

1 **Working Title: Cognitive mechanisms underpinning successful perception of different**
2 **speech distortions.**

3 **Authors:**

4 Dan Kennedy-Higgins^{a,b},

5 *Department of Speech, Hearing and Phonetic Sciences, University College London, Chandler*
6 *House, 2 Wakefield Street, London, United Kingdom, WC1N 1PF*

7

8 Joseph T. Devlin,

9 *Department of Experimental Psychology, University College London, 26 Bedford Way,*
10 *London, United Kingdom, WC1H 0AP*

11

12 Patti Adank

13 *Department of Speech, Hearing and Phonetic Sciences, University College London, Chandler*
14 *House, 2 Wakefield Street, London, United Kingdom, WC1N 1PF*

15

16

17 ^a Electronic mail: daniel.kennedy-higgins@kcl.ac.uk.

18 ^b Current address: *Department of Psychology, King's College London, Guy's Campus,*
19 *London, United Kingdom, SE1 1UL*

20

21

1 **Abstract**

2 Few studies thus far have investigated whether perception of distorted speech is consistent
3 across different types of distortion. We investigated whether participants show a consistent
4 perceptual profile across three speech distortions: time-compressed, noise-vocoded and
5 speech in noise. **Additionally, we investigated whether/how individual differences in**
6 **performance on a battery of audiological and cognitive tasks links to perception.**

7 **Eighty-eight** participants completed a speeded sentence-verification task with increases in
8 accuracy and reductions in response times used to indicate performance. **Audiological and**
9 **cognitive task measures include pure tone audiometry, speech recognition threshold,**
10 **working memory, vocabulary knowledge, attention switching, and pattern analysis.**

11 Despite previous studies suggesting that temporal and spectral/environmental perception
12 require different lexical or phonological mechanisms, we show significant positive
13 correlations in accuracy and response time performance across all distortions. Results of a
14 principal component analysis and multiple linear regressions suggest that a **component**
15 **based on vocabulary knowledge and working memory** predicted performance in the
16 speech in quiet, time-compressed and speech in noise conditions. These results suggest that
17 listeners employ a similar cognitive strategy to perceive different temporal and
18 spectral/environmental speech distortions and that this mechanism is supported by vocabulary
19 knowledge and working memory.

20
21 **Keywords:** Time-compressed speech; noise-vocoded speech; speech in noise; individual
22 differences.

1 they cannot speak compared to a language they can comprehend. However, more adaptation
2 occurs if the training language shares an isochrony with the test language. For example,
3 training on a syllable-timed language such as Italian will benefit participants only if the test
4 sentences are also in syllable-timed languages, such as Spanish. Such a benefit will not occur
5 for stress-timed languages such as English. This phenomenon suggests that changes may be
6 occurring at phonological as opposed to lower acoustic levels of processing, as simple
7 exposure to time-compression does not result in equivalent benefits across training languages.
8 Furthermore, as comprehension does not appear to be necessary for adaptation, changes in
9 attention at the phonological level appear most important for temporal distortions.

10 In comparison for spectral manipulations such as noise-vocoded speech or environmental
11 distortions such as speech in noise, these changes in attention are believed to occur at
12 lexical/semantic levels (Bradlow & Alexander, 2007; Burk, Humes, Amos, & Strauser, 2006;
13 Cainer et al., 2008; Davis et al., 2005; Hervais-Adelman, Davis, Johnsrude, & Carlyon, 2008;
14 Hervais-Adelman, Davis, Johnsrude, Taylor, & Carlyon, 2011; Loizou et al., 1999; Mayo,
15 Florentine, & Buus, 1997). The importance of lexical/semantic level changes was
16 demonstrated by Davis et al. (2005). Two groups of participants were trained, via passive
17 listening, on 20 sentences containing words or 20 sentences containing non-word noise-
18 vocoded sentences and then were tested on 20 noise-vocoded English words. Overall, the
19 group trained with words performed significantly better than the group trained on non-words.
20 Additionally, the group trained on non-words performed at a level that was equivalent to
21 subjects that were completely naïve, i.e., subjects that had no prior history of exposure to
22 vocoded speech. This result suggests that adaptation to noise-vocoded speech is dependent on
23 either lexical, semantic and/or syntactic information with phonological information being less
24 important, i.e., the exact opposite of time-compressed speech. For a full review see
25 (Kennedy-Higgins, 2019; Mattys, Davis, Bradlow, & Scott, 2012; Samuel & Kraljic, 2009)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

A. Perception of different speech distortions

Most studies thus far have investigated how listeners perceive a single distortion (Bradlow & Bent, 2008; Cainer et al., 2008; Clarke & Garrett, 2004; Davis et al., 2005; Hervais-Adelman et al., 2008; Mehler et al., 1993; Pallier et al., 1998; Zaballos, Plasencia, Gonzalez, de Miguel, & Macias, 2016). Where adaptation to multiple forms of distortion has been investigated, the study either used a between-group design (Davis & Johnsrude, 2003) or a within-group design without reporting any transfer of learning effects (Peelle & Wingfield, 2005). Yet, a between-group design does not provide any insight into whether individuals are capable of perceptually learning equally to multiple distortions or are more skilled at perceiving certain distortions more than others whilst a within-group design without reporting transfer of learning (i.e. order) effects assumes that the adaptation process for the different types of distortion is independent of each other. Thus far only one study has investigated whether the order of presentation influences the degree of subsequent perception to another speech distortion. Adank and Janse (2009) investigated the degree to which exposure and adaptation to one form of compressed speech transferred to the learning of a secondary form of compression. Participants were either presented with a block of artificially time-compressed sentences followed by naturally fast spoken sentences or vice versa. Adapting to the easier, artificially compressed sentences before the more difficult, naturally fast sentences gave participants an advantage over those who adapted to the natural speech before the artificially compressed sentences. This result indicates that learning transferred from the artificial distortion to the natural distortion, but not vice versa suggesting that adaptation to different types of distortion is not independent of each other. Adank and Janse (2009) argue, with reference to Reverse Hierarchy Theory (Ahissar, Nahum, Nelken, & Hochstein, 2009), that this difference in transfer of learning is due to the fact that the artificially time-

1 compressed sentences posed less of a challenge to the perceptual system and therefore
2 participants were able to process this stimuli at a higher level. For artificially time-
3 compressed sentences, adjustments to the timings of expected word boundaries need to occur
4 (cognitive level changes), whereas for **naturally fast** speech, participants need to adapt to the
5 temporal compression as well as additional spectral variability (cognitive and perceptual
6 changes). Consequently, when subsequently faced with the, more difficult, naturally fast
7 speech sentences, participants were better able to process the higher-level features (e.g.
8 temporal compression) and focus on lower-level distortion-specific cues (e.g. spectral
9 variability), as the cognitive temporal adjustments have already been made, but the spectral
10 variability has not been encountered in the previous artificial compression condition. This
11 relationship fits with RHT's prediction that transfer of learning occurs when an easy
12 condition is followed by a more difficult condition. The lack of transfer from naturally fast to
13 artificially fast sentences is believed to be due to the need to immediately focus attention to
14 lower-level properties for the naturally fast stimuli, resulting in learning of stimulus-specific
15 information that does not transfer easily to alternative stimuli.

16 Only three other studies have specifically investigated whether participants'
17 perceptual profile is consistent across distortions. Bent, Baese-Berk, Borrie, and McKee
18 (2016) investigated recognition of words in phrases across a nonnative accent, a regional
19 dialect and ataxic dysarthric speech within the same group of participants. Results show a
20 significant correlation between performances in the nonnative accented speech condition and
21 both the regional dialect and dysarthric speech conditions, suggesting that individuals who
22 were able to successfully perceive nonnative speech were also more successful in the other
23 conditions. However, no correlation was found between performance in a regional dialect and
24 dysarthric speech conditions. The authors conclude that these results suggest that listeners are
25 not "globally skilled" at perceiving distorted speech. Instead, different individuals can map

1 acoustic-phonetic features found only in certain types of distortions onto words in their
2 mental lexicons. However, in a follow-up study, Borrie, Baese-Berk, Engen, and Bent (2017)
3 investigated the overlap in ability to report words spoken by an individual with dysarthria or
4 words presented in noise. The authors found a significant positive correlation between
5 performances in the two conditions and concluded that similar cognitive-perceptual processes
6 aid comprehension in both conditions, i.e., it appears that participants do possess a global
7 skill that allows them to adapt to a relatively equal level when the speech signal is distorted in
8 an array of forms.

9 Of particular relevance to the current experiment, McLaughlin et al. (2018)
10 investigated if individuals possess a global skill to perceive multiple different distortions.
11 McLaughlin et al. (2018) used four different speech conditions, (1) a native speaker masked
12 in speech shaped noise (energetic masking), (2) a native speaker masked by a single-talker
13 masker (informational masking), (3) a non-native speaker in quiet and (4) a non-native
14 speaker masked by speech-shaped noise. Similarly, to Bent et al. (2016), McLaughlin et al.
15 (2018) found that performance in one condition did not always correlate with performance in
16 all other conditions, but instead performance appeared to correlate for conditions with shared
17 characteristics. For instance, conditions with energetic masking correlated with each other
18 and conditions with a non-native speaker correlated whilst the informational and energetic
19 masker conditions did not correlate. In addition to the previous research however,
20 McLaughlin et al. (2018) also investigated the underlying cognitive mechanism supporting
21 perception in each speech condition and found that greater receptive vocabulary performance
22 was linked to better performance in all conditions. Furthermore, their measure of working
23 memory positively predicted performance for the non-native accented speech. McLaughlin et
24 al. (2018) suggest that vocabulary knowledge may act as a global predictor of individual

1 differences in the perception of any form of difficult speech, while other measures e.g.
2 working memory may be engaged only in certain listening environments.

