

# **Improving the understanding and diagnosis of Auditory Processing Disorder (APD) in Children**

Shiran Koifman

A thesis submitted in partial fulfilment of the requirements for the  
degree of Doctor of Philosophy at University College London

March 15 2021



For Yihui Xie



# Acknowledgements

This is where you will normally thank your advisor, colleagues, family and friends, as well as funding and institutional support. In our case, we will give our praises to the people who developed the ideas and tools that allow us to push open science a little step forward by writing plain-text, transparent, and reproducible theses in R Markdown.

We must be grateful to John Gruber for inventing the original version of Markdown, to John MacFarlane for creating Pandoc (<http://pandoc.org>) which converts Markdown to a large number of output formats, and to Yihui Xie for creating `knitr` which introduced R Markdown as a way of embedding code in Markdown documents, and `bookdown` which added tools for technical and longer-form writing.

Special thanks to Chester Ismay, who created the `thesisdown` package that helped many a PhD student write their theses in R Markdown. And a very special tahnks to John McManigle, whose adaption of Sam Evans' adaptation of Keith Gillow's original maths template for writing an Oxford University DPhil thesis in  $\text{\LaTeX}$  provided the template that I adapted for R Markdown.

Finally, profuse thanks to JJ Allaire, the founder and CEO of RStudio, and Hadley Wickham, the mastermind of the tidyverse without whom we'd all just given up and done data science in Python instead. Thanks for making data science easier, more accessible, and more fun for us all.

Ulrik Lyngs  
Linacre College, Oxford  
2 December 2018



# Abstract

This *R Markdown* template is for writing an Oxford University thesis. The template is built using Yihui Xie's `bookdown` package, with heavy inspiration from Chester Ismay's `thesisdown` and the `OxThesis` L<sup>A</sup>T<sub>E</sub>X template (most recently adapted by John McManigle).

This template's sample content include illustrations of how to write a thesis in R Markdown, and largely follows the structure from this R Markdown workshop.

Congratulations for taking a step further into the lands of open, reproducible science by writing your thesis using a tool that allows you to transparently include tables and dynamically generated plots directly from the underlying data. Hip hooray!





# Contents

<b>List of Figures</b>	<b>xiii</b>
<b>List of Tables</b>	<b>xvii</b>
<b>List of Abbreviations</b>	<b>xxi</b>
<b>Introduction</b>	<b>1</b>
Speech-in-noise in children . . . . .	1
APD definition . . . . .	2
Diagnosis . . . . .	2
Binaural and spatial listening in APD . . . . .	2
Summary . . . . .	3
<b>1 Binaural listening: interrupted and alternated speech-in-noise in adults</b>	<b>5</b>
1.1 Influence of distractor type on IM . . . . .	6
1.1.1 Introduction . . . . .	6
1.1.2 Experiment I: speech vs. non-speech distractors . . . . .	6
1.1.3 Experiment II: speech distractors spoken in a familiar vs. unfamiliar language . . . . .	6
1.1.4 General discussion and conclusion . . . . .	8
1.2 Dichotic vs. monotic presentation and the influence of speech material	8
1.2.1 Introduction . . . . .	9
1.2.2 Methods . . . . .	9
1.2.3 Results . . . . .	9
1.2.4 Discussion . . . . .	10
1.2.5 Conclusion . . . . .	10

<b>2</b>	<b>Spatial listening: development and normalisation of a children’s spatialised speech-in-noise test</b>	<b>11</b>
2.1	Introduction . . . . .	12
2.2	Methods . . . . .	12
2.3	Discussion . . . . .	12
2.4	Conclusion . . . . .	12
<b>3</b>	<b>APD study</b>	<b>13</b>
3.1	Introduction . . . . .	14
3.2	Methods . . . . .	16
3.2.1	Participants . . . . .	16
3.2.2	Measurements . . . . .	19
3.2.3	Procedure . . . . .	29
3.2.4	Data Analysis . . . . .	31
3.3	Results . . . . .	33
3.3.1	Standard audiometry . . . . .	33
3.3.2	EHF audiometry . . . . .	37
3.3.3	ST . . . . .	40
3.3.4	LiSNS-UK . . . . .	52
3.3.5	ENVASA . . . . .	56
3.3.6	CELF-RS . . . . .	60
3.3.7	Questionnaires . . . . .	62
3.4	Overall performance . . . . .	66
3.4.1	Unsupervised machine learning (UML) . . . . .	68
3.4.2	Interaction between measures . . . . .	68
3.5	Discussion . . . . .	80
3.5.1	EHF . . . . .	80
3.5.2	ST . . . . .	81
3.5.3	CCC-2 . . . . .	82
3.5.4	ECLiPS . . . . .	82
3.6	Conclusion . . . . .	82
	<b>General discussion</b>	<b>83</b>
	Summary of main findings . . . . .	83
	Conclusion . . . . .	83

## Appendices

<i>Contents</i>	<i>xi</i>
<b>A The First Appendix</b>	<b>87</b>
<b>B The Second Appendix</b>	<b>89</b>
<b>C The Third Appendix</b>	<b>91</b>
<b>Works Cited</b>	<b>95</b>



# List of Figures

2.1	Code chunk syntax . . . . .	12
3.1	Schematic of the ENVASA experimental paradigm (taken from Leech et al., 2009) . . . . .	26
3.2	Standard audiometry: APD participants pure-tone detection thresholds plotted separately for the left and the right ear (black lines). The shaded grey area represents the TD group thresholds range and the white line represents the TD group mean at each frequency. The dashed line represents the threshold criteria of hearing level $\leq 25$ dB HL. . . . .	34
3.3	Standard audiometry: Pure-tone detection thresholds by frequency bands between 0.25 to 8 kHz (A), and averaged thresholds (B). Individual scores are indicated by circles. The boxes show the data interquartile range (25th-75th percentile) and the horizontal line indicate the median (i.e., 50th percentile). Values that fall within 1.5 times the interquartile range are indicated by the whiskers. . . . .	34
3.4	EHF audiometry: Pure-tone detection thresholds for extended high-frequency bands measured in the left and the right ear. The thin black lines represents the individual thresholds in the APD group and the group mean is marked by the bold black line. The shaded grey area represents the TD group threshold range and the white line represents the TD group mean at each frequency. . . . .	38
3.5	EHF audiometry: Boxplots for pure-tone detection thresholds measured at the extended high-frequency bands split by ear and groups (A). Boxplots of the groups averaged PTAs and better-ear BE thresholds are depicted in figure B. Individual scores are indicated by circles. The boxes show the data interquartile range (25th-75th percentile) and the horizontal line indicate the median (i.e., 50th percentile). Values that fall within 1.5 times the interquartile range are indicated by the whiskers. . . . .	38

- 3.6 ST raw data: Frequency of potential outliers with  $\text{LevsPC} \leq 35\%$ . LevsPC denotes the proportion of correct keywords within the final test trials. . . . . 42
- 3.7 ST: Scatterplot and linear regression lines for the listeners SRdTs measured with the ASL (A) and CCRM speech material (B) as a function of age. Corresponding regression coefficients and statistics is provided for TD group only. Red indicates data from the APD group and cyan indicates data from the TD control group. Data for normal hearing adults taken from Chapter 2 is shown in the boxplots as a reference. . . . . 43
- 3.8 ST: Age effect - a comparison between the regression lines slopes fitted for the CCRM (x-axis) and ASL speech material (y-axis). Test conditions are represented by the different symbols. The diagonal line represents an optimal agreement between the speech materials. Observations falling below the line indicate a steeper slope for the ASL material than for the CCRM material. . . . . 47
- 3.9 ST: Boxplots of the listeners age-independent standardised residuals for data measured with the ASL (A) and the CCRM speech material (B). Residuals were calculated separately for each condition and are based on a model prediction for TD group only. The grey area represents the deviance cut-off for abnormal score ( $\text{SD} \pm 1.96$  below and above the TD mean), where about 95% of the normal population is expected to lay within. The dashed line represents the theoretical TD group mean ( $z = 0$ ). Individual scores are indicated by circles. The boxes show the data interquartile range (25th-75th percentile) and the horizontal line indicate the median (i.e., 50th percentile). Values that fall within 1.5 times the interquartile range are indicated by the whiskers. . . . . 49
- 3.10 LiSNS-UK: Age-effect - scatterplot and linear regression lines for SRTs obtained for SSN and the spatialised conditions S0N0 (collocated) and S0N90 (separated) (A) and the derived measure SRM (B) as a function of the listeners age. Corresponding regression coefficients and statistics is provided for TD group only. Red indicates data from the APD group and cyan indicates data from the TD control group. 53

3.11	LiSNS-UK: Boxplots of the listeners age-independent standardised residuals (open circles) for data measured with LiSNS-UK task (A) and the derived measure SRM (B). Residuals were calculated separately for each condition and are based on a model prediction for TD group only. The grey area represents the deviance cut-off for abnormal score ( $SD \pm 1.96$ below and above the TD mean), where about 95% of the normal population is expected to lay within. The dashed line represents the theoretical TD group mean ( $z = 0$ ). The boxes show the data interquartile range (25th-75th percentile) and the horizontal lines indicate the median (i.e., 50th percentile). Values that fall within 1.5 times the interquartile range are indicated by the whiskers. . . . .	56
3.12	ENVASA: Scatterplot and linear regression lines for the listeners' PC (%-correct) as a function of age for single background, dual backgrounds and the combined measure. Red indicates data from the APD group and cyan indicates data from the TD control group. . . . .	59
3.13	ENVASA: Listeners' age-independent standardised residuals for single background, dual backgrounds & the combined measure. Residuals were calculated separately for each condition and are based on a model prediction for TD group only. The grey area represents the deviance cut-off for abnormal score ( $SD \pm 1.96$ below and above the TD mean), where about 95% of the normal population is expected to lay within. The dashed line represents the theoretical TD group mean ( $z = 0$ ). . . . .	61
3.14	CELF-RS: Boxplots for CELF-5 UK Recall Sentences subtest scaled scores by groups. The dashed line represents the norms mean and the grey area indicates the upper and lower limit average performance in the normal population ( $\pm 1$ SD). . . . .	62
3.15	CCC-2 parental reports for the APD (red) and TD group (cyan). (A) Boxplots for scaled scores in the ten sub-scales. (B) Scatterplot for General Communication Composite (GCC) as a function of Social-Interaction Deviance Composite, (SIDC). APD children with diagnosed high-functioning Autism (HF-ASD) are denoted with open circles. APD children with undergoing ASD assessment on the day of testing are marked with open rectangles. The lines indicates the GCC cut-off criteria for typically developing children (TD) SIDC scores indicative of predominantly structural developmental language disorder (DLD) and more social communication deficits (cf. Norbury, 2013). . . . .	64

3.16	ECLiPS parental report scaled scores split by groups and sub-scales.	65
3.17	Overall performance: Abnormal (black cells) and normal (empty cells) performance in the present study test battery of individuals from the APD group (n=20) and the TD group (n=23). Missing data is marked by the grey cells. . . . .	67
3.18	Overall performance: Proportion of abnormal score per measure or task split by groups. . . . .	67
3.19	Switching task PCA: Scatterplot for the input variables as a function of PCA components: PC1.ST vs. PC2.Material (A), PC1.ST vs. PC3.Nz (B). Loadings for ASL conditions are indicated by circles and loadings for CCRM conditions are indicated by rectangles. Filled shapes denotes conditions with speech distractors (Spch) and non-filled shapes denote nonspeech conditions (No-Spch). . . . .	70
3.20	Switching task PCA: Listeners weighted scores split by components and group. . . . .	71
3.21	Switching task PCA: Comparison between PCA weighted scores and calculated measures: (1) ST = mean score across all ST data, (2) Material = $\overline{ASL} - \overline{CCRM}$ , (3) Nz = $\overline{NoSpch} - \overline{Spch}$ . . . . .	72
3.22	Language measures PCA: Listeners weighted scores split by components and group . . . . .	74
3.23	Language measures PCA: Individual scores split by groups for loadings in PC1.Lang as a function of scores for PC2.Lang (A), and PC3.Lang (B). . . . .	75
3.24	Language measures PCA: Comparison between the listeners weighted scores by components, PC1.Lang - PC3.Lang (A), and calculated measures, Lang1 - Lang3 (B). . . . .	75
3.25	Association between predictors and performance in the APD group for the switching task composite (PC1.ST), language composite (PC1.Lang), SRM, standard and EHF pure-tone PTA. Predictors included: 1. APD diagnosis (APD vs. LiD), 2. SPD diagnosis (SPD vs. non-SPD), 3. Regular use of FM-device (FM vs. No FM), 4. History of OME (OME vs. No OME), 5. Pressure equalisation tube history (PET vs. No PET), and 6. Auditory training (Training vs. No training). Individual observations are marked in circles. Observations of children diagnosed with APD are filled in dark blue, and LiD observations are filled in light blue. Significant p-values for independent t-test comparison are marked with asterisk ( $p < 0.05$ ). . . . .	80
C.1	LiSNS-UK: Correlations for listeners SRTs (dB SNR). . . . .	93
C.2	LiSNS-UK: Correlations for listeners age-independent z-scores. . . . .	94



# List of Tables

3.1	APD group demographics and APD-related history background. . .	18
3.2	Summary of the study test battery. . . . .	20
3.3	Experimental design and measurements order. . . . .	30
3.4	Standard audiometry: Descriptives for pure-tone detection thresholds (dB HL) by frequency bands (kHz) and ear split by the two groups.	35
3.5	Standard audiometry: Statistical analysis for the effects of Frequency, Ear and Group and their interaction (6x2x2 factorial design with repeated measures) tested with a robust rank-based method for analysis of nonparametric data using nparLD package (Noguchi et al., 2012). Analysis was based on a f1-ld-f2 design ANOVA-type statistic (ATS) test, whereby f1 refers to an experimental design with a single between-subjects factor (Group) and f2 refers to two within-subjects factors (Frequency and Ear). . . . .	36
3.6	Post-hoc paired comparison t-test for PTA x Group. The test was performed on the fitted LMEM model and included adjusted least-squared-mean for the random intercept (subjects) using lsmeans package (Lenth, 2020). . . . .	36
3.7	EHF audiometry: Descriptive for pure-tone detection thresholds (dB HL) by extended-high frequency bands (kHz) split by ear and group.	39
3.8	EHF audiometry: statistical analysis for the effects of Frequency, Ear and Group as well as their interaction (3x2x2 factorial design with repeated measures) tested with a robust rank-based method for analysis of nonparametric data using nparLD package (Noguchi et al., 2012). Analysis was based on a f1-ld-f2 design ANOVA-type statistic (ATS) test, whereby f1 refers to an experimental design with a single between-subjects factor (Group) and f2 refers to two within-subjects factors (Frequency and Ear). . . . .	39

3.9	EHF audiometry: Statistical analysis for the effects of the listeners calculated measures ( $PTA_{Right}$ , $PTA_{Left}$ , PTA, and BE) and Group as well as their interaction (4x2 factorial design with repeated measures) tested with a robust rank-based method for analysis of nonparametric data using nparLD package (Noguchi et al., 2012). Analysis was based on a f1-ld-f1 design ANOVA-type statistic (ATS) test, whereby f1 refers to an experimental design with a single between-subjects factor (Group) and a single within-subjects factor (Measure). . . . .	40
3.10	ST: Age effect analysis using LMEM for SRdTs measured across condition, speech material and age as fixed factors and a random intercept for subjects. Reference levels: Condition = Quiet-NoAlt, Group = APD, Material = ASL). Note: only data measured with the control group following outliers trimming was included (trimmed TD). . . . .	46
3.11	ST: Age-effect - post-hoc paired comparison t-test for Condition x Material two-way interaction. The test was performed on the fitted LMEM model and included adjusted least-squared-mean for the random intercept (subjects) using lsmeans package (emmeans package; Lenth, 2020). . . . .	47
3.12	ST: Descriptives for standardised residuals (z-scores) calculated for data measured with the ASL and CCRM speech material. . . . .	48
3.13	ST: Statistical analysis for the effects of Group, Material, and Condition as well as their interaction (2x2x5 factorial design with repeated measures) tested with a robust rank-based method for analysis of nonparametric data using nparLD package (Noguchi et al., 2012). Analysis was based on a f2-ld-f1 design ANOVA-type statistic (ATS) test, whereby f2 refers to an experimental design with two between-subjects factors (Group and Material) and f1 refers to a single within-subjects factor (Condition). . . . .	52
3.14	LiSNS-UK: Age effect - LMEM model for SRT with condition (reference level: SSN) and age as fixed factors and a random intercept for subjects. Note: only data measured with the control group following outliers trimming was included. . . . .	54
3.15	LiSNS-UK standard residuals (z-scores) descriptives by group. abnormal: defined as the percentage of abnormal z-score $> 1.96$ (SSN, S0N0, & S0N90) and z-score $< 1.96$ (SRM). . . . .	55
3.16	LiSNS-UK: Group difference - LMEM model for age-independent z-scores with condition and group as fixed factors (reference levels: Condition = SSN, Group = APD) and a random intercept for subjects. . . . .	57

3.17 ENVASA: Age effect - LMEM model for PC (%-correct) in the three background measures single, dual, & combine background/s (reference levels: Background=single-background, APDsibling=1), and age as fixed factors and random intercept for subjects. Note: only data measured with the control group following outliers trimming was included. . . . .	58
3.18 ENVASA: Descriptive and statistics of the listeners age-independent standard residuals (z-scores) split by groups and test measures. . . .	60
3.19 CCC-2 subscales descriptives split by groups. . . . .	63
3.20 ECLiPS descriptives split by groups and sub-scales. . . . .	65
3.21 Switching task PCA: Input variables loading. . . . .	70
3.22 Language measures PCA: Input variables loading. . . . .	73
3.23 Correlation matrix (Spearman) between the study test measures for aggregated data across the two groups. . . . .	77



# List of Abbreviations

- 1-D, 2-D** . . . One- or two-dimensional, referring in this thesis to spatial dimensions in an image.
- Otter** . . . . . One of the finest of water mammals.
- Hedgehog** . . . Quite a nice prickly friend.



# Introduction

Welcome to the *R Markdown* Oxford University thesis template. This sample content is adapted from **thesisdown** and the formatting of PDF output is adapted from the OxThesis LaTeX template. Hopefully, writing your thesis in R Markdown will provide a nicer interface to the OxThesis template if you haven't used TeX or LaTeX before. More importantly, using *R Markdown* allows you to embed chunks of code directly into your thesis and generate plots and tables directly from the underlying data, avoiding copy-paste steps. This will get you into the habit of doing reproducible research, which benefits you long-term as a researcher, but also will greatly help anyone that is trying to reproduce or build upon your results down the road.

Using LaTeX together with *Markdown* is more consistent than the output of a word processor, much less prone to corruption or crashing, and the resulting file is smaller than a Word file. While you may never have had problems using Word in the past, your thesis is likely going to be about twice as large and complex as anything you've written before, taxing Word's capabilities.

## Speech-in-noise in children

*R Markdown* creates a simple and straightforward way to interface with the beauty of LaTeX. Packages have been written in **R** to work directly with LaTeX to produce nicely formatting tables and paragraphs. In addition to creating a user friendly interface to LaTeX, *R Markdown* allows you to read in your data, analyze it and to visualize it using **R**, **Python** or other languages, and provide documentation and commentary on the results of your project.

Further, it allows for results of code output to be passed inline to the commentary of your results. You'll see more on this later, focusing on **R**. If you are more into

**Python** or something else, you can still use *R Markdown* - see ‘Other language engines’ in Yihui Xie’s *R Markdown: The Definitive Guide*.

## APD definition

Anyone who needs to use data analysis, math, tables, a lot of figures, complex cross-references, or who just cares about reproducibility in research can benefit from using *R Markdown*. If you are working in ‘softer’ fields, the user-friendly nature of the *Markdown* syntax and its ability to keep track of and easily include figures, automatically generate a table of contents, index, references, table of figures, etc. should still make it of great benefit to your thesis project.

## Diagnosis

Anyone who needs to use data analysis, math, tables, a lot of figures, complex cross-references, or who just cares about reproducibility in research can benefit from using *R Markdown*. If you are working in ‘softer’ fields, the user-friendly nature of the *Markdown* syntax and its ability to keep track of and easily include figures, automatically generate a table of contents, index, references, table of figures, etc. should still make it of great benefit to your thesis project.

## Binaural and spatial listening in APD

Anyone who needs to use data analysis, math, tables, a lot of figures, complex cross-references, or who just cares about reproducibility in research can benefit from using *R Markdown*. If you are working in ‘softer’ fields, the user-friendly nature of the *Markdown* syntax and its ability to keep track of and easily include figures, automatically generate a table of contents, index, references, table of figures, etc. should still make it of great benefit to your thesis project.



## Summary

Anyone who needs to use data analysis, math, tables, a lot of figures, complex cross-references, or who just cares about reproducibility in research can benefit from using *R Markdown*. If you are working in ‘softer’ fields, the user-friendly nature of the *Markdown* syntax and its ability to keep track of and easily include figures, automatically generate a table of contents, index, references, table of figures, etc. should still make it of great benefit to your thesis project.



*Neque porro quisquam est qui dolorem ipsum quia dolor sit amet, consectetur, adipisci velit...*

*There is no one who loves pain itself, who seeks after it and wants to have it, simply because it is pain...*

— Cicero's *de Finibus Bonorum et Malorum*.

# 1

## Binaural listening: interrupted and alternated speech-in-noise in adults

### Contents

---

<b>1.1</b>	<b>Influence of distractor type on IM</b>	<b>6</b>
1.1.1	Introduction	6
1.1.2	Experiment I: speech vs. non-speech distractors	6
1.1.3	Experiment II: speech distractors spoken in a familiar vs. unfamiliar language	6
1.1.4	General discussion and conclusion	8
<b>1.2</b>	<b>Dichotic vs. monotic presentation and the influence of speech material</b>	<b>8</b>
1.2.1	Introduction	9
1.2.2	Methods	9
1.2.3	Results	9
1.2.4	Discussion	10
1.2.5	Conclusion	10

---

Here is a brief introduction to using *R Markdown*. *Markdown* is a simple formatting syntax for authoring HTML, PDF, and MS Word documents and much, much more. *R Markdown* provides the flexibility of *Markdown* with the implementation of **R** input and output. For more details on using *R Markdown* see <http://rmarkdown.rstudio.com>.

Be careful with your spacing in *Markdown* documents. While whitespace largely

is ignored, it does at times give *Markdown* signals as to how to proceed. As a habit, try to keep everything left aligned whenever possible, especially as you type a new paragraph. In other words, there is no need to indent basic text in the Rmd document (in fact, it might cause your text to do funny things if you do).

## 1.1 Influence of distractor type on IM

### 1.1.1 Introduction

- *Italics* are done like `*this*` or `_this_`
- **Bold** is done like `**this**` or `__this__`
- ***Bold and italics*** is done like `***this***`, `____this____`, or (the most transparent solution, in my opinion) `**_this_**`

### 1.1.2 Experiment I: speech vs. non-speech distractors

- Inline code is created with backticks like ``this``

#### Methods

Sub<sub>2</sub> and super<sup>2</sup> script is created like `this~2~` and `this^2^`

#### Results

- ~~Strikethrough~~ is done `~~like this~~`

#### Discussion

- To include an actual `*`, `_` or `\`, add another `\` in front of them: `\*`, `\_`, `\\`

### 1.1.3 Experiment II: speech distractors spoken in a familiar vs. unfamiliar language

- `--` and `---` with `--` and `---`

## Methods

Do like this:

Put a > in front of the line.

## Results

- are done with #’s of increasing number, i.e.
  - # First-level heading
  - ## Second-level heading
  - ### Etc.

In PDF output, a level-five heading will turn into a paragraph heading, i.e. `\paragraph{My level-five heading}`, which appears as bold text on the same line as the subsequent paragraph.

## Discussion

Unordered list by starting a line with an \* or a -:

- Item 1
- Item 2

Ordered lists by starting a line with a number:

1. Item 1
2. Item 2

Notice that you can mislabel the numbers and *Markdown* will still make the order right in the output.

To create a sublist, indent the values a bit (at least four spaces or a tab):

1. Item 1
2. Item 2

3. Item 3

- Item 3a
- Item 3b

### 1.1.4 General discussion and conclusion

The official *Markdown* way to create line breaks is by ending a line with more than two spaces.

Roses are red.   Violets are blue.

This appears on the same line in the output, because we didn't add spaces after red.

Roses are red.

Violets are blue.

This appears with a line break because I added spaces after red.

I find this is confusing, so I recommend the alternative way: Ending a line with a backslash will also create a linebreak:

Roses are red.

Violets are blue.

To create a new paragraph, you put a blank line.

Therefore, this line starts its own paragraph.

## 1.2 Dichotic vs. monotonic presentation and the influence of speech material

- This is a hyperlink created by writing the text you want turned into a clickable link in [square brackets followed by a](https://hyperlink-in-parentheses)

### 1.2.1 Introduction

- Are created<sup>1</sup> by writing either `^[my footnote text]` for supplying the footnote content inline, or something like `[^a-random-footnote-label]` and supplying the text elsewhere in the format shown below <sup>2</sup>:

`[^a-random-footnote-label]: This is a random test.`

### 1.2.2 Methods

To write comments within your text that won't actually be included in the output, use the same syntax as for writing comments in HTML. That is, `<!-- this will not be included in the output -->`.

### 1.2.3 Results

The syntax for writing math is stolen from LaTeX. To write a math expression that will be shown **inline**, enclose it in dollar signs. - This: `$A = \pi*r^{2}$`  
Becomes:  $A = \pi * r^2$

To write a math expression that will be shown in a block, enclose it in two dollar signs.

This: `$$A = \pi*r^{2}$$`

Becomes:

$$A = \pi * r^2$$

To create numbered equations, put them in an 'equation' environment and give them a label with the syntax `(\#eq:label)`, like this:

```
\begin{equation}
f\left(k\right) = \binom{n}{k} p^k\left(1-p\right)^{n-k}
(\#eq:binom)
\end{equation}
```

---

<sup>1</sup>my footnote text

<sup>2</sup>This is a random test.

