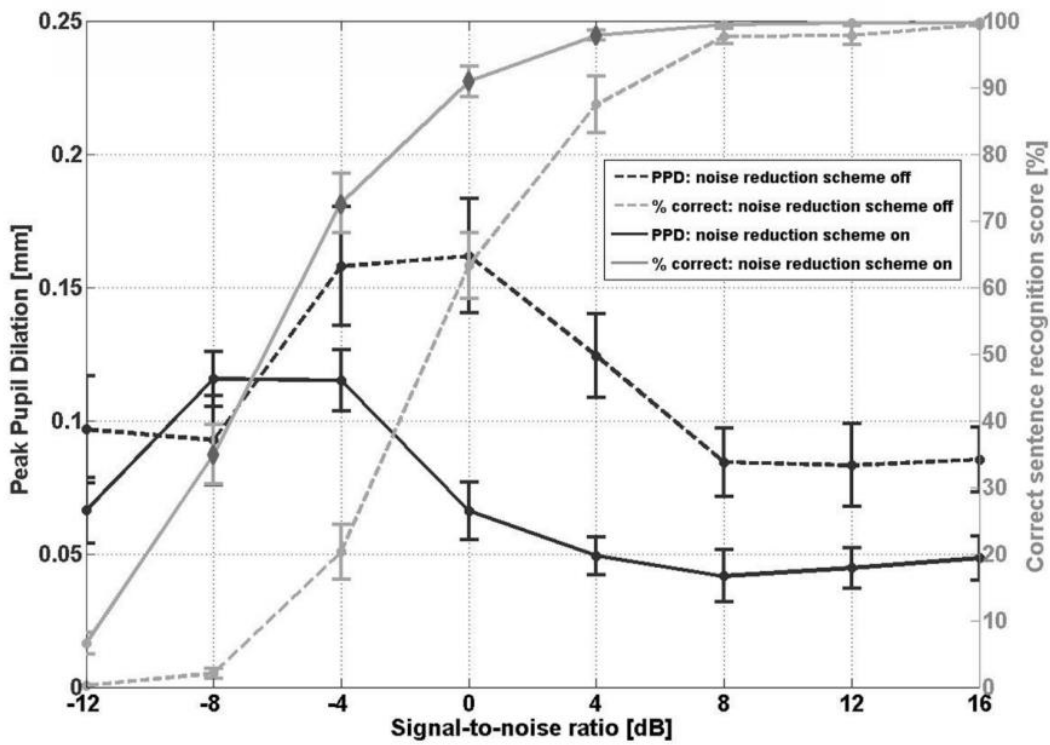


Fig. 4: Peak pupil dilation (PPD) (black color) and the percentage of correct sentence recognition scores (gray color) are shown on the right y-axis across the signal-to-noise ratios (SNRs) with the 4-talker masker and noise reduction scheme on or off. Error bars represent the standard error of the mean. Dark gray diamonds at -8, -4, 0 and +4 dB SNR represent significant differences in sentence recognition performance between the active and inactive noise reduction in the pairwise comparison at each SNR level ($p < 0.01$).