3 Thus, only Adank and Janse (2009), Bent et al. (2016) Borrie et al. (2017) and
4 McLaughlin et al. (2018) have so far investigated how individuals adapt to different
5 distortions. The aim of this study was to further elucidate individual differences in the
6 successful perception of different speech distortions and to systematically evaluate links
7 between sensory/cognitive abilities and perception of temporal (time-compressed); spectral
8 (noise-vocoded speech) and environmental (speech in noise) distortions. Based on previous
9 research, where participants have been shown to be capable of adapting fully to time-
10 compressed speech in up to 20 sentences (Adank & Devlin, 2010; Dupoux & Green, 1997;
11 Sebastián-Gallés et al., 2000), while adaptation to noise-vocoded and speech in noise can
12 occur within 30 sentences but can take many hours' worth of training before full adaptation
13 occurs (Cainer et al., 2008; Davis et al., 2005; Zaballos et al., 2016), participants in the
14 current experiment were predicted to adapt rapidly to time-compressed speech and slower
15 and less extensively to the noise-vocoded and speech in noise conditions. However, it is
16 unclear whether performance in one condition will equate to a relatively similar performance
17 in all conditions. Research by Borrie et al. (2017) suggests that participants may possess (or
18 lack) a global skill to perceive speech in any difficult listening environment. Moreover,
19 research by Bent et al. (2016) and McLaughlin et al. (2018) suggest that participants may be
20 better at perceiving speech in some – but not all - listening environments. If participants are
21 better at perceiving speech in specific conditions, performance in the noise-vocoded and
22 speech in noise conditions may correlate as perception of these conditions depends on similar
23 changes at the lexical/semantic level. In contrast, a weaker correlation between noise-
24 vocoded/speech in noise and time-compressed speech perception performance would be
25 expected as this condition is dependent to a greater extent on phonological level changes.

1 Alternatively, as Borrie et al. (2017) suggest, **participants may possess (or indeed lack) a**
2 **‘global skill’ for perceiving speech that deviates from the norm** and thus performance
3 across all three manipulations will be correlated.

4

5 **B. Individual differences in perceptual adaptation**

6 A key aim of this research is to systematically evaluate links between a battery of
7 audiological and cognitive tests and successful perception of three distortions. The ability to
8 perceive distorted speech has been related to a range of cognitive factors, yet no
9 comprehensive model currently exists that explains which factors are most important and
10 how these factors interact with the type of adverse condition. It is not known, for example,
11 whether different distortions depend on a common cognitive mechanism or whether different
12 distortions require different mechanisms to underpin adaptation. Thus far, associations
13 between four audiological/cognitive abilities and perception of distorted speech have been
14 investigated most: individual hearing thresholds; working memory; selective
15 attention/inhibition and vocabulary knowledge. The impact of individual differences in
16 hearing ability has predominantly been investigated in older populations where difficulty
17 perceiving speech, especially in the presence of background noise, is a common trait.
18 Although overall hearing thresholds have been associated with poorer overall performance on
19 distorted speech tasks (Adank & Janse, 2010; Akeroyd, 2008; Janse & Adank, 2012),
20 research that adjusts for differences in auditory sensitivities suggests that it is not just the
21 decline of the auditory periphery causing the speech in noise deficit; effective listening also
22 relies upon general cognitive processes (Bilodeau-Mercure, Lortie, Sato, Guitton, &
23 Tremblay, 2015; Golomb et al., 2007; Moore, Peters, & Stone, 1999; Tun, 1998; Tun &
24 Wingfield, 1999; Wong et al., 2009).

1 The Ease of Language Understanding (ELU) model (Rönnberg et al., 2013;
2 Rönnberg, Rudner, Foo, & Lunner, 2008) emphasises the role of working memory capacity
3 in suboptimal conditions where the incoming perceived signal is distorted and does not match
4 any internal phonological representations. Working memory is required to initially keep track
5 of the incoming signal before subsequently assisting in inferring meaning from the
6 incomplete information gained from the distorted signal. In support of this model, working
7 memory has been shown to be an important cognitive mechanism when perceiving speech in
8 noisy environments (Akeroyd, 2008; Rönnberg et al., 2013; Zekveld, Rudner, Johnsrude, &
9 Rönnberg, 2013). Akeroyd (2008) suggests that after hearing thresholds, working memory is
10 the most effective cognitive mechanism in explaining individual differences in performance
11 on tasks requiring perception of speech in noise. This conclusion is in agreement with
12 research investigating perception of an artificial accent (Banks, Gowen, Munro, & Adank,
13 2015; Janse & Adank, 2012). **However, working memory has not consistently been found**
14 **to be a significant predictor of distorted speech perception. For instance,** no relationship
15 was found between individual working memory capabilities and performance requiring
16 perception of foreign-accented (Gordon-Salant, Yeni-Komshian, Fitzgibbons, Cohen, &
17 Waldroup, 2013), frequency compressed (Ellis & Munro, 2013), noise-vocoded (Erb, Henry,
18 Eisner, & Obleser, 2012; Neger, Rietveld, & Janse, 2014), or speech in an array of noise
19 backgrounds (Boebinger et al., 2015). Finally, a meta-analysis from Füllgrabe and Rosen
20 (2016) concluded that for young listeners with normal hearing, differences in working
21 memory account for less than two percent of the variance in speech in noise perception.

22 A similar inconclusive relationship has also been found between individual
23 differences in attention switching/inhibition. For example, Huyck and Johnsrude (2012)
24 simultaneously exposed their participants to noise-vocoded sentences, auditory distractors,
25 and visual distractors. One group of participants were asked to attend to the vocoded

1 sentences, whilst two other groups performed a target detection task for either the visual or
2 auditory distractors. Huyck and Johnsrude (2012) found that in order to effectively learn to
3 perceive the noise-vocoded speech, simple exposure to the stimuli (as occurred in the two
4 distractor groups) is not sufficient, participants must also attend to the noise-vocoded speech.
5 Additionally, attention switching/inhibition has also been linked with greater overall
6 performance for foreign (Tao & Taft, 2017) and novel accented speech (Adank & Janse,
7 2010; Banks et al., 2015) with a mediating effect in the perception of noise-vocoded speech
8 (Erb et al., 2012). However, Bent et al. (2016) found no relationship with foreign-accented or
9 regionally-accented speech; Ellis and Munro (2013) found no relationship with frequency
10 compressed speech; and Boebinger et al. (2015) found no relationship between attention
11 switching/inhibition and speech in noise. Finally, whilst individual differences in vocabulary
12 knowledge have mainly only been investigated for adaptation to accented speech, the results
13 thus far have been more consistent, with greater vocabulary knowledge associated with
14 greater adaptation to accented speech across numerous studies (**Adank & Janse, 2010; Janse**
15 **& Adank, 2012; McLaughlin et al., 2018; Neger et al., 2014**). The current experiment
16 aimed to establish how individual differences in a single battery of audiological and cognitive
17 assessments including pure tone audiometry (PTA), speech recognition thresholds (SRT),
18 working memory, vocabulary knowledge, attention-switching and pattern analysis associate
19 with performance across three distortions in a within-subject design, with particular focus on
20 the degree of overlap or divergence in how each cognitive measure relates to each separate
21 speech condition.

22

23 **C. Summary of research aims**

- 24 1. Determine the extent to which perceptual performance for one condition correlates
25 with performance on other forms of distortion.

- 1 2. Establish whether transfer of learning effects occur between temporal, spectral or
- 2 environmental distortions.
- 3 3. Determine the extent to which individual differences in a battery of audiological and
- 4 cognitive assessments relate to performance for each distortion and whether the
- 5 pattern of associations is consistent across distortions.

6

7 **I. Methods**

8 **A. Participants**

9 Ninety participants took part in this experiment (mean age 21.4 ± 2.74 *SD*; range 18-30; 25

10 males). All participants were native British English speakers, had normal or corrected to

11 normal vision and were highly educated ($15.8\text{yrs} \pm 1.67$). No participants reported a history

12 of speech, language, neurological or psychiatric disorder. **The data for two participants**

13 **were excluded due to their performance on a preliminary test of general cognition**

14 **falling below a standardised cut-off score (both participants scored below 26 on the**

15 **Montreal Cognitive Assessment, Nasreddine et al. 2005).** All participants gave informed

16 consent and were compensated with monetary payment or course credit. The study was

17 approved by the UCL research ethics committee (#0599/001).

18

19 **B. Procedure**

20 Participants underwent audiological and cognitive testing in addition to the main speech

21 adaptation task. All testing was performed in a double-walled soundproof room and lasted up

22 to 90 minutes.

23

24 **C. Audiological assessments**

1 Two audiological assessments were performed: (1) Pure Tone Audiometry (PTA) using a
2 clinical audiometer (Maico, MA 41) with each ear tested separately at octave frequencies
3 between 250 and 8000Hz. For each participant, a PTA (average threshold across all measured
4 frequencies) was computed for both ears. (2) Speech Recognition Threshold (SRT) was used
5 to assess the lowest level at which participants could comprehend 50% of an auditorily
6 presented sentence (Plomp & Mimpen, 1979). Each test started at +20dB and varied
7 systematically thereafter. Each sentence had five key words, if participants repeated three or
8 more of the key words then the SNR value would decrease on the subsequent trial, initially in
9 steps of -10dB and subsequently in steps of -2dB, thus making the following trials harder to
10 perceive. The SNR value decreases until participants were only able to comprehend two or
11 fewer of the key words at which point the SNR value would initially increase in steps of
12 +6dB and subsequently in steps of +2dB. The first six lists of the IEEE Harvard Sentences
13 (IEEE, 1969) were used (60 sentences). On average 36 trials/sentences were required to
14 establish each individual speech recognition threshold. Sentences were presented in the same
15 order to all participants.

16

17 **D. Cognitive assessments**

18 *Handedness* was assessed using the 10 point Edinburgh Handedness Inventory (Oldfield,
19 1971). With scores between 50-100 indicative of right-hand dominance and scores from zero-
20 50 indicative of left-hand dominance.

21

22 *Working Memory* was assessed using a forward digit span task. Participants initially heard a
23 set of three numbers and were asked to repeat them back in the same order, for six lists. If
24 participants correctly recalled five or six lists correctly, then the list size increased to four

1 numbers and so on until more than one list was incorrectly recalled for a list size. The last
2 correctly recalled list size was taken as the working memory threshold.

3

4 *Vocabulary Knowledge* was assessed using the auditory version of the spot-the-word section
5 of the Speed and Capacity of Language Processing (SCOLP) test (Baddeley, Emslie, &
6 Nimmo-Smith, 1993; Baddley, Emslie, & Nimmo-Smith, 1992). Participants were presented
7 with 60 pairs of letter strings and had to indicate which one of the letter strings per pair
8 spelled out a real British English word. Reported scores are number of correct identifications
9 out of 60.