Becomes:

$$f(k) = \binom{n}{k} p^k (1-p)^{n-k} \quad (1.1)$$

For more (e.g. how to theorems), see e.g. the documentation on [bookdown.org](https://bookdown.org)

#### 1.2.4 Discussion

- *R Markdown: The Definitive Guide* - <https://bookdown.org/yihui/rmarkdown/>
- *R for Data Science* - <https://r4ds.had.co.nz>

#### 1.2.5 Conclusion



# 2

## Spatial listening: development and normalisation of a children's spatialised speech-in-noise test

### Contents

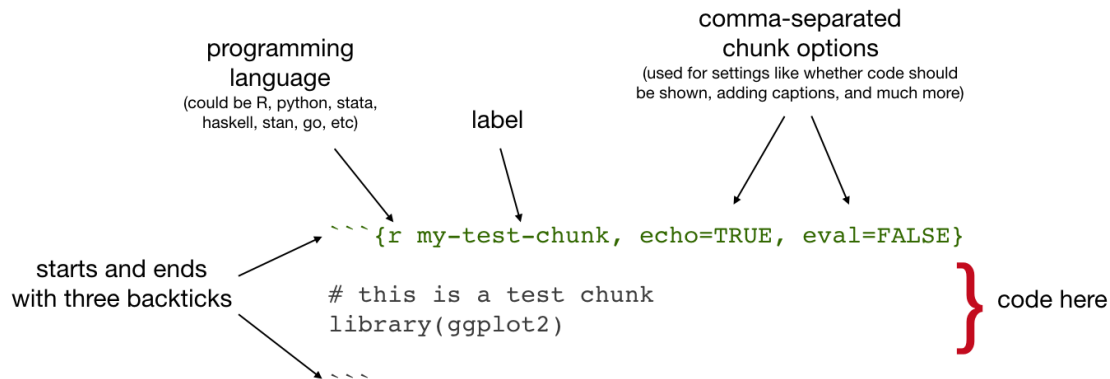
---

<b>2.1</b>	<b>Introduction</b>	<b>12</b>
<b>2.2</b>	<b>Methods</b>	<b>12</b>
<b>2.3</b>	<b>Discussion</b>	<b>12</b>
<b>2.4</b>	<b>Conclusion</b>	<b>12</b>

---

The magic of R Markdown is that we can add code within our document to make it dynamic.

We do this either as *code chunks* (generally used for loading libraries and data, performing calculations, and adding images, plots, and tables), or *inline code* (generally used for dynamically reporting results within our text).



**Figure 2.1:** Code chunk syntax

## 2.1 Introduction

## 2.2 Methods

## 2.3 Discussion

## 2.4 Conclusion

The syntax of a code chunk is shown in Figure 2.1.

Common chunk options include (see e.g. [bookdown.org](http://bookdown.org)):

- `echo`: whether or not to display code in knitted output
- `eval`: whether or to to run the code in the chunk when knitting
- `include`: wheter to include anything from the from a code chunk in the output document
- `fig.cap`: figure caption
- `fig.scap`: short figure caption, which will be used in the ‘List of Figures’ in the PDF front matter

**IMPORTANT:** Do *not* use underscoores in your chunk labels - if you do, you are likely to get an error in PDF output saying something like “! Package caption Error: \caption outside float”.

# 3

## APD study

### Contents

---

<b>3.1</b>	<b>Introduction . . . . .</b>	<b>14</b>
<b>3.2</b>	<b>Methods . . . . .</b>	<b>16</b>
3.2.1	Participants . . . . .	16
3.2.2	Measurements . . . . .	19
3.2.3	Procedure . . . . .	29
3.2.4	Data Analysis . . . . .	31
<b>3.3</b>	<b>Results . . . . .</b>	<b>33</b>
3.3.1	Standard audiometry . . . . .	33
3.3.2	EHF audiometry . . . . .	37
3.3.3	ST . . . . .	40
3.3.4	LiSNS-UK . . . . .	52
3.3.5	ENVASA . . . . .	56
3.3.6	CELF-RS . . . . .	60
3.3.7	Questionnaires . . . . .	62
<b>3.4</b>	<b>Overall performance . . . . .</b>	<b>66</b>
3.4.1	Unsupervised machine learning (UML) . . . . .	68
3.4.2	Interaction between measures . . . . .	68
<b>3.5</b>	<b>Discussion . . . . .</b>	<b>80</b>
3.5.1	EHF . . . . .	80
3.5.2	ST . . . . .	81
3.5.3	CCC-2 . . . . .	82
3.5.4	ECLIPS . . . . .	82
<b>3.6</b>	<b>Conclusion . . . . .</b>	<b>82</b>

---

### 3.1 Introduction

- APD definition: “unexplained idiopathic (spontaneous) listening difficulty (LiD) is often termed auditory processing disorder (APD) in children who have symptoms of difficulty hearing and understanding speech, and abnormal results on more complex auditory tests, despite having normal pure-tone hearing sensitivity (Jerger & Musiek 2000; Musiek et al., 2017)” [Hunter et al., 2020]
- Prevalance of LiD  $\sim 10\%$  (Sharma et al., 2009). Prevalence of LiD complaints with measured NH, complying with APD definition is estimated at  $\sim 0.5$  to  $1\%$  of the general population (Hind et al., 2011; Halliday et al., 2017)
- Association with other developmental disorders and lack of understanding of the underlying auditory deficits of APD.
- “Hearing involves both”bottom-up" (ear to brain) and “top-down” (cortical to subcortical) pathways through simultaneous and sequatial processing (Moore & Hunter, 2013)" [Hunter et al., 2020]
- Two general mechanistic hypotheses of APD:
  - (1) **Sensory processing difficulties (bottom-up)**: involving the central auditory nervous system, are based on animal and human lesion studies (Snow et al., 1997). Supporters of this hypothesis suggested this can be assessed using low-redundancy (simple) speech tests (e.g., using added noise, filtering, rapid speech,..) to “stress” the highly redundant central auditory pathways to reveal deficits (Keith 1995, 2000; Cameron et al., 2014).
  - (2) **Higher-level cognition or attention (top-down)**: especially in children with language disorders (Rees, 1973; Moore et al., 2010).

Individuals may have a combination of both.

- There is no accepted consensus or gold standard diagnosis of APD (Wilson & Arnott 2013)
- Possible link between OME (+ grommets) or EHF HL and APD in a sub-group of children.
- OME related HL has been shown to persist after recovery at frequencies above 4 kHz (Hunter et al., 2020; REFs..)
- OME or EHF HL can potentially be a basis for poorer speech perception, especially in noise. Findings are not consistent. Studies that tested both TD and APD with OME or EHF HL found that they are predictors of measurable peripheral damage in both groups.
  - Besser et al. (2015) and Levy et al. (2015) found that better thresholds between 6 to 12.5 kHz were associated with better reception of speech in noise (adult studies).
  - Motlagh Zadeh et al. (2019): impairment in higher frequency regions could negatively impact speech perception.

Conductive loss results in impaired spatial processing (Cameron et al., 2014) and binaural interaction (Hall et al., 1995; Hogan et al., 1996)

## 3.2 Methods

### 3.2.1 Participants

Forty-four primary school children who are native British English speakers with normal hearing acuity participated in the study. Amongst them 21 belonged to the APD clinical group (5 females) with an average age of  $11.04 \pm 1.42$  years (range: 7.8 - 12.9 years). The remaining 23 (12 females) comprised of typically developing control children (TD) with no reported concerns or diagnosis of an auditory, language or other cognitive developmental disorder. The TD group average age was  $9.47 \pm 1.58$  years and ranged between 7 to 12.1 years. Since not all the measurement equipment was easily portable and in order to maintain the same environment during the assessment across the complete sample, the children and their caregivers were required to travel to central London for the testing. In order to maximise the number of children taking part in the study, 8 out of the 23 TD children (35%) had APD siblings which took part in a parallel study on the same day of testing. All the children who participated in the study were required to have normal hearing acuity, defined as thresholds  $\leq 25$  dB HL at the octave frequency bands between 0.25 to 8 kHz and their eardrum had to be visible, healthy and intact in both ears following otoscopic inspection. One APD participant was excluded from the analysis due to raised thresholds predominantly in the right ear, ranging between 30 to 45 dB HL ( $PTA_{Right} = 36.25$  dB HL;  $PTA_{Left} = 13.75$  dB HL), thus resulting in a final APD group size of twenty. Otoloscopic inspection of the child's ear canal revealed a large accumulation of cerumen in both ears with an occluded right ear. Two additional children (x1 APD, x1 TD) had slightly raised thresholds at 8 kHz in one ear of 35 and 30 dB HL, respectively. However, since thresholds at all other frequency bands were well within the  $\leq 25$  dB HL criteria they were not excluded.

APD children were recruited in two ways. Children diagnosed with APD at Great Ormond Street Hospital (GOSH) or at the London Hearing and Balance Centre (LHBC), London, UK, were identified based on their clinical records and were

contacted by a clinical team member. The caregivers were provided with information about the study and means of contact to express interest in participation. Others, including the TD group were recruited by advertisements on social networks (e.g., APD Support UK Facebook group), science events, local information boards and UCL staff newsletter email, where parents were requested to fill-out an online interest form with short screening questions to ensure that the child met the participation requirements. Most of the children in the APD group (85%, 17/20) were reported to undergo an APD assessment at GOSH, about a third were directly recruited from the clinic. The remaining three were reported to be assessed at LHBC, at the University of Southampton Auditory Implant Service and Chime Audiology Royal Devon & Exeter Hospital (screening only).

Our initial aim was to take a conservative stance on inclusion criteria by including only those who met a clinical APD criteria (2 SD below the norms on two or more tests during the assessment). Moreover, being aware of the high morbidity with other developmental disorders (PERCENTAGE + REF), we strived to recruit children who display a “pure” form of APD without reported diagnosis or concerns for additional developmental disorder/s. However, very few APD children met these strict criteria, only 75% (15/20) met the clinical criteria of APD, out of which 60% (9/15) were diagnosed with spatial processing disorder (SPD) due to abnormal SRM in the LiSN-S task (see Table 3.1 for descriptives of the APD group). Of the remaining children in the APD group, four did not meet the diagnosis criteria for various reasons (e.g., young age, lack of psychological educational evaluation report and the need to exclude other deficits), however their assessment report acknowledged some “auditory processing difficulties”, whereas the fifth child awaited an APD assessment following an APD screening. Furthermore, half of the APD group (10/20) were reported for being diagnosed with one or more secondary developmental disorder/s (x6 Dyslexia, x3 HF-ASD, x3 DLD, x1 ADHD, x1 ADD, x1 Dyspraxia, x1 visual stress, x1 sensory integration disorder, and x1 poor short-term WM). Nonetheless, several caregivers reported that their motivation for seeking additional diagnosis

**Table 3.1:** APD group demographics and APD-related history background.

<b>School type</b>	85% (17/20) Mainstream, (1 child in a special ASD unit, 2 in a private school), 15% (3/20) non-mainstream school
<b>Assessment location</b>	85% (17/20) GOSH, 15% (3/20) Other
<b>APD Diagnosis</b>	75% (15/20) APD, 25% (5/20) LiD/susAPD
<b>SPD subtype</b>	60% (9/15) SPD
<b>Additional disorder</b>	50% (10/20) secondary developmental disorder/s
<b>Additional disorder (undergoing assessment)</b>	25% (5/20)
<b>OME history</b>	60% (12/20)
<b>PET history</b>	25% (5/20)
<b>FM-device usage</b>	55% (11/20)
<b>Auditory training</b>	35% (7/20)
OME: Otitis media with effusion	
PET: Pressure equalisation tube	

was to get more help from the school, rather than a real concern, after feeling that their support for their child's APD was lacking.

Caregivers from both groups completed a comprehensive background questionnaire, similar to the one that is typically given prior to an APD assessment, concerning the guardian/s educational level, child and family history of hearing, listening problems and developmental disorders, child history of otitis media, grommets, pregnancy-related questions (e.g., complications, prematurity, etc.), APD-related (e.g., date of diagnosis, location, use of FM device and auditory training), any diagnosis or concerns regarding the child's speech, language, educational and/or cognitive skills, speech and language therapy, medication taken, musical training and the type of school the child attends.

Children in the APD group were on average 1.5 years older than children in the TD group. Difference in age between the two groups was tested using t-test with



Welch degrees of freedom correction for uneven sample-size (independent-samples with bootstrapping  $n=9999$ ; MKinfer package; Kohl, 2020), showing a significant difference in age between the groups [ $t(40.95) = 3.43$ ,  $p = 0.001$ ]. Nonetheless, since age is often reported as a strong indicator for performance in other similar behavioural studies, analysis of the results obtained in the current study was conducted for age-independent scaled scores and should not affect the comparison between the two groups. The project was approved by the UCL Research Ethics Committee (Project ID Number 0544/006) and the NHS Health Research Authority (REC reference: 18/LO/0250). The testing commenced once an informed consent was given by both the caregiver and the child.

### 3.2.2 Measurements

The test battery used in the present study is described in the following section and summarised in Table 3.2.

#### Auditory evaluation

##### Standard & extended high-frequency (EHF) audiometry

Otoscopic inspection was performed prior to the audiometric test to ensure the ear was clear from cerumen and to avoid harming the eardrum when inserting the ear probe. Both standard and extended high-frequency (EHF) audiometry thresholds were measured using the Hughson-Westlake manual procedure, starting from 1 kHz. Standard air conduction pure-tone audiometry was carried out at six octave frequency bands ranging between 0.25 to 8 kHz using ??? audiometer and ??? headphones.

Extended high-frequency pure-tone detection thresholds were measured at four octave frequency bands 8, 11, 16, & 20 kHz using locally written MATLAB based software which generated the stimulus and collected the data. Target tones were pulsed (3 repetitions) with a duration time of 700 ms and 50 ms rise/fall time. EHF measurements took place in a designated sound attenuated chamber with the

**Table 3.2:** Summary of the study test battery.

Task	Information	Measure
<b>Standard &amp; extended high-frequency (EHF) audiometry</b>	Pure-tones detection thresholds measured at the octave frequency bands between 0.25 and 8 kHz (standard), and 8 to 20 kHz (EHF).	Detection threshold in dB HL
<b>Switching task (ST)</b>	Adaptive speech-on-speech listening task that involves perception of interrupted and periodically segmented speech that is switched between the two ears out-of-phase with an interrupted distractor. ST assesses the ability to switch attention and integration of binaural information.	Proportion of speech required to understand 50% of the keywords, Speech Reception duty cycle Threshold (SRdT)
<b>Listening in Spatialised Noise Sentences UK (LiSNS-UK)</b>	Locally developed version of the LiSN-S (Cameron & Dillon, 2007), an adaptive speech-on-speech listening task that assesses the ability to use spatial release from masking (SRM), measured as the difference in perception between collocated and separated speech distractors.	Signal-to-noise-ratio (SNR) yielding 50% speech intelligibility, Speech Reception Threshold (SRT)
<b>Speech-shaped-noise (SSN)</b>	Conventional adaptive speech in noise task that assesses speech perception of ASL sentences (MacLeod & Summerfield, 1990) in a speech-shaped-noise with a spectrum matched to the ASL material.	SRT
<b>The Environmental Auditory Scene Analysis task, ENVASA</b> (Leech et al., 2009)	Non-linguistic self-administered task involves detection of everyday environmental sounds presented in naturalistic auditory scenes and can be used to assess IM effects as well as sustained selective auditory attention skills.	%-correct
<b>Recalling sentences, CELF-RS</b> (Wiig et al., 2017)	A subtest from the Clinical Evaluation of Language Fundamentals UK 5 <sup>th</sup> edition (CELF-5-UK), assess expressive language skills, measured by the ability to repeat in verbatim sentences with varying length and complexity. Standardised for children aged 5 to 16 years.	Age-corrected scaled scores
<b>The Evaluation of Children's Listening and Processing Skills, ECLiPS</b> (Barry & Moore, 2014)	Standardised questionnaire comprises of 38 statements grouped into five categories designed to identify listening and communication difficulties in children aged 6 to 11 years. Respondents agreement is expressed using a five-point Likert scale (" <i>strongly agree</i> " - " <i>strongly disagree</i> ").	Age-corrected scaled scores
<b>The Children's Communication Checklist 2<sup>nd</sup> edition, CCC-2</b> (Bishop, 2003)	Standardised questionnaire comprising of 70 items designed to screen language and/or communication problems in children aged 4 to 16 years. Item comprises of a behaviour statement (e.g., " <i>Mixes up words of similar meaning</i> ") with respondents asked to judge how often the behaviours occur using a four-point Likert scale (0-3).	Age-corrected scaled scores

child sitting in the centre of the chamber while the examiner was situated outside. Communication during the testing was carried out via a video-audio intercom system. The child was instructed to raise his/hers hand each time s/he heard a tone. The MATLAB script was executed using a Windows PC which which was connected via USB to an RME sound card (Audio AG, Haimhausen Germany) and an ER10X Extended-Bandwidth Acoustic Probe System (Etymotic Research, Elk Grove Village, IL). Stimulus was presented via an otoacoustic emission probe with silicon tips in variable sizes (10, 12 or 13 mm), depending on the size of

the child's ear.

Standing waves in the ear canal produces spatially non-uniform sound pressure at frequencies above 2-3 kHz, introducing calibration errors when estimating the sound pressure level arriving at the eardrum (Lee et al., 2012; Richmond et al., 2011; Siegel, 1994). Together with other factors such as individual variations in the ear canal length and differences in depth in which the ear probe is inserted into the ear canal, these factors can introduce up to 20 dB calibration error (Siegel, 1994). To account for that, a sound pressure level calibration procedure was used (chirp noise) using a similar technique as described by Lee et al. (2012). For each frequency, the first half-wave resonance of the ear canal was measured, estimating the distance between the ear probe and the eardrum. The target stimulus was then scaled to the desired output level that corresponds to 0 dB HL using the in-situ calibration forward-pressure level data (FPL) and EHF-specific weighting thresholds (in dB SPL) measured across 84 NH listeners aged 10 to 21 years (see Table 1 in Lee et al., 2012).

### **Switching task (ST)**

Estimating the effect of IM while minimising peripheral EM on speech perception was measured using the switching task (ST) which is believed to assess the listeners ability to switch attention and integration of binaural information. The exact same test procedure and equipment was used as described in Chapter ???. Listeners were presented with both test versions using the ASL and the CCRM speech material. As for the stimuli, the ASL target sentences, spoken by a single male talker, were taken from the final sentences selected following the normalisation study. In addition, a level correction was applied to each sentence using the sentence-specific weighing factors estimated in the normalisation study (see Chapter ???). The first five test lists out of the eight phonetically-balanced normalised test lists (à 25 sentences each) were used, whereby their order was quasi-randomised to account for order, masker combinations, and fatigue effect. The target CCRM sentences were the same as described in Chapter ??, spoken by three different male talkers, were selected at

random every trial and always began with the priming animal ‘dog’. The target speech material was presented either without a distractor (Quiet) with and without switching (NoAlt / Alt) or with a distractor. A selection of four distractors were used (see Chapter ??? for detailed description): English (ENG\_F) and Mandarin (MDR\_F) unrelated connected-speech, each spoken by ten different female talkers, and a non-speech amplitude-modulated speech-spectrum-noise (AMSSN) with the envelope of a single talker out of 40 talkers (20 females). The fourth distractor was presented only with the CCRM speech material and comprised of CCRM target-like sentences (CCRM\_F) with a different priming animal, colour and digit, spoken by ten different female talkers. Each participant was presented with a total of 11 runs, one for each test condition, with 5 conditions for the ASL (Quiet-NoAlt, Quiet-Alt, MDR\_F-Alt, ENG\_F-Alt), and 6 for the CCRM (with the additional CCRM\_F-Alt condition). Testing started following a practice phase, where four trials of each of the eleven test conditions were presented. Practice runs started at an easy-to-moderate DC rate of 0.8 in order to expose the listeners to the adaptive procedure. In addition, every test run started with two practice sentences (DC start = 0.97) to orient the listeners to the test condition that is about to be presented.

### **Listening in Spatialised Noise Sentences UK (LiSNS-UK)**

The locally developed Listening in Spatialised Noise Sentences UK (LiSNS-UK) assesses the ability to use binaural cues in speech-on-speech listening conditions. The test development, speech material normalisation, and norms standardisation followed Cameron and Dillon (2007) development steps and are described in detail in Chapter ???. The test uses virtualisation techniques to create spatial distribution of sound sources in space for headphones presentation where target sentences (ASL; MacLeod & Summerfield, 1990) are presented in two simultaneous speech distractors (unrelated children’s stories spoken by the target talker). The LiSNS-UK comprises of two main listening conditions, differing in their availability of spatial cues. The target sentences are configured to always appear in front of the listener’s head,

at  $0^\circ$  azimuth on the horizontal plane, with the two streams of speech distractors either collocated in space with the target (S0N0), resulting in relatively poor speech perception, or offset in space, with one distractor to either side of the target at  $\pm 90^\circ$ . The spatial separation in the later condition results in an improvement in speech perception of circa 13 dB (Cameron et al., 2011), typically termed as spatial release from masking (SRM). This SRM advantage is calculated by taking the difference between performance in the collocated and the separated condition.

Speech distractors were presented continuously throughout a run at a fixed 65 dB SPL output level and comprised of a combination of two out of three available passages. A 1-up/1-down adaptive procedure was used, varying the level of the target talker relative to the distractors depending on the listener’s response to measure their speech reception threshold (SRT), i.e., the signal-to-noise-ratio (SNR) yielding 50% speech intelligibility. A 2 ms long reference cue (1 kHz pure-tone) was presented 500 ms before the target sentence onset at 65 dB SPL. The initial target output level was 75 dB SPL for the collocated condition and 70 dB SPL for the separated condition with an initial step-size of 4 dB SNR. The step-size was reduced after every reversal, reaching a minimum step-size of 2 dB SNR after three practice reversals. The adaptive procedure ended once all 25 test trials were presented and stopped in case a maximal output level of 89 dB SPL was reached more than three times. Nonetheless, such event did not occur in the present study. Since each listener was only presented once with each condition, it was decided not to introduce any other stopping rules that could have expedited the testing time but may as well introduced an estimation error for the SRTs in some cases. The SRT was calculated by averaging the test reversals SNRs, whereby test reversals were defined as any reversals following three practice reversals.

The order of the listening condition, test lists, sentences within a run, and distractors combinations was fixed across all the participants and started with the collocated condition. Each test list consisted of 25 sentences taken from

the 8-phonetically-balanced ASL test lists which were constructed following the normalisation study and a sentence-specific level correction was applied (see Chapter. ???). Spatialisation was applied by convolving each stimuli with head-related transfer functions (HRTFs) at the corresponding azimuthal direction separately for the left and the right channel. The HRTFs were measured with a Knowles Electronics Manikin for Acoustic Research (KEMAR) with a small pinnae taken from the CIPIC HRTF database<sup>1</sup> (see Algazi et al., 2001, “special” HRTF data). A post-equalisation step was applied in order to flatten the magnitude of the headphones frequency response. Headphone-to-ear Transfer Functions (HpTFs) measured with KEMAR manikin for HD-25 supraaural headphones were extracted from Wierstorf et al. (2011) HRTF database. The final mixed stimulus was filtered with the inverse HpTFs separately for the left and the right channel before being combined together as a final step. Every participant was presented with two runs, one for each listening condition (collocated / separated). Testing started following a practice phase of two runs, one for each of the test conditions with five BKB sentences each (Bench et al., 1979). Listeners were instructed to verbally repeat the target sentences to the experimenter who was situated alongside in a sound treated chamber. The experimenter scored the response by selecting the correctly repeated keywords on the screen. Listeners were encouraged to guess if unsure while no feedback was given at any time. A loose keyword scoring method was used, whereby errors of case or declension were considered as correct responses, e.g., a repetition of the keywords ‘<clowns> <funny> <faces>’ to the stimulus ‘The <clown> had a <funny> <face>’.

### Speech-shaped-noise (SSN)

A speech-in-noise test was used as a more conventional listening task that is widely used in the clinic as opposed to the more complex listening conditions measured by the ST or the LiSNS-UK. The normalised ASL sentences were presented in a speech-spectrum-noise (SSN) with spectrum matched to the ASL corpus. The

---

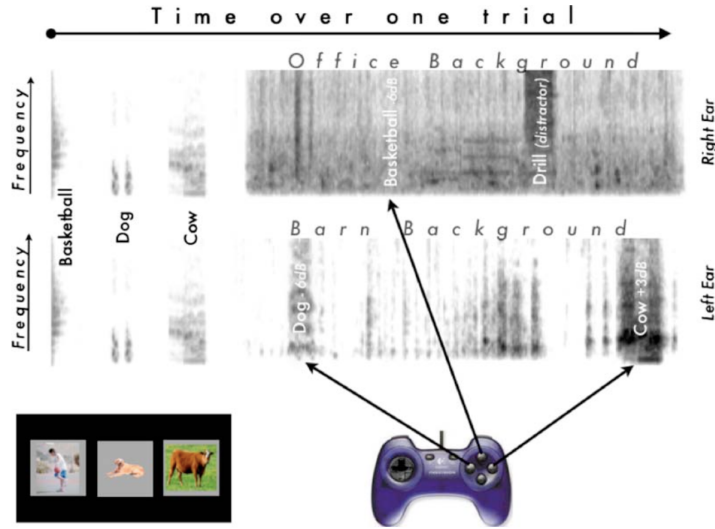
<sup>1</sup>The database is available online in: <https://www.ece.ucdavis.edu/cipic/spatial-sound/hrtf-data/>

SSN onset was 500 ms before the target sentence began. The exact same adaptive procedure as for the LiSNS-UK was used with the same stopping-rules and SRT calculation. Each listener was presented with a single run of 25 sentences following a practice phase with seven BKB sentences. The same test list and sentences order was used across all the listeners.

