

The agreement between parent-reported and directly measured child language and parenting behaviors

Shannon K. Bennetts^{1, 2, 3*}, Fiona K. Mensah^{1, 2, 4}, Elizabeth M. Westrupp^{1, 2, 3}, Naomi J. Hackworth^{2, 3}, Sheena Reilly^{1, 2, 4, 5}

¹Department of Paediatrics, The University of Melbourne, Australia, ²Murdoch Childrens Research Institute, Australia, ³Judith Lumley Centre, La Trobe University, Australia, ⁴The Royal Children's Hospital, Australia, ⁵Menzies Health Institute, Griffith University, Australia

Submitted to Journal:
Frontiers in Psychology

Specialty Section:
Developmental Psychology

ISSN:
1664-1078

Article type:
Original Research Article

Received on:
03 Jun 2016

Accepted on:
17 Oct 2016

Provisional PDF published on:
17 Oct 2016

Frontiers website link:
www.frontiersin.org

Citation:

Bennetts SK, Mensah FK, Westrupp EM, Hackworth NJ and Reilly S(2016) The agreement between parent-reported and directly measured child language and parenting behaviors. *Front. Psychol.* 7:1710. doi:10.3389/fpsyg.2016.01710

Copyright statement:

© 2016 Bennetts, Mensah, Westrupp, Hackworth and Reilly. This is an open-access article distributed under the terms of the [Creative Commons Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use, distribution and reproduction in other forums is permitted, provided the original author(s) or licensor are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

Provisional

The agreement between parent-reported and directly measured child language and parenting behaviors

Bennetts, S.K.^{1,2,3*}, Mensah, F.K.^{1,2,4}, Westrupp, E.M.^{1,2,3}, Hackworth, N.J.^{2,3}, & Reilly, S.^{1,2,4,5}.

¹Department of Paediatrics, The University of Melbourne, Parkville, Victoria, Australia

²Murdoch Childrens Research Institute, Parkville, Victoria, Australia

³Judith Lumley Centre, La Trobe University, Melbourne, Victoria, Australia

⁴The Royal Children's Hospital, Parkville, Victoria, Australia

⁵Menzies Health Institute, Griffith University, Gold Coast, Queensland, Australia

*Correspondence: Shannon Bennetts, Murdoch Childrens Research Institute Hearing, Language & Literacy Group, The Royal Children's Hospital, 50 Flemington Rd, Parkville, 3052, Victoria, Australia, shannon.bennetts@mcri.edu.au

Keywords: agreement, bias, Bland-Altman Method, Reduced Major Axis regression, measurement, parent-report, child language, parenting

Word count: 9,452

Number of figures: 6

Abstract

Parenting behaviors are commonly targeted in early interventions to improve children's language development. Accurate measurement of both parenting behaviors and children's language outcomes is thus crucial for sensitive assessment of intervention outcomes. To date, only a small number of studies have compared parent-reported and directly measured behaviors, and these have been hampered by small sample sizes and inaccurate statistical techniques, such as correlations. The Bland-Altman Method and Reduced Major Axis regression represent more reliable alternatives because they allow us to quantify fixed and proportional bias between measures. In this study, we draw on data from two Australian early childhood cohorts ($N=201$ parents and slow-to-talk toddlers aged 24 months; and $N=218$ parents and children aged 6-36 months experiencing social adversity) to (1) examine agreement and quantify bias between parent-reported and direct measures, and (2) to determine socio-demographic predictors of the differences between parent-reported and direct measures. Measures of child language and parenting behaviors were collected from parents and their children. Our findings support the utility of the Bland-Altman Method and Reduced Major Axis regression in comparing measurement methods. Results indicated stronger agreement between parent-reported and directly measured child language, and poorer agreement between measures of parenting behaviors. Child age was associated with difference scores for child language; however the direction varied for each cohort. Parents who rated their child's temperament as more difficult tended to report lower language scores on the parent questionnaire, compared to the directly measured scores. Older parents tended to report lower parenting responsiveness on the parent questionnaire, compared to directly measured scores. Finally, speaking a language other than English was associated with less responsive parenting behaviors on the videotaped observation. Variation in patterns of agreement across the distribution of scores highlighted the importance of assessing agreement

49 comprehensively, providing strong evidence that simple correlations are grossly insufficient
50 for method comparisons. We discuss implications for researchers and clinicians, including
51 guidance for measurement selection, and the potential to reduce financial and time-related
52 expenses and improve data quality. Further research is required to determine whether
53 findings described here are reflected in more representative populations.
54

55 **1 Introduction**

56 The success of early intervention programs relies on accurate and sensitive measurement of
57 intervention processes and outcomes. It is surprising then, that research comparing agreement
58 between different types of measurement methods has been extremely limited. There has been
59 increasing attention over the past decade on early intervention programs targeting parenting
60 behaviors in order to improve children's language outcomes. Language delay affects around
61 one in five children at age four (e.g. Reilly et al., 2010) and persistent difficulties can impact
62 upon future academic success, employment prospects and socio-emotional functioning
63 (Campbell & Ramey, 1994; Stothard, Snowling, Bishop, Chipchase, & Kaplan, 1998).
64 Parenting characterized by warm, positive and responsive interactions can facilitate language
65 development in the early years (Cartmill et al., 2013; Sim, 2012), serving as a buffer against
66 the above-mentioned risks. Understanding how to identify language concerns and intervene
67 early relies upon accurate and reliable measurement. This paper uses existing data from two
68 early childhood cohorts to examine agreement between parent-reported and directly measured
69 child language and parenting behaviors, and the socio-demographic predictors of the
70 difference between measures.
71