10

11 *Attention-Switching* was assessed using the trail-making test (Batterey, 1944; Tombaugh,
12 2004). This task consists of two parts, in part A, participants must draw a line to connect 25
13 ascending numbers in ascending numerical order (1-2-3-4 etc.) as quickly as possible. In part
14 B, participants have to draw a line to connect 24 circles - 12 of which contain numbers and
15 12 of which contain letters of the alphabet – in an alternating numerical and alphabetic
16 sequence (1-A-2-B-3-C etc.) again they were required to do this as quickly as possible. We
17 took ratio scores of the two parts as the main outcome statistic (part B/part A).

18

19 *Pattern/Rule Analysis* was assessed using the Wisconsin Card Sorting test (WCST; Grant &
20 Berg, 1948). In this test participants are required to sort a deck of 128 cards into stacks
21 depending on how they correspond to one of four reference cards. Each card (playing and
22 reference) contains a symbol of a certain shape, color and size. The participant must sort the
23 cards depending on one of these features. Critically, participants are initially unaware of how
24 the playing cards and reference cards correspond, with the researcher simply informing them
25 whether each placement is correct or incorrect. After 10 correct placements (for example

1 matching ten playing cards in front of the corresponding color matched reference card) the
2 correspondence rule changes and participants must first notice the rule has changed and then
3 find the new rule. Each of the correspondence rules are repeated twice per test (making six
4 rules), the outcome measure reported here is the number of trials required to complete each
5 rule (i.e. two sets of ten correct placements). Perfect performance would be completing this
6 task in 60 trials.

7

8 **E. Speech perception task**

9 *Task:* a computerized version of the SCOLP speed of comprehension test. Participants
10 listened to simple sentences in each of the four conditions outlined below and had to decide
11 whether the sentence was true or false, indicating their response by pressing either the left
12 (true) or right (false) key of a standard PC keyboard. All sentences were clearly true
13 (*'Admirals are people'*) or false (*'Admirals have fins'*). Accuracy and RTs were recorded per
14 trial with adaptation to each condition adjudged via improvements in speed and accuracy of
15 sentence verification. Each set of 48 sentences per condition were retrospectively divided into
16 four sub-blocks of 12 so that the time course of adaptation could be fully assessed.

17

18 *Stimuli:* The auditory sentences were recordings of 192 SCOLP sentences, 96 of which were
19 true and 96 false, with 48 sentences presented per speech condition. All sentences were
20 recorded by four different male speakers of standard Southern British English. At time of
21 recording all speakers were between 30-32 years of age and all were born, raised and
22 educated in South East England (see Adank, Evans, Stuart-Smith, & Scott, 2009 for more
23 details on the recording parameters). Sentences from each speaker were used 12 times per
24 condition with the order of speaker randomized (results related to the effect of multiple
25 speakers will not be discussed here). All sentences were saved to separate files with the

1 beginning and end trimmed to zero crossings as closely as possible to the onset/offset of the
2 initial/final speech sounds, resampled to 22050 Hz, peak normalized to 99% of maximum
3 amplitude and scaled to 70dB SPL using Praat (Boersma & Weenink, 2011). Stimulus
4 presentation was performed using a custom-made MATLAB 2014a program (The
5 MathWorks Inc., Natick, MA, 2000) and Sennheiser headphones, with all stimuli delivered at
6 a comfortable listening level (preset at 74dB SPL but where necessary this was adjusted to fit
7 individual participant preference).

8 Participants' ability to perceive different types of speech was tested using four
9 different conditions. (1) Time compressed sentences shortened to 40% of their original length
10 using PSOLA implemented in Praat (Charpentier & Stella, 1986). (2) Noise-vocoded
11 sentences were filtered into four logarithmically spaced frequency bands from 50 to 5000Hz
12 (50-528; 528-1248; 1248-2541; 2541-5000Hz). (3) Speech in noise sentences were embedded
13 in a stream of speech-shaped noise at a signal to noise ratio of -4dB. The spectrum of the
14 speech-shaped noise was derived from the 192 sentences used in the adaptation task. (4)
15 Speech-in-quiet sentences were presented without any manipulation (beyond the zero
16 trimmings, peak normalization etc. outlined above). The speech-in-quiet sentences condition
17 was always the first condition that all participants heard. **Presenting the speech-in-quiet**
18 **sentences consistently as the first condition** ensured that any task practice effects were
19 overcome before the distorted speech stimuli were encountered. Theoretically, therefore any
20 improvement in task performance for the time-compressed, noise-vocoded and speech in
21 noise conditions comes purely from the participants adapting to the specific manipulation and
22 was not due to greater familiarity with the task. Order of presentation of the three non-speech
23 in quiet conditions was fully randomized between participants. Half of the participants were
24 tested on the audiological and cognitive measures first followed by the adaptation task; the
25 other half of participants had the opposite order to ensure results were not due to fatigue. No

1 significant effect of this procedural manipulation was found ($p > .25$) and therefore in all
2 subsequent analyses the data are collapsed across this variable.

3

4 **II. Results**

5 Differences in performance on the speech perception task across the different speech
6 conditions and sub-blocks of the experiment were assessed using two separate mixed analysis
7 of variance (mixed ANOVA) performed in SPSS 25. One ANOVA was for accuracy and one
8 was for response time (RT) data with condition (speech-in-quiet, time-compressed, noise-
9 vocoded, speech in noise) and sub-block (each set of 12 sentences) as within-subjects factors
10 and order of distortion presentation (speech-in-quiet (Q)-speech in noise (N)-time-
11 compressed (T)-noise-vocoded (V); Q-N-V-T; Q-T-N-V; Q-T-V-N; Q-V-N-T; Q-V-T-N) as
12 a between-subjects factor. For the accuracy descriptive statistics and figures, raw percentage
13 scores are reported, however for the mixed ANOVA the data were transformed to
14 rationalized arcsine units (RAU) to ensure consistent variance over the range of scores
15 obtained (Studebaker, 1985).

16

17 **A. Accuracy**

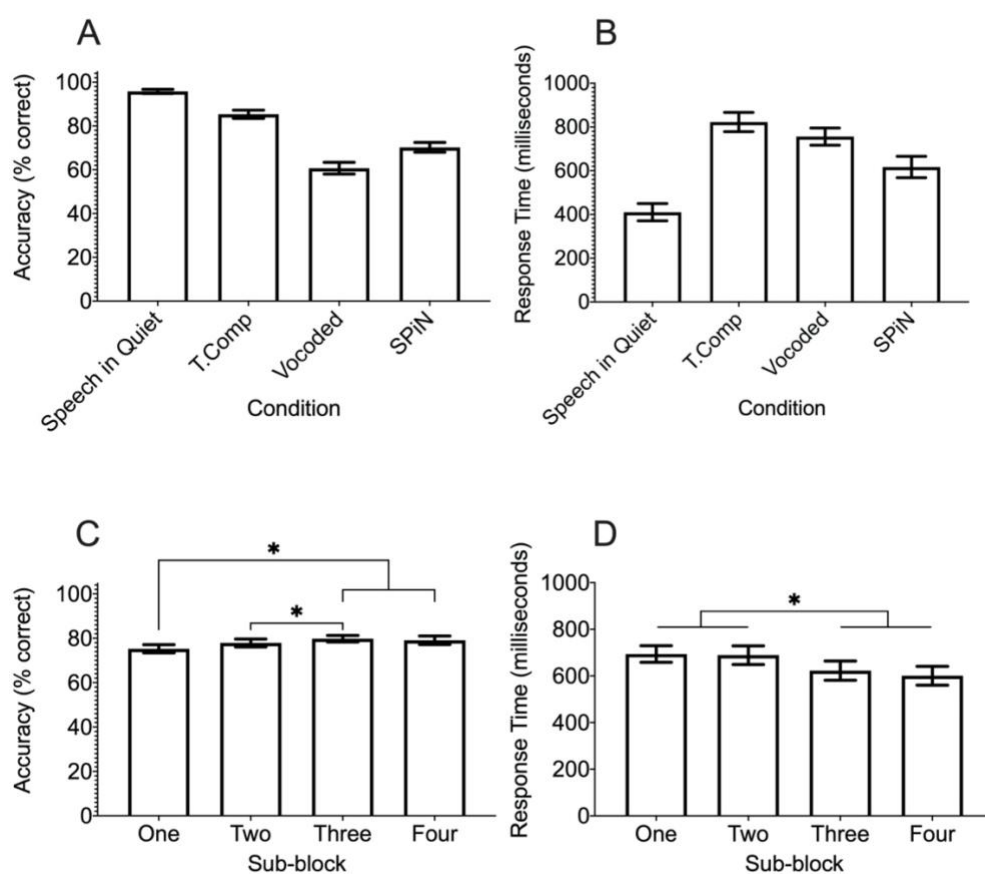
18 Accuracy was highest for speech in quiet ($M = 95.8$, $SD = 4.45$), followed by the time-
19 compressed ($M = 85.3$, $SD = 8.86$), speech in noise ($M = 70.2$, $SD = 10.7$) and noise-vocoded
20 speech ($M = 60.7$, $SD = 12.6$). The mixed ANOVA found a significant main effect of
21 condition $F(3, 246) = 527$, $p < .001$, $\eta_p^2 = .865$ and post-hoc paired samples t -tests revealed
22 significant differences in percentage correct between all four speech conditions at a
23 Bonferroni corrected alpha level ($p = .05/6 = .008$; see Fig. 1 and Supplementary table I₁). In
24 addition, the mixed ANOVA also revealed a significant main effect of sub-block
25 $F(3, 246) = 11.3$, $p < .001$, $\eta_p^2 = .121$, reflecting an improvement in performance from the first to

1 the last sub-block (see Fig. 1 and Supplementary table Ii). Post-hoc paired samples *t*-tests
2 revealed significant differences in the average percentage correct for the first 12 sentences in
3 a block of 48 compared to the third set of 12 (sentences 25-36) $t(87) = -5.52, p < .001$, Cohen's
4 $d = -0.59$; a significant difference between the first and final sub-block of 12 sentences
5 $t(87) = -4.03, p < .001$, Cohen's $d = -0.43$; and a significant difference between the average
6 percent correct for the second and third sub-blocks $t(87) = -3.68, p < .001$, Cohen's $d = -0.39$.
7 The comparisons between the first and second sub-blocks $t(87) = -2.1, p = .038$, Cohen's $d =$
8 -0.22 and the second and fourth sub-blocks $t(87) = -2, p = .049$, Cohen's $d = -0.21$ did not
9 survive Bonferroni correction ($p = .05/6 = .008$).