### **The Environmental Auditory Scene Analysis task (ENVASA)**

In analogy to the classic ‘cocktail-party’ scenario, ENVASA is a non-linguistic paradigm (Leech et al., 2009) that measures detection of everyday environmental sounds presented in naturalistic auditory scenes and can be used to assess IM effects as well as sustained selective auditory attention skills. In the task, short environmental target sounds (e.g., a dog’s bark, a door knock, or a bouncing ball) were presented in a dichotic background scene (i.e., the target sound is presented only in one ear), consisting of either a single background scene, presented in both ears, or two background scenes, each presented in a different ear. The number of targets, the onset time and the ear of presentation varied across trials. Four SNRs were employed split into two categories ‘low’ (-6 and -3 dB) and ‘high’ (0 and +3 dB). Target-background contextual agreement was manipulated by embedding the target sound in a *congruent* background scene that is in agreement with the listener’s expectations (e.g., a cow’s ‘moo’ in a farmyard scene) or in an *incongruent* background scene which violate these expectations (e.g., a cow’s ‘moo’ in a traffic scene). A schematic illustration of a single test sequence is shown in Figure 3.1.

The experiment was carried out using the original setup as described by Leech et al. (2009). Sounds were presented via Sennheiser HD-25 headphones (Wedemark, Germany) and the participants response was recorded using a gamepad. The output level was adjusted to a comfortable level before the test started. The participants were situated in front of a laptop and were instructed to hold the gamepad. Prior to the test, the listeners were presented with a short child-friendly demonstration video with audio instructions. Next, a short recap was given verbally



**Figure 3.1:** Schematic of the ENVASA experimental paradigm (taken from Leech et al., 2009)

by the examiner and an exemplary trial was simulated together with the child to ensure that the child fully understood the task’s instructions. The task began with three short practice trials with provided feedback, while no further feedback was given during the test phase.

Every trial was made of two parts, starting with a target audio and visual familiarisation phase before the main target detection phase. Target identification was recorded by pressing one of the three buttons on the gamepad which corresponded to the location of the target objects on the screen. A response was counted as correct only if the participants pushed the corresponding button within 2 seconds, 300 ms after the target onset. The outcome measure was calculated as the percentage of target sounds correctly identified within a condition (%-correct). In total there were 115 target sounds presented over 40 trials, where 46 target sounds were presented in a single background condition and another 46 in a dual-background condition. The 23 remaining target sounds served as foil items which were played at 0 dB SNR without a corresponding picture on the screen. The order of the foil items was quasi-randomised and was used to estimate the quality of the participants performance.



### **CELF-RS**

The Recalling Sentences (RS) sub-test of the Clinical Evaluation of Language Fundamentals UK fifth edition (CELF-5-UK Wiig et al., 2017) was administered to assess the listeners expressive language skills, measuring the ability to repeat in verbatim sentences with varying length and complexity. Standardised norms are available for children aged 5 to 16 years. The CELF-RS is simple and quick to administer and has been shown to be a good psycholinguistic marker for children with Developmental Language Disorder (DLD) and to provide high levels of sensitivity and specificity (Conti-Ramsden et al., 2001), thus making it a good screening tool. Scoring were marked by hand by the examiner as instructed by the test manual. The sentences were presented using a local MATLAB program via headphones using the same experimental equipment as listed above at a comfortable output level of 70 dB HL. The sentences were spoken by a female speaker with a standard southern British English accent and were recorded in a sound-treated recording booth at the Speech, Hearing and Phonetics Sciences (SHaPS) UCL laboratory, London. The task began with two practice sentences while the number of test items varied depending on the child's age and performance. No repetitions or feedback was given during the testing and the test was discontinued in case the child failed to score any points for four consecutive items. Age-scaled score were calculated based on the test norms with a mean score of 10 and SD of 3. Scaled scores within  $\pm 1$  SD from the norms mean (between 8 to 12) are classified as average scores, whereas performance beyond  $\pm 1$  SD are classified as above / below the average score, with scaled-scores  $< 7$  considered as abnormally poor.

### **Questionnaires**

#### **The Evaluation of Children's Listening and Processing Skills (ECLiPS)**

The ECLiPS questionnaire (Barry & Moore, 2014) comprises of 38 items, where the respondents are asked to express their agreement on simple statements about the child's listening and other related skills or behaviours using a five-point Likert

scale (from “strongly agree” to “strongly disagree”). The ECLiPS was design to identify listening and communication difficulties in children aged 6 to 11 years. Nonetheless, in their evaluation study, Barry and Moore (2014) found little to no age effect in many of the scale items, suggesting that testing age could be extended below and beyond the population used for the development. Based on factor analysis the items were grouped into five subcategories: 1. Speech & Auditory Processing (SAP), assessing ability to interpret speech and non-speech input, 2. Environmental & Auditory Sensitivity (EAS), estimating the ability to cope with environmentally challenging conditions, 3. Language, literacy & laterality (L/L/L), assessing different abilities that are known to be coupled with language and literacy difficulties, 4. Memory & Attention (M&A), covering short-term and serial memory as well as attention, 5. Pragmatic & Social skills (PSS), assessing pragmatic language or non-normative social behaviours. Aggregated measures were calculated for *Listening* (SAP, M&A, & PSS), *Language* (L/L/L & M&A), *Social* (PSS & EAS), and a *Total* aggregate, calculated by taking the mean of scores across all the sub-scales. Individual age- and sex-scaled scores were computed using the test excel scorer. A score below the 10<sup>th</sup> percentile (corresponding to a scale score of circa 6) is generally considered clinically significant.

### **The Children’s Communication Checklist 2<sup>nd</sup> edition (CCC-2)**

Communication abilities were assessed using the Children’s Communication Checklist second edition questionnaire (CCC-2; Bishop, 2003) which is designed to screen communication problems in children aged 4 to 16 years and comprises of 70 checklist items each comprising of a behaviour statement, like “*Mixes up words of similar meaning*”. The respondents are asked to judge how often the behaviours occur using a four-point Likert rating scale: 0. *less than once a week (or never)*, 1. *at least once a week, but not every day*, 2. *once or twice a day*, 3. *several times (more than twice) a day (or always)*. The items are grouped into ten sub-scales of behaviours tapping into different skills (A. Speech, B. Syntax, C. Semantics, D. Coherence, E.

Inappropriate initiation, F. Stereotyped language, G. Use of context, H. Non-verbal communication, I. Social relations, J. Interests). Taking the sum of scores for the sub-scales A to H are used to derive the General Communication Composite (GCC) which is used to identify clinically abnormal communication competence. A GCC score  $< 55$  was found to well separate between control and clinical groups, identifying children with scores at the bottom 10% (Norbury & Bishop, 2005). Another proposed composite, is the SIDC (Social-Interaction Deviance Composite) which is calculated by taking the difference in sum of subscales E, H, I, and J (tapping into pragmatic language and social skills) from the sum of scales of A to D (describes structural language skills). Abnormal GCC ( $< 55$ ) combined with a negative SIDC score has been shown to be indicative of an autistic spectrum disorder profile (Bishop, 2003). The CCC-2 scaled and composite scores were computed using the test scorer.

### 3.2.3 Procedure

Testing took place at the SHaPS laboratory (UCL, London) in a sound-attenuated chamber. Unfortunately, since many of the APD children had to travel from outside London and because of difficulties in recruitment, all the testing had to be made in a single session, lasting in total circa 2.5 to 3 hours (including breaks). To minimise possible fatigue effect, the session was carefully designed to ensure several planned and unplanned breaks. The participants were encouraged to request for a break between test runs whenever they required and were observed for any signs of fatigue by the examiner. The different tasks were gathered into short blocks and different measures were scattered throughout the session to keep the session fun and engaging for the child. At the end of the session, each child received a certificate and an Amazon voucher as a token of appreciation for taking part in the study and travel costs of the family were reimbursed.

**Table 3.3:** Experimental design and measurements order.

Order	Group A	Group B	Group C	Group D
1	Standard audiogram	Standard audiogram	Standard audiogram	Standard audiogram
2	Otoscopy	Otoscopy	Otoscopy	Otoscopy
3	ST-ASL	ST-ASL	ST-CCRM	ST-CCRM
4	CELF-RS	SSN	CELF-RS	SSN
5	ST-CCRM	ST-CCRM	ST-ASL	ST-ASL
6	SSN	CELF-RS	SSN	CELF-RS
7	EHF	EHF	EHF	EHF
8	ENVASA	ENVASA	ENVASA	ENVASA
9	LiSNS-UK	LiSNS-UK	LiSNS-UK	LiSNS-UK

Participants from both the TD and the APD group completed the same test battery in the below listed order (see Table 3.3). The ECLiPS, CCC-2 and the locally compiled background questionnaire were completed by the caregiver during the testing day. The session started with a standard pure-tone audiogram and otoscopy to ensure that detection thresholds fulfil the study criteria and that there are no abnormalities in the ear canal and the eardrum. Next, the switching task was conducted. Since performance in the task was one of the main focuses in the study, and because little is known about any possible learning effect in the task, presentation of the two speech materials (ASL and CCRM) was counterbalanced within each group, where about half of the children started with the ASL and the other half with the CCRM speech material. In between the two ST versions, each child completed the CELF-RS and the SSN task, whereby again, the order of presentation was counterbalanced within each group. Since both CELF-RS and SSN test duration are relatively short, they served as a short informal break between the ST test versions and kept the child engaged. Next, about half-way through the session, with a fixed order, all the participants were presented with the EHF audiometry, and the ENVASA task. The session was concluded with the LiSNS-UK, in-line with typical clinical assessment where the test is often presented last.

### 3.2.4 Data Analysis

All the data extraction, management and analysis in the present study was computed in R environment (Version 4.0.3; R Core Team, 2020) using RStudio (Version 1.4.938; RStudio Team, 2020).

#### Age scaled scores

Age-independent scores were estimated using a linear regression model. The model was fitted per condition separately for each measure (ST-ASL, ST-CCRM, LiSNS-UK, SSN, & ENVASA) and was based on the control group data only with the respective test raw scores (e.g., SRdT, SRT or %-correct) as an dependent variable and age as a predictor. A two-steps model comparison was performed to test the assumption that performance displays a monotonic linear relationship with age versus a non-monotonic (segmented) linear relationship. Extreme outliers were initially trimmed from the TD group to reduce noise in the data and to improve the models fit. In the first step, both models were computed and the best model was selected based on F-statistic model comparison using analysis of variance using *anova()* function. Standard residuals were next calculated for each TD listener, based on the selected model prediction. Since age was included in the model, the standardised residuals are age-independent and are comparable to z-scores for data with normal distribution, with a mean and SD of approximately 0 and 1. Since the main goal of the study was to find a measure that is able to well separate between the APD group and the typically developed control group, individual differences and group differences were explored using a deviance analysis procedure proposed by Ramus et al. (2003). Abnormal scores were defined by a two-tailed deviance cut-off of  $\pm 1.96$  SD from the TD group mean. Thus, circa 95% of the normal population residuals are expected to be within the deviance range of  $\pm 1.96$ . Occasional occurrence of abnormal scores in the normal population is not unusual in behavioural measures. Therefore, since the prediction of the residuals is based on the control data, such outliers may skew the TD group true mean or

SD and thus may introduce an error in the model prediction. Therefore, in the second step, additional TD outliers (with standardised residuals below/above TD mean  $\pm 1.96$ ) were trimmed from the data and the two models were refitted and compared again. Finally, the model with the best fit was selected and was used to calculate the standardised residuals for all the listeners, including the trimmed TD observations and the APD group.

### Statistical analyses

Residual analysis was performed separately for each measure to determine whether the data fulfils parametric methods assumptions of normal distribution using Shapiro-Wilk test (*shapiro.test()*, R Core Team, 2020) and homogeneity of variance using Levene’s test (*leveneTest()*; Fox & Weisberg, 2019). Consequently, statistical analyses for factorial design data that met these requirement was performed using linear mixed-effects regression models (LMEMs). LMEM was fitted using the *lmer()* function (lme4 package; Bates et al., 2015). Backward model selection procedure was applied to find the model that gives the best fit using a likelihood ratio test ( $\chi^2$ ). Main effects and interaction terms were tested by comparing predictions of the full model to a reduced model where each fixed term was separately removed, starting with the interaction terms. When applicable, post-hoc paired comparison t-test was performed on the fitted model and included adjusted least-squared-mean for the random intercepts (subjects) using the *lsmeans()* function from the emmeans R package (Lenth, 2020). In addition, group differences for single parametric measures such as CELF-RS and CCC-2 total score were examined using *boot.t.test()* function (MKinfer package; Kohl, 2020) which performs independent-samples t-test with bootstrapping (n=9999).

Nonparametric data was analysed using *nparLD()* function (nparLD package; Noguchi et al., 2012) which is a robust rank-based method for analysis of skewed data or for data with outliers or from a small sample size (see Feys, 2016, for a good introduction on robust nonparametric techniques). The function enables different

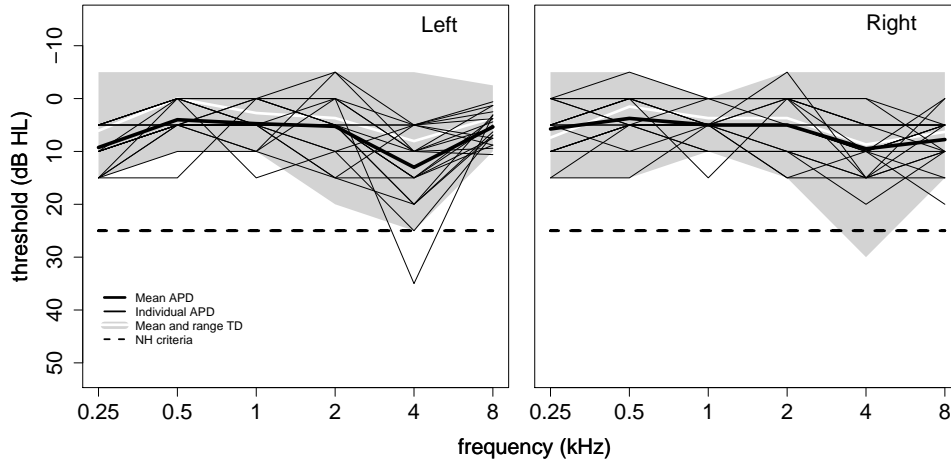
types of nonparametric tests for factorial design data with repeated measures with variable between-/within-subjects factors. The results reported in the present study were based on the ANOVA-type statistic test (ATS) output. Inspection of the ENVASA task age-independent z-scores revealed that the assumption of sphericity (Mauchly's test) was violated. Therefore, analysis was performed using *npIntFactRep* package (Feys, 2015), which is another robust aligned rank technique that enables sphericity correction (Greenhouse-Geisser). When applicable, post hoc comparison and group difference were examined using Wilcoxon rank-sum test with permutation (N=999999, independent two samples) which is a t-test equivalent for non-parametric data (*coin::wilcox\_test()*; Hothorn et al., 2006).

### 3.3 Results

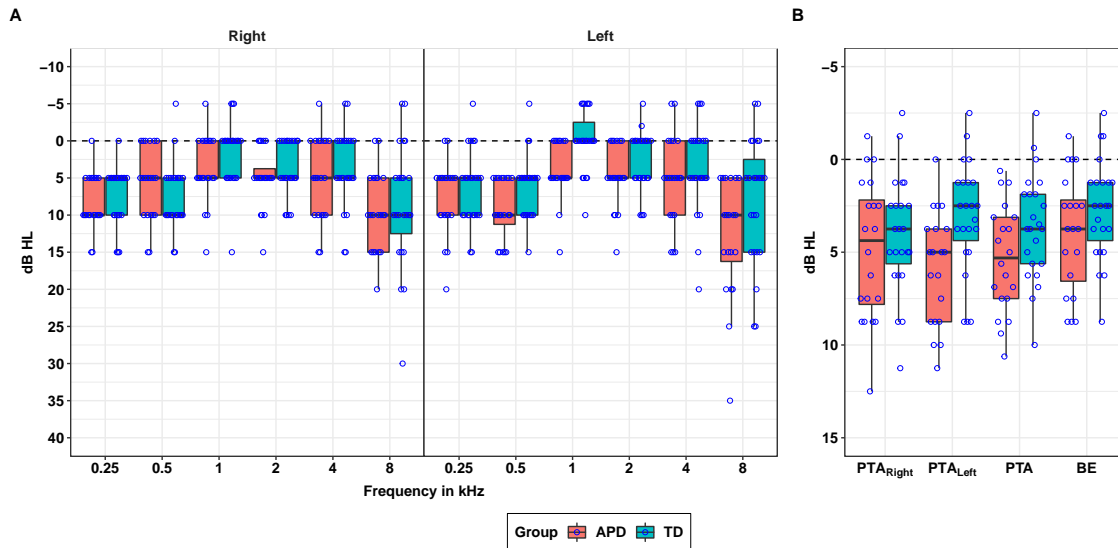
#### 3.3.1 Standard audiometry

The listeners' detection thresholds for the left and the right ear are plotted in Figure 3.2. The shaded grey area represents the TD group thresholds range and the white line represents the group mean at each frequency. The black lines marks the individual thresholds in the APD group and the group mean is marked by the bold black line. The dashed line indicates the maximal thresholds criteria of  $\leq 25$  dB HL for participation in the study.

Boxplots of listeners pure-tones detection thresholds measured at six octave frequency bands between 0.25 to 8 kHz and their corresponding pure-tone-average (PTA) are shown in Figure 3.3 A-B. Individuals PTAs were calculated by averaging thresholds at the frequency bands 0.5, 1, 2 and 4 kHz separately for the right and left ear ( $PTA_{Right}$ ,  $PTA_{Left}$ ) and by taking the grand mean for thresholds in both ears (denoted as PTA), whereas the listeners' PTA at the better-ear is denoted as BE. Thresholds descriptives by frequency bands and ear split by the two groups is given in Table 3.4, as well as Table 3.6 for PTAs and BE with additional statistics.



**Figure 3.2:** Standard audiometry: APD participants pure-tone detection thresholds plotted separately for the left and the right ear (black lines). The shaded grey area represents the TD group thresholds range and the white line represents the TD group mean at each frequency. The dashed line represents the threshold criteria of hearing level  $\leq 25$  dB HL.



**Figure 3.3:** Standard audiometry: Pure-tone detection thresholds by frequency bands between 0.25 to 8 kHz (A), and averaged thresholds (B). Individual scores are indicated by circles. The boxes show the data interquartile range (25th-75th percentile) and the horizontal line indicate the median (i.e., 50th percentile). Values that fall within 1.5 times the interquartile range are indicated by the whiskers.



**Table 3.4:** Standard audiometry: Descriptives for pure-tone detection thresholds (dB HL) by frequency bands (kHz) and ear split by the two groups.

Frequency	Ear	APD					TD				
		N	median	sd	min	max	N	median	sd	min	max
0.25	R	20	10	3.73	0	15	23	5	3.95	0	15
0.5	R	20	5	4.94	0	15	23	10	4.49	-5	15
1	R	20	5	4.55	-5	15	23	0	4.38	-5	15
2	R	20	5	3.97	0	15	23	5	3.76	0	10
4	R	20	5	5.62	-5	15	23	5	5.27	-5	15
8	R	20	10	5.36	0	20	23	10	8.29	-5	30
0.25	L	20	5	4.99	0	20	23	5	5.05	-5	15
0.5	L	20	10	4.06	5	15	23	5	4.25	-5	15
1	L	20	5	3.84	0	15	23	0	3.99	-5	10
2	L	20	5	4.13	0	15	23	0	4.03	-5	10
4	L	20	5	5.95	-5	15	23	5	6.07	-5	20
8	L	20	10	8.01	5	35	23	5	8.49	-5	25

Group differences for detection thresholds across frequency bands and ears was statistically tested with a three-way 6 x 2 x 2 factorial design with repeated measures. Inspection of the data by groups revealed that the assumption of normality and homoscedasticity were violated. Therefore, a non-parametric approach was adopted, using an rank-based ANOVA-type statistic test (ATS) with the *nparLD()* function (nparLD package; Noguchi et al., 2012). The ATS test results are given in Table 3.5. There was no significant three-way or two-way interaction between the three factors Group, Frequency and Ear ( $p > 0.05$ ), albeit Group x Ear interaction reached significance level ( $p = 0.062$ ). There was a significant main effect for group difference ( $p < 0.05$ ) as well as a highly significant difference in detection thresholds across frequencies ( $p < 0.001$ ), however no significant main effect for Ear was found ( $p = 0.827$ ).

Group differences for PTAs and BE were examined using a 4 x 2 LMEM model (parametric model assumptions were met). Detection measures (PTA<sub>Right</sub>, PTA<sub>Left</sub>, PTA and BE) and Group were set as fixed factors (reference levels: PTA<sub>Right</sub> and TD group) and detection threshold (in dB HL) as dependent variable, as well as a

**Table 3.5:** Standard audiometry: Statistical analysis for the effects of Frequency, Ear and Group and their interaction (6x2x2 factorial design with repeated measures) tested with a robust rank-based method for analysis of nonparametric data using nparLD package (Noguchi et al., 2012). Analysis was based on a f1-ld-f2 design ANOVA-type statistic (ATS) test, whereby f1 refers to an experimental design with a single between-subjects factor (Group) and f2 refers to two within-subjects factors (Frequency and Ear).

	Statistic	df	p-value
Group	4.350	1.000	<b>0.037</b>
Frequency	26.856	3.970	<b>0.000</b>
Ear	0.048	1.000	0.827
Group:Frequency	0.990	3.970	0.411
Ear:Frequency	1.234	4.121	0.294
Group:Ear	3.493	1.000	0.062
Group:Frequency:Ear	1.716	4.121	0.141

\* significant p-values ( $p < 0.05$ ) are shown in bold.

**Table 3.6:** Post-hoc paired comparison t-test for PTA x Group. The test was performed on the fitted LMEM model and included adjusted least-squared-mean for the random intercept (subjects) using lsmeans package (Lenth, 2020).

	APD					TD					post-hoc paired t-test					
	N	median	sd	min	max	N	median	sd	min	max	Estimate	SE	Df	t-value	p-value	95%-CI
PTA <sub>Right</sub>	20	4.38	3.78	-1.25	12.50	23	3.75	3.16	-2.5	11.25	0.799	0.942	59.762	0.848	0.4	-1.09 - 2.68
PTA <sub>Left</sub>	20	5.00	3.04	0.00	11.25	23	2.50	3.01	-2.5	8.75	2.682	0.942	59.762	2.846	<b>0.006</b>	0.8 - 4.57
PTA	20	5.31	2.92	0.62	10.62	23	3.75	2.87	-2.5	10.00	1.740	0.942	59.762	1.846	0.07	-0.15 - 3.62
BE	20	3.75	3.17	-1.25	8.75	23	2.50	2.68	-2.5	8.75	1.242	0.942	59.762	1.318	0.193	-0.64 - 3.13

\* significant p-values ( $p < 0.05$ ) are shown in bold.

PTA: average detection threshold (dB HL) at 0.5, 1, 2, & 4 kHz.

BE: PTA at the better ear.

random intercept for subjects. A model with interaction term was found to give the best fit, showing a significant interaction between the calculated detection measures and Group [ $\chi^2(3) = 12.27$ ,  $p < 0.05$ ]. Post-hoc paired comparison t-test based on the fitted model was computed using lsmeans function (*emmeans* package, Lenth (2020)); see Table 3.6) revealed a significant difference between the groups for PTA measured in the left ear ( $p < 0.05$ ), nonetheless, the magnitude of the difference between the two groups of 2.5 dB is rather small and clinically negligible, and is likely to occur due to sampling error. No significant difference was found in the remaining measures (all  $p$ 's  $> 0.05$ ).

**Age effect?**

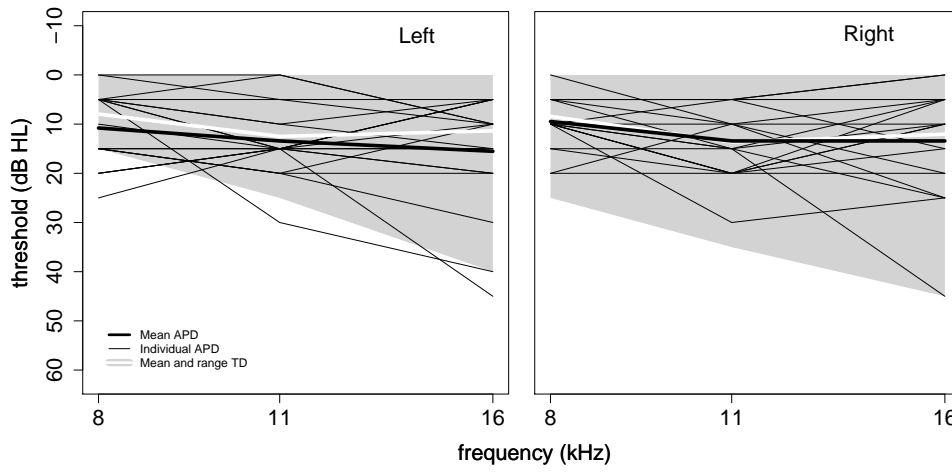
Regression lines for PTAs are not sig. nparLD with Age was not executable. LMEM  $dBHL \sim Freq + Ear + Group + Age + Freq:Ear + Freq:Group + (1 | listener)$  (although parametric assumptions were not met) resulted in a p-value of 0.0594.