72 Research focused on understanding complex child developmental and family processes
73 requires highly sensitive assessment tools. Two primary options for researchers seeking to
74 quantify constructs related to child language and parenting behaviors are parent-reported
75 measures and direct (observational or standardized) measures. Both possess notable strengths
76 and limitations, yet there is a lack of comparable data to help researchers identify the
77 circumstances in which one or both methods should be employed. Parents are uniquely
78 positioned to report on their children's behavior retrospectively and across multiple settings
79 (Gartstein & Marmion, 2008). Parent-reported data is relatively straightforward and
80 inexpensive to collect and analyze (Hawes & Dadds, 2006), making it an appealing
81 measurement approach in large-scale trials where time and cost are significant considerations.
82 However parents' unique set of experiences, opinions and attitudes (both explicit and
83 implicit) can contribute to response bias. For example, parents may vary in their
84 interpretation of key terms (Aspland & Gardner, 2003); psychological difficulties can color
85 parents' perceptions of their children's behavior (Hayden, Durbin, Klein, & Olino, 2010); or
86 responses can be influenced by social desirability (Law & Roy, 2008). In contrast, direct
87 measures permit the collection of data which is more objective (Wysocki, 2015). For this
88 reason, direct measures are often considered the "gold standard" for assessing both parenting
89 behaviors (Hawes & Dadds, 2006) and child language (Sachse & Von Suchodoletz, 2008).
90 However collection of such measures requires considerable time and financial resources
91 (Gardner, 2000), and generalizability to other time points and settings has been questioned
92 (Gardner, 1997).
93

94 Myriad factors can affect observed behavior or parent-reported responses. Direct measures
95 can be influenced by the presence of the observer or assessor, illness, tiredness, or
96 distractions. Parent-reported measures may be biased by factors associated with a parent's

97 background. Parents from low socio-economic backgrounds (e.g. low income, low education)
98 have been shown to over- or under-estimate children's vocabulary on the Communicative
99 Development Inventory (Feldman et al., 2000; Roberts, Burchinal, & Durham, 1999),
100 suggesting that caution in interpretation is required. Furthermore, acquiescence or "yea-
101 saying" (i.e. the tendency to agree with items irrespective of their content) may be a
102 particularly important consideration when administering parent-reported measures with
103 socially disadvantaged populations (Meisenberg & Williams, 2008). It has also been
104 suggested that less educated parents may be less able than well-educated parents to
105 discriminate between expressive and receptive items on a vocabulary checklist, thus
106 providing an inflated estimate of their child's language abilities (Reese & Read, 2000). Child
107 characteristics such as temperament and gender have similarly been shown to affect parent
108 responses or parent behaviors (Hayden et al., 2010; Olino, Durbin, Klein, Hayden, & Dyson,
109 2013).

110

111 In light of these relative strengths and limitations of parent-reported and direct measures, it is
112 important to establish the extent to which these measurement methods concur and for whom.
113 This information will allow researchers to make more informed decisions about the most
114 appropriate and cost-effective measurement option, given the specific context of their study
115 and finite study resources. For example, given evidence of strong agreement between parent-
116 reported and direct measures, researchers may opt to administer only parent-report; whereas
117 evidence suggesting weak agreement may require researchers to administer multiple methods
118 or only the agreed "gold standard" method.

119

120 Few studies have investigated agreement between parent-reported and directly measured
121 behaviors. Of those that have, two primary limitations can be identified. Firstly, these studies
122 tend to employ small sample sizes (in the range of $N=50-70$). While understandable given the
123 expense associated with using direct measures, small samples reduce the power of the study
124 to identify the limits of agreement with precision. Secondly, these studies typically employ
125 correlational analyses to quantify agreement between measures. For example, moderate
126 correlations have been reported between parent-reported and directly measured child
127 language (e.g. Ring & Fenson, 2000; Sachse & Von Suchodoletz, 2008) and weak to
128 negligible correlations between parent-reported and directly measured parenting behaviors
129 (e.g. Arney, 2004). The use of correlations is problematic because correlations provide a
130 single figure representing the strength of the association between two related variables; they
131 do not assess agreement (Eadie et al., 2014). That is, correlations do not allow for differences
132 in agreement to be examined *across* the spectrum, and they do not account for bias which
133 may be present between two measures, including fixed bias (i.e. bias which is constant across
134 the distribution) or proportional bias (i.e. bias which varies proportionally across the
135 distribution) (See Bennetts et al., 2016; Bland & Altman, 1986; Carstensen, 2010). We agree
136 with Stolarova and colleagues (2014) that greater awareness of the difference between
137 agreement and correlation will lead to the use of more appropriate statistical methods.

138

139 Methods such as the Bland-Altman Method (Bland & Altman, 1986) or Reduced Major Axis
140 (RMA) regression (Ludbrook, 2010) represent appropriate alternatives for assessing
141 agreement, allowing researchers to quantify fixed and proportional bias, respectively. These
142 techniques are commonly used for method comparisons in fields such as medicine and
143 chemistry, but are seldom applied in psychology due to a lack of awareness and paucity of
144 literature in the field (Miles & Banyard, 2007). The Bland-Altman Method involves plotting
145 the mean of two measures against the difference between two measures (Altman & Bland,
146 1983). This provides a visual means of examining the variation in agreement across the

147 spectrum of scores. RMA regression is particularly helpful for identifying proportional bias
 148 between measures (Ludbrook, 1997). Execution of this technique involves minimizing the
 149 sum of the vertical and horizontal residuals. RMA is suitable for contexts in which
 150 measurement error is present in both x and y , as would be expected in the current study.

151
 152 This study used data from two cohorts of parents and their children aged 6-36 months to: (1)
 153 quantify the agreement between parent-reported and directly measured child language and
 154 parenting behaviors, and to (2) determine the association between a range of socio-
 155 demographic factors and the difference between parent-reported and direct measures.

156 **2 Materials and method**

157 **2.1 Participants**

158 Data were drawn from two randomized controlled trials of early childhood parenting
 159 interventions; (1) a community-based sample of parent-child dyads participating in the
 160 Language for Learning program for slow-to-talk toddlers aged 24 months ($N=201$), and (2)
 161 parent-child dyads participating in the Early Home Learning Study; an evaluation of a
 162 community-based program to support disadvantaged parents to provide a rich home learning
 163 environment for their children aged 6-36 months ($N=218$). Parents and children completed a
 164 suite of assessments, including parent-reported and direct measures of child language and
 165 parenting behaviors.