10 The mixed ANOVA revealed a significant three-way interaction of condition, sub-
11 block and order $F(45, 738) = 372, p < .001, \eta_p^2 = .108$, reflecting differential rates of adaptation
12 between conditions depending on the order in which the participants were exposed to the
13 different distortions. This interaction was investigated with separate two-way repeated
14 measures ANOVAs for each order. This follow-up analysis found a significant condition by
15 sub-block interaction for the Q-N-T-V order $F(9, 117) = 4.3, p < .001, \eta_p^2 = .249$ with post-hoc
16 paired samples *t*-tests showing a significant difference between the first and third sub-blocks
17 in the time-compressed condition $t(13) = -4.28, p = .001$, Cohen's $d = 1.14$ and between the first
18 and last sub-blocks for the noise condition $t(13) = -4.17, p = .001$, Cohen's $d = 1.12$. All
19 remaining comparisons had a significance level greater than the Bonferroni-corrected alpha
20 level ($.05/24 = .002$). In addition, a significant condition \times sub-block interaction was found in
21 the Q-T-N-V order $F(9, 117) = 1.99, p = .046, \eta_p^2 = .13$, but no follow-up comparison survived
22 Bonferroni correction (Supplementary Fig. 1a2).

23 The mixed ANOVA also revealed a significant condition \times order interaction. When
24 investigating the main effect of condition separately for each order, the speech-in-quiet, time-
25 compressed and speech in noise conditions all differed significantly from each other

1 irrespective of their position in the order. Whilst performance in the noise-vocoded condition
 2 always differed significantly from the speech-in-quiet and time-compressed condition
 3 (always significantly poorer), differences between the noise-vocoded and speech in noise
 4 conditions only occurred when the noise-vocoded condition was second in the order (Q-V-N-
 5 T and Q-V-T-N). In both cases, performance was significantly poorer in the noise-vocoded
 6 relative to speech in noise condition. This interaction suggests that participants found the
 7 noise-vocoded condition especially difficult when they had not yet encountered any of the
 8 other distortions (as speech in quiet always came first in the order). The interaction between
 9 speech condition and sub-block approached significance $F(9, 738)=1.87, p=.054, \eta_p^2 = .022$.
 10 All other main effects and interactions were non-significant (p 's $>.05$).



11
 12 **FIG. 1.** Average accuracy and response time data for each condition (A, B) and sub-block (C,
 13 D). Errors bars represent 95% confidence intervals of the mean. For A and C all comparisons
 14 are significantly different at Bonferroni corrected alpha level of .008. For C and D a *
 15 represents all comparisons that are significantly different at the Bonferroni corrected alpha
 16 level of .008.

1

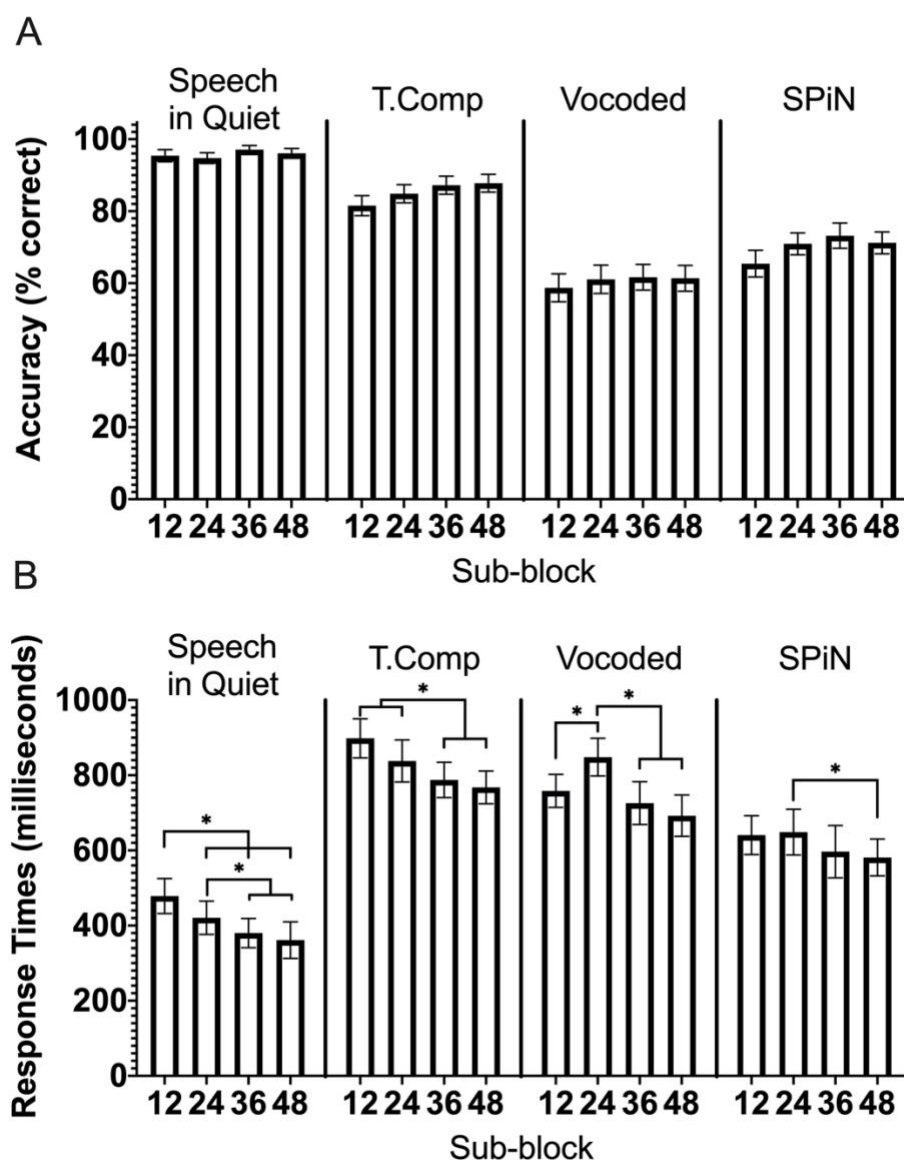
2 **B. Response times**

3 Response times (RTs) in milliseconds were analysed for correct responses only. RTs were
4 measured relative to the end of each sentence, therefore allowing for negative RTs. Overall
5 participants were quickest to respond in the speech in quiet condition ($M = 410$, $SD = 186$)
6 followed by speech in noise ($M = 617$, $SD = 230$), then noise-vocoded ($M = 756$, $SD = 185$)
7 and finally the slowest overall RTs were in response to time-compressed speech ($M = 823$,
8 $SD = 206$). The mixed ANOVA revealed a significant main effect of condition $F(2.69,$
9 $220)=200$, $p<.001$, $\eta_p^2=.709$, post-hoc comparisons revealed significant differences in RTs
10 between all four conditions at the corrected alpha level ($p =.008$; Fig. 1, Supplementary table
11 I₁). Additionally, the mixed ANOVA revealed a significant effect of sub-block $F(2.49,$
12 $246)=31.5$, $p<.001$, $\eta_p^2=.278$, post-hoc paired samples t -tests revealed significant RT
13 differences between the first two and final two sub-blocks. The comparison between the third
14 and fourth sub-blocks $t(87)=-2.03$, $p=.046$, Cohen's $d = -0.22$ did not survive Bonferroni
15 correction ($p = .05/6 = .008$), no difference was found between the first and second sub-
16 blocks (Fig. 1, Supplementary table I₁).

17 The mixed ANOVA result revealed a significant condition \times sub-block \times order
18 interaction $F(36.3, 596)=1.48$, $p=.037$, $\eta_p^2=.083$. When investigating the condition by sub-
19 block interaction separately for each order a significant two-way interaction was only found
20 for the Q-N-V-T order $F(9,126)=3.15$, $p=.002$, $\eta_p^2=.184$, within this order, the only post-hoc
21 paired samples t -test to survive Bonferroni correction was the comparison between the first
22 and third sub-block for the speech in quiet condition $t(14)=4.8$, $p<.001$, Cohen's $d = 1.24$,
23 reflecting the quicker response of participants in the third set of 12 sentences relative to the
24 first set of 12 sentences in the speech in quiet condition (Supplementary Fig. 1b₃).

1 The mixed ANOVA also revealed a significant condition \times sub-block interaction,
2 $F(7.27, 596)=3.61, p=.001, \eta_p^2 = .042$. When analysing the four conditions in separate RM
3 ANOVAs, a significant effect of sub-block was found for all conditions: speech in quiet
4 $F(2.48, 216)=17.8, p<.001, \eta_p^2 = .17$; time-compressed speech $F(2.65, 230)=17.8, p<.001,$
5 $\eta_p^2 =.170$; noise-vocoded $F(3, 261)=12.5, p=.001, \eta_p^2 =.125$; and speech in noise $F(3,$
6 $261)=3.14, p=.026, \eta_p^2 =.035$. In the speech in quiet condition, a significant difference was
7 found in RTs between all sub-blocks except the third and fourth where reductions in RTs
8 levelled off. In the time compressed condition, the first two sub-blocks differed significantly
9 from the final two sub-blocks and for the speech in noise condition a significant difference
10 was found between the second and fourth sub-block. In all condition's participants became
11 quicker to make a correct response as the number of trials increased. For the noise-vocoded
12 speech, however, a significant difference was observed between the second sub-block and all
13 other sub-blocks, this difference was due to participants on average becoming slower to
14 respond in the second sub-block of 12 sentences compared to the other sub-blocks (see Fig.
15 2). All other main effects and interactions were non-significant (p 's $>.05$).

16 In summary, overall adaptation was most evident in the response time data where
17 participants became quicker to make a correct response as the number of trials increased for
18 the speech-in-quiet, time-compressed and speech in noise conditions. Whilst, transfer of
19 learning effects were most noticeable in the accuracy data where participants found the noise-
20 vocoded condition especially difficult when they had not yet encountered any of the other
21 distortions.



1

2 **FIG. 2. A.** Accuracy of responses for each condition across the four sub-blocks per condition.
 3 **B.** Average RTs for each condition across the four sub-blocks per condition. Error bars
 4 represent 95% confidence intervals of the mean. All significant differences are represented by
 5 a *. T.Comp = Time compressed, Vocoded = Noise Vocoded, SPiN = Speech in Noise.