**3.3.2 EHF audiometry**

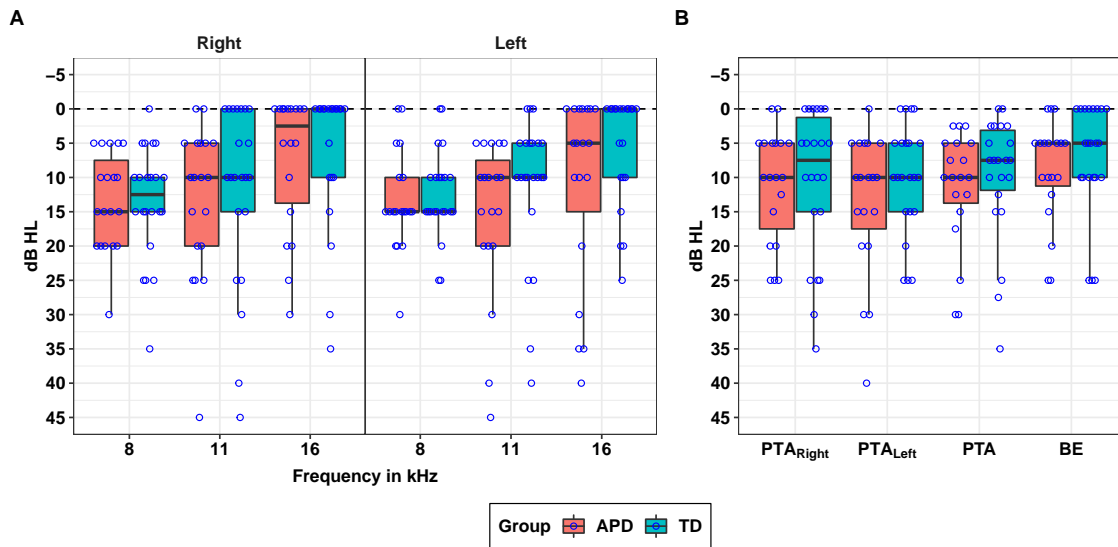
The listeners pure-tone detection thresholds measured at the octave frequency bands 8, 11 and 16 kHz are plotted in Figure 3.4 separately for the left and the right ear. Again, the thin black lines represents individuals' thresholds in the APD group and the group mean is marked by the bold black line. The shaded grey area represents the TD group thresholds range and the white line represents their mean at each frequency. It is worth noting that in many children from both groups it was not possible to record a reliable response for thresholds measured at 20 kHz, resulting in a large portion of missing data points. Therefore, thresholds measured at 20 kHz were not included in the analysis. A comparison of the group means reveals relatively small differences in thresholds between the groups, with a relatively larger difference in the left ear, where APD thresholds at 11 and 16 kHz were on average 5 dB higher (i.e., poorer).

Boxplots of the listeners thresholds by frequency and ear as well as their calculated PTA's and BE are shown in Figure 3.5 A-B. Descriptives of the groups detection thresholds is given in Table 3.7.

Group difference for thresholds across frequencies (8, 11, & 16 kHz) and ears (left/right) were examined for a 3 x 2 x 2 repeated measures factorial design. Inspection of parametric model assumptions revealed that the groups data violated the assumption of normality and homoscedasticity. Therefore, the exact same nonparametric procedure as used for the standard audiometry was performed using nparLD package. The ATS ANOVA-type test given in Table 3.8 found no significant



**Figure 3.4:** EHF audiometry: Pure-tone detection thresholds for extended high-frequency bands measured in the left and the right ear. The thin black lines represents the individual thresholds in the APD group and the group mean is marked by the bold black line. The shaded grey area represents the TD group threshold range and the white line represents the TD group mean at each frequency.



**Figure 3.5:** EHF audiometry: Boxplots for pure-tone detection thresholds measured at the extended high-frequency bands split by ear and groups (A). Boxplots of the groups averaged PTAs and better-ear BE thresholds are depicted in figure B. Individual scores are indicated by circles. The boxes show the data interquartile range (25th-75th percentile) and the horizontal line indicate the median (i.e., 50th percentile). Values that fall within 1.5 times the interquartile range are indicated by the whiskers.

**Table 3.7:** EHF audiometry: Descriptive for pure-tone detection thresholds (dB HL) by extended-high frequency bands (kHz) split by ear and group.

		APD					TD				
	Ear	N	median	sd	min	max	N	median	sd	min	max
<b>Octave frequency bands</b>											
8	R	19	15.0	7.08	5.0	30	22	12.5	8.34	0	35
11	R	19	10.0	11.19	0.0	45	22	10.0	13.24	0	45
16	R	19	2.5	10.04	0.0	30	22	0.0	10.51	0	35
8	L	19	15.0	7.27	0.0	30	22	15.0	6.50	0	25
11	L	19	10.0	11.65	5.0	45	22	10.0	10.71	0	40
16	L	19	5.0	13.80	0.0	40	22	0.0	8.11	0	25
<b>PTAs and better-ear</b>											
PTA <sub>Right</sub>	R	19	10.0	8.27	0.0	25	22	7.5	10.83	0	35
PTA <sub>Left</sub>	L	19	10.0	10.39	0.0	40	22	10.0	8.09	0	25
PTA		19	10.0	8.59	2.5	30	22	7.5	9.05	0	35
BE		19	5.0	7.60	0.0	25	22	5.0	8.27	0	25

PTA: average detection threshold at 8, 11, &amp; 16 kHz.

BE: PTA at the better ear.

**Table 3.8:** EHF audiometry: statistical analysis for the effects of Frequency, Ear and Group as well as their interaction (3x2x2 factorial design with repeated measures) tested with a robust rank-based method for analysis of nonparametric data using nparLD package (Noguchi et al., 2012). Analysis was based on a f1-ld-f2 design ANOVA-type statistic (ATS) test, whereby f1 refers to an experimental design with a single between-subjects factor (Group) and f2 refers to two within-subjects factors (Frequency and Ear).

	Statistic	df	p-value
Group	1.323	1.000	<i>0.25</i>
Frequency	22.019	1.717	<b>0.000</b>
Ear	0.395	1.000	<i>0.53</i>
Group:Frequency	1.036	1.717	<i>0.346</i>
Ear:Frequency	0.427	1.969	<i>0.649</i>
Group:Ear	0.348	1.000	<i>0.555</i>
Group:Frequency:Ear	0.220	1.969	<i>0.799</i>

\* significant p-values ( $p < 0.05$ ) are shown in bold.

three-way nor two way interaction between the different factors. There was however a highly significant difference in thresholds between the three frequency bands ( $p < 0.001$ ), whereas no significant main effect for Group or Ear was found.

Similarly, additional nonparametric 4 x 2 factorial design model was used to

**Table 3.9:** EHF audiometry: Statistical analysis for the effects of the listeners calculated measures ( $PTA_{Right}$ ,  $PTA_{Left}$ , PTA, and BE) and Group as well as their interaction (4x2 factorial design with repeated measures) tested with a robust rank-based method for analysis of nonparametric data using nparLD package (Noguchi et al., 2012). Analysis was based on a fl-ld-fl design ANOVA-type statistic (ATS) test, whereby fl refers to an experimental design with a single between-subjects factor (Group) and a single within-subjects factor (Measure).

	Statistic	df	p-value
Group	0.907	1.000	<i>0.341</i>
Measure	7.695	1.389	<b>0.002</b>
Group:Measure	0.154	1.389	<i>0.777</i>

\* significant p-values ( $p < 0.05$ ) are shown in bold.

examine the difference between the two groups for the four combined threshold measures ( $PTA_{Right}$ ,  $PTA_{Left}$ , PTA, and BE). The nparLD ATS test found no significant two-way interaction between Group and Measure nor a main effect of groups, while there was a significant main effect of measure (see Table 3.9).

### 3.3.3 ST

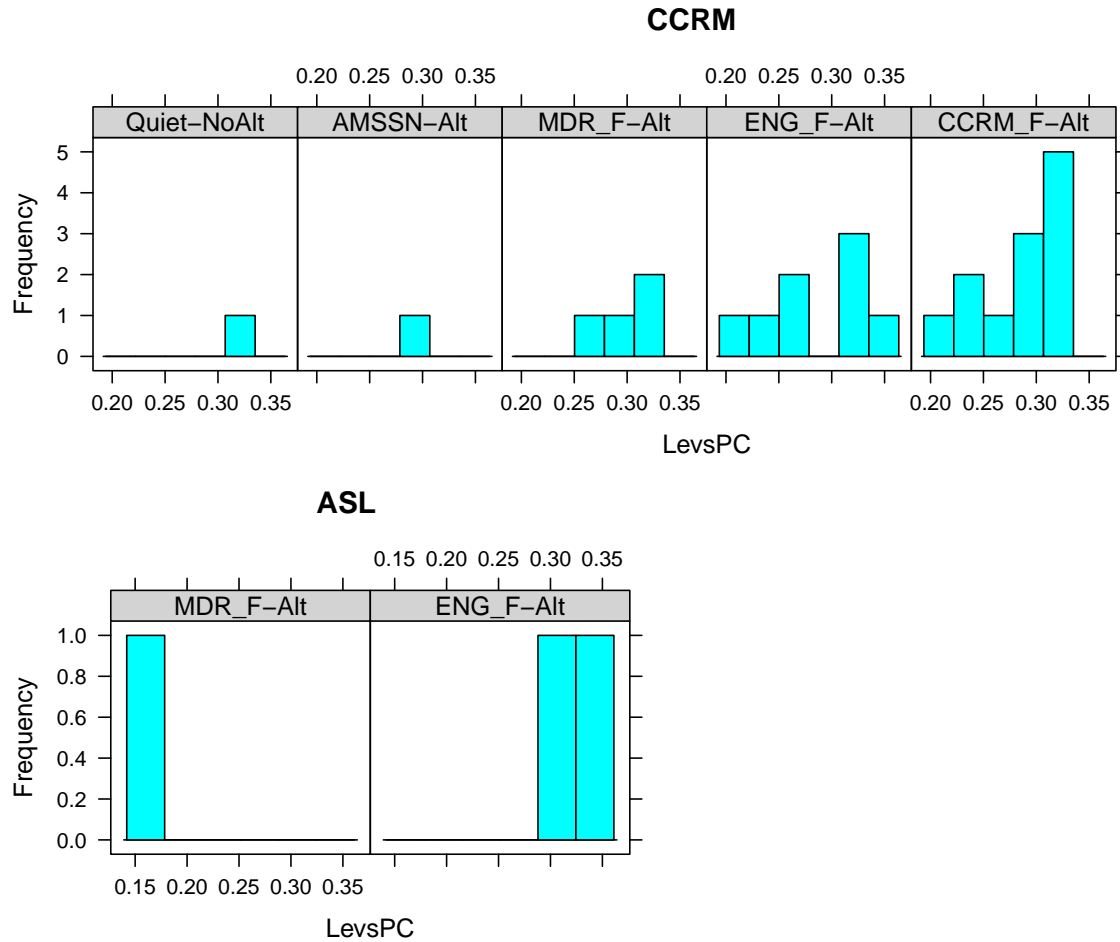
#### Outliers & missing data

As a first step, the listeners adaptive track and psychometric functions were manually inspected for abnormalities. The proportion of correct keywords within the final test trials (LevsPC) was calculated as a measure describing the success of the adaptive procedure. Since the adaptive procedure was set to yield 50%-correct, a successful procedure is expected to have a LevsPC at approximately 50% range. A binomial statistical test was applied to identify observations that significantly differ from 50%. Observations with  $LevsPC \leq 35\%$  were flagged as possible outliers and were further inspected (see Figure 3.6). Interestingly, most of the flagged outliers belonged to the CCRM material with 29 observations out of 258 (6 conditions x 43 listeners), whereas only 3 observations out of 215 (5 conditions x 43 listeners) were flagged for data measured with the ASL speech material.

As expected, most of the identified cases in both materials were for observations measured with the more demanding conditions with speech distractors. In five cases (2 ASL; 3 CCRM) we were able to confidently determine that the listener's true score was near to ceiling, and thus these observations were set to the maximal DC in the task (0.97). In other cases it was not possible to confidently determine the true SRdT, either because the procedure ended after reaching the maximum number of trials before a minimum number of test reversals was obtained (x1 CCRM, x2 ASL), or due to aberrant adaptive tracks (x5 CCRM). Since all these cases belonged to more challenging test conditions with speech distractors, it is very likely that the children's true score is at or beyond the upper DC limit (i.e., at ceiling). Thus, to account for that, rather than removing these observations, which will consequently reduce the statistical power and may not represent the true performance in the group, they were set to a DC of 1, which is above the task's upper DC limit of 0.97.

### **SRdT by age**

Since the present study sample comprised of young children from different age groups from circa 7 to 13 years, developmental age effect was expected, whereby performance was expected to improve with an increasing age due to different maturity effects. This is illustrated by the scatterplots and linear regression lines plotted in Figure 3.7 A-B split by groups for the listeners SRdT obtained across the different test conditions and speech material (ASL / CCRM) as a function of age. Note that smaller SRdT indicate better performance. Age effect was tested against the TD group alone because this group is more heterogeneous and thus expected to display smaller variability than the APD group. Nonetheless, despite the larger spread in the APD group, the group showed similar trend in performance, albeit shifted towards higher SRdT (i.e., poorer performance). The TD regression lines were determined based on a model comparison and outliers trimming procedure described in Section 3.2.4 to improve model prediction. Regular regression lines were found to be the most suitable in describing the relationship between the TD

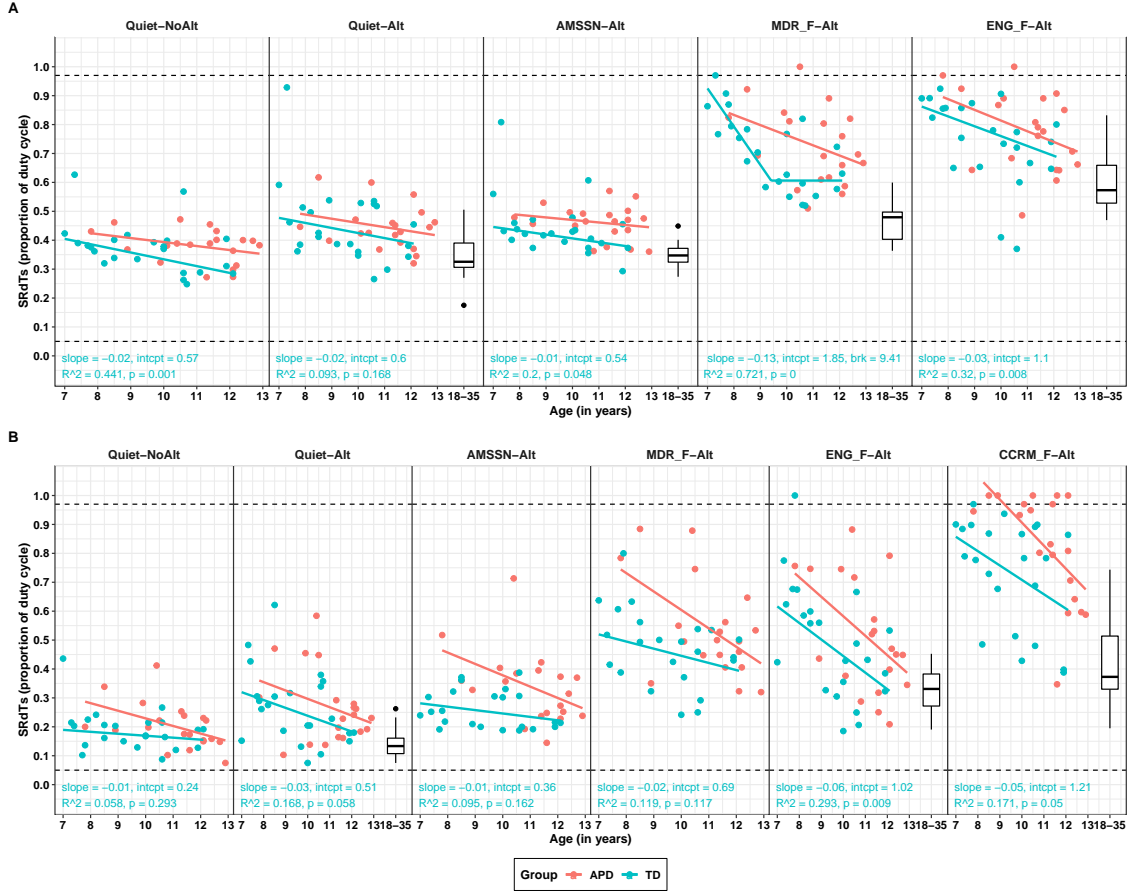


**Figure 3.6:** ST raw data: Frequency of potential outliers with  $\text{LevsPC} \leq 35\%$ .  $\text{LevsPC}$  denotes the proportion of correct keywords within the final test trials.

children performance and age in all test conditions but the MDR\_F condition for the ASL material, where a segmented line was found to give the best fit. MDR\_F segmented line indicated that DC improved with age by circa 0.1 per year until reaching a plateau at the age of 9.5 years.

Looking at Figure 3.7 A-B, it is noticeable that children in both groups showed a larger decrement in performance when presented with speech distractors. The regression lines indicates that the improvement in performance by age was more prominent for speech distractors, with relatively steeper slopes (at least twice as steep) than for the non-speech distractor (AMSSN) or for conditions without a distractor. Furthermore, as expected, CCRM sentences were more intelligible, with





**Figure 3.7:** ST: Scatterplot and linear regression lines for the listeners SRdTs measured with the ASL (A) and CCRM speech material (B) as a function of age. Corresponding regression coefficients and statistics is provided for TD group only. Red indicates data from the APD group and cyan indicates data from the TD control group. Data for normal hearing adults taken from Chapter 2 is shown in the boxplots as a reference.

performance shifted towards lower DC range relative to performance for the ASL speech material. The lower DC meant that the children were able to understand 50% of the sentences with larger portions of the speech information missing.

A closer look at the linear lines shows several interesting trends. The non-speech AMSSN distractor showed to have little-to-no effect on performance, at least in the TD group, where performance was fairly similar to performance in the Quiet conditions. Introducing alternations (as in Quiet-Alt vs. Quiet-NoAlt), seems to hinder intelligibility in both groups, however the effect is relatively small and may not be significant due to the large spread in the APD group. Furthermore, when

comparing the regression lines, there appears to be a relatively larger separation between the groups for data measured with the CCRM material, especially for AMSSN, but also for the speech distractors. However, it is possible that the APD regression lines do not reflect the true population due to the large spread in performance and the small sample size and thus any interpretation should be taken with caution. Another interesting observation is that the children showed little-to-no *masking-release* for speech spoken in an unfamiliar language (MDR\_F) when compared with a distractor spoken in English (ENG\_F). This is in agreement with findings in the adults' study in Chapter 2. Lastly, it is apparent from the figure that performance for CCRM\_F distractor was near-to-ceiling for some children, mostly among the APD group.

#### *Belongs to discussion?*

An exploratory comparison between the children's data measured in the present study with data measured across young NH adults collected in Chapter 2 further highlight the strong developmental trend, with SRdTs still not entirely "adult-like" even at the age of 13 years, especially for speech distractors (see boxplots in Figure 3.7 A-B). The children in both groups seems to be markedly susceptible to competing CCRM sentences and for familiar- or unfamiliar-speech presented with ASL sentences, with performance at the age of 12 years still largely differing from those obtained by the adults. On the other hand, by the age of 12 years, the TD children reached near to "adult-like" performance when CCRM target sentences were presented with ENG\_F speech distractor or when ASL sentences were presented with AMSSN distractor.

Next, age effect was evaluated using LMEM model, with Condition (Quiet-NoAlt, Quiet-Alt, AMSSN, MDR\_F, & ENG\_F), Material (ASL / CCRM) and Age as fixed factors, SRdT as dependent variable and a random intercept for subjects (reference levels: Condition = Quiet-NoAlt; Material = ASL). Note that data for CCRM\_F was excluded from the model since it was only measured for

the CCRM material<sup>2</sup>. Parametric assumption inspection of normal distribution was met, whereas the assumption of homogeneity of variance was marginal ( $p = 0.04$ ). Inspect a model with a fix effect for APDsibling.. A model without three-way interaction between Condition, Material and Age was found to give the best fit (see Table 3.10). Model comparison revealed a highly significant two-way interaction between Condition x Age ( $p < 0.001$ ) as well as a marginal two-way interaction between Condition x Material and Material x Age ( $p = 0.052$ ) which further support the trends seen in the scatterplots.

The significant Material x Age interaction indicates that the developmental trend is different between the two speech materials, with a larger age effect (i.e., steeper slopes) for the ASL sentences, with an average improvement of 0.014 DC per 1 year, which is approximately 8% higher than for the CCRM sentences across the age span. Furthermore, the significant Condition x Material interaction implies that performance in the different test conditions differed between the two speech materials. A post-hoc t-test comparison based on the fitted model given in Table 3.11, revealed a highly significant difference in performance between the speech materials across all five test conditions (all p-values  $< 0.001$ ). The estimated mean difference between the contrast pairs ranged between +0.18 to +0.27, hence, the CCRM speech material was significantly more intelligible than the ASL material, across all test conditions.

Lastly, the highly significant Condition x Age interaction supports the observation in Figure 3.7 A-B, that the magnitude of the age effect was different across the test conditions. These findings raises the following questions – do all the conditions show a significant age effect? Moreover, since the effect of age is not the same across the test conditions, which conditions showed the largest age effect? One possible way to tackle these questions is to compare the separate regression models

---

<sup>2</sup>A separate model for CCRM data only, also gave a significant Condition x Age interaction ( $p < 0.001$ ).

**Table 3.10:** ST: Age effect analysis using LMEM for SRdT<sub>s</sub> measured across condition, speech material and age as fixed factors and a random intercept for subjects. Reference levels: Condition = Quiet-NoAlt, Group = APD, Material = ASL). Note: only data measured with the control group following outliers trimming was included (trimmed TD).

SRdT $\sim$ Condition + Material + Age + Condition:Material + Condition:Age + Material:Age + (1   Subjects)			
Main effects	Df	$\chi^2$	p
Condition:Material	4	9.385	0.052
Condition:Age	4	14.919	<b>0.005</b>
Material:Age	1	3.786	0.052

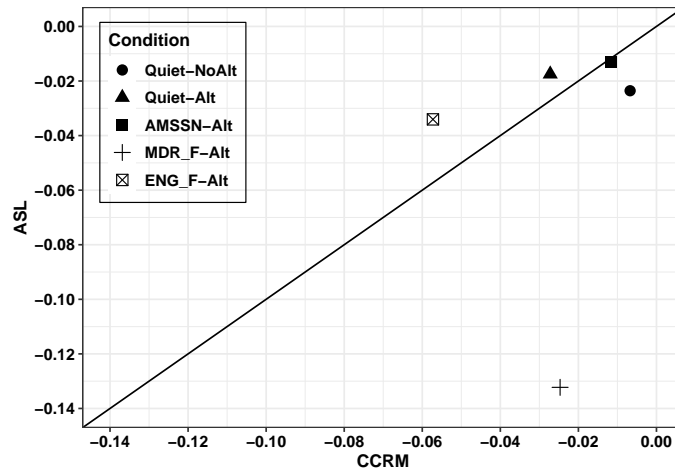
\* significant p-values ( $p < 0.05$ ) are shown in bold.

using F-statistics. Nonetheless, due to the small sample-size and the large number of paired comparisons, such test lacks a statistical power and the results may not reflect the true effect in a larger sample. The TD group regression model's  $R^2$  and p-values are given at the bottom part of Figure A and B. The ASL models p-values indicated a highly significant age effect for ENG\_F, MDR\_F and Quiet-NoAlt condition as well as a marginal effect for AMSSN ( $p = 0.048$ ), whereas no significant age effect was found for Quiet-Alt ( $p = 0.168$ ). As for the CCRM material, there was a highly significant age effect for ENG\_F and a marginal effect for the Quiet-Alt condition ( $p = 0.058$ ) and for CCRM\_F condition ( $p = 0.05$ ) which was not included in the LMEM model, whilst there was no significant age effect found for Quiet-NoAlt, AMSSN and MDR\_F conditions. Furthermore, age was found to be a better predictor (i.e., accounting for larger variance in SRdT) for conditions with speech distractors, with  $R^2$  ranging between 32% to 72% for the ASL material and about 12% to 29% for the CCRM. A comparison between the test conditions regression line slopes fitted for the CCRM (x-axis) and ASL speech material (y-axis) is depicted in Figure 3.8. A possible pattern emerges from the figure, where slopes for the quiet and non-speech conditions are fairly similar across the two speech material (indicated by their proximity to the diagonal line), while, differences between the slopes are relatively larger for speech distractors, in particular for MDR\_F where the slope for the ASL material (-0.13) is about six times steeper than the slope for the CCRM material (-0.02).

**Table 3.11:** ST: Age-effect - post-hoc paired comparison t-test for Condition x Material two-way interaction. The test was performed on the fitted LMEM model and included adjusted least-squared-mean for the random intercept (subjects) using lsmeans package (emmeans package; Lenth, 2020).

ASL - CCRM	Estimate	SE	Df	t-value	p-value	95%-CI
Quiet-NoAlt	0.19	0.03	216.45	7.62	< <b>0.001</b>	0.14 - 0.24
Quiet-Alt	0.20	0.03	216.21	8.01	< <b>0.001</b>	0.15 - 0.25
AMSSN-Alt	0.18	0.03	216.25	7.37	< <b>0.001</b>	0.14 - 0.23
MDR_F-Alt	0.24	0.03	216.21	9.66	< <b>0.001</b>	0.19 - 0.29
ENG_F-Alt	0.27	0.03	216.21	10.90	< <b>0.001</b>	0.22 - 0.32

\* significant p-values ( $p < 0.05$ ) are shown in bold.



**Figure 3.8:** ST: Age effect - a comparison between the regression lines slopes fitted for the CCRM (x-axis) and ASL speech material (y-axis). Test conditions are represented by the different symbols. The diagonal line represents an optimal agreement between the speech materials. Observations falling below the line indicate a steeper slope for the ASL material than for the CCRM material.

### Age-independent z-scores

Age-independent standardised residuals (z-scores) were calculated based on a model prediction for the TD group data using a multiple-case study approach (Ramus et al. (2003); or see section 3.2.4 for more details). Descriptive statistics for the listeners z-scores is given in Table 3.12. Additional boxplots are shown in Figure 3.9 A-B, for the ASL and CCRM speech material respectively. Scores were calculated separately for each test condition, with better performance indicated by lower z-score. The grey area marks the two-tailed 1.96 deviance cut-off for abnormal

**Table 3.12:** ST: Descriptives for standardised residuals (z-scores) calculated for data measured with the ASL and CCRM speech material.