166
 167 Language for Learning participants were recruited by maternal and child health nurses in
 168 three local government areas in Victoria, Australia. All children residing in these areas were
 169 recruited at 12 months of age. Children were excluded if there was a known cognitive delay,
 170 a major medical condition, or if parents unable to complete written questionnaires. At child
 171 age 18 months, parents completed the Sure Start Language Measure. Children falling below
 172 the 20th percentile were invited to participate in the current study of slow-to-talk toddlers.

173
 174 Early Home Learning Study participants were recruited by child and family service workers
 175 and maternal and child health nurses within twenty local government areas in Victoria,
 176 Australia. Eligibility criteria included: living within the geographical boundaries of a trial
 177 locality; having at least one child aged 6-36 months; and evidence of at least one risk
 178 indicator for social disadvantage including: low family income; receipt of government
 179 benefits (e.g. Health Care Card for low income families); single, socially isolated or young
 180 parent (<25 years); and culturally and linguistically diverse background. Parents were not
 181 eligible if they were aged less than 18 years, did not speak English, or were receiving
 182 intensive support or child protection services.

183 **2.2 Measures**

184 A summary of parent-reported and direct measures administered for each cohort is provided
 185 in Table 1. Both cohorts completed parent-reported and direct measures of child language.
 186 Direct measures included a standardized language assessment for the Language for Learning
 187 cohort, and a videotaped observation for the Early Home Learning Study cohort. Participants
 188 in the Early Home Learning Study also completed a videotaped observation of parent-child
 189 interaction, as well as parent-reported measures of parenting behaviors.

190

191 Table 1. Parent-reported and direct measures.

Parent-reported measures	Direct measures
--------------------------	-----------------

Child language	<ul style="list-style-type: none"> ▪ Sure Start Language Measure (SSLM)^a ▪ Ages & Stages Questionnaire (ASQ)^{a,b} communication subscale ▪ MacArthur-Bates Communicative Development Inventory (Short-Form, CDI)^b 	<ul style="list-style-type: none"> ▪ Preschool Language Scale (PLS-4)^a ▪ Early Communication Indicator (ECI)^b
Parenting behaviors	<ul style="list-style-type: none"> ▪ Parental Verbal Responsivity (PVR)^b ▪ Home Activities with Child (HAC)^b 	<ul style="list-style-type: none"> ▪ Indicator of Parent-Child Interactions (IPCI)^b

192 ^aLanguage for Learning; ^bEarly Home Learning Study

193 2.3.1 Parent-reported measures

194 *MacArthur-Bates Communicative Development Inventories.* The CDI is a brief, reliable and
 195 commonly used measure of children’s language skills (Fenson et al., 2000). One of three
 196 versions was used depending on the child’s age in months. The CDI Short-Form Level I was
 197 used for children up to 18 months, consisting of an 89-word list, resulting in a total score
 198 from 0-89. Parents were asked to indicate if their child “understands” or “understands and
 199 says” each word. Parents of children aged 19-30 months completed the Short-Form Level II.
 200 Parents were asked to report whether their child ‘says’ 100 listed words resulting in a total
 201 score from 0-100 for word production and a single item assessing word combinations.
 202 Parents of children aged 31 months and above completed the CDI III, consisting of a 100-
 203 word vocabulary checklist, 12 sentence pairs to evaluate complexity of language use, and 12
 204 yes/no items assessing language comprehension, resulting in a total score from 0-124. Minor
 205 changes in word items were made for the Australian context, in-line with other Australian
 206 studies (Reilly et al., 2009; Skeat, Eadie, Ukoumunne, & Reilly, 2010). Scores were
 207 standardized for each of the three age-appropriate versions.
 208

209 *Sure Start Language Measure.* Children’s expressive vocabulary was assessed with the Sure
 210 Start Language Measure (SSLM) 100-word checklist (Roy, Kersley, & Law, 2005). The
 211 SSLM was developed based on the commonly used MacArthur-Bates Communicative
 212 Development Inventory, with some items adjusted for the United Kingdom, rather than
 213 American context. Parents were asked to indicate whether their child could say 100 words,
 214 (e.g., “meow”, “finish” or “happy”) and whether their child was combining words “not yet”,
 215 “sometimes” or “often” to produce a total score out of 100.
 216

217 *Ages & Stages Questionnaire (ASQ-3) communication subscale.* The ASQ allows for
 218 developmental and social-emotional screening of children, aged between 1-66 months
 219 (Squires, Twombly, Bricker, & Potter, 2009). Questionnaires comprise five sub-scales:
 220 communication, gross motor, fine motor, problem solving, and personal-social, with 6 items
 221 in each subscale, plus an additional 8 open-ended questions addressing overall child
 222 development. Only the communication subscale is reported here. Parents were asked to
 223 indicate whether their child performs a specific activity using three response categories:
 224 ‘yes’, ‘sometimes’ or ‘not yet’ across six items, each scored as 10, 5 or 0 for ‘yes’,
 225 ‘sometimes’ or ‘not yet’ respectively (e.g. “Does your child correctly use at least two words
 226 like “me”, “I”, “mine” and “you”?). Scores were summed to give a total score ranging from 0
 227 to 60. Higher scores indicated stronger communicative abilities. Fourteen age-appropriate
 228 versions were administered; therefore scores were standardized within age bands to derive z-
 229 scores.
 230

231 *Parental Verbal Responsivity*. The 4-item PVR subscale from the StimQ-Toddler (Dreyer,
 232 Mendelsohn, & Tamis-LeMonda, 1996) measures how verbally responsive the parent is in
 233 interactions with their child on a dichotomous “yes”/“no” scale. To detect greater variability,
 234 an alternative 4-point Likert scale was used, where 1 =not at all and 4 =every day. (e.g. “I
 235 talk about the day while my child is eating”). Scores were summed to produce a total score
 236 between 4 and 16, with higher scores indicating greater parental verbal responsivity.