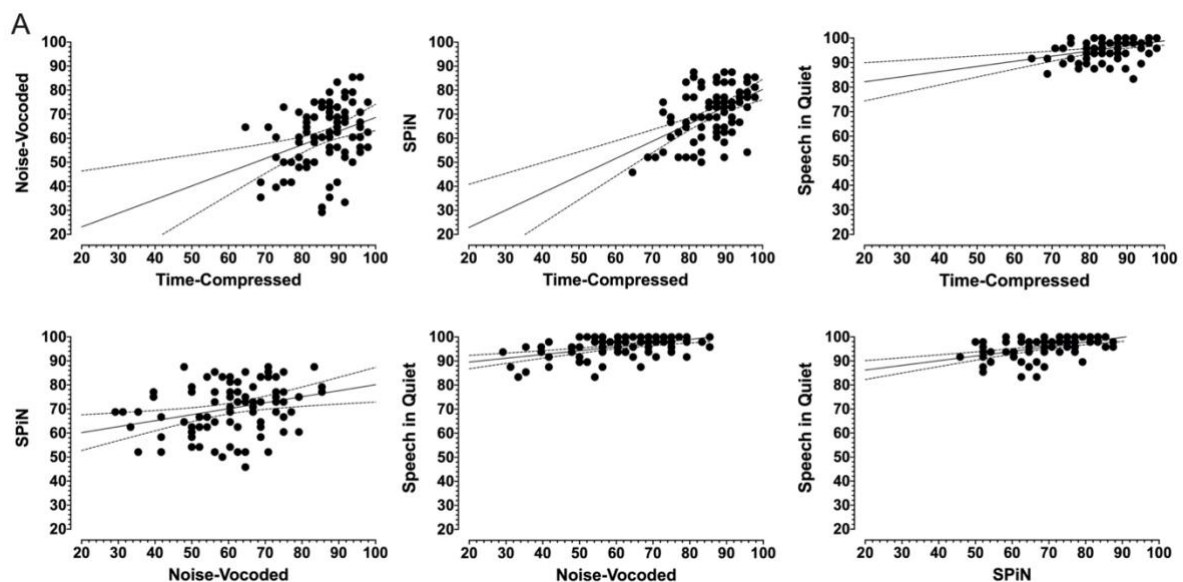
6

7 **C. Correlation in performance between conditions**

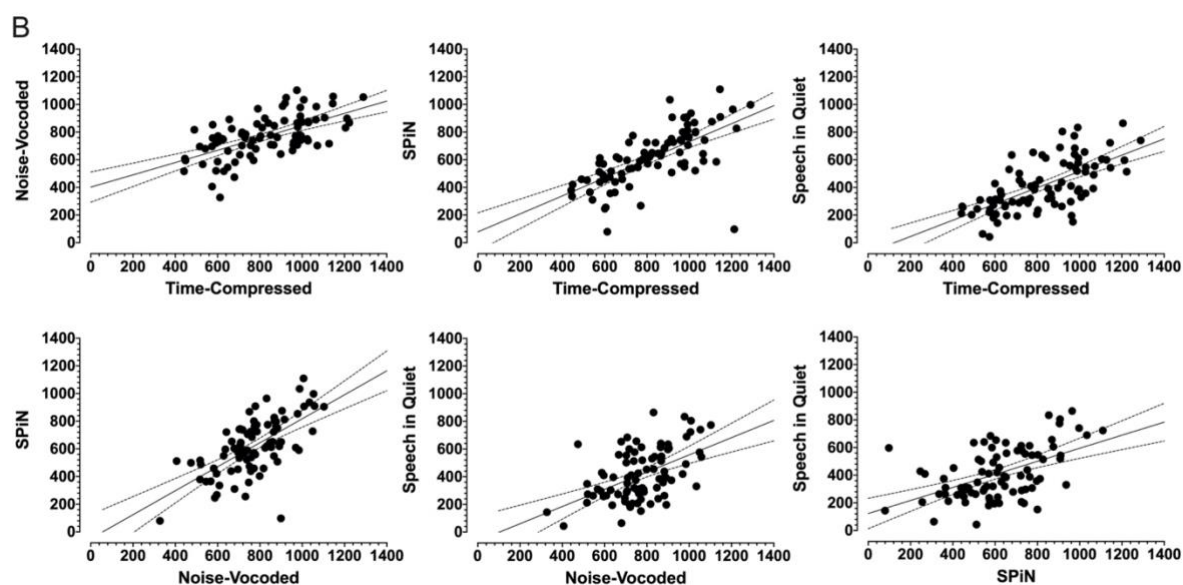
8 Results for both accuracy and RT data reveal significant positive Pearson correlations
 9 between all conditions ($p < .008$; Fig. 3, Supplementary table II₄). This result suggests that
 10 participants appear to possess (or lack) a general ability to perform relatively equally in

1 different adverse listening conditions, irrespective of the type of distortion (spectral, temporal
 2 or environmental).

3



5



6

7 **FIG. 3.** Scatterplot matrices representing the correlations between overall accuracy (A) and
 8 RTs (B) for each participant across the 4 speech conditions. Dotted lines represent 95%
 9 confidence intervals of the regression line. All subplots represent significant positive
 10 correlations.

10

1 **D. Relationship between cognitive assessments and performance in different speech**
 2 **conditions**

3 Overall performance on each of the audiological and cognitive assessments was high, as
 4 would be expected from a homogenous young, highly educated population of participants
 5 (Table I). Prior to establishing the relationship between individual differences in
 6 audiological/cognitive ability and performance on each of the speech distortions a Principal
 7 Component Analysis (PCA) was run on the audiological and cognitive assessment scores to
 8 establish correlations between variables. Initially, average PTA thresholds for both ears,
 9 SRTs, working memory, vocabulary knowledge, attention-switching (trail-making test ratio)
 10 and pattern analysis (number of trials required to complete all rules on WCST) were included
 11 in the analysis. Scores on the SRT, attention-switching or pattern analysis tests showed no
 12 correlation coefficients $>.3$ and were therefore excluded from the subsequent PCA. Bartlett's
 13 test of sphericity was found to be significant indicating sufficient between variables
 14 correlations in the remaining measures to be suitable for a PCA ($X^2(21) = 54.6, p < .001$).

15 **Table I.** Descriptive statistics for battery of audiological and cognitive tests.

Assessment	N	<i>M</i>	<i>SD</i>	Min.	Max.
PTA Left Ear	88	3.90	4.74	-5	25
PTA Right Ear	88	5.30	5.58	-2.50	35
SRT	87	-3.78	1.11	-7	0
Working Memory	88	6.31	1.22	4	9
Vocabulary Knowledge	88	50.1	3.77	42	58
Attention Switching	88	2.14	0.72	1.01	5.65
Pattern Analysis	88	77	11.7	63	129

16

17 Two PCA components showed eigenvalues >1 and together explained 70.5% of all
 18 variance. A Varimax orthogonal rotation was employed to aid interpretability. Component 1
 19 loads most strongly onto vocabulary knowledge and working memory. Component 2 reflects
 20 general hearing ability loading most strongly onto the PTA thresholds of the two ears.

1 **Table II.** Summary of Principal component analysis loadings. Bolded font represents the
 2 strongest loadings for each item.

Items	Rotated Component Coefficients	
	1	2
Vocabulary Knowledge	.834	-.109
Working Memory	.839	.077
PTA Left Ear	.077	.844
PTA Right Ear	-.180	.811

3

4 Multiple linear regressions were conducted to establish if performance in each speech
 5 distortion condition could be predicted based on individual differences in both of the two
 6 components and the SRT, attention switching and pattern analysis tasks (separate analyses
 7 were conducted for accuracy and response time data). The regression model significantly
 8 predicted overall accuracy performance in the speech-in-quiet, time-compressed and speech
 9 in noise conditions (Table III) but not in the noise-vocoded condition ($p=.541$). Across all
 10 three significant models, component 1 significantly predicted performance, suggesting a link
 11 between vocabulary knowledge, working memory and performance in the different
 12 distortions.

13

14

1 **Table III.** Summary of multiple regression analyses for RAU accuracy data across conditions.

Clear: $F(5,79)=3.42, p=.008, \text{adj. } R_2=.126$				
Variable	B	SE_B	β	p
Intercept	108	7.75		<.001
Component One	2.01	0.96	.226	.038
Component Two	-0.6	0.93	-.066	.525
SRT	-1.86	0.88	-.232	.037
Attention Switching	-0.52	1.3	-.042	.69
Pattern Analysis	-0.09	0.08	-.112	.304

Time-Compressed: $F(5,79)=4.13, p=.002, \text{adj. } R_2=.157$				
Variable	B	SE_B	B	p
Intercept	85.4	9.18		<.001
Component One	3.51	1.13	.327	.003
Component Two	-0.3	1.1	-.027	.789
SRT	-1.73	1.04	-.179	.1
Attention Switching	1.78	1.54	.119	.25
Pattern Analysis	-0.08	0.1	-.085	.426

Noise-Vocoded: $F(5,79)=0.82, p=.541, \text{adj. } R_2=-.011$				
Variable	B	SE_B	B	p
Intercept	52.7	11.8		<.001
Component One	0.97	1.46	.077	.509
Component Two	-0.02	1.42	-.002	.988
SRT	-2.07	1.34	-.181	.127
Attention Switching	0.5	1.99	.028	.802
Pattern Analysis	-0.02	0.13	-.014	.904

SPiN: $F(5,79)=2.64, p=.029, \text{adj. } R_2=.089$				
Variable	B	SE_B	B	p
Intercept	61.3	10.3		<.001
Component One	3.59	1.27	.311	.006
Component Two	-0.29	1.24	-.025	.813
SRT	-1.09	1.16	-.105	.351
Attention Switching	2.38	1.72	.148	.172
Pattern Analysis	-.001	0.11	-.001	.994

Notes: B = standardised regression coefficient; SE_B = Standard error of the coefficient; β = standardised coefficient. All p -values highlighted in bold are significant at an alpha level of <0.05

2

3 For the RT data, the multiple linear regression model was non-significant for all of the

4 distorted speech conditions (all p 's>.05; Table IV).

1

2 **Table IV.** Summary of multiple regression analyses for response time data across conditions.