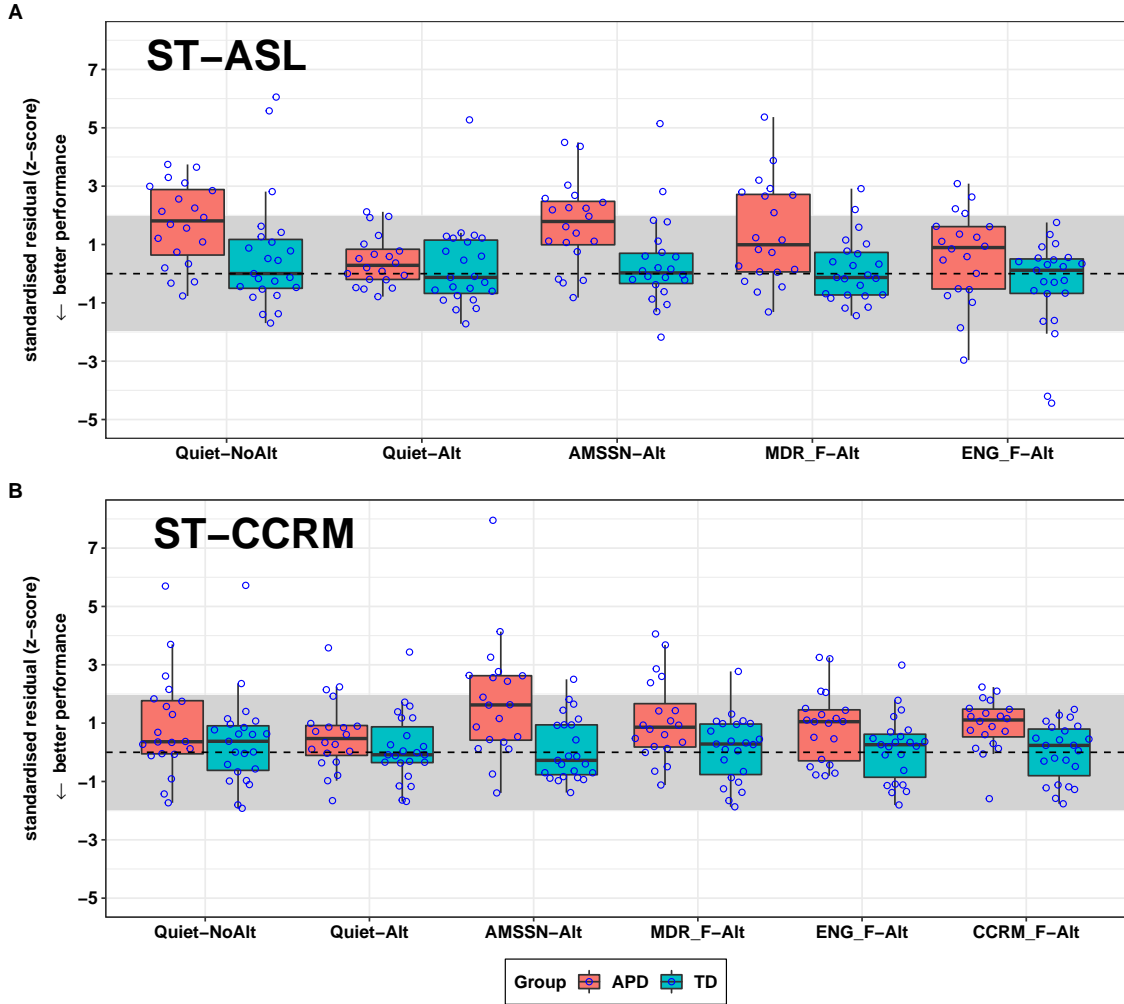
	APD						TD					
	N	median	sd	min	max	abnormal	N	median	sd	min	max	abnormal
<b>ASL</b>												
Quiet-NoAlt	20	1.81	1.39	-0.76	3.74	45.00%	23	0.00	1.96	-1.69	6.05	13.04%
Quiet-Alt	20	0.29	0.87	-0.79	2.12	10.00%	23	-0.13	1.46	-1.72	5.27	4.35%
AMSSN-Alt	20	1.79	1.45	-0.82	4.50	50.00%	23	0.10	2.35	-2.18	9.04	13.04%
MDR_F-Alt	20	0.99	1.75	-1.31	5.37	40.00%	23	-0.13	1.11	-1.44	2.91	8.70%
ENG_F-Alt	20	0.90	1.53	-2.96	3.09	20.00%	23	0.12	1.55	-4.44	1.75	0.00%
<b>CCRM</b>												
Quiet-NoAlt	20	0.36	1.75	-1.73	5.70	20.00%	23	0.38	1.57	-1.92	5.72	8.70%
Quiet-Alt	20	0.47	1.23	-1.66	3.58	15.00%	23	-0.08	1.19	-1.68	3.44	4.35%
AMSSN-Alt	20	1.62	2.03	-1.39	7.95	40.00%	23	-0.28	1.09	-1.38	2.50	4.35%
MDR_F-Alt	20	0.86	1.40	-1.12	4.06	25.00%	23	0.28	1.11	-1.87	2.77	4.35%
ENG_F-Alt	20	1.05	1.22	-0.80	3.25	20.00%	23	0.26	1.14	-1.80	2.99	4.35%
CCRM_F-Alt	20	1.11	0.89	-1.59	2.24	10.00%	23	0.24	0.98	-1.76	1.47	0.00%

abnormal: defined as the percentage of abnormal z-score > 1.96.

score from the theoretical control group mean ( $z = 0$ ), where only about 5% of the normal population is expected to score below and above it. Overall, APD children performance was noticeably poorer in both test material, with higher median z-scores than compared with the TD children. The next paragraphs will cover the examination and statistical analysis of the individuals and group differences separately for each speech material.

### ASL speech material

Surprisingly, a comparison of the groups averaged z-score reveals that the non-switched quiet condition (Quiet-NoAlt) and the switched condition with the nonspeech distractor (AMSSN) yielded the largest separation between the groups, with APD median z-score of 1.81 and 1.79, respectively, laying just within the norms upper limit. Performance of the APD children was also noticeably poorer for conditions with speech distractors (MDR\_F and ENG\_F), each with a median z-score of circa 1, whereas performance for Quiet-Alt condition was fairly similar between the groups.



**Figure 3.9:** ST: Boxplots of the listeners age-independent standardised residuals for data measured with the ASL (A) and the CCRM speech material (B). Residuals were calculated separately for each condition and are based on a model prediction for TD group only. The grey area represents the deviance cut-off for abnormal score ( $SD \pm 1.96$  below and above the TD mean), where about 95% of the normal population is expected to lay within. The dashed line represents the theoretical TD group mean ( $z = 0$ ). Individual scores are indicated by circles. The boxes show the data interquartile range (25th-75th percentile) and the horizontal line indicate the median (i.e., 50th percentile). Values that fall within 1.5 times the interquartile range are indicated by the whiskers.

Within the APD group AMSSN, Quiet-NoAlt and MDR\_F resulted in the highest proportion of abnormal scores<sup>3</sup>. Surprisingly, AMSSN distractor yielded the highest proportion of abnormal scores, where half of the APD children fell outside the norm (20/10, 50%). Followed by the non-switched condition Quiet-NoAlt, where paradoxically and against our expectation 45% of the APD group (9/20) had abnormally poor score, whereas only 10% (2/20) had abnormal score in the switched condition Quiet-Alt. Another interesting finding was that the APD children did not benefit from a release from masking for a speech distractor spoken in an unfamiliar language (MDR\_F) as opposed to a familiar speech spoken in English (ENG\_F), with median scores very similar in both conditions. This sits well with our previous findings with adults where adults showed no benefit for MDR\_F speech masker (see chapter ???). Moreover, while the overall performance was similar in the two conditions, the percentage of abnormal score was twice as large for MDR\_F condition (8/20, 40%) than for ENG\_F condition (4/20, 20%). The proportion of abnormal scores amongst the TD group ranged between 0% to 13% (mean = 7.8%), which is relatively higher than expected in the normal population.

### CCRM speech material

Figure 3.9 B reveals a similar trend for the CCRM sentences, nonetheless with more modest differences between the two groups. Again, AMSSN yielded the largest separation between the groups, where 40% (8/20) of the APD children obtained an abnormal score and with a median score of 1.62, which is relatively close to the +1.96 upper deviance cut-off. In comparison, only 4.3% of the TD children (1/23) had abnormal performance for AMSSN condition. The APD group median score for the speech distractors was approximately 1 (range: 0.86-1.11), however the proportion of abnormal APD children was noticeably smaller than seen for AMSSN, with 25% (5/20) for MDR\_F, 20% (4/20) for ENG\_F, and only 10% (2/20) for

---

<sup>3</sup>with the aim to develop a clinically applicable test that exhibits good sensitivity and specificity rates, we were only interested in identifying children with clinically poor performance. Thus, abnormal score was defined as a one-tailed deviance cut-off of z-score > 1.96, within which circa 97.5% of the normal population is expected to lay.



CCRM\_F distractor. Lastly, in contrast to the ASL material, performance for the CCRM sentences presented in quiet were relatively better without switching (NoAlt) than with switching (Alt). Nonetheless, the spread in performance for the non-switched condition was larger. The percentage of abnormal scores in the TD group were relatively low, ranging between 0 to 8.7% (mean = 4.3%).

A three-way 2 x 2 x 5 factorial design model with repeated measures was used to test the main effects of Group, Material (ASL / CCRM) and Condition as well as their interaction on performance in the task with z-scores as a dependent variable. Note that the model did not include the CCRM test condition with CCRM-type sentences as distractor (CCRM\_F) since there was no comparable condition in the ASL speech material. Inspection of the data revealed that the assumption of normal distribution was rejected for data measured with both speech material. Furthermore, homogeneity of the variance in the APD group was rejected for the ASL corpus. Thus, due to the small sample size and the incomplete fulfilment of parametric statistical methods assumptions, a non-parametric approach was adopted. This was tested with a rank-based ANOVA-type statistic test (ATS) using the *nparLD()* function (*nparLD* package; Noguchi et al., 2012). The analysis was based on a f2-ld-f1 design ATS test, whereby f2 refers to an experimental design with two between-subjects factors (Group & Material) and f1 refers to a single within-subjects factor (Condition). The test results are given in Table 3.13. No significant three-way or two-way interactions were found (Group x Material x Condition, Group x Material, Material x Condition, and Group x Condition; all  $p$ 's  $> 0.05$ ), while there was a highly significant main effect of Group ( $p < 0.0001$ ) and a strong main effect of Condition ( $p < 0.001$ ). Despite some evidence for better performance in the CCRM material, the main effect of Material was not significant ( $p = 0.62$ ). Furthermore, in spite of some apparent differences in performance between the two groups across the different test conditions, these differences were found to be insignificant based on the Group x Condition two-way interaction.

**Table 3.13:** ST: Statistical analysis for the effects of Group, Material, and Condition as well as their interaction (2x2x5 factorial design with repeated measures) tested with a robust rank-based method for analysis of nonparametric data using nparLD package (Noguchi et al., 2012). Analysis was based on a f2-ld-f1 design ANOVA-type statistic (ATS) test, whereby f2 refers to an experimental design with two between-subjects factors (Group and Material) and f1 refers to a single within-subjects factor (Condition).

	Statistic	df	p-value
Group	13.555	1.000	<b>0.000</b>
Material	0.246	1.000	0.62
Condition	3.730	3.450	<b>0.008</b>
Group:Material	0.181	1.000	0.67
Condition:Material	1.473	3.472	0.214
Group:Condition	1.957	3.450	0.109
Group:Material:Condition	0.688	3.472	0.58

\* significant p-values ( $p < 0.05$ ) are shown in bold.

Since data for the CCRM material with CCRM-type distractor was not included in the aforementioned model, a separate 2 x 6 model was computed for the CCRM speech material only. The model included Group and Condition as between- and within-subjects predictors, respectively, with z-scores as the dependent variable using nparLD ATS test (f1.ld.f1 design). The ATS test found a strong significant difference between the groups (Statistic = 10.980, df = 1.000,  $p < 0.001$ ), while no significant main effect was found for Condition (Statistic = 0.819, df = 4.618,  $p = 0.527$ ) nor for Group x Condition interaction (Statistic = 1.215, df = 4.618,  $p = 0.301$ ).

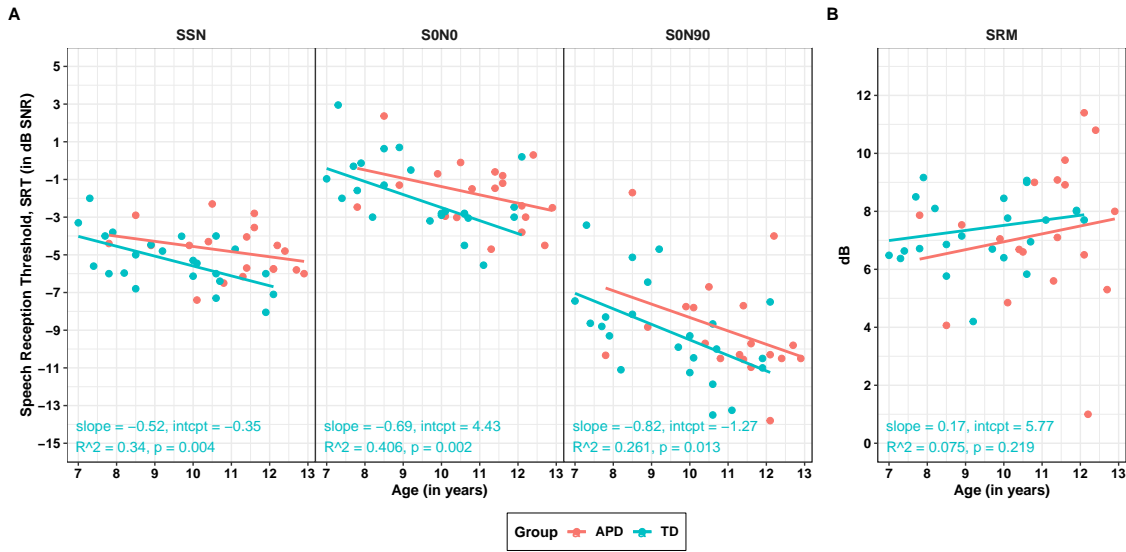
*ROC curves table?*

*Switching effect (derived measures)*

### 3.3.4 LiSNS-UK

#### SRTs by age

The distribution of the listeners SRTs and their corresponding regression lines split by group is shown in Figure 3.10 A for the spatially- collocated (S0N0) and separated condition (S0N90), as well as for the non-spatialised condition where the ASL sentences were presented with a speech-shaped-noise (SSN). The listeners binaural advantage, calculated as the difference between the collocated and separated spatial



**Figure 3.10:** LiSNS-UK: Age-effect - scatterplot and linear regression lines for SRTs obtained for SSN and the spatialised conditions S0N0 (collocated) and S0N90 (separated) (A) and the derived measure SRM (B) as a function of the listeners age. Corresponding regression coefficients and statistics is provided for TD group only. Red indicates data from the APD group and cyan indicates data from the TD control group.

conditions ( $\text{SRM} = \text{S0N0} - \text{S0N90}$ ) is shown in Figure 3.10 B. As in the switching task, age effect was tested against the TD group only, where the regression lines for the TD group were estimated based on a model comparison and outliers trimming procedure to improve the model's fit (model coefficients and statistic are given at the bottom of the figures).

As previously reported by other researchers that used similar test paradigm in children from a similar age group (e.g., Cameron & Dillon, 2007; Murphy et al., 2019), the scatterplots shows a clear developmental trend, with an overall improvement in performance with an increase in age. S0N90 and S0N0 conditions showed the largest age effect, with near to 1 dB improvement in performance per 1 year increase (TD slope: -0.82 & -0.69, respectively). The regression lines slope for SSN conditions was shallower, with roughly half a dB improvement in performance per 1 year increase, with a TD slope of -0.52. Difference in performance with age for the SRM was negligible, with a predicted improvement of circa 1 dB between the age of 7 to 13 years. There was a significant effect of age in all three test

**Table 3.14:** LiSNS-UK: Age effect - LMEM model for SRT with condition (reference level: SSN) and age as fixed factors and a random intercept for subjects. Note: only data measured with the control group following outliers trimming was included.

SRT ~ Condition + Age + (1   Subjects)			
Main effects	Df	$\chi^2$	p
Condition	2	100.356	<b>&lt;0.001</b>
Age	1	13.364	<b>&lt;0.001</b>

\* significant p-values ( $p < 0.05$ ) are shown in bold.

conditions (moderate effect size), with the largest effect for S0N0, accounting for circa 40% of variability in performance, followed by SSN with 34% and about 26% for S0N90. The linear regression fit for SRM showed no significant age effect for SRM ( $R^2 = 0.075$ ,  $p = 0.219$ ).

A two-way factorial design model with repeated measures was used to test the main effects for Condition (SSN, S0N0, & S0N90) and Age with TD group SRTs as a dependent variable. Interaction terms were included as well as a random intercept for subjects. Note that also here the model included only data for the control group. Assumptions of normal distribution and homogeneity were met, and thus a parametric approach was chosen using LMEM. The model with the best fit and main effects are given in Table 3.14. A model with a fixed factor for TD children with APD siblings did not improve the model's fit and was thus removed, suggesting that having an APD sibling did not affect performance in the present study sample. The final model did not include interaction terms and thus indicates that age affected performance in a similar way in the three test conditions. The model revealed a highly significant main effect of Age and Condition ( $p < 0.001$ ).

### Age-independent z-scores

Boxplots of the listeners age-independent standardised residuals z-scores (blue circles) collapsed across the different test conditions are shown in Figure 3.11,

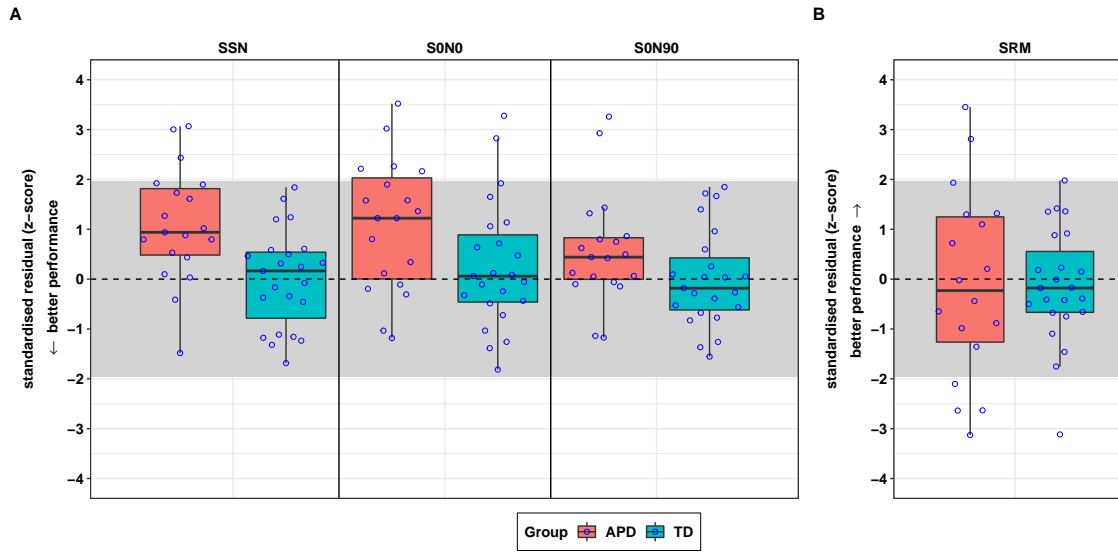
**Table 3.15:** LiSNS-UK standard residuals (z-scores) descriptives by group. abnormal: defined as the percentage of abnormal z-score  $> 1.96$  (SSN, S0N0, & S0N90) and z-score  $< -1.96$  (SRM).

	APD						TD					
	N	median	sd	min	max	abnormal	N	median	sd	min	max	abnormal
SSN	19	0.94	1.14	-1.48	3.07	15.79%	23	0.16	0.98	-1.68	1.84	0.00%
S0N0	19	1.22	1.31	-1.18	3.52	26.32%	23	0.06	1.28	-1.81	3.28	8.70%
S0N90	19	0.44	1.11	-1.17	3.26	10.53%	23	-0.18	0.98	-1.55	1.85	0.00%
SRM	19	-0.44	2.39	-6.77	3.45	26.32%	23	-0.18	1.15	-3.12	1.98	4.35%

separately for the APD group (red) and TD group (cyan). The z-scores were calculated in the exact same way as for ST. Again, the dashed line indicates the theoretical TD group mean of zero, and the grey area indicates the lower and upper limit of the normal population (TD mean  $\pm 1.96$ ). Descriptive statistics collapsed by group and test conditions are given in Table 3.15. Overall, when compared with the control group, the APD children exhibited poorer performance across all three test conditions (i.e., higher z-score) as well as for the derived SRM measure (i.e., lower z-score).

S0N0 and SRM yielded the largest separation between the groups, however the spread in scores was relatively large and the percentage of abnormal performance in the APD group was rather small, with only circa 26% (5/19) in each condition. Whereas only about 16% (3/19) and 10% (2/19) of the APD children had abnormal score for SSN and S0N90, respectively. No abnormal performance was obtained in the TD group for SSN and S0N90, while two TD children (~9%) had abnormal score for S0N0 condition and one child for SRM. Nonetheless, when excluding the TD outliers that were trimmed during the z-score calculation procedure, all the TD observations were within the norms.

Group differences for the SSN and the spatialised conditions S0N0 and S0N90 were tested with a 3 x 2 factorial design LMEM model with z-score as a dependent variable and random intercept for subjects (reference levels: Condition = SSN, Group = APD). Model assumptions for normal distribution and homogeneity of



**Figure 3.11:** LiSNS-UK: Boxplots of the listeners age-independent standardised residuals (open circles) for data measured with LiSNS-UK task (A) and the derived measure SRM (B). Residuals were calculated separately for each condition and are based on a model prediction for TD group only. The grey area represents the deviance cut-off for abnormal score ( $SD \pm 1.96$  below and above the TD mean), where about 95% of the normal population is expected to lay within. The dashed line represents the theoretical TD group mean ( $z = 0$ ). The boxes show the data interquartile range (25th-75th percentile) and the horizontal lines indicate the median (i.e., 50th percentile). Values that fall within 1.5 times the interquartile range are indicated by the whiskers.

variance were verified. The model that was found to give the best fit did not include a two-way interaction term between the main effects Group and Condition, thus suggesting again that the two groups behaved in a similar way in the three test conditions (see Table 3.16). Comparison of the full model with a simplified model without the term of interest revealed a significant effect of Group ( $p < 0.05$ ), while no significant difference in performance between the three conditions was found ( $p = 0.149$ ).

### 3.3.5 ENVASA

Due to technical problems, observations for six listeners are missing (TD: 2; APD: 4), resulting in a total sample-size of 21 and 17 for the TD and APD group, respectively. Initial inspection of the individuals performance was performed to ensure that the task instructions were followed and well understood. Performance

**Table 3.16:** LiSNS-UK: Group difference - LMEM model for age-independent z-scores with condition and group as fixed factors (reference levels: Condition = SSN, Group = APD) and a random intercept for subjects.

$z \sim \text{Condition} + \text{Group} + (1 \mid \text{Subjects})$			
Main effects	Df	$\chi^2$	p
Condition	2	3.809	<i>0.149</i>
Group	1	8.673	<b>0.003</b>

\* significant p-values ( $p < 0.05$ ) are shown in bold.

for the reference condition (single incongruent background at a high SNR), which is expected to least impact performance, was compared with a cut-off criterion of 56%, calculated as 2 SD from the TD group mean ( $84\% \pm 14\%$ ). Individuals with performance below the cut-off criterion were excluded from the analysis. One TD listener aged 7 years old scored 45 % and was thus excluded, resulting in a total of 20 listeners in the TD group.

### %-correct by age

The ENVASA measurements followed the same factorial design as used by Leech et al. (2009), with 2 background types (single/dual) x 4 SNRs (low: -6, -3 dB; high: 0 +3 dB), resulting in a total of 92 responses (%-correct, PC) per listener or between 10 to 11 test items per background-SNR combination. Because of the small number of test items per condition, responses were averaged into three measures: 1) *single background*, 2) *dual backgrounds*, and 3) *combined background* which reflects the overall performance across the two background types.

The relationship between performance and age was inspected in the same way as carried out for the other auditory tasks, with the listeners average response plotted as a function age, with linear regression lines, model coefficients and statistics for the trimmed TD group (see Figure 3.12). The regression lines revealed a noticeable developmental trend in all three measures, where performance improved with increasing age. A single linear regression line with a monotonic increase in

**Table 3.17:** ENVASA: Age effect - LMEM model for PC (%-correct) in the three background measures single, dual, & combined background/s (reference levels: Background=single-background, APDsibling=1), and age as fixed factors and random intercept for subjects. Note: only data measured with the control group following outliers trimming was included.

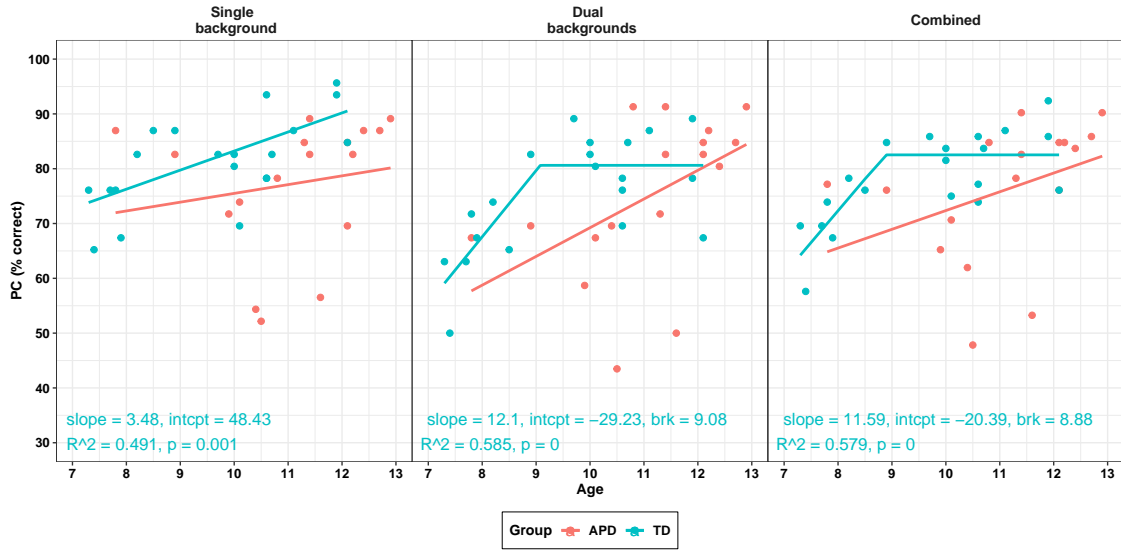
PC ~ Background + Age + APDsibling + (1   Subjects)			
Main effects	Df	$\chi^2$	p
Background	2	21.173	<b>&lt;0.001</b>
Age	1	17.516	<b>&lt;0.001</b>
APDsibling	1	3.834	0.050

\* significant p-values ( $p < 0.05$ ) are shown in bold.

performance by age was found to best fit performance for a single background, with an increase of circa 3.5% in PC per year. Performance for dual backgrounds and the combined score on the other hand were best described using segmented linear regression models, with an increase of PC by circa 12% per year until the age of 9 years, where PC plateaued thereafter. The effect of age was statistically tested using an LMEM model with the three background measures, age, and APDsibling, which indicates whether a child has an APD sibling (1) or not (0) as fixed factors and a random intercept for subjects with PC as dependent variable (reference levels: Background = single-background, APDsibling = 1). A model without an interaction term was found to give the best fit (see Table 3.17). Model comparison revealed a highly significant main effect for age and condition ( $p < 0.001$ ). This is in agreement with the Krishnan et al. (2013) study where they found a strong developmental effect across normal-hearing typically developing children in a similar age range to those measured in the present study.