237

238 *Home Activities with Child*. The 5-item “Home Activities with Child” scale (Nicholson,
 239 Berthelsen, Abad, Williams, & Bradley, 2008) assessed the frequency with which parents
 240 engage in developmentally important activities with their child in a typical week. The scale is
 241 administered on 4-point Likert scale, where 1=not at all, and 4=every day (e.g. “How often
 242 do you involve your child in everyday activities at home, such as cooking or caring for
 243 pets?”). Item scores were summed to produce a total score between 5 and 20, with higher
 244 scores indicating greater frequency of home activities between the parent and child.

245 **2.3.2 Direct measures**

246 *Preschool Language Scale, Fourth Edition*. The PLS-4 (Zimmerman, Steiner, & Pond,
 247 2002a) is a standardized and norm-referenced instrument to evaluate children’s receptive and
 248 expressive language skills from birth to six years and 11 months. This assessment can be used
 249 as a screening tool for a range of developmental delays such as problems with language,
 250 articulation, connected speech, social communication skills, stuttering, or voice disorders.
 251 Although this measure is normed on a US, rather than Australian sample, ($n=1564$)
 252 (Zimmerman, Steiner, & Pond, 2009), it is one of the most widely used, directly-assessed,
 253 standardized tools for assessing language ability in very young children. The PLS-4 has been
 254 used in other Australian studies with young children (e.g. Ching, Leigh, & Dillon, 2013).
 255 This study reports only on the PLS standard score for expressive language.

256

257 *Early Communication Indicator*. The ECI (Carta, Greenwood, Walker, & Buzhardt, 2010)
 258 aims to assess early communicative development of children aged 6-36 months across four
 259 key domains: vocalizations; single words; multiple words; and gestures. Parents were asked
 260 to play with their child with a standardized set of toys for six minutes while being videotaped.
 261 Accredited expert coders scored video data according to standardized protocols. Frequencies
 262 for each of the four domains were recorded at one-minute intervals. A total communication
 263 composite score was generated by weighting single words by two and multiple words by
 264 three, before summing all four domain scores. Inter-rater agreement on 20% of observations
 265 independently coded by both assessors was 93.9%, consistent with previously reported figures
 266 (Greenwood, Walker, & Buzhardt, 2010). Families from a non-English speaking background
 267 were not instructed to speak English. Rather, all families were asked to “do what they
 268 normally do”. Videos featuring families who spoke a language other than English could not
 269 be coded due to the need to employ interpreters; only families who chose to interact in
 270 English are included in this analysis.

271

272 *Indicator of Parent-Child Interactions*. The IPCI (Carta et al., 2010) was used to quantify the
 273 frequency of specific parent and child behaviors during a set of four common early childhood
 274 activities: free play (4 minutes); looking at books (2 minutes); distraction (2 minutes); and
 275 getting dressed (2 minutes). The distraction task required parents to keep their child on a
 276 small blanket without the child touching a small musical device which was placed within
 277 reach. This activity was not administered to children less than 12 months of age. The
 278 activities are designed to elicit natural interactions which would typically occur between the

279 parent and child. Activities were videotaped, resulting in a total of 8-10 minutes' footage.
 280 Accredited expert coders scored video data according to standardized protocols by counting
 281 the frequency of interactions for each activity across six parent domains: conveys acceptance
 282 and warmth; uses descriptive language; follows child's lead; maintains or extends child's
 283 focus; uses criticism or harsh voice; uses restrictions or intrusions. For each activity, a
 284 relative frequency was allocated to each domain based on a 4-point scale where 0=never;
 285 1=rarely; 2=sometimes or inconsistently; 3=often or consistently. After each activity was
 286 rated, a domain percentage score was calculated by summing all activity scores and dividing
 287 by the total number of possible points for that domain. This study reports on the total positive
 288 caregiver score only, which captures the frequency of responsive parenting behaviors
 289 occurring during the videotaped observation. This total score was generated by summing the
 290 percentage scores for the first four domains listed above. Inter-rater agreement on 20% of
 291 observations independently coded by both assessors was 87.4%, consistent with previously
 292 reported figures (Baggett & Carta, 2006). As described above, all families were asked to "do
 293 what they normally do". However, videos featuring families who spoke a language other than
 294 English could not be coded due to the need to employ interpreters.

295

296 *Socio-demographic Factors*

297 Variables available for both cohorts included: parent age, child age, child gender, parent
 298 education, household income, household unemployment, language other than English and
 299 socio-economic disadvantage. Socio-economic disadvantage was assessed with the Socio-
 300 Economic Indexes for Areas Disadvantage indicator (Australian Bureau of Statistics, 2011),
 301 which summarizes the economic and social circumstances for people and households in a
 302 particular area ($m=1000$; $sd=100$). Lower scores indicate greater disadvantage. A single-item
 303 indicator of child temperament was included for both cohorts (higher scores indicated more
 304 difficult temperament). Additional variables were included in the Early Home Learning Study
 305 analysis due to availability of data, and evidence that these factors may affect parent
 306 responses or behavior: global parenting self-efficacy, assessed using a single-item indicator
 307 ("Overall, as a parent, do you feel that you are ..." not very good at being a parent; a person
 308 who has some trouble being a parent; an average parent; a better than average parent; a very
 309 good parent); psychosocial distress assessed with the K6 (Kessler et al., 2002); and health-
 310 related quality of life evaluated with the SF-12 UK version (Jenkinson & Layte, 1997).

311 **2.3 Procedure**

312 *Language for Learning:* Children identified as slow-to-talk at 18 months were assessed at 24
 313 months by a trained research assistant. Researchers visited families at home to collect parent-
 314 reported data and to administer a standardized child language assessment. *Early Home*
 315 *Learning Study:* Prior to intervention, trained research assistants videotaped parents and
 316 children at home during play activities to examine child language development and parent-
 317 child interactions. Parents also completed a brief measure of child language during the visit.
 318 A 30-minute parent questionnaire was administered via computer-assisted telephone
 319 interview.