Clear: $F(5,79)=0.73$, $p=.607$, adj. $R_2=-.017$

Variable	B	SE_B	β	p
Intercept	389	172		<.001
Component One	-10.4	21.3	-.056	.626
Component Two	17.8	20.5	.098	.388
SRT	13.03	19.4	.079	.503
Attention Switching	-17	29.1	-.066	.561
Pattern Analysis	1.46	1.82	.093	.424

Time-Compressed: $F(5,79)=2.25$, $p=.057$, adj. $R_2=.069$

Variable	B	SE_B	β	p
Intercept	903	182		<.001
Component One	-22.4	22.5	-.11	.323
Component Two	53.7	21.7	.266	.016
SRT	30.1	20.5	.165	.146
Attention Switching	13.3	30.8	.047	.667
Pattern Analysis	0.15	1.92	.008	.939

Noise-Vocoded: $F(5,79)=0.9$, $p=.486$, adj. $R_2=-.006$

Variable	B	SE_B	β	p
Intercept	827	141		<.001
Component One	17.4	17.4	.142	.22
Component Two	18.5	16.8	.123	.274
SRT	19.6	15.9	.144	.22
Attention Switching	27.8	23.9	.131	.247
Pattern Analysis	-0.55	1.49	-.042	.715

SPiN: $F(5,79)=.868$, $p=.507$, adj. $R_2=-.008$

Variable	B	SE_B	β	p
Intercept	871	190		<.001
Component One	5.65	23.4	.028	.81
Component Two	32.4	22.6	.161	.156
SRT	27.9	21.3	.153	.194
Attention Switching	15.2	32	.053	.637
Pattern Analysis	-2.33	2	-.135	.248

Notes: B = standardised regression coefficient; SE_B = Standard error of the coefficient; β = standardised coefficient. All p -values highlighted in bold are significant at an alpha level of <0.05

3

4III. Discussion

5 The present study investigated whether participants show a consistent perceptual profile

6 across three speech distortions: time-compressed, noise-vocoded and speech in noise.

1 Additionally, we investigated whether/how individual differences in performance on a battery
2 of audiological and cognitive tasks links to perceptual performance to uncover the underlying
3 cognitive mechanisms that underpin the perceptual process.

4 **A. Speech perception in different listening conditions**

5 Speech perception performance was highest in the speech in quiet, followed by time-
6 compressed, speech in noise, and noise-vocoded speech condition. Second, individual
7 participants' performance in one condition was highly predictable from their performance in
8 the three other conditions. This supports the notion that participants possess (or lack) a
9 general ability to successfully perceive multiple forms of speech, even when the distortions
10 differ in the degree of spectral, temporal or environmental manipulation. It should be noted
11 that the significant correlations in both the current experiment and those of Bent et al. (2016)
12 and Borrie et al. (2017) do not imply a causal link and thus do not definitively illustrate a
13 common cognitive mechanism underpinning general auditory perceptual adaptation to any
14 form of distorted speech. It is possible, even likely, that listeners can reach a similar level of
15 perception for *different* distortions using *different* cognitive strategies with individuals who
16 are skilled on one strategy also being more skilled on other strategies. As highlighted
17 previously, we think the perceptual system is stressed differently by the speech distortions
18 used in this experiment. The temporal distortion likely resulted in changes at phonological
19 processing levels (Sebastián-Gallés et al., 2000), whereas the spectral and environmental
20 distortions likely resulted in changes at lexical/semantic levels (Davis et al., 2005). Future
21 research could investigate this further by using functional imaging to establish whether the
22 neural patterns that occur during the perceptual/adaptation processes are similar or different
23 across distortions, **this research** may elucidate whether one common mechanism underpins
24 this process or whether multiple mechanisms are required/responsible (Adank, Davis, &
25 Hagoort, 2012; Adank & Devlin, 2010). Such work could consequently inform research using

1 neuromodulatory techniques such as Transcranial Magnetic Stimulation, which could be used
2 to demonstrate causality between the observed neural activation and the associated cognitive
3 mechanisms during tasks requiring perceptual adaptation to distorted speech.

4 In the current experiment, individual differences in accuracy performance in the
5 speech in quiet, time-compressed and speech in noise conditions were all associated with
6 individual differences in performance on tests of working memory and vocabulary
7 knowledge. These results are strikingly similar to the results of McLaughlin et al. (2018) who
8 also found a link between (receptive) vocabulary knowledge, working memory, and
9 performance on a speech perception task, using accented speech. The association between
10 greater vocabulary knowledge and perception of accented speech has been shown previously
11 (Adank & Janse, 2010; Banks et al., 2015; Bent et al., 2016). The findings from the present
12 experiment, extend those from the accented speech literature to show that individual
13 differences in vocabulary knowledge and working memory are also associated with
14 individual differences in the ability to perceive speech in quiet and artificially time-
15 compressed speech. McLaughlin et al. (2018) suggest that individuals with greater receptive
16 vocabularies are able to perceive distorted speech to a greater extent “*because they have*
17 *stronger lexical mappings that allow them to access semantic representations from input even*
18 *when it is environmentally degraded*” (McLaughlin et al, 2018, page 1567). It is expected that
19 the increased working memory capability can be used to retain larger chunks of the degraded
20 speech signal for longer. This retention of larger chunks of information in turn allows the
21 listener to analyse and compare the stimulus with pre-existing speech templates. Thus, a
22 larger vocabulary may increase the chance of correctly identifying words and/or phrases in
23 the distorted stimuli, with improved speech perception as a result.

24 Future research could build on the current research to investigate which other
25 mechanisms, if any, enable successful perception of speech in different adverse listening

1 conditions in an attempt to build a more comprehensive model of the supporting cognitive
2 and neural mechanisms. One potential mechanism to explore is statistical learning. Research
3 in infant language learning suggests that statistical learning is critical for early language
4 acquisition and development (Romberg & Saffran, 2010; Saffran, 2003). Similar statistical
5 learning abilities may be utilised in older children and adults when faced with the challenge
6 of adapting to distorted speech. For example, Neger et al. (2014) found that perceptual
7 learning of distorted speech was modified by statistical learning ability, with participants who
8 showed better performance on a statistical learning task also showing greater perceptual
9 learning of noise-vocoded speech (in their younger group of participants, but not in the older
10 group). Additionally, Neger et al. (2014) found a significant effect of vocabulary knowledge
11 on perceptual learning of noise-vocoded speech. The authors argue that perceptual learning
12 abilities may rely directly on sensitivity to the probabilistic information inherent in all speech
13 (i.e., statistical learning). It is believed that individuals who are more capable, and faster, to
14 identify subunits of the distorted speech signal are able to transfer this information to higher
15 level processors thus facilitating faster access to their larger store of lexical representations
16 and greater overall perception and potential adaptation to the distortion. Future research
17 should extend this research and investigate the role of statistical learning in perception of
18 distorted speech, along with other potential cognitive mechanisms and investigate how each
19 of the different mechanisms interact to underpin successful perception of speech in adverse
20 listening conditions.

21 We found that vocabulary knowledge and working memory did not predict accuracy
22 of performance in the noise-vocoded condition. It is possible that this result is indicative of an
23 alternative mechanism sub serving perception of vocoded speech. Alternatively, the
24 parameters we used to vocode the speech were more stringent than previous studies (Davis et
25 al., 2005; Hervais-Adelman et al., 2012; Hervais-Adelman et al., 2008; Huyck & Johnsrude,

1 2012). Indeed Loizou et al. (1999) noted that adaptation to noise-vocoded speech drops
2 rapidly below five channels. Future research could establish whether perception of noise-
3 vocoded speech is underpinned by a separate cognitive mechanism compared to other speech
4 distortions or whether the lack of a relationship found in the current experiment is the result
5 of the parameters used during vocoding.

6

7 **B. Transfer of learning effects**

8 With reference to Reverse Hierarchy Theory (Ahissar et al., 2009), Adank and Janse
9 (2009) explain their transfer of learning result on the basis of the easier to comprehend
10 artificial condition providing the listener with a training signal (Hervais-Adelman et al.,
11 2008) for the harder naturally fast spoken sentences. In the current experiment, overall
12 performance for the noise-vocoded condition was poorest of all four conditions, however
13 performance improved for the vocoded condition when listeners had encountered either or
14 both time-compressed and speech in noise conditions before the noise-vocoded condition.
15 When encountered third or fourth in the order, participants would have experience of either
16 96 or 144 occurrences of the noun plus predicate structure of the target sentences. The
17 superior knowledge of the sentence structure may have helped the listeners to perform better
18 on the noise-vocoded condition when it occurred later in the order. These advantages would
19 be absent when noise-vocoding occurred immediately after the speech in quiet condition and
20 thus this knowledge could not be transferred to assist in performance on this distortion
21 resulting in poorer performance when noise-vocoding was encountered early in the order.
22 Furthermore, of the three distorted speech conditions, the noise-vocoded condition represents
23 the condition that was designed to sound most different from anything experienced regularly
24 by individuals with hearing in the normal range. Listeners can be expected to hear speech in
25 quiet and fast speech on a daily basis but cannot be expected to hear the specific changes to

1 the spectral composition of the sentences in the noise-vocoded condition. Therefore, it is
2 possible that the poorer performance in the noise-vocoded condition, when presented first,
3 may in part be attributable to attentional effects. Indeed, Floccia, Butler, Goslin, and Ellis
4 (2009) found that the initial perturbation in performance (accuracy and RTs) that is often
5 observed when a stimulus changes from one condition to another e.g., from speech in quiet to
6 noise-vocoded speech, can be dependent on the task instructions. Floccia et al. (2009) found
7 that individuals who were aware that the accent changed during testing - and who were
8 exposed to this accent during training – showed less perturbation when the accent changed.
9 **The reduced perturbation** was the case for those who were not told to specifically focus on
10 the differences between accents only. Future studies could explore whether the effect of
11 providing the participants with 1 or 2 “training” sentences before the test phase of the
12 experiment would reduce any initial attentional effects (Peelle & Wingfield, 2005). In turn,
13 this reduction could then change the nature of any transfer of learning effects, as found in the
14 current experiment.

15

16 **C. Limitations and directions for future research**

17 There are a few limitations of the current experiment that could be addressed in future
18 research. For example, the working memory task that was chosen required participants to
19 listen as a list of numbers were read to them, internally rehearse the list and subsequently
20 repeat the list in the exact (forward) order as the list was presented to them. While this design
21 fits within the definition of working memory as the ‘temporary storage of information in
22 connection with performance of other cognitive tasks’ (Baddeley, 1983), it might be useful
23 for future research to use a working memory task that has been more closely related to the

1 more traditional definition of working memory as an online memory manipulation process,
2 such as an operation span or reading span task (Banks et al., 2015).