The main effect of APDsibling was marginal ( $p = 0.05$ ) with an estimated mean difference of 5.15 (SE = 2.50, CI = 0.24 - 10.05). Although there is no significant interaction between Background and APDsibling, a post-hoc comparison finds no significant difference between APD-siblings and non-siblings within a condition...





**Figure 3.12:** ENVASA: Scatterplot and linear regression lines for the listeners' PC (%-correct) as a function of age for single background, dual backgrounds and the combined measure. Red indicates data from the APD group and cyan indicates data from the TD control group.

### Age-independent z-scores

For further analysis, age was controlled for using the same multiple-case approach method described in Section 3.2.4. Boxplots of the age-independent z-scores for the three ENVASA measures are shown in Figure 3.13, with larger z-score indicating better performance. The grey area indicates the upper and lower cut-off  $\pm 1.96$  for normal score, where scores of about 95% of the normal population are expected to lay within. Surprisingly, the less demanding condition with single competing background yielded the largest separation between the group with a median z-score of roughly -1, while the median performance for dual backgrounds and the combined score was relatively similar to those in the control group, albeit with larger spread. The percentage of abnormal APD scores was relatively low, with circa 29% (5/17) for the combined score, 24% (4/17) for single background and 18% (3/17) for dual backgrounds condition. There was only one case of abnormal score in the TD group for single background (5%, 1/20) when trimmed TD outliers are included.

A two-way interaction between Group and Condition (2 x 3 factorial design

**Table 3.18:** ENVASA: Descriptive and statistics of the listeners age-independent standard residuals (z-scores) split by groups and test measures.

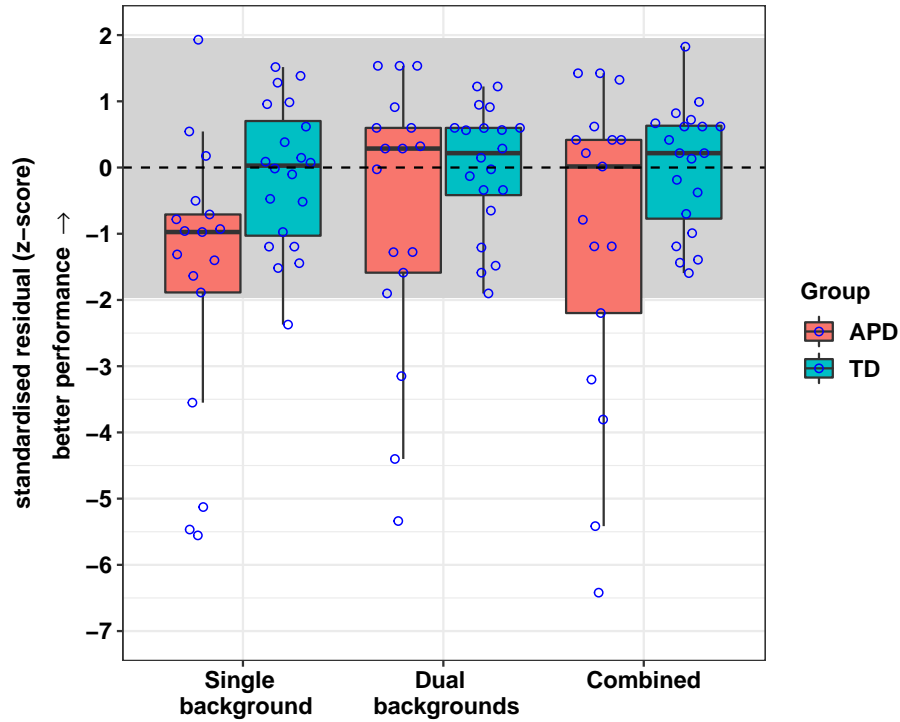
background	APD						TD						Wilcoxon rank-sum test			
	N	median	sd	min	max	abnormal	N	median	sd	min	max	abnormal	95%-CI	p	r	magnitude
Single	17	-0.97	2.11	-5.56	1.93	23.53%	20	0.03	1.08	-2.37	1.52	5.00%	-2.27 - -0.34	<b>0.02</b>	0.39	moderate
Dual	17	0.29	2.07	-5.34	1.54	17.65%	20	0.22	0.95	-1.90	1.22	0.00%	-1.56 - 0.62	<i>0.66</i>	0.08	small
Combined	17	0.02	2.39	-6.42	1.42	29.41%	20	0.22	0.95	-1.59	1.83	0.00%	-1.81 - 0.4	<i>0.29</i>	0.18	small

\* significant p-values ( $p < 0.05$ ) are shown in bold.

data with repeated measures) was tested with a non-parametric robust aligned rank test using *npIntFactRep* package (Feys, 2015). Mauchly's test indicated that the assumption of sphericity for the two-way interaction term had been violated ( $p < 0.001$ ), therefore the degrees of freedom was corrected using Greenhouse-Geisser estimate of sphericity ( $\varepsilon = 0.55$ ). The test showed a significant two-way interaction between Group and Condition [ $F(1.64, 57.57) = 10.82$ ,  $p < 0.001$ ]. Difference between the groups were examined using unpaired two samples Wilcoxon rank-sum test with permutation ( $N=999999$ ) which is a t-test equivalent for non-parametric data (*coin::wilcox\_test()*; Hothorn et al., 2006). Groups descriptives collapsed by the three test measures as well as p statistics and effect size r are given in Table 3.18. Performance of the APD children was significantly poorer than of the TD children in the single background condition ( $p < 0.05$ , moderate effect), whereas there was no significant difference between the groups in the dual backgrounds or the combined background measure (both p's  $> 0.05$ ).

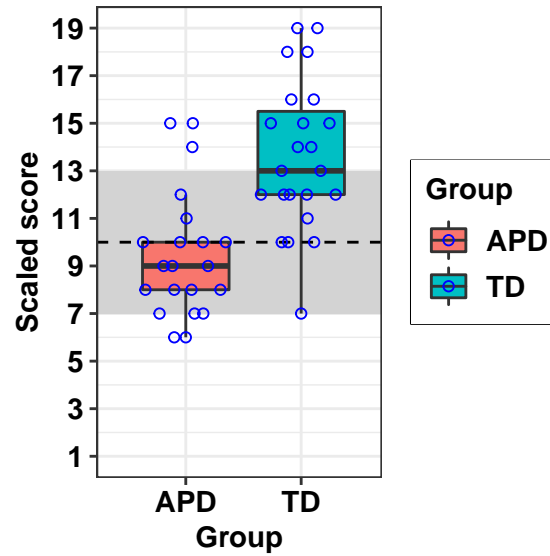
### 3.3.6 CELF-RS

The children's raw scores were converted into age-corrected scaled scores using the CELF-5 UK Recalling Sentences subtest standardised norms ( $M = 10$ ,  $SD = 3$ ). Boxplots of the children's scaled scores split by groups are given in Figure 3.14. The grey area indicates the upper and lower limit among the normal population ( $\pm 1$  SD). On average, performance was within the norms range in both the APD group ( $Mdn = 9$ ) and the TD group, albeit laying within the upper limit ( $Mdn = 13$ ). Thus, although the majority of the APD children expressive language skills were within the norms, the figure shows a clear difference in performance between the



**Figure 3.13:** ENVASA: Listeners’ age-independent standardised residuals for single background, dual backgrounds & the combined measure. Residuals were calculated separately for each condition and are based on a model prediction for TD group only. The grey area represents the deviance cut-off for abnormal score ( $SD \pm 1.96$  below and above the TD mean), where about 95% of the normal population is expected to lay within. The dashed line represents the theoretical TD group mean ( $z = 0$ ).

group, where the TD children expressive language skills are noticeably better. Almost half of the TD children obtained a scaled score above the average and none exhibited abnormal scores. On the other hand, only three APD children performed above the average and performance of two children was considered abnormal (scaled score  $< 7$ ). An independent-samples t-test with bootstrapping ( $n=9999$ ) was computed using *boot.t.test()* function (MKinfer package; Kohl, 2020) to compare the listeners scaled scores in the two groups (parametric data assumptions were met). There was a highly significant difference in scaled scores between the APD ( $Mdn = 9.0$ ,  $SD = 2.7$ ) and TD group ( $Mdn = 13.0$ ,  $SD = 3.1$ ) [ $t(41.81) = -4.71$ ,  $p < 0.001$ ].



**Figure 3.14:** CELF-RS: Boxplots for CELF-5 UK Recall Sentences subtest scaled scores by groups. The dashed line represents the norms mean and the grey area indicates the upper and lower limit average performance in the normal population ( $\pm 1$  SD).

### 3.3.7 Questionnaires

#### CCC-2

Data for one TD listener was flagged as inconsistent using the test scorer and was thus removed from the analysis. The groups descriptives for the parental reports in the different sub-scales as well as the GCC and SIDC composites are given in Table 3.19. GCC stands for general communication composite, calculated by taking the sum for scaled scores A to H. It is used to clinically identify abnormal communication skills, defined by a GCC  $< 55$  (10<sup>th</sup> percentile). The SIDC stands for social-interaction deviance composite [ $\text{sum}(\text{E}+\text{H}+\text{I}+\text{J})-\text{sum}(\text{A}+\text{B}+\text{C}+\text{D})$ ], where in combination with abnormal GCC score, the SIDC can be used to identify the child's primary difficulty, whereby, a positive SIDC is indicative of a predominantly structural language deficit (referred here as DLD), and a negative SIDC reflects social communication problems and is indicative of autistic spectrum disorder (ASD) traits (Bishop, 2003; Norbury, 2014).

Boxplots of the groups scaled scores in the ten sub-scales and a scatterplot depicting the relationship between GCC and SIDC are shown in Figure 3.15 A-B,

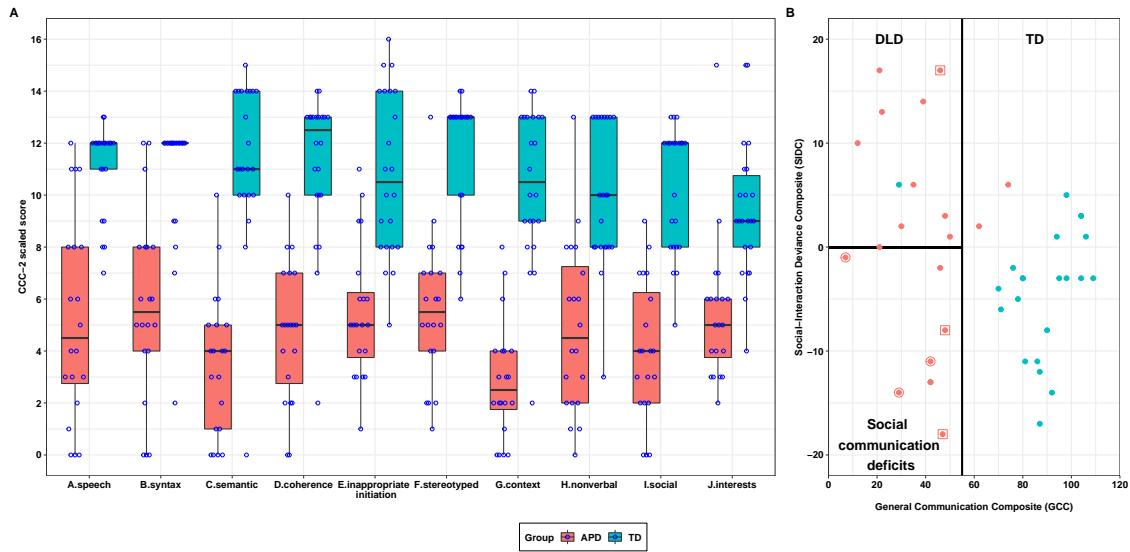
**Table 3.19:** CCC-2 subscales descriptives split by groups.

Measure	APD					TD				
	N	median	sd	min	max	N	median	sd	min	max
A.speech	20	4.5	3.96	0	12	22	12.0	1.72	7	13
B.syntax	20	5.5	3.61	0	12	22	12.0	2.49	2	12
C.semantic	20	4.0	2.78	0	10	22	11.0	3.23	0	15
D.coherence	20	5.0	2.68	0	10	22	12.5	2.87	2	14
E.inappropriate.initiation	20	5.0	2.61	1	11	22	10.5	3.17	5	16
F.stereotyped	20	5.5	2.82	1	13	22	13.0	2.52	6	14
G.use.of.context	20	2.5	2.28	0	8	22	10.5	2.97	2	14
H.nonverbal	20	4.5	3.31	0	13	22	10.0	2.75	3	13
I.social	20	4.0	2.68	0	9	22	12.0	2.41	5	13
J.interests	20	5.0	2.84	2	15	22	9.0	2.63	4	15
GCC	20	42.0	16.38	7	74	22	88.5	17.38	29	109
SIDC	20	1.5	10.70	-18	17	22	-3.0	6.14	-17	6

GCC, General Communication Composite  $\text{sum}(A+B+C+D+E+F+G+H)$ ;

SIDC, Social Interaction Deviance Composite  $\text{sum}(E+H+I+J) - \text{sum}(A+B+C+D)$

respectively. A striking 90% of the APD children (18/20) obtained a scaled score below the 5th percentile two or more times, which has been found to indicate clinically significant communication problems (Bishop, 2003), whereas, only one such case (out of 22) was found in the TD group. The single-value GCC composite showed the exact same proportion of abnormal scores in both groups when a cut-off value of 55 was used, where only one TD child had abnormal communication skills (see Figure 3.15 B). Half of the APD children with abnormal GCC score (45%, 9/20) exhibited a score pattern that is indicative of DLD, whereas the other half exhibited a negative SIDC, indicating social communication deficits as the primary difficulty. Interestingly, out of the nine APD children who fell within the later category, three were reported by their parents to have HF-ASD diagnosis, and an additional two children were undergoing an ASD assessment at the time of testing (see scores marked with open circles and rectangles in Figure 3.15 B). Difference in GCC between the two groups was tested using an independent-samples t-test with bootstrapping (*MKinfer::boot.t.test()*,  $n=9999$ ; Kohl, 2020). There was a highly significant difference in GCC between the APD ( $Mdn = 42.0$ ,  $SD = 16.4$ ) and TD group [ $Mdn = 88.5$ ,  $SD = 17.4$ ;  $t(39.80) = -9.42$ ,  $p < 0.001$ ].



**Figure 3.15:** CCC-2 parental reports for the APD (red) and TD group (cyan). (A) Boxplots for scaled scores in the ten sub-scales. (B) Scatterplot for General Communication Composite (GCC) as a function of Social-Interaction Deviance Composite, (SIDC). APD children with diagnosed high-functioning Autism (HF-ASD) are denoted with open circles. APD children with undergoing ASD assessment on the day of testing are marked with open rectangles. The lines indicates the GCC cut-off criteria for typically developing children (TD) SIDC scores indicative of predominantly structural developmental language disorder (DLD) and more social communication deficits (cf. Norbury, 2013).

## ECLIPS

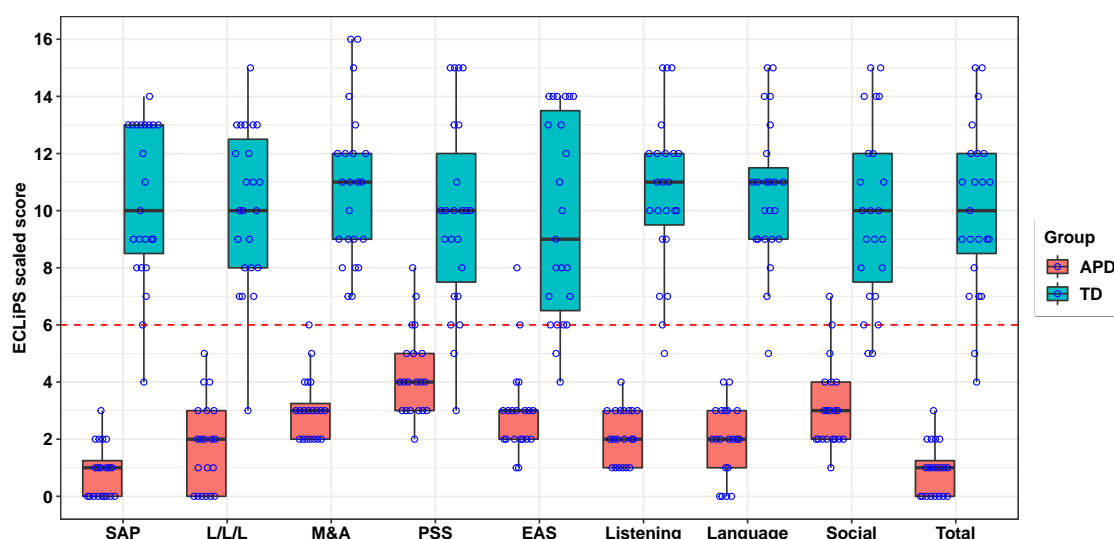
Descriptives of the ECLiPS parental report scaled scores for the different subscales and composite measures split by groups is given in Table 3.20 and depicted in Figure 3.16. A score below the 10<sup>th</sup> percentile (corresponding to a scale score of circa 6) is generally considered as clinically significant listening and processing difficulties (Barry & Moore, 2014). Overall, the ECLiPS was able to well separate between the two groups across all the different sub-scales. All APD children exhibited abnormal Total score, whereas only two TD children (out of 22) obtained abnormal Total score.

A closer look at the boxplots in Figure 3.16 reveals a clear difference in the distribution of the scaled scores across the two groups, with relatively larger spread for the TD group. Inspection of the groups Total score revealed that data for the APD group did not follow a normal distribution and that the assumption of homoscedasticity was violated ( $p < 0.001$ ). Thus, group difference for the listeners

**Table 3.20:** ECLiPS descriptives split by groups and sub-scales.

Measure	APD					TD				
	N	median	sd	min	max	N	median	sd	min	max
SAP	20	1	0.93	0	3	23	10	2.77	4	14
L/L/L	20	2	1.55	0	5	23	10	2.78	3	15
M&A	20	3	1.10	2	6	23	11	2.69	7	16
PSS	20	4	1.53	2	8	23	10	3.37	3	15
EAS	20	3	1.64	1	8	23	9	3.52	4	14
Listening	20	2	0.93	1	4	23	11	2.69	5	15
Language	20	2	1.28	0	4	23	11	2.48	5	15
Social	20	3	1.52	1	7	23	10	3.15	5	15
Total	20	1	0.91	0	3	23	10	2.92	4	15

SAP = Speech & Auditory Processing; L/L/L = Language, Literacy & Laterality;  
M&A = Memory & Attention; PSS = Pragmatic & Social skills; EAS = Environmental  
& Auditory sensitivity; Listening = (SAP + PSS) / 2; Language = (L/L/L + M&A) / 2;  
Social = (PSS + EAS) / 2; Total = mean of all sub-scales

**Figure 3.16:** ECLiPS parental report scaled scores split by groups and sub-scales.

Total score was examined using a non-parametric independent two samples Wilcoxon rank-sum test with permutation ( $N=999999$ , `coin::wilcox_test()`; Hothorn et al., 2006), showing a highly significant difference between the groups with a large effect-size ( $p < 0.001$ ,  $r = 0.86$ ).

## 3.4 Overall performance

An overview of the children's performance split by group is given in Figure 3.17 providing an overlook at individuals that performed outside the norm in one or more tasks (filled black cells). Abnormally poor performance for the listeners age-independent scores was defined using standardised norms for the CELF-RS, ECLiPS and the CCC-2 data or was defined as a one-tailed cut-off of  $\pm 1.96$  (where circa 97.5% of the normal population is expected to lay within) for the rest of the tasks.

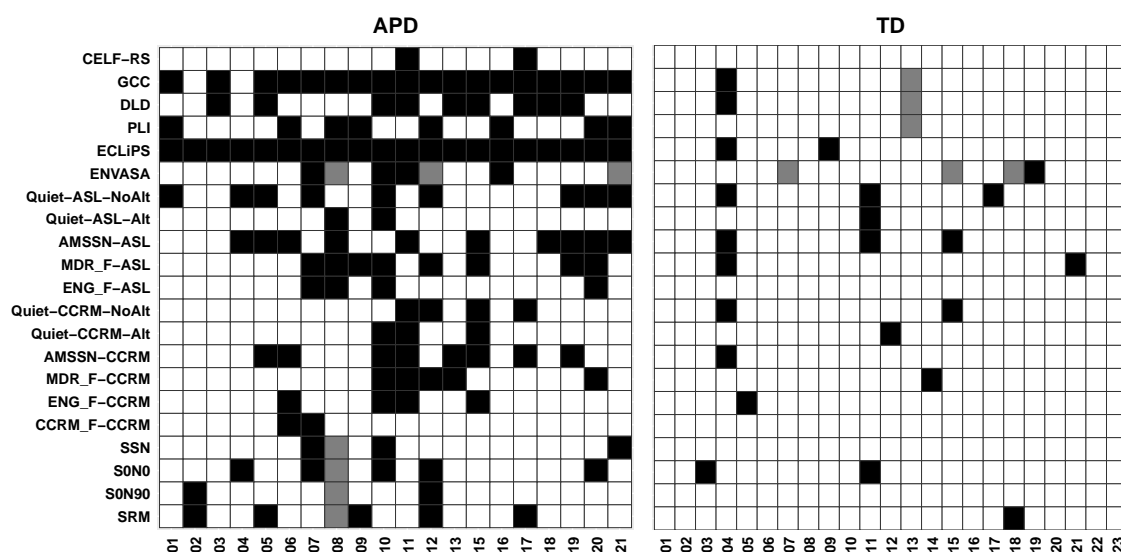
As seen in the figure, the proportion of abnormal scores across the APD group is substantially higher than in the TD group. The majority of the APD children (80%, 16/20) performed abnormally in at least two test conditions either in the ST or LiSNS-UK task, whereas there were only three cases (13%, 3/23) in the TD children.

Another interesting observation is that apart from one TD child, which experienced difficulties in various measures including the CCC-2, none of the other TD children experienced language difficulties. This is in contrast to the APD group where 90% (18/20) of the children experienced some kind of language deficit. The CELF-RS has been reported to be a good marker for children with DLD, nevertheless, the results of the present study suggests otherwise. While performance in the APD group was noticeably poorer than in the TD group, only two APD children obtained abnormally poor CELF-RS score, whereas nearly half of the APD children (45%, 9/20) exhibited a CCC-2 score indicative of DLD, and the about the remaining half (40% 8/20) obtained a CCC-2 score indicative of pragmatic language and social communication deficit (PLI).

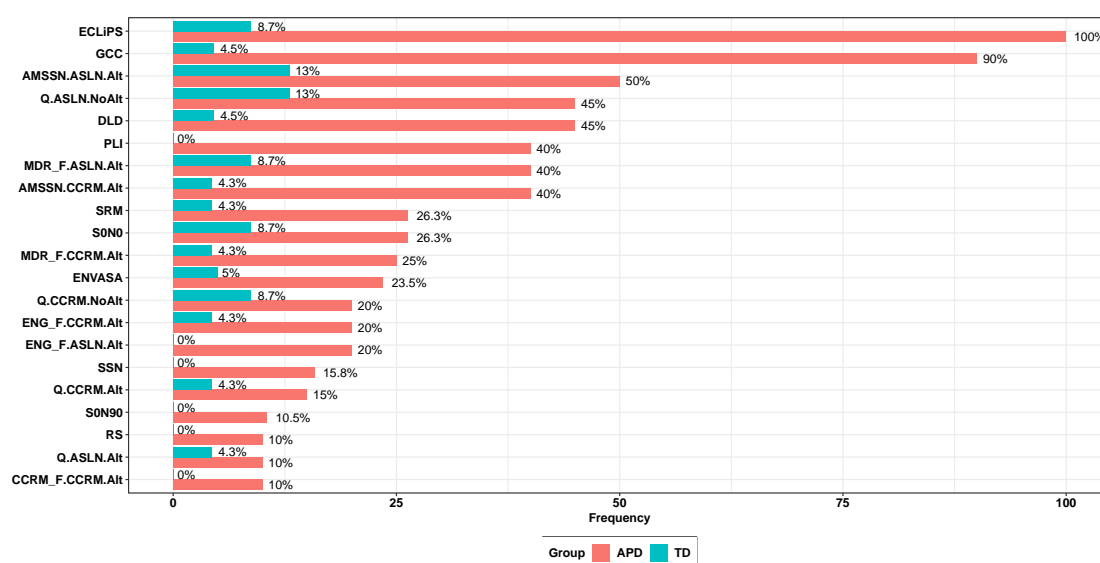
Potential experimental bias of reporters when recruited due to an informed group affiliation? email Courtenay!

The proportion of abnormal scores by measure or task split by group is shown in Figure 3.18. Both the ECLiPS total score and the CCC-2 GCC sum score





**Figure 3.17:** Overall performance: Abnormal (black cells) and normal (empty cells) performance in the present study test battery of individuals from the APD group (n=20) and the TD group (n=23). Missing data is marked by the grey cells.



**Figure 3.18:** Overall performance: Proportion of abnormal score per measure or task split by groups.

resulted in the largest separation between the groups. Out of the auditory tasks, the tests conditions that resulted in the highest proportion of abnormal scores in the APD group were AMSSN (ASL: 50%, CCRM: 40/%), Quiet-ASL-NoAlt (45%) and MDR\_F-ASL (40%), whereas only 26% of the APD children had abnormal SRM score.

### 3.4.1 Unsupervised machine learning (UML)

see Bradshaw paper..

Idea: separate PCA per material (CCRM/ASL) and d' comparison or calculation of d' per condition..

### 3.4.2 Interaction between measures

The present study involved a large number of test conditions and various measures assessing different skills. For example, the ST data alone comprises of 11 different conditions (x5 ASL, x6 CCRM speech material). Another set of measures consisting of the CELF-RS, ECLiPS and the CCC-2 taps into language and communication related skills, whereby the latter two consists of a sum of 15 different sub-scales and have been shown to strongly correlate with one another (Barry & Moore, 2014). Examining the extent to which the groups performance is explained by such a large number of measures will result in a very conservative significance level in order to minimise Type-I error (false positive), and could increase Type II error rate (false negative) (McDonald, 2014). Since the measures within the ST and within the language dataset are expected to strongly correlate, it was decided to use an exploratory data analysis technique using Principal Components Analysis (PCA). PCA is a technique used to reduce a large number of correlated parameters into a smaller set of components that together explain a considerable amount of the variability in the large dataset. Whereby, each of the PCA components is composed of a linear combination of the input parameters (James et al., 2013). PCA was performed separately for the ST and language data set using FactoMineR package (Lê et al., 2008) with scaled units and will be discussed separately below.

## ST

PCA for the ST z-scores comprised of 11 input variables and a sample size of 43. Sample size adequacy for PCA was verified using Kaiser-Meyer-Olkin test (psych::KMO; Revelle, 2020), with an overall KMO of 0.76 ('good'; Field et al., 2012),

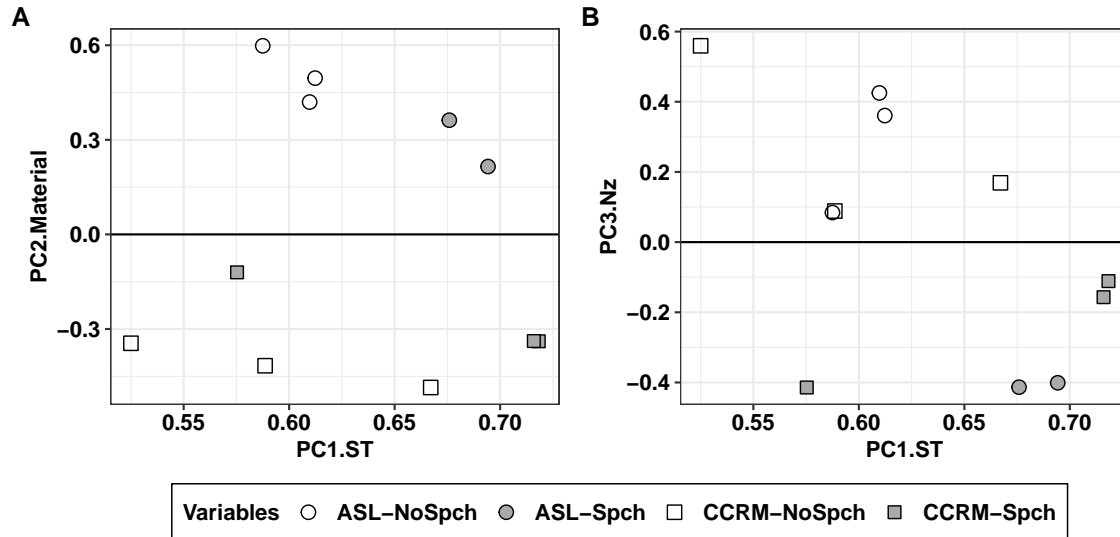
and a KMO range between 0.66 to 0.85 across the conditions. Bartlett's sphericity test was significant [ $\chi^2(55) = 190.36$ ,  $p < 0.001$ ], indicating that the correlations between the different items were large enough for a PCA. Table 3.21 shows the variables loadings (no rotation was applied), their eigenvalues and percentage of variance explained. Loadings are indicators of substantive importance of a given variable to a given component (Field et al., 2012). The first three components were used, yielding eigenvalues  $> 1$  (Kaiser's criterion), explaining together circa 67% of the variance in the data. The first component (PC1.ST) accounted for the largest portion of spread in the data of 40.6% and was interpreted as an overall measure for performance in the switching task with relatively high loadings across all input variables. The remaining components explained each circa 16% and 11% of the variance (ascending order). Figure 3.19 illustrates the different dimensions in the data captured by the three PCA components. Clustering in the second component (PC2.Material) reflected differences in performance across the two speech materials (ASL & CCRM). The third component (PC3.Nz) reflected the degree of distractibility introduced by speech distractors (MDR\_F, ENG\_F, CCRM\_F) irrespective of the speech material used, resulting in decrement in performance when compared with non-speech distractors or target-only conditions (Quiet and AMSSN). Boxplots of the listeners weighted scores for the PCA components split by group is shown in Figure 3.20. PC1.ST shows to separate very well between the two groups, with very little overlap in scores between the TD group and the majority of the APD children. Whereas separation between the two groups in the remaining components are noticeably smaller.

Figure 3.21 illustrates the relationship between the listeners weighted scores based on the three PCA components (PC1.ST, PC2.Material and PC3.Nz) and three calculated composites composed from the listeners z-scores based on the interpretation stated above; where *ST* denoted the listeners' aggregated overall score across all ST conditions, and the two calculated discrepancy composites denoted as *Material* and *Nz*. The Material composite was calculated by subtracting the mean

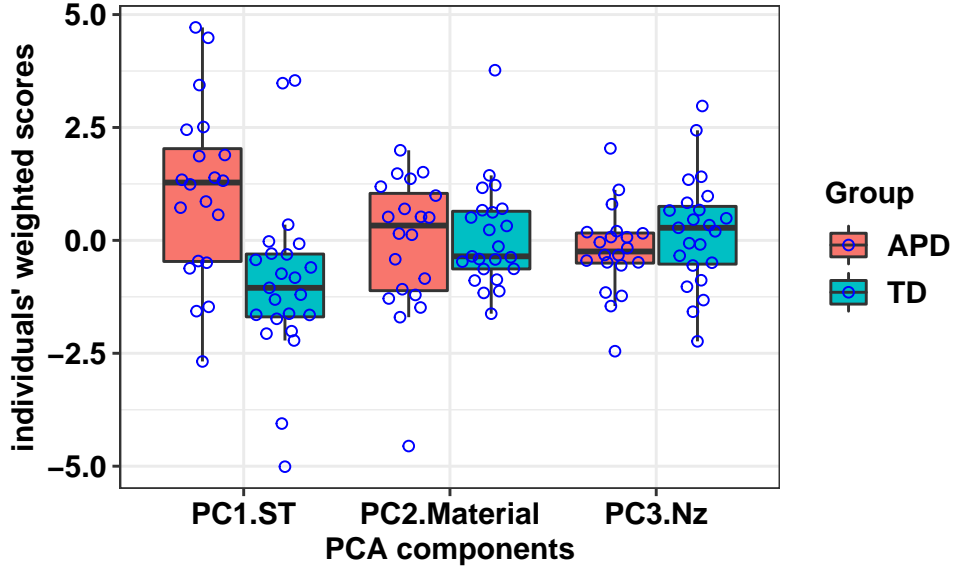
**Table 3.21:** Switching task PCA: Input variables loading.

Item	PC1.ST	PC2.Material	PC3.Nz
Q-ASLN-NoAlt	<b>0.59</b>	<b>0.60</b>	0.08
Q-ASLN-Alt	<b>0.61</b>	<b>0.42</b>	<b>0.43</b>
AMSSN-ASLN.Alt	<b>0.61</b>	<b>0.50</b>	<b>0.36</b>
MDR_F-ASLN-Alt	<b>0.68</b>	<b>0.36</b>	<b>-0.41</b>
ENG_F-ASLN-Alt	<b>0.69</b>	0.22	<b>-0.40</b>
Q-CCRM-NoAlt	<b>0.52</b>	<b>-0.35</b>	<b>0.56</b>
Q-CCRM-Alt	<b>0.59</b>	<b>-0.42</b>	0.09
AMSSN-CCRM-Alt	<b>0.67</b>	<b>-0.49</b>	0.17
MDR_F-CCRM-Alt	<b>0.72</b>	<b>-0.34</b>	-0.11
ENG_F-CCRM-Alt	<b>0.72</b>	<b>-0.34</b>	-0.16
CCRM_F-CCRM-Alt	<b>0.58</b>	-0.12	<b>-0.41</b>
eigenvalue	4.46	1.73	1.21
variance (%)	40.52	15.72	10.98
cumulative variance (%)	40.52	56.24	67.22

|loading| >0.3 are highlighted in bold.



**Figure 3.19:** Switching task PCA: Scatterplot for the input variables as a function of PCA components: PC1.ST vs. PC2.Material (A), PC1.ST vs. PC3.Nz (B). Loadings for ASL conditions are indicated by circles and loadings for CCRM conditions are indicated by rectangles. Filled shapes denotes conditions with speech distractors (Spch) and non-filled shapes denote nonspeech conditions (No-Spch).

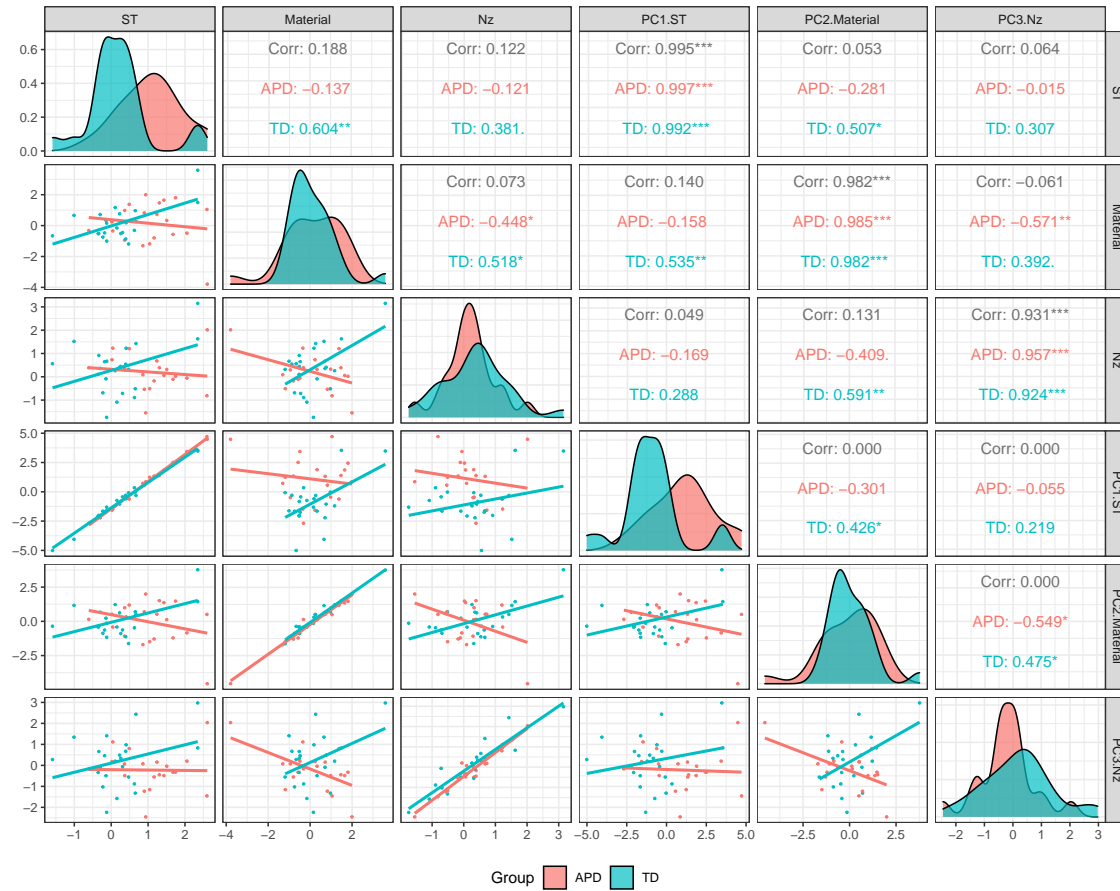


**Figure 3.20:** Switching task PCA: Listeners weighted scores split by components and group.

score of all CCRM conditions ( $\overline{CCRM}$ ) from the mean score of all ASL conditions ( $\overline{ASL}$ ), i.e.,  $\text{Material} = \overline{ASL} - \overline{CCRM}$ . The remaining composite, Nz, was calculated by subtracting the listeners performance averaged across conditions with speech distractors ( $\overline{Spch}$ ) from the average performance taken across the nonspeech and quiet conditions ( $\overline{NoSpch}$ ), i.e.,  $\text{Nz} = \overline{NoSpch} - \overline{Spch}$ . As can be seen in the figure, the PCA components highly correlated with the respective calculated composites (PC1.ST - ST, PC2.Material - Material, PC3.Nz - Nz), whereas none correlated with another composite, thus indicating that the components are independent from one another and that each describe different dimensions within the data.

### Language measures

PCA with three components was computed for the listeners scaled scores obtained in the different language measures, comprising of 16 input variables (x1 CELF-RS, x5 ECLiPS, x10 CCC-2) with a sample size of 42. Data for one TD child was excluded from the analysis due to inconsistent CCC-2 responses. Kaiser-Meyer-Olkin test for sample-size adequacy was ‘superb’ (Field et al., 2012) with an overall KMO of 0.93 (range: 0.86 - 0.97) and the assumption of sphericity was verified using



**Figure 3.21:** Switching task PCA: Comparison between PCA weighted scores and calculated measures: (1) ST = mean score across all ST data, (2) Material =  $\overline{ASL} - \overline{CCRM}$ , (3) Nz =  $\overline{NoSpch} - \overline{Spch}$ .

Bartlett's sphericity test [ $\chi^2(120) = 787.52, p < 0.0001$ ]. The PCA variables loadings, eigenvalues and percentage of variance explained split by components is given in Table 3.22. The first component (PC1.Lang) yielded eigenvalue  $> 1$ , explaining circa 73% of the variance, reflecting an overall performance averaged across all the language measures. The remaining components had eigenvalue of just under 1 (0.95 & 0.85, respectively), each explaining circa 6% and 5% of the variance. The second component (PC2.Lang) reflected discrepancy between expressive language skills, measured by the CELF-RS and listening and communication skills measured by the ECLiPS subscales. Interestingly, the third component (PC3.Lang) reflected once again a discrepancy, clustering together variables that taps onto pragmatic language and social interaction skills such as the ECLiPS subscale PSS (pragmatic & social skills) and the CCC-2 subscales E, H, I & J, separating them from other

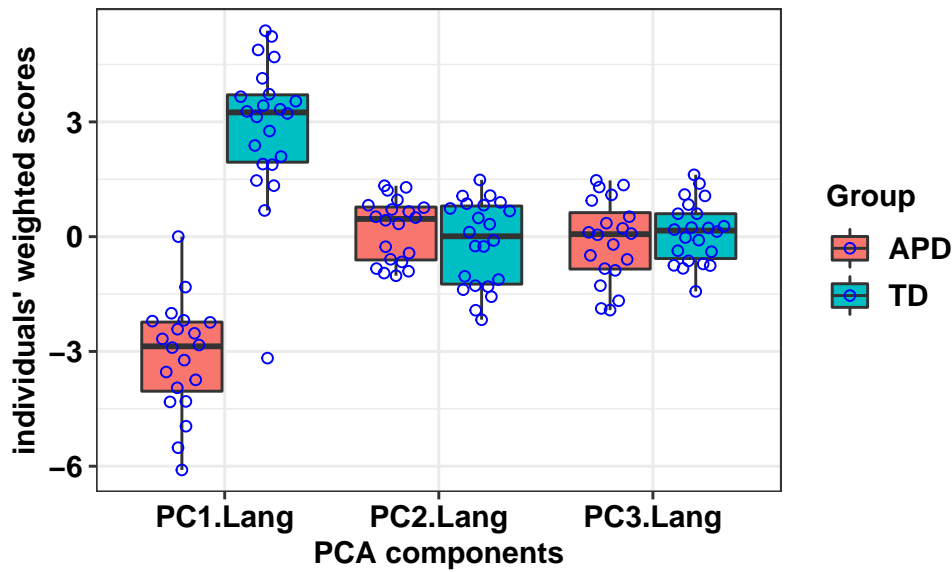
**Table 3.22:** Language measures PCA: Input variables loading.

Item	PC1.Lang	PC2.Lang	PC3.Lang
CELF-RS	<b>0.69</b>	<b>0.40</b>	<b>0.37</b>
ECLIPS.SAP	<b>0.91</b>	<b>-0.32</b>	0.14
ECLIPS.LLL	<b>0.92</b>	-0.14	0.11
ECLIPS.M.A	<b>0.88</b>	<b>-0.30</b>	0.14
ECLIPS.PSS	<b>0.83</b>	<b>-0.36</b>	-0.17
ECLIPS.EAS	<b>0.79</b>	<b>-0.52</b>	0.07
CCC2.A.speech	<b>0.78</b>	0.08	<b>0.35</b>
CCC2.B.syntax	<b>0.82</b>	0.19	0.27
CCC2.C.semantic	<b>0.92</b>	0.14	0.05
CCC2.D.coherence	<b>0.92</b>	0.09	0.05
CCC2.E.inappropriate.initiation	<b>0.82</b>	0.13	<b>-0.41</b>
CCC2.F.stereotyped	<b>0.89</b>	0.22	-0.03
CCC2.G.use.of.context	<b>0.93</b>	0.04	-0.08
CCC2.H.nonverbal	<b>0.84</b>	0.13	-0.26
CCC2.I.social	<b>0.88</b>	0.15	-0.16
CCC2.J.interests	<b>0.80</b>	0.11	<b>-0.40</b>
eigenvalue	11.67	0.95	0.85
variance (%)	72.96	5.95	5.3
cumulative variance (%)	72.96	78.91	84.21

|loading| >0.3 are highlighted in bold.

variables that assess more structural language skills such as the CELF-RS and the CCC-2 subscales speech (A) and Syntax (B). Boxplots of the listeners weighted scores for the PCA components split by group is shown in Figure 3.22. As seen in the ST data, the first component (PC1.Lang) best separated between the two groups, whereas separation between the two groups in the remaining components were noticeably smaller.

Despite the small proportion of variance explained by the later two principal components, they yet capture other aspects of language and communication skills that may be relevant in explaining the individual and group differences in the auditory tasks and were therefore included in the analysis. Nevertheless, interpretation of the relationship between these components with performance in the auditory tasks should be viewed with caution. Inspection of the individuals' scaled

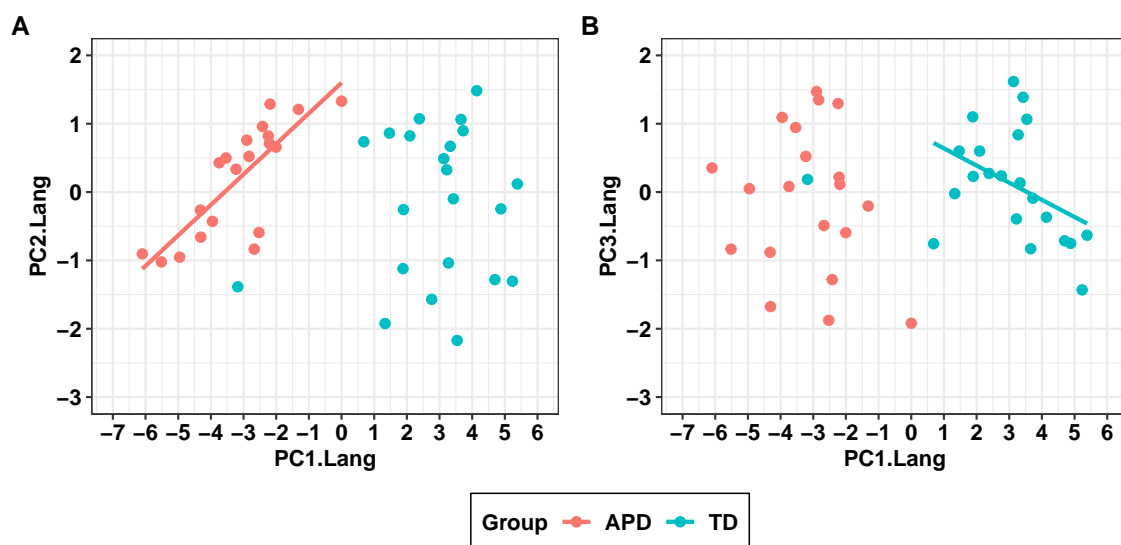


**Figure 3.22:** Language measures PCA: Listeners weighted scores split by components and group

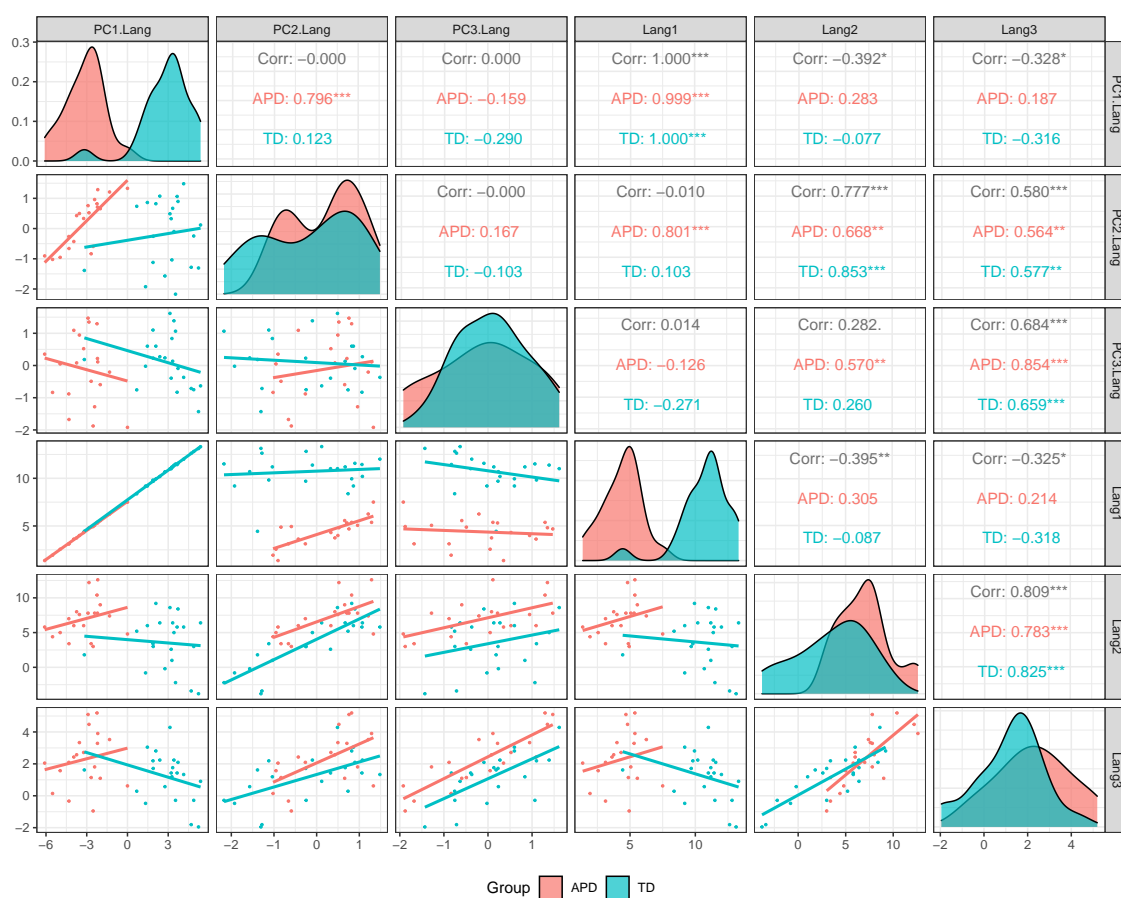
scores split by groups for loadings in PC1.Lang as a function of loadings in PC2.Lang and PC3.Lang shown in Figure 3.23 A-B revealed a linear relationship between PC1.Lang and PC2.Lang (APD group) and between PC1.Lang and PC3.Lang (TD group), thus indicating that they are not entirely independent from one another. The partial lack of independence may be in part explained by the large polarity in scores between the groups across the different input variables.

Again, the relationship between the PCA components (PC1.Lang, PC2.Lang and PC3.Lang) and the three calculated composites that reflects the components interpretations is illustrated in Figure 3.24. The calculated components were based on the listeners scaled scores, where *Lang1* represents the overall performance aggregated across all the language scores, *Lang2* represents discrepancy between expressive language skills (CELF-RS) and listening and communication skills (all ECLiPS subscales), and lastly, *Lang3* stands for discrepancy between structural and pragmatic & social skills. As seen in the figure, correlations were high between each PCA component and the corresponding calculated composites (range: 1 to 0.58).





**Figure 3.23:** Language measures PCA: Individual scores split by groups for loadings in PC1.Lang as a function of scores for PC2.Lang (A), and PC3.Lang (B).



**Figure 3.24:** Language measures PCA: Comparison between the listeners weighted scores by components, PC1.Lang - PC3.Lang (A), and calculated measures, Lang1 - Lang3 (B).

*Discussion(?): Which measures were best described by the PCA?*

All the measures showed strong correlation with PC1.Lang, whereas the CCC-2 GCC score showed the largest correlation ( $\rho = 0.98$ ,  $p < .0001$ ). This was true not only for the data aggregated across groups, but also when correlations were examined separately in each group. Therefore, taking into account the short administration time and simplicity, the CCC-2 alone provides a good screening tool for children's language and communication skills with high levels of sensitivity and specificity. Nonetheless, children in the present study knowingly consented to take part in the study either as part of the clinical APD group or the control group, which may introduced bias in the reporters response, and may resulted in a larger separation between the two groups than one would expect across the true population.

## Correlations

Next, the extent to which individual differences in speech perception could be explained by other measures was examined for the aggregated data across the two groups with multiple Spearman's rho correlations using *rcorr* function (Hmisc R package; Harrell Jr, 2020) between SSN scores, LiSNS-UK scores for the spatialised conditions and the derived score for spatial release from masking (S0N0, S0N90 & SRM), the principal components for the switching task PC1.ST, PC2.Material and PC3.NZ, and for the language measures PC1.Lang, PC2.Lang and PC3.Lang, average PTA at standard audiometry frequency bands (0.5-4 kHz), average PTA at high-frequency bands ( $PTA_{EHF}$ , at 8, 11 and 16 kHz), and ENVASA total score as a measure for sustained and selective-attention skills. Age effect was accounted for either by using standardised norms when available or by a regression model based z-score transformation. The correlation matrix outcomes are given in Table 3.23.

There was a significant correlation between the listeners overall performance in the switching task (PC1.ST) and their language skills (PC1.Lang;  $\rho = -0.55$ ,  $p < 0.001$ ), PTA ( $\rho = 0.46$ ,  $p < 0.01$ ), speech perception in noise (SSN;  $\rho = 0.46$ ,

**Table 3.23:** Correlation matrix (Spearman) between the study test measures for aggregated data across the two groups.

	PTA	PTA <sub>EHF</sub>	ENVASA	SSN	S0N0	S0N90	SRM	PC1.ST	PC2.Material	PC3.Nz	PC1.Lang	PC2.Lang
PTA												
PTA <sub>EHF</sub>	0.31											
ENVASA	-0.10	-0.13										
SSN	0.26	0.01	-0.40*									
S0N0	0.26	0.02	-0.19	0.39*								
S0N90	0.45**	0.34*	-0.23	0.30	0.64****							
SRM	-0.39*	-0.36*	0.12	-0.07	0.08	-0.67****						
PC1.ST	0.46**	0.09	-0.27	0.46**	0.34*	0.44**	-0.23					
PC2.Material	-0.17	-0.14	0.01	0.30	0.34*	0.20	0.12	0.06				
PC3.Nz	-0.03	0.03	0.05	0.09	-0.10	-0.10	-0.03	-0.11	0.01			
PC1.Lang	-0.16	-0.07	0.46**	-0.51***	-0.19	-0.15	-0.02	-0.55***	-0.03	0.16		
PC2.Lang	0.07	0.07	0.12	-0.02	0.21	0.23	-0.14	-0.01	0.16	-0.04	0.08	
PC3.Lang	-0.10	-0.26	0.00	0.09	-0.03	-0.12	0.08	-0.02	0.09	-0.14	-0.05	-0.02

significant p-values: \*\*\*\* p < .0001, \*\*\* p < .001, \*\* p < .01, \* p < .05

p < 0.01), and the spatialised LiSNS-UK test conditions S0N0 ( $\rho = 0.35$ , p < 0.05) and S0N90 ( $\rho = 0.45$ , p < 0.01). The second ST principal component, PC2.Material, significantly correlated with S0N0 ( $\rho = 0.33$ , p < 0.05) and SRM ( $\rho = 0.33$ , p < 0.05), whereas no relationship was found between the third PC3.Nz and any of the study measures.

Performance in the LiSNS-UK exhibited the highest correlation coefficients, with highly significant correlation between S0N0 and S0N90, where better performance in one condition was highly associated with better performance in the other ( $\rho = 0.64$ , p < 0.0001), and between S0N90 and SRM ( $\rho = -0.67$ , p < 0.0001), where better SRM was predicted by better performance for S0N90, whereas correlation between S0N0 and SRM was not significant ( $\rho = 0.08$ , p = 0.62). Note that lower z-score in the spatialised conditions denotes better performance, whereas the opposite holds for SRM with higher z-scores marking better performance, which explains the negative correlation between SRM and S0N90. A separate group-wise analysis gave similar results for correlation between S0N90 and SRM, whereas correlations in the APD group between S0N0 and S0N90, and between S0N0 and SRM were smaller and not significant ( $\rho$ : 0.35 and 0.30, respectively). The non-significant correlation between SRM and S0N0 stands in contrasts to our expectations, which are supported by the norming study in Chapter ???, for a positive correlation, where listeners with poorer (i.e., higher) S0N0 score were expected to have a larger (i.e., better) SRM. The insignificant and reduced correlation in the APD group is likely

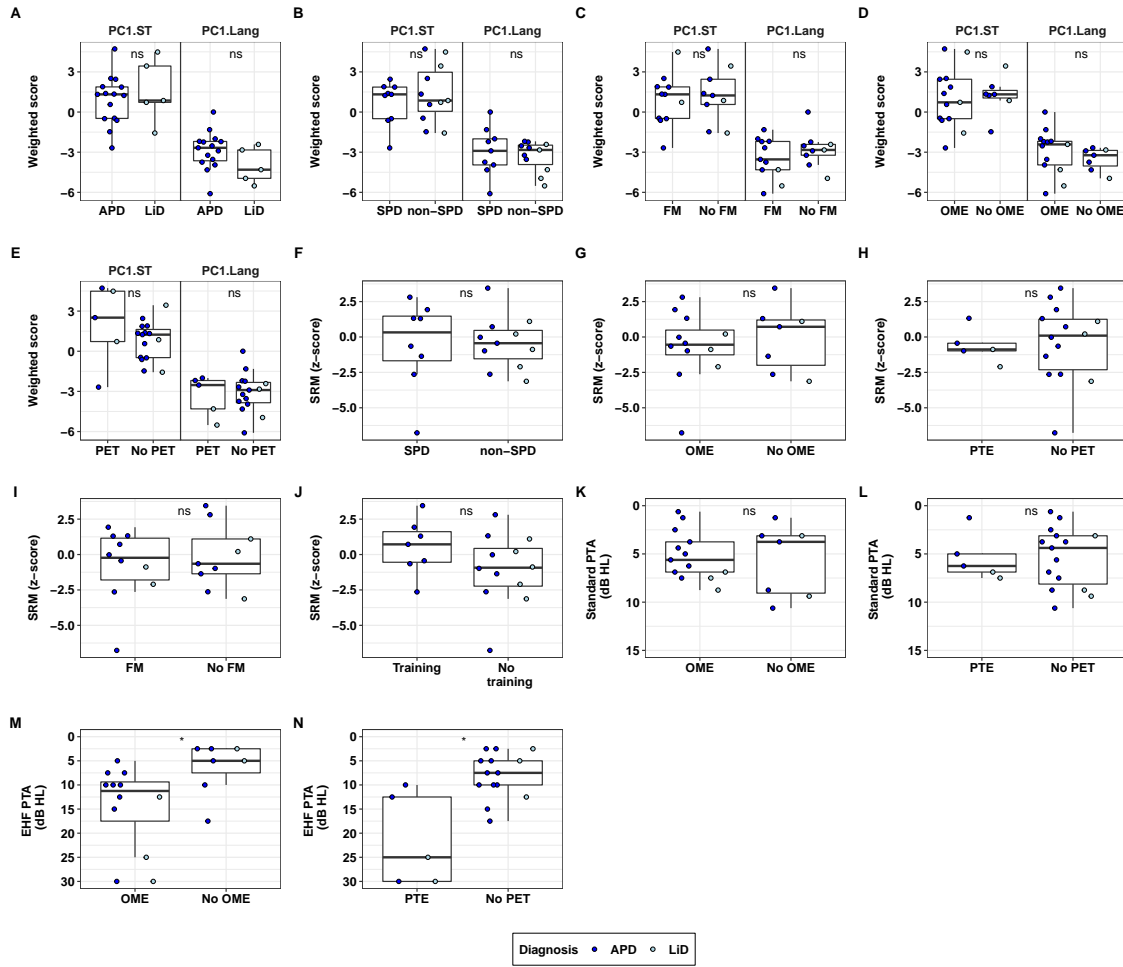
due to sampling error and the small sample size in the present study (correlation between the LiSNS-UK condition for the listeners SRT and z-scores are given in appendices in Figures C.1 and C.2).

SSN score was found to be related to performance in the two spatialised LiSNS-UK test conditions with correlation coefficients of 0.30 (S0N90) and 0.39 (S0N0), however only correlation for S0N0 was significant ( $p < 0.05$ ), while p-value for correlation with S0N90 was just above the significance level ( $p = 0.055$ ). The listeners S0N90 score significantly correlated with hearing sensitivity thresholds measured at both standard (PTA;  $\rho = 0.45$ ,  $p < 0.01$ ) and extended frequency bands ( $PTA_{EHF}$ ;  $\rho = 0.34$ ,  $p < 0.05$ ). Moreover, none of the LiSNS-UK measures significantly correlated with the language principle components PC1.Lang - PC3.Lang or the attention measure ENVASA. Additional significant correlations were found between PC1.Lang and SSN ( $\rho = -0.51$ ,  $p < 0.0001$ ) and between PC1.Lang and the ENVASA task ( $\rho = 0.46$ ,  $p < 0.001$ ). No p-value Bonferroni correction for multiple comparisons was applied.

### Exploratory predictors in the APD group

Association of potential predictors with performance in the APD group was examined in the following section. Nevertheless, it is important to emphasise that this is an exploratory examination across a small sample size and thus its outcomes may not be generalised in a larger sample. Predictors were selected based on the caregivers response in the background questionnaire, where the APD children were subdivided into the following group-pairs: 1. APD diagnosis (APD vs. LiD), 2. SPD diagnosis (SPD vs. non-SPD), 3. Regular use of FM-device (FM vs. No FM), 4. History of OME (OME vs. No OME), 5. Pressure equalisation tube history (PET vs. No PET), and 6. Auditory training (Training vs. No training). The listeners performance subdivided by predictors is shown in Figure 3.25 for data measured with the ST task (PCA1.ST), the language composite (PCA1.Lang), SRM, and thresholds for standard audiometry (PTA) and EHF audiometry (EHF PTA).

Individual observations are marked in circles, whereby observations of children diagnosed with APD are filled in dark blue, while LiD observations are filled in light blue. From the boxplots, PTE and OME emerges as the best predictors, explaining the largest portion of the within group differences. History of PET showed the highest association with poorer EHF PTA thresholds, and to a relatively smaller extent with PC1.ST (higher score indicates poorer performance) and with the SRM score (higher score indicates better performance). Consequently, it is not surprising that a related predictor – history of OME was also highly related to poorer EHF PTA thresholds, nevertheless, association between OME and the other measures was weak. Interestingly, there was no association between SRM performance and a diagnosis of SPD, with only a small difference between APD children with or without an SPD diagnosis.



**Figure 3.25:** Association between predictors and performance in the APD group for the switching task composite (PC1.ST), language composite (PC1.Lang), SRM, standard and EHF pure-tone PTA. Predictors included: 1. APD diagnosis (APD vs. LiD), 2. SPD diagnosis (SPD vs. non-SPD), 3. Regular use of FM-device (FM vs. No FM), 4. History of OME (OME vs. No OME), 5. Pressure equalisation tube history (PET vs. No PET), and 6. Auditory training (Training vs. No training). Individual observations are marked in circles. Observations of children diagnosed with APD are filled in dark blue, and LiD observations are filled in light blue. Significant p-values for independent t-test comparison are marked with asterisk ( $p < 0.05$ ).

## 3.5 Discussion

### 3.5.1 EHF

Lee's thresholds for 10-21 yrs group: 8=16.35 (1.46-29.33); 11=22.99, 16=48 (20.01-91.35); 20=93.07 (48.57-105.00) all dB SPL

EHF in children: Read Schechter et al., 1986:

- 6-10 yrs: 10k=23, 12k=20, 16k=39 dB SPL
- 11-15 yrs: 10k=21, 12k=22, 16k=51 dB SPL

### 3.5.2 ST

#### *Why CCRM performance is better*

The improved intelligibility in the CCRM material is amongst others due to the more simple speech material, the reduced confusion between the target sentences and the connected speech distractors as well as the restricted alternative responses of the CCRM matrix-based sentences.

z- scores by material: proportion of abnormal TD kids:

ASL: The proportion of abnormal scores amongst the TD group ranged between 0% to 13% (mean = 7.8%), which is relatively higher than expected in the normal population. Nonetheless, when taking into account TD observations that were trimmed during the z-score calculation procedure, the proportion of abnormal scores are smaller, ranging between 0% to 9.5% (mean = 3.8%), which corroborate fairly well with the theoretical probability of 2.5% (one-tailed).

CCRM: The percentage of abnormal scores in the TD group were relatively low ranging between 0 to 8.7% (mean = 4.3%) and were at 0% across all conditions when TD observations that were trimmed as part of the z-score calculation procedure were accounted for.

Why there was no interaction between Group x Condition x Material? *Discussion or here?*

The lack of significant interaction (Group x Condition or Group x Condition x Material), is somewhat surprising and do not reflect some of the differences seen in Figure 3.7 A-B between the two groups in some conditions or the overall difference in performance between the speech materials and may suggest that the model was under-powered to test these questions.

*Points for age effect:*

- Goldsworthy et al. 2018 found that age explained only a small portion of variability

in speech perception performance (n.s.) for Quiet, SSN and 2-talker connected-speech distractors (children aged 5-17). See table 3.

*Points for SSN:*

- “Despite mature peripheral encoding, school-children have more difficulty understanding speech in noise compared with adults. For example, 5-7 year-old children require 3 to 6 dB more favourable SNR than adults to achieve comparable speech detection, word identification, or sentence recognition performance in a speech-shaped noise maker (e.g., Corbin et al., 2016)” [Leibold, Buss and Calandruccio, 2019, Acoustics today]. - “Speech recognition gradually improves until 9-10 years of age , after which mature performance is generally observed” [Leibold, Buss and Calandruccio, 2019, Acoustics today].

- SSN age effect in other studies are smaller

### 3.5.3 CCC-2

### 3.5.4 ECLiPS

**Discussion:** Correlation with CCC-2 sub-scales (Barry & Moore, 2014): Overall, all the ECLiPS sub scales shows strong correlation with most of the CCRM 10 sub-scales. Interestingly, PSS strongly correlates with all 10 CCC-2 sub-scales, suggesting that both tests taps into similar abilities.

**In the results:** compare scores with scores obtained by: <https://www.nature.com/articles/s41598-018-25316-9.pdf> and Moore et al. 2020 (Listening Difficulties in Children: Behaviour and Brain Activation Produced by Dichotic Listening of CV Syllables)

Discussion: - Compare data with Ferguson et al. 2011

## 3.6 Conclusion



*Alles Gescheite ist schon gedacht worden.  
Man muss nur versuchen, es noch einmal zu denken.*

*All intelligent thoughts have already been thought;  
what is necessary is only to try to think them again.*

— Johann Wolfgang von Goethe (von Goethe, 1829) **General discussion**

If we don't want Conclusion to have a chapter number next to it, we can add the {-} attribute.

### **More info**

And here's some other random info: the first paragraph after a chapter title or section head *shouldn't be* indented, because indents are to tell the reader that you're starting a new paragraph. Since that's obvious after a chapter or section title, proper typesetting doesn't add an indent there.

## **Summary of main findings**

## **Conclusion**



# Appendices





## The First Appendix

This first appendix includes an R chunk that was hidden in the document (using `echo = FALSE`) to help with readability:

**In 02-rmd-basics-code.Rmd**

```
library(tidyverse)
knitr::include_graphics("figures/chunk-parts.png")
```

**And here's another one from the same chapter, i.e. Chapter ??:**



# B

## The Second Appendix





# C

## The Third Appendix

```
Loadings <- as.data.frame(PCA.Lang$var$coord)
colnames(Loadings)[1:3] <- c("PC1.Lang", "PC2.Lang", "PC3.Lang")
# arrange variables into groups
Loadings$CondCode <- rownames(Loadings)
Loadings$CondCode <- ifelse(str_detect(Loadings$CondCode, "CCC2.")==TRUE, sprintf(
Loadings$Variables <- factor(c("CELF-RS", "ECLIPS", "ECLIPS", "ECLIPS", "ECLIPS", "ECLIPS", "ECLIPS", "ECLIPS",
                                "CCC2", "CCC2", "CCC2", "CCC2", "CCC2", "CCC2", "CCC2", "CCC2", "CCC2", "CCC2",
                                "CCC2", "CCC2"), levels=c("CELF-RS", "ECLIPS", "CCC2"), ordered=TRUE)
pragmatic <- c("ECLIPS.PSS", "CCC2.E", "CCC2.H", "CCC2.I", "CCC2.J")
Loadings$Variables2 <- ifelse(Loadings$CondCode %in% pragmatic, "pragmatic", "structural")

t1 <- ggplot(Loadings, aes(PC1.Lang, PC2.Lang)) +
  geom_hline(yintercept=0, color = "black", size=0.5) +
  geom_point(aes(shape = Variables, fill=Variables), size = 3) +
  # geom_text(label=Loadings$CondCode, size=2) +
  scale_fill_manual(values=c("white", "darkgrey", "white", "darkgrey")) +
  scale_shape_manual(values=c(0, 1, 2)) +
```

```

theme_bw()+
theme(axis.text = element_text(size = 11, face="bold",colour = "black"),
      axis.title.x = element_text(size=11, face="bold"),
      axis.title.y = element_text(size=11, face="bold"),
      legend.title = element_text(size=11, face="bold"),
      legend.text = element_text(size=11, face="bold"))

t2 <- ggplot(Loadings, aes(PC1.Lang, PC3.Lang)) +
  geom_hline(yintercept=0,color = "black", size=0.5)+
  geom_point(aes(shape = Variables2,fill=Variables2), size = 3) +
  # geom_point(aes(color=CondCode), size = 3) +
  geom_text(label=Loadings$CondCode,size=3) +
  scale_fill_manual(values=c("white","darkgrey")) +
  scale_shape_manual(values=c(1,21)) +
  # scale_fill_manual(values=c("white","darkgrey","white","darkgrey")) +
  # scale_shape_manual(values=c(1,2,3)) +
  theme_bw()+
  theme(axis.text = element_text(size = 11, face="bold",colour = "black"),
        axis.title.x = element_text(size=11, face="bold"),
        axis.title.y = element_text(size=11, face="bold"),
        legend.title = element_text(size=11, face="bold"),
        legend.text = element_text(size=11, face="bold"))

library(patchwork)
(t1 + t2) + plot_layout(ncol = 2, nrow = 1, heights = c(1, 1), widths = c(1,1)) + plo
  theme(legend.position='bottom',
        legend.direction = "horizontal",
        legend.background = element_blank(),
        legend.box.background = element_rect(colour = "black"),
        plot.tag = element_text(face = 'bold'))

```



**Figure C.1:** LiSNS-UK: Correlations for listeners SRTs (dB SNR).



**Figure C.2:** LiSNS-UK: Correlations for listeners age-independent z-scores.

## Works Cited

- Algazi, V. R., Duda, R. O., Thompson, D. M., & Avendano, C. (2001). The CIPIC HRTF database. *IEEE ASSP Workshop on Applications of Signal Processing to Audio and Acoustics*. <https://doi.org/10.1109/aspaa.2001.969552>
- Barry, J. G., & Moore, D. R. (2014). *Evaluation of Children's Listening and Processing Skills (ECLiPS)* (tech. rep.). MRC-T. London, United Kingdom.
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48. <https://doi.org/10.18637/jss.v067.i01>
- Bench, J., Kowal, Å., & Bamford, J. (1979). The Bkb (Bamford-Kowal-Bench) Sentence Lists for Partially-Hearing Children. *British Journal of Audiology*, 13(3), 108–112. <https://doi.org/10.3109/03005367909078884>  
doi: 10.3109/03005367909078884
- Bishop, D. V. M. (2003). *The Children's Communication Checklist, Version 2 (CCC-2)* (tech. rep.). The Psychological Corporation. London, United Kingdom.
- Cameron, S., & Dillon, H. (2007). Development of the Listening in Spatialized Noise-Sentences Test (LISN-S). *Ear and Hearing*, 28(2), 196–211. <https://doi.org/10.1097/AUD.0b013e318031267f>
- Cameron, S., Glyde, H., & Dillon, H. (2011). Listening in Spatialized Noise—Sentences Test (LiSN-S): Normative and Retest Reliability Data for Adolescents and Adults up to 60 Years of Age. *Journal of the American Academy of Audiology*, 22(10), 697–709. <https://doi.org/10.3766/jaaa.22.10.7>
- Conti-Ramsden, G., Botting, N., & Faragher, B. (2001). Psycholinguistic markers for specific language impairment (SLI). *Journal of Child Psychology and Psychiatry and Allied Disciplines*. <https://doi.org/10.1111/1469-7610.00770>
- Feys, J. (2015). *Npintfactrep: Nonparametric interaction tests for factorial designs with repeated measures* [R package version 1.5]. <https://CRAN.R-project.org/package=npIntFactRep>
- Feys, J. (2016). Nonparametric tests for the interaction in two-way factorial designs using R. *R Journal*. <https://doi.org/10.32614/rj-2016-027>
- Field, A., Miles, J., & Field, Z. (2012). Discovering Statistics Using R - 17 Exploratory factor analysis. *Discovering statistics using r*.
- Fox, J., & Weisberg, S. (2019). *An R companion to applied regression* (Third). Sage. <https://socialsciences.mcmaster.ca/jfox/Books/Companion/>
- Harrell Jr, F. E. (2020). *Hmisc: Harrell miscellaneous* [R package version 4.4-2]. <https://CRAN.R-project.org/package=Hmisc>
- Hothorn, T., Hornik, K., van de Wiel, M. A., & Zeileis, A. (2006). A Lego system for conditional inference. *The American Statistician*, 60(3), 257–263. <https://doi.org/10.1198/000313006X118430>
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning: With applications in r* (Vol. 103). Springer.

- Kohl, M. (2020). *MKinfer: Inferential Statistics* [R package version 0.6].  
<http://www.stamats.de>
- Krishnan, S., Leech, R., Aydelott, J., & Dick, F. (2013). School-age children's environmental object identification in natural auditory scenes: Effects of masking and contextual congruence. *Hearing Research*, 300, 46–55.  
<https://doi.org/10.1016/j.heares.2013.03.003>
- Lê, S., Josse, J., & Husson, F. (2008). FactoMineR: A package for multivariate analysis. *Journal of Statistical Software*, 25(1), 1–18. <https://doi.org/10.18637/jss.v025.i01>
- Lee, J., Dhar, S., Abel, R., Banakis, R., Grolley, E., Lee, J., Zecker, S., & Siegel, J. (2012). Behavioral Hearing Thresholds between 0.125 and 20 kHz Using Depth-Compensated Ear Simulator Calibration. *Ear and Hearing*.  
<https://doi.org/10.1097/AUD.0b013e31823d7917>
- Leech, R., Gygi, B., Aydelott, J., & Dick, F. (2009). Informational factors in identifying environmental sounds in natural auditory scenes. *The Journal of the Acoustical Society of America*, 126(6), 3147–3155. <https://doi.org/10.1121/1.3238160>
- Lenth, R. V. (2020). *Emmeans: Estimated marginal means, aka least-squares means* [R package version 1.5.3]. <https://CRAN.R-project.org/package=emmeans>
- MacLeod, A., & Summerfield, Q. (1990). A procedure for measuring auditory and audiovisual speech-reception thresholds for sentences in noise: Rationale, evaluation, and recommendations for use. *British Journal of Audiology*, 24(October), 29–43. <https://doi.org/10.3109/030053690009077840>
- McDonald, J. (2014). Multiple comparisons. *Handbook of biological statistics* (3rd ed., pp. 254–260). Sparky House Publishing.
- Murphy, C. F., Hashim, E., Dillon, H., & Bamiau, D. E. (2019). British children's performance on the listening in spatialised noise-sentences test (LISN-S). *International Journal of Audiology*.  
<https://doi.org/10.1080/14992027.2019.1627592>
- Noguchi, K., Gel, Y. R., Brunner, E., & Konietzschke, F. (2012). nparLD: An R software package for the nonparametric analysis of longitudinal data in factorial experiments. *Journal of Statistical Software*, 50(12), 1–23.  
<http://www.jstatsoft.org/v50/i12/>
- Norbury, C. F. (2014). Practitioner Review: Social (pragmatic) communication disorder conceptualization, evidence and clinical implications. *Journal of Child Psychology and Psychiatry and Allied Disciplines*. <https://doi.org/10.1111/jcpp.12154>
- Norbury, C. F., & Bishop, D. V. M. (2005). Children's Communication Checklist - 2 : a validation study. *Publie dans Revue Tranel*, 42, 53–63.
- R Core Team. (2020). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. Vienna, Austria.  
<https://www.R-project.org/>
- Ramus, F., Rosen, S., Dakin, S. C., Day, B. L., Castellote, J. M., White, S., & Frith, U. (2003). Theories of developmental dyslexia: Insights from a multiple case study of dyslexic adults. *Brain*, 126(4), 841–865. <https://doi.org/10.1093/brain/awg076>
- Revelle, W. (2020). *Psych: Procedures for psychological, psychometric, and personality research* [R package version 2.0.12]. Northwestern University. Evanston, Illinois.  
<https://CRAN.R-project.org/package=psych>
- Richmond, S. A., Kopun, J. G., Neely, S. T., Tan, H., & Gorga, M. P. (2011). Distribution of standing-wave errors in real-ear sound-level measurements. *The Journal of the Acoustical Society of America*. <https://doi.org/10.1121/1.3569726>

- RStudio Team. (2020). *Rstudio: Integrated development environment for r*. RStudio, PBC. Boston, MA. <http://www.rstudio.com/>
- Siegel, J. H. (1994). Ear-canal standing waves and high-frequency sound calibration using otoacoustic emission probes. *Journal of the Acoustical Society of America*.  
<https://doi.org/10.1121/1.409829>
- von Goethe, J. W. (1829). *Wilhelm Meisters Wanderjahre oder die Entsagenden*. Cotta.
- Wierstorf, H., Geier, M., Raake, A., & Spors, S. (2011). A Free Database of Head-Related Impulse Response Measurements in the Horizontal Plane with Multiple Distances. *AES130*.
- Wiig, E., H, Semel, E., & Secord, W. (2017). *Clinical Evaluation of Language Fundamentals - Fifth Edition UK (CELF-5UK)* (tech. rep.). PsychCorp, Pearson Clinical Assessment.