320

321 Ethical approval for the Language for Learning study was granted by the Royal Children's
 322 Hospital Human Research Ethics Committee (EHRC #26028) and The University of
 323 Melbourne (#0829736). All parents provided written informed consent. Ethical approval to
 324 access existing Language for Learning data for the current study was covered under the
 325 Centre for Excellence in Child Language and approved by the Royal Children's Hospital

326 Human Research Ethics Committee (HREC #32261 B). Ethical approval for the Early Home
327 Learning Study was granted by the Victorian Government Department of Health (HREC
328 08/10). All parents provided written informed consent. Ethical approval to access existing
329 Early Home Learning Study data for the current study was granted by The University of
330 Melbourne Human Research Ethics Committee (ID 1543863.1).

331 **2.4 Statistical analyses**

332 All analysis was conducted using Stata/IC Version 13.0 (StataCorp, 2013). Prior to analyses,
333 two fathers were excluded from the *Language for Learning* dataset and nine from the Early
334 Home Learning Study dataset, given that parent gender has been found to contribute to
335 differences in data collection (Olino et al., 2013) and the inclusion of such small numbers of
336 fathers was considered insufficient to identify differences between mothers and fathers. A
337 total of nine measures were examined across the two cohorts. Between these measures, nine
338 comparisons were conducted: six compared parent-reported and directly measured behaviors,
339 and three compared parent-reported and parent-reported behaviors. Histograms of the
340 differences were examined for all nine comparisons, followed by scatterplots with a line of
341 best fit to determine linearity. Both Pearson's Correlation Coefficients and Spearman Rank
342 Correlation Coefficients were calculated for each comparison. Pearson's is reported here to
343 enable cross-study comparisons with existing literature, and Spearman's is also reported to
344 account for non-normality of distributions. The Concordance Correlation Coefficient (CCC)
345 was also computed using the Stata "-concord" command. Developed by Lin (1989) as a
346 measure of agreement, the CCC quantifies the degree to which pairs of observations fall on
347 the 45° line through the origin. It contains a measure of precision using the Pearson's
348 Correlation Coefficient, as well as a bias correction for accuracy.

349
350 Z-scores were derived for each of the outcome variables to enable cross-measure
351 comparisons on the same scale. Bland-Altman plots were then generated using the Stata "-
352 concord" command (Cox & Steichen, 2007) for all nine comparisons. This plots the mean of
353 the measures against the difference between the measures, as well as the line of mean
354 difference and the 95% limits of agreement. RMA regression (or "ordinary least products"
355 regression) was conducted using the Stata "-concord" command.

356
357 The associations between a range of socio-demographic factors and the difference between z-
358 scores were estimated using unadjusted and adjusted linear regression. Difference scores
359 were calculated by subtracting one z-score from the other, and these were then used as the
360 outcome variables for the regressions. Unadjusted associations were examined, before the
361 adjusted models were tested. Only variables associated with the outcome at $p \leq .1$ were
362 included in the adjusted models. All continuous variables were screened for evidence of
363 multicollinearity ($r \geq .70$); none were excluded. Factors included in the adjusted models for
364 both cohorts included parent age, child age, child gender, parental education, household
365 income, household unemployment, SEIFA disadvantage score, language other than English
366 and a single-item indicator of child temperament. Additional variables included in the
367 analyses from the Early Home Learning Study dataset were parenting self-efficacy,
368 psychosocial distress, and health-related quality of life. The Stata "-mixed" command was
369 used for this cohort, to account for the cluster-RCT study design and Intraclass Correlation
370 Coefficients were examined.

371
372 Finally, quantile regressions were conducted to determine whether the association between
373 the socio-demographic factors and the difference scores varied across the distribution of the

374 difference scores. Associations were examined across the 25th, 50th and 75th quantiles. Each
 375 model was compared to the standard ordinary least squares output and a test for
 376 heteroscedascity was used to determine whether there was evidence against the null
 377 hypothesis of constant variance across the quantiles.

378 2.5 Sample size

379 Bland (2004) provides a formula to evaluate the precision of the sample size to accurately
 380 assess agreement between measures. Bland advises that the 95% confidence interval around
 381 the limits for agreement may be estimated as $\pm 1.96\sqrt{\left(\frac{3}{n}\right)}s$ where s is the standard deviation of
 382 the differences between measurements by the two methods, and n is the sample size. Bland
 383 recommends that a sample size of 100 is adequate for method comparisons. Applying this
 384 formula provides excellent precision for the Language for Learning cohort of $N=201$ (+/
 385 0.24s). For the Early Home Learning Study, direct measures were only available for a subset
 386 of the cohort (Early Communication Indicator, $N=100$; and Indicator of Parent-Child
 387 Interactions, $N=163$) providing adequate precision for comparisons involving these measures
 388 (+/- 0.34s and +/- 0.27s, respectively).

389 3 Results

390 3.1 Sample

391 Sample characteristics for each study are summarized in Table 2. *Language for Learning*:
 392 Nearly half of the parents had completed higher education and fewer than one in ten families
 393 spoke a non-English language. There were approximately equal proportions of male and
 394 female children, and more than three-quarters of parents were married. Most parents reported
 395 earning a mid to high range household income, with one in five reporting a low income.
 396 *Early Home Learning Study*: Similar characteristics were observed in terms of education,
 397 marital status and child gender compared to families in the Language for Learning study.
 398 However, Early Home Learning Study parents were more likely to be younger, to speak a
 399 language other than English, and to live in a household without an employed person.
 400 Language for Learning participants were on average, less disadvantaged compared with the
 401 Australian mean ($m=1026.6$) and Early Home Learning Study participants were slightly more
 402 disadvantaged ($m=984.2$); however there was also a large degree of variation in scores
 403 (ranges: 888.2–1117.5 and 816.7-1105.9, respectively)
 404

405 Table 2. Sample characteristics for participants in each cohort.

Variable	Language for Learning ($n=201$)	Early Home Learning Study ($n=218$)
Parent age, years, mean (SD)	35.3 (4.4)	32.6 (5.1)
Child age, months, mean (SD)	24.4 (1.1)	16.2 (9.3)
Child female, n (%)	95 (47.0)	113 (51.8)
Parent marital status n , (%)		
Single/separated/divorced	11 (5.5)	17 (7.8)
Married/de facto	190 (94.5)	201 (92.2)
Household unemployment n (%)`	10 (5.0)	18 (8.3)
Parent education, n (%)		
Higher education	93 (46.7)	112 (51.4)
No higher education	106 (53.3)	106 (48.6)
LOTE, n (%) [^]	19 (9.5)	46 (21.1)

Household income p/a, <i>n</i> (%) [*]		
<\$46,800	38 (19.3)	-
\$46,800-\$70,200	69 (35.0)	-
>\$70,200	90 (45.7)	-
<\$36,400	-	26 (12.0)
\$36,400-51,999	-	36 (16.6)
>=\$52,000	-	147 (67.7)
SEIFA [#] , <i>mean</i> (<i>SD</i>)	1026.6 (54.1)	984.2 (57.9)

406 Notes: [`]Single parent unemployed or both parents unemployed; [^]Language other than
 407 English; ^{*}Different categories of income were administered for each sample; [#]Socio-
 408 Economic Index for Areas (SEIFA) Disadvantage score is an indicator of relative
 409 disadvantage, based on postcode of residence, accounting for low income, low educational
 410 attainment and high unemployment. Lower index scores indicate greater disadvantage.

411 3.2 Descriptive statistics

412 The means, standard deviations and ranges for the parent-reported and directly measured
 413 behaviors are presented in Table 3. Alpha coefficients indicate excellent internal consistency
 414 for the Sure Start Language Measure and Preschool Language Scale, consistent with figures
 415 reported elsewhere (Roy et al., 2005; Zimmerman, Steiner, & Pond, 2002b; Zubrick, Taylor,
 416 & Rice, 2007). There was poorer internal consistency for Parental Verbal Responsivity and
 417 the Home Activities with Child scales, which is typically expected for measures with few
 418 items (Gliem & Gliem, 2003). Children's expressive language standard scores on the
 419 Preschool Language Scale (*m* = 91.2, *sd* = 12.3) indicate that, on average, children were
 420 performing below the 50th percentile. As shown, direct measures were only available for a
 421 sub-sample of participants in the Early Home Learning Study cohort. Due to the financial
 422 expenses associated with video coding, the data used in this paper represents a sub-sample of
 423 a larger dataset; this sub-sample was selected at random.

424
 425 Table 3. Descriptives for parent-reported and directly measured behaviors.

	<i>M</i> (<i>SD</i>)	<i>Range</i>	<i>α</i>	<i>N Missing from Total Sample N</i>
<i>Child language</i>				
Sure Start Language Measure ^a	35.0 (22.7)	0 - 98	.97	7/201
Ages & Stages Questionnaire ^{a*}	0 (1)	-2.8 - 1.1	n/a	1/201
Ages & Stages Questionnaire ^{b*}	0 (1)	-3.01 - 1.7	n/a	1/218
Communicative Development Inventory ^b	100.39 (9.7)	80.99 - 160.7	n/a	5/218
Preschool Language Scale ^a	91.2 (12.3)	64 - 135	.86	2/201
Early Communication Indicator ^b	10.1 (7.3)	.3 - 32.3	n/a	118/218
<i>Parenting behaviors</i>				
Parental Verbal Responsivity ^b	12.94 (2.21)	6 - 16	.40	0/218
Home Activities with Child ^b	17.11 (2.52)	9 - 20	.49	0/218
Indicator of Parent-Child Interactions ^b	200.15 (55.3)	50 - 370	n/a	55/218

426 ^aLanguage for Learning; ^bEarly Home Learning Study; ^{*}Z-Scores were derived to account for
 427 different age-appropriate versions.

428 3.3 Correlations

429 The strongest correlations were obtained for comparisons involving two parent-reported
 430 measures, with moderate positive associations (see Table 4). The strongest correlation for any

431 parent-reported and direct comparison was between the Ages & Stages Questionnaire
 432 (communication subscale) and Preschool Language Scale (expressive language), with a
 433 moderate, positive correlation. Weaker associations were obtained for the remaining
 434 comparisons, with a moderate positive correlation between the Communicative Development
 435 Inventory and the Early Communication Indicator, and a weak non-significant correlation
 436 between the Ages & Stages Questionnaire and the Early Communication Indicator.
 437 Associations between measures of parenting behaviors were much weaker than the child
 438 language comparisons, with near-negligible associations between parent-reported and direct
 439 measures. In contrast to the Pearson's and Spearman's coefficients, which produced similar
 440 coefficients for each comparison, the Lin's Concordance Correlation Coefficient produced
 441 markedly smaller correlations for several comparisons. This suggests that, although the
 442 measures are associated, the level of *agreement* is much poorer. Lin's correlation line passes
 443 through the origin, with a slope of one. Thus, it provides a measure of correspondence
 444 between measures, rather than association. As shown in Table 4, the Lin's coefficient for two
 445 comparisons was close to zero, yet highly significant. The confidence intervals for these
 446 comparisons were very narrow, hence the significant p-values.

447
 448 Table 4. Pearson's (r), Spearman's Rank (ρ) and Lin's Concordance (ρ_c) correlation
 449 coefficients for each of the nine comparisons.

	r	ρ	ρ_c
<i>Child language</i>			
1. Ages & Stages Questionnaire vs. Sure Start Language Measure	.70***	.73***	0.70***
2. Ages & Stages Questionnaire vs. Preschool Language Scale	.61***	.59***	0.61***
3. Sure Start Language Measure vs. Preschool Language Measure	.56***	.60***	0.56***
4. Ages & Stages Questionnaire vs. Communicative Development Inventory	.44***	.48***	.001***
5. Ages & Stages Questionnaire vs. Early Communication Indicator	.12	.16	.01
6. Communicative Development Inventory vs. Early Communication Indicator	.32**	.33**	.01**
<i>Parenting behaviors</i>			
7. Parental Verbal Responsivity vs Home Activities with Child	.45***	.45***	.17***
8. Parental Verbal Responsivity vs Indicator of Parent-Child Interaction	-.03	-.04	.00
9. Home Activities with Child vs Indicator of Parent-Child Interaction	.06	.08	.00

450 *** $p \leq .001$; ** $p \leq .01$; * $p < .05$

451 3.4 Agreement between methods

452 Application of the Bland-Altman Method requires the differences between measures to be
 453 approximately normally distributed (Bland & Altman, 2003). Histograms of the differences
 454 showed approximate normality, with slight negative skewness evident for the ASQ-ECI and
 455 ASQ-CDI, and positive skewness for the PVR-HAC. The association between each of the
 456 nine comparisons was examined using scatterplots with a fitted line of equality. Scatterplots
 457 suggested a positive, approximately linear relationship whereby higher scores on one measure
 458 correspond with increasing scores on the other measure. Exceptions were the PVR-IPCI and
 459 HAC-IPCI scatterplots, which did not provide evidence of a linear association. Bland-Altman
 460 Plots were generated for each of the nine comparisons (Figures 1-3). In each plot, the solid
 461 horizontal line represents the mean difference between the measures and the dotted lines
 462 represent the 'limits of agreement' within which 95% of data points lie. The overall bias
 463 (mean difference) was close to zero for most comparisons, reflecting the scaling of the

464 measures to z-scores; we therefore focus on the limits of agreement and patterns of agreement
 465 across the range of scores.

466

467 *Child language.* Figure 1 shows the three Language for Learning comparisons. Plot A shows
 468 the agreement between the parent-reported ASQ and the standardized language measure,
 469 PLS. The points are widely dispersed around the mid-section, indicating poorer agreement for
 470 children with average language abilities. Points are closest to $y=0$ at the lower end, indicating
 471 the strongest agreement for children with the poorest language abilities. Vertical reference
 472 lines at $x = -1$ and $x = 1$ have been included for ease of interpretation. At the upper end of the
 473 spectrum (scores >1 on the x axis), points are all below $y=0$, indicating that the parent-
 474 reported ASQ is systematically underestimating children's language, compared to the direct
 475 measure (PLS). The limits of agreement tell us that 95% of the points lie between -1.75 and
 476 1.73 standard deviations. Plot B shows the agreement between the parent-reported SSLM and
 477 the standardized language measure, the PLS. At the lower end (scores below $x = -1$) points
 478 are more tightly clustered around the line $y=0$, indicating stronger agreement between these
 479 measures for children with poorer language abilities. Agreement then appears to deteriorate
 480 across the spectrum, as children's average language abilities increase. This is shown by the
 481 much wider dispersion of points from $x=0$ and above. The limits of agreement tell us that
 482 95% of the points lie between -1.80 and 1.80 standard deviations. The strongest agreement of
 483 all nine comparisons was found for two parent-reported language measures, the ASQ and the
 484 SSLM (Plot C), with the narrowest limits of agreement (-1.53 to 1.53 standard deviations).
 485 For children with average language abilities (scores between -1 and 1 on the x axis) parents
 486 both underestimate and overestimate on the ASQ, compared to the SSLM. For children with
 487 poorer language abilities (below -1 on the x axis) and higher average language abilities
 488 (above 1 on the x axis), the ASQ produces lower scores than the SSLM.

489

490 Figure 2 shows the three Early Home Learning Study comparisons. Plot A shows the
 491 agreement between the parent-reported ASQ and the direct videotaped observation, the ECI.
 492 For children with average language abilities, parents are over- and under-estimating their
 493 children's language abilities on the ASQ, compared to scores from the directly measured
 494 ECI. For children with poorer language abilities (scores below -1 on the x axis) and stronger
 495 language abilities (scores above 1 on the x axis), most points are positioned below the line
 496 $y=0$. This suggests that parents of children with very poor or very strong average language
 497 abilities are underestimating on the ASQ, compared to the directly measured ECI. The limits
 498 of agreement tell us that 95% of the points lie between -2.55 and 2.48 standard deviations. A
 499 different pattern of agreement is evident between the parent-reported CDI and the ECI (Plot
 500 B), whereby the strongest agreement occurred for children with the poorest language ability,
 501 and agreement progressively deteriorated as children's language ability improved (95% limits
 502 of agreement: -2.21 to 2.33 standard deviations). Not surprisingly, the strongest agreement of
 503 the six Early Home Learning Study comparisons was between the two parent-reported
 504 measures, the ASQ and the CDI (Plot C) (95% limits of agreement: -2.10 to 2.06 standard
 505 deviations). However, the distribution of points suggests that the poorest agreement between
 506 the measures is for children with average language abilities (scores between -1 and 1 on the x
 507 axis). For children with poorer average language abilities (scores < -1 on the x axis) and
 508 stronger average language abilities (scores > 1 on the x axis), the ASQ is underestimating,
 509 compared to the SSLM.

510

511 *Parenting behaviors.* Figure 3 shows poorer agreement between measures of parenting
 512 behaviors compared to the child language measures. Plot A presents agreement between the
 513 parent-reported PVR and the direct videotaped observation, the IPCI. The more dispersed

514 scatter of points around the mid-section reveals that the poorest agreement is for parents of
 515 average responsiveness (95% limits of agreement: -2.78 to 2.80 standard deviations). Parents
 516 with poorer average responsiveness (scores < -1 on the x axis) and parents with stronger
 517 average responsiveness (scores > 1 on the x axis) tend to underestimate their responsiveness
 518 on the PVR, compared to scores on the IPCI. A similar pattern can be seen between the
 519 parent-reported HAC and the IPCI (Plot B), with slightly stronger agreement indicated by
 520 narrower 95% limits of agreement (-2.52 to 2.70 standard deviations). As shown with the
 521 child language comparisons, the strongest agreement between measures of parenting
 522 behaviors was between the two parent-reported measures, the PVR and the HAC (Plot C),
 523 whereby 95% of the points lie between -2.06 and 2.06 standard deviations. The horizontal
 524 scatter of points indicates that the bias between these measures is relatively fixed across the
 525 distribution of scores.

526 **3.5 Identification of proportional bias**

527 Figures 4-6 present the RMA regression plots to identify the presence of proportional bias.
 528 As shown in Figure 4, the Language for Learning language measures show very minimal
 529 proportional bias, evidenced by the slopes which are close to one and the intercepts which are
 530 close to zero. The three Early Home Learning Study child language comparisons also show
 531 minimal proportional bias; however Figure 5A shows a degree of bias between the parent-
 532 reported ASQ and the directly measured videotaped observation, the ECI, indicated by the
 533 slight divergence of lines. Figure 6 shows the three parenting behavior comparisons.
 534 Substantial proportional bias is evident between the parent-reported PVR and the directly
 535 measured videotaped observation, the IPCI (6A). This is shown by the strong divergence of
 536 lines in the plot. The slope of around -1 indicates that for lower PVR scores, IPCI scores are
 537 relatively higher, and for lower IPCI scores, PVR scores are relatively higher. Only slight
 538 proportional bias can be seen between the parent-reported HAC and the IPCI (6B). Figure 6C
 539 indicates the absence of proportional bias between the parent-reported PVR and the parent-
 540 reported HAC.

541 **3.6 Socio-demographic factors and agreement**

542 The results of the adjusted linear regressions are presented in Table 5 for the Language for
 543 Learning cohort and Tables 6 and 7 for the Early Home Learning cohort (See supplementary
 544 tables for unadjusted models). Non-significant variables at the unadjusted level ($p > .1$) were
 545 excluded from the adjusted analyses. The outcome variables in all regression analyses are
 546 difference scores, calculated by subtracting one z-score from another. The intraclass
 547 correlations from the multilevel mixed-effects linear regression for each outcome measure
 548 (Early Home Learning Study cohort) ranged from 0.00 to 0.22. This reflected the cluster
 549 randomized controlled trial design and was accounted for in the regression models.

550

551 *Child language (Language for Learning cohort)*

552 Child age was a significant predictor of difference scores for this cohort. Parents of older
 553 children tended to report higher child language scores on the parent-reported ASQ and
 554 SSLM, compared to scores generated by the directly measured PLS. Older child age was also
 555 associated with lower scores on the ASQ, compared with the SSLM. The included predictors
 556 explained nearly twice the amount of variance in difference scores for the CDI and PLS
 557 (19%), compared to the ASQ and PLS (9%) and the ASQ and CDI (10%).

558

559 *Child language (Early Home Learning Study cohort)*

560 Child age and temperament predicted the difference scores between the parent-reported ASQ
 561 and the directly measured ECI, as well as the parent-reported CDI and ECI. For both
 562 comparisons, older child age was associated with lower scores on the ASQ and CDI,
 563 compared to the ECI. Parents who perceived their child as more difficult also tended to report
 564 lower scores on the ASQ and CDI, compared with the ECI. The included predictors explained
 565 negligible variance in difference scores between the two parent-reported measures, the ASQ
 566 and the CDI (1%), but explained substantial variance between the ASQ and ECI (40%) and
 567 the CDI and ECI (33%).

568

569 *Parenting behaviors (Early Home Learning Study cohort)*

570 The differences between measures of parenting behaviors were associated with parent age
 571 and English language status. Parents who spoke a language other than English were more
 572 likely than native English speakers to report greater parental responsiveness on a parent
 573 questionnaire (PVR or HAC), compared to scores generated from the directly measured IPCI.
 574 Older parents were also more likely to report less parent responsiveness on the parent-
 575 reported PVR and HAC, compared with scores on the IPCI. The included predictors
 576 explained minimal variance in difference scores: PVR and IPCI (18%); HAC and IPCI:
 577 (14%); PVR and HAC (5%).

578 **3.7 Socio-demographic factors across quantiles of agreement**

579 Quantile regression analyses provided scant evidence that the association between the socio-
 580 demographic factors and the difference scores varied across the distribution of the difference
 581 scores. The Breusch-Pagen/Cook-Weisberg test for heteroscedasticity provided non-
 582 significant p -values for eight of the nine comparisons (Language for Learning: ASQ-PLS,
 583 $p=.40$; SSLM-PLS, $p=0.16$; ASQ-SSLM, $p=.31$ and Early Home Learning Study: ASQ-ECI,
 584 $p=.20$; ASQ-CDI, $p=.12$; CDI-ECI, $p=.87$; PVR-IPCI, $p=.87$; PVR-HAC, $p=.11$). This
 585 suggests that the standard Ordinary Least Squares regression is sufficient for quantifying
 586 these associations. However, associations did vary across the distribution of difference scores
 587 for the HAC-IPCI comparison ($p=.03$). Closer inspection of the HAC-IPCI comparison
 588 revealed that income (low vs mid income), varied across the quantiles of difference (25th
 589 quantile: coefficient = .27, $p=.68$; 50th quantile: coefficient = -.64, $p=.30$; 75th quantile:
 590 coefficient = -1.44, $p=.01$). That is, participants with a low income were more likely to have a
 591 large difference between HAC and IPCI scores, compared to participants with a mid-range
 592 income. This finding should be interpreted with caution; given the number of comparisons
 593 made, it is potentially attributable to chance.

594 **4 Discussion**

595 This is the first study to specifically examine agreement between parent-reported and directly
 596 measured behaviors using the Bland-Altman Method and RMA regression. Nine comparisons
 597 were conducted using