3 Future studies could consider task-related effects on performance. We used a speeded
4 sentence-verification task with accuracy and RTs of responses as a measure of perception
5 differences between conditions. While accuracy performance was above chance, sentence
6 verification is a somewhat crude measure of successful perception of speech, especially in
7 difficult listening environments. A task that requires word identification and recall, e.g., a
8 transcription task, might have been better. First, a transcription task requires participants to
9 explicitly recognise, recall, and reproduce each word of the sentence, which is arguably
10 associated with more focused and extensive linguistic processing. Second, this task could also
11 provide a more fine-grained accuracy measure. However, the downside of such a task would
12 be that it would not be feasible to measure RTs, so the speed of processing could not be
13 assessed. Nevertheless, as it is still possible that the type of task interacts with processing of
14 different types of distortions of the speech signal, different tasks should be considered in
15 future research.

16

17 IV. Conclusions

18 In conclusion, results from the current experiment suggest that individual listeners possess a
19 general skill to adapt to various speech distortions. This skill allows participants to perform to
20 a relatively equivalent level across different conditions irrespective of whether the distortion
21 is temporal, spectral or environmental in nature. The results of this experiment suggest that
22 measures of vocabulary knowledge and working memory could underpin the perceptual
23 learning process. **Overall, the current research adds to and extends the work of Adank
24 and Janse (2009), Bent et al. (2016), Borrie et al. (2017) and McLaughlin et al. (2018).**

1 **The present work supports previous work by providing a more comprehensive overview**
2 **of the potential underlying cognitive mechanism(s) that underpins perception across a**
3 **range of different speech distortions, an area of the field that is currently under**
4 **researched.** Future research should try to establish whether equivalent perception across
5 distortions is dependent on the use of a singular neural-cognitive-perceptual mechanism in all
6 adverse listening conditions as suggested here or whether different mechanisms are employed
7 dependent on the distortion.

8

1 See supplementary material at [URL will be inserted by AIP] for a summary of the pairwise comparison statistics for accuracy and response times.

2 See supplementary material at [URL will be inserted by AIP] for a figure showing the average accuracy data for each condition across the four sub-blocks and each order of presentation.

3 See supplementary material at [URL will be inserted by AIP] for a figure showing the average response time for each condition across the four sub-blocks and each order of presentation.

4 See supplementary material at [URL will be inserted by AIP] for a table of the Pearson correlation coefficients for average performance between conditions for RAU accuracy and response time data.

.

9 **References**

- 10 Adank, P., Davis, M. H., & Hagoort, P. (2012). Neural dissociation in processing noise and
11 accent in spoken language comprehension. *Neuropsychologia*, *50*(1), 77-84.
12 doi:<https://doi.org/10.1016/j.neuropsychologia.2011.10.024>
- 13 Adank, P., & Devlin, J. T. (2010). On-line plasticity in spoken sentence comprehension:
14 Adapting to time-compressed speech. *NeuroImage*, *49*(1), 1124-1132.
15 doi:<https://doi.org/10.1016/j.neuroimage.2009.07.032>
- 16 Adank, P., Evans, B. G., Stuart-Smith, J., & Scott, S. K. (2009). Comprehension of familiar
17 and unfamiliar native accents under adverse listening conditions. *J Exp Psychol Hum*
18 *Percept Perform*, *35*(2), 520-529. doi:10.1037/a0013552
- 19 Adank, P., & Janse, E. (2009). Perceptual learning of time-compressed and natural fast speech.
20 *J Acoust Soc Am*, *126*(5), 2649-2659. doi:10.1121/1.3216914

- 1 Adank, P., & Janse, E. (2010). Comprehension of a novel accent by young and older listeners.
2 *Psychol Aging*, 25(3), 736-740. doi:10.1037/a0020054
- 3 Ahissar, M., Nahum, M., Nelken, I., & Hochstein, S. (2009). Reverse hierarchies and sensory
4 learning. *Philosophical Transactions of the Royal Society B: Biological Sciences*,
5 364(1515), 285-299. doi:10.1098/rstb.2008.0253
- 6 Akeroyd, M. A. (2008). Are individual differences in speech reception related to individual
7 differences in cognitive ability? A survey of twenty experimental studies with normal
8 and hearing-impaired adults. *Int J Audiol*, 47 Suppl 2, S53-71.
9 doi:10.1080/14992020802301142
- 10 Baddeley, A. (1983). Working memory. *Philosophical Transactions of the Royal Society B:*
11 *Biological Sciences*, 302(1110), 311-324. doi:doi:10.1098/rstb.1983.0057
- 12 Baddeley, A., Emslie, H., & Nimmo-Smith, I. (1993). The Spot-the-Word test: a robust
13 estimate of verbal intelligence based on lexical decision. *Br J Clin Psychol*, 32 (Pt 1),
14 55-65.
- 15 Baddley, A., Emslie, H., & Nimmo-Smith, I. (1992). *The Speech and Capacity of Language*
16 *Processing Test manual*. Suffolk, UK: Thames Valley Test Company.
- 17 Banks, B., Gowen, E., Munro, K. J., & Adank, P. (2015). Cognitive predictors of perceptual
18 adaptation to accented speech. *J Acoust Soc Am*, 137(4), 2015-2024.
19 doi:10.1121/1.4916265
- 20 Battery, A. I. T. (1944). *Manual of Directions and Scoring*. Washington D.C., USA: War
21 Department, Adjutant General's Office.
- 22 Bent, T., Baese-Berk, M., Borrie, S. A., & McKee, M. (2016). Individual differences in the
23 perception of regional, nonnative, and disordered speech varieties. *The Journal of the*
24 *Acoustical Society of America*, 140(5), 3775-3786. doi:10.1121/1.4966677
- 25 Bilodeau-Mercure, M., Lortie, C. L., Sato, M., Guitton, M. J., & Tremblay, P. (2015). The
26 neurobiology of speech perception decline in aging. *Brain Struct Funct*, 220(2), 979-
27 997. doi:10.1007/s00429-013-0695-3
- 28 Boebinger, D., Evans, S., Rosen, S., Lima, C. F., Manly, T., & Scott, S. K. (2015). Musicians
29 and non-musicians are equally adept at perceiving masked speech. *J Acoust Soc Am*,
30 137(1), 378-387. doi:10.1121/1.4904537
- 31 Boersma, P., & Weenink, D. (2011). Praat: doing phonetics by computer (Version 5.4.02).
- 32 Borrie, S. A., Baese-Berk, M., Engen, K. V., & Bent, T. (2017). A relationship between
33 processing speech in noise and dysarthric speech. *The Journal of the Acoustical Society*
34 *of America*, 141(6), 4660-4667. doi:10.1121/1.4986746
- 35 Bradlow, A. R., & Alexander, J. A. (2007). Semantic and phonetic enhancements for speech-
36 in-noise recognition by native and non-native listeners. *The Journal of the Acoustical*
37 *Society of America*, 121(4), 2339-2349. doi:10.1121/1.2642103
- 38 Bradlow, A. R., & Bent, T. (2008). Perceptual adaptation to non-native speech. *Cognition*,
39 106(2), 707-729. doi:10.1016/j.cognition.2007.04.005
- 40 Burk, M. H., Humes, L. E., Amos, N. E., & Strauser, L. E. (2006). Effect of training on word-
41 recognition performance in noise for young normal-hearing and older hearing-impaired
42 listeners. *Ear Hear*, 27(3), 263-278. doi:10.1097/01.aud.0000215980.21158.a2
- 43 Cainer, K. E., James, C., & Rajan, R. (2008). Learning speech-in-noise discrimination in adult
44 humans. *Hearing Research*, 238(1), 155-164.
45 doi:<https://doi.org/10.1016/j.heares.2007.10.001>
- 46 Charpentier, F., & Stella, M. (1986, Apr 1986). *Diphone synthesis using an overlap-add*
47 *technique for speech waveforms concatenation*. Paper presented at the ICASSP '86.
48 IEEE International Conference on Acoustics, Speech, and Signal Processing.
- 49 Clarke, C. M., & Garrett, M. F. (2004). Rapid adaptation to foreign-accented English. *J Acoust*
50 *Soc Am*, 116(6), 3647-3658.

- 1 Davis, M. H., & Johnsrude, I. S. (2003). Hierarchical Processing in Spoken Language
2 Comprehension. *The Journal of Neuroscience*, 23(8), 3423-3431.
- 3 Davis, M. H., Johnsrude, I. S., Hervais-Adelman, A. G., Taylor, K., & McGettigan, C. (2005).
4 Lexical Information Drives Perceptual Learning of Distorted Speech: Evidence From
5 the Comprehension of Noise-Vocoded Sentences. *Journal of Experimental*
6 *Psychology: General*, 134(2), 222-241. doi:10.1037/0096-3445.134.2.222
- 7 Dupoux, E., & Green, K. (1997). Perceptual adjustment to highly compressed speech: effects
8 of talker and rate changes. *J Exp Psychol Hum Percept Perform*, 23(3), 914-927.
- 9 Ellis, R. J., & Munro, K. J. (2013). Does cognitive function predict frequency compressed
10 speech recognition in listeners with normal hearing and normal cognition? *Int J Audiol*,
11 52(1), 14-22. doi:10.3109/14992027.2012.721013
- 12 Erb, J., Henry, M. J., Eisner, F., & Obleser, J. (2012). Auditory skills and brain morphology
13 predict individual differences in adaptation to degraded speech. *Neuropsychologia*,
14 50(9), 2154-2164. doi:http://dx.doi.org/10.1016/j.neuropsychologia.2012.05.013
- 15 Fairbanks, G., & Jr., F. K. (1957). Word Intelligibility as a Function of Time Compression. *The*
16 *Journal of the Acoustical Society of America*, 29(5), 636-641. doi:10.1121/1.1908992
- 17 Füllgrabe, C., & Rosen, S. (2016). On The (Un)importance of Working Memory in Speech-in-
18 Noise Processing for Listeners with Normal Hearing Thresholds. *Frontiers in*
19 *Psychology*, 7(1268). doi:10.3389/fpsyg.2016.01268
- 20 Goldstone, R. L. (1998). Perceptual learning. *Annu Rev Psychol*, 49, 585-612.
21 doi:10.1146/annurev.psych.49.1.585
- 22 Golomb, J. D., Peelle, J. E., & Wingfield, A. (2007). Effects of stimulus variability and adult
23 aging on adaptation to time-compressed speech. *J Acoust Soc Am*, 121(3), 1701-1708.
- 24 Gordon-Salant, S., Yeni-Komshian, G. H., Fitzgibbons, P. J., Cohen, J. I., & Waldroup, C.
25 (2013). Recognition of accented and unaccented speech in different maskers by
26 younger and older listeners. *The Journal of the Acoustical Society of America*, 134(1),
27 618-627. doi:10.1121/1.4807817
- 28 Grant, D. A., & Berg, E. (1948). A behavioral analysis of degree of reinforcement and ease of
29 shifting to new responses in a Weigl-type card-sorting problem. *Journal of*
30 *Experimental Psychology*, 38(4), 404-411. doi:10.1037/h0059831
- 31 Hervais-Adelman, A. G., Davis, M. H., Johnsrude, I. S., & Carlyon, R. P. (2008). Perceptual
32 learning of noise vocoded words: Effects of feedback and lexicality. *Journal of*
33 *Experimental Psychology: Human Perception and Performance*, 34(2), 460-474.
34 doi:10.1037/0096-1523.34.2.460
- 35 Hervais-Adelman, A. G., Davis, M. H., Johnsrude, I. S., Taylor, K. J., & Carlyon, R. P. (2011).
36 Generalization of perceptual learning of vocoded speech. *J Exp Psychol Hum Percept*
37 *Perform*, 37(1), 283-295. doi:10.1037/a0020772
- 38 Huyck, J. J., & Johnsrude, I. S. (2012). Rapid perceptual learning of noise-vocoded speech
39 requires attention. *The Journal of the Acoustical Society of America*, 131(3), EL236-
40 EL242. doi:10.1121/1.3685511
- 41 IEEE. (1969). IEEE Recommended Practice for Speech Quality Measurements. *IEEE No 297-*
42 *1969*, 1-24. doi:10.1109/IEEESTD.1969.7405210
- 43 Janse, E., & Adank, P. (2012). Predicting foreign-accent adaptation in older adults. *The*
44 *Quarterly Journal of Experimental Psychology*, 65(8), 1563-1585.
45 doi:10.1080/17470218.2012.658822
- 46 Kennedy-Higgins, D. (2019). *Neural and cognitive mechanisms affecting perceptual*
47 *adaptation to distorted speech*. UCL (University College London),
- 48 Loizou, P. C., Dorman, M., & Tu, Z. (1999). On the number of channels needed to understand
49 speech. *J Acoust Soc Am*, 106(4 Pt 1), 2097-2103.

- 1 Mattys, S. L., Davis, M. H., Bradlow, A. R., & Scott, S. K. (2012). Speech recognition in
2 adverse conditions: A review. *Language and Cognitive Processes*, 27(7-8), 953-978.
3 doi:10.1080/01690965.2012.705006
- 4 Mayo, L. H., Florentine, M., & Buus, S. (1997). Age of Second-Language Acquisition and
5 Perception of Speech in Noise. *Journal of Speech, Language, and Hearing Research*,
6 40(3), 686-693. doi:10.1044/jslhr.4003.686
- 7 McLaughlin, D. J., Baese-Berk, M. M., Bent, T., Borrie, S. A., Van Engen, K. J. J. A.,
8 Perception,, & Psychophysics. (2018). Coping with adversity: Individual differences in
9 the perception of noisy and accented speech. 80(6), 1559-1570. doi:10.3758/s13414-
10 018-1537-4
- 11 Mehler, J., Sebastian, N., Altmann, G., Dupoux, E., Christophe, A., & Pallier, C. (1993).
12 Understanding Compressed Sentences: The Role of Rhythm and Meaning a. *Ann N Y*
13 *Acad Sci*, 682(1), 272-282. doi:10.1111/j.1749-6632.1993.tb22975.x
- 14 Moore, B. C., Peters, R. W., & Stone, M. A. (1999). Benefits of linear amplification and
15 multichannel compression for speech comprehension in backgrounds with spectral and
16 temporal dips. *J Acoust Soc Am*, 105(1), 400-411.
- 17 Nasreddine, Z. S., Phillips, N. A., Bedirian, V., Charbonneau, S., Whitehead, V., Collin, I., . .
18 . Chertkow, H. (2005). The Montreal Cognitive Assessment, MoCA: a brief screening
19 tool for mild cognitive impairment. *J Am Geriatr Soc*, 53(4), 695-699.
20 doi:10.1111/j.1532-5415.2005.53221.x
- 21 Neger, T. M., Rietveld, T., & Janse, E. (2014). Relationship between perceptual learning in
22 speech and statistical learning in younger and older adults. *Frontiers in Human*
23 *Neuroscience*, 8(628). doi:10.3389/fnhum.2014.00628
- 24 Oldfield, R. C. (1971). The assessment and analysis of handedness: The Edinburgh Inventory.
25 *Neuropsychologia*, 9, 97-112.
- 26 Pallier, C., Sebastian-Gallés, N., Dupoux, E., Christophe, A., & Mehler, J. (1998). Perceptual
27 adjustment to time-compressed speech: A cross-linguistic study. *Memory & Cognition*,
28 26(4), 844-851. doi:10.3758/bf03211403
- 29 Peelle, J. E., & Wingfield, A. (2005). Dissociations in perceptual learning revealed by adult
30 age differences in adaptation to time-compressed speech. *J Exp Psychol Hum Percept*
31 *Perform*, 31(6), 1315-1330. doi:10.1037/0096-1523.31.6.1315
- 32 Plomp, R., & Mimpen, A. M. (1979). Speech-reception threshold for sentences as a function
33 of age and noise level. *The Journal of the Acoustical Society of America*, 66(5), 1333-
34 1342. doi:10.1121/1.383554
- 35 Romberg, A. R., & Saffran, J. R. (2010). Statistical learning and language acquisition. 1(6),
36 906-914. doi:10.1002/wcs.78
- 37 Rönnerberg, J., Lunner, T., Zekveld, A., Sörqvist, P., Danielsson, H., Lyxell, B., . . . Rudner, M.
38 (2013). The Ease of Language Understanding (ELU) model: theoretical, empirical, and
39 clinical advances. *Frontiers in Systems Neuroscience*, 7(31).
40 doi:10.3389/fnsys.2013.00031
- 41 Rönnerberg, J., Rudner, M., Foo, C., & Lunner, T. (2008). Cognition counts: a working memory
42 system for ease of language understanding (ELU). *Int J Audiol*, 47 Suppl 2, S99-105.
43 doi:10.1080/14992020802301167
- 44 Saffran, J. R. (2003). Statistical Language Learning: Mechanisms and Constraints. 12(4), 110-
45 114. doi:10.1111/1467-8721.01243
- 46 Samuel, A. G., & Kraljic, T. (2009). Perceptual learning for speech. *Atten Percept Psychophys*,
47 71(6), 1207-1218. doi:10.3758/app.71.6.1207
- 48 Sebastián-Gallés, N., Dupoux, E., Costa, A., & Mehler, J. (2000). Adaptation to time-
49 compressed speech: Phonological determinants. *Perception & Psychophysics*, 62(4),
50 834-842. doi:10.3758/bf03206926

- 1 Shannon, R. V., Zeng, F. G., Kamath, V., Wygonski, J., & Ekelid, M. (1995). Speech
2 recognition with primarily temporal cues. *Science*, *270*(5234), 303-304.
- 3 Song, J. H., Skoe, E., Banai, K., & Kraus, N. (2012). Training to Improve Hearing Speech in
4 Noise: Biological Mechanisms. *Cerebral Cortex*, *22*(5), 1180-1190.
5 doi:10.1093/cercor/bhr196
- 6 Studebaker, G. A. (1985). A "Rationalized" Arcsine Transform. *Journal of Speech, Language,
7 and Hearing Research*, *28*(3), 455-462. doi:doi:10.1044/jshr.2803.455
- 8 Tao, L., & Taft, M. (2017). Influences of Cognitive Processing Capacities on Speech
9 Perception in Young Adults. *Frontiers in Psychology*, *8*(266).
10 doi:10.3389/fpsyg.2017.00266
- 11 Tombaugh, T. N. (2004). Trail Making Test A and B: Normative data stratified by age and
12 education. *Archives of Clinical Neuropsychology*, *19*(2), 203-214.
13 doi:https://doi.org/10.1016/S0887-6177(03)00039-8
- 14 Tun, P. A. (1998). Fast noisy speech: Age differences in processing rapid speech with
15 background noise. *Psychol Aging*, *13*(3), 424-434. doi:10.1037/0882-7974.13.3.424
- 16 Tun, P. A., & Wingfield, A. (1999). One voice too many: adult age differences in language
17 processing with different types of distracting sounds. *J Gerontol B Psychol Sci Soc Sci*,
18 *54*(5), P317-327.
- 19 Voor, J. B., & Miller, J. M. (1965). The effect of practice upon the comprehension of time-
20 compressed speech. *Communications Monographs*, *32*(4), 452-454.
- 21 Watson, C. S. (1980). Time Course of Auditory Perceptual Learning. *Annals of Otology,
22 Rhinology & Laryngology*, *89*(5_suppl), 96-102. doi:10.1177/00034894800890s525
- 23 Wong, P. C., Jin, J. X., Gunasekera, G. M., Abel, R., Lee, E. R., & Dhar, S. (2009). Aging and
24 cortical mechanisms of speech perception in noise. *Neuropsychologia*, *47*(3), 693-703.
25 doi:10.1016/j.neuropsychologia.2008.11.032
- 26 Zaballos, M. T., Plasencia, D. P., Gonzalez, M. L., de Miguel, A. R., & Macias, A. R. (2016).
27 Air traffic controllers' long-term speech-in-noise training effects: A control group
28 study. *Noise Health*, *18*(85), 376-381. doi:10.4103/1463-1741.195804
- 29 Zekveld, A. A., Rudner, M., Johnsrude, I. S., & Rönnberg, J. (2013). The effects of working
30 memory capacity and semantic cues on the intelligibility of speech in noise. *134*(3),
31 2225-2234. doi:10.1121/1.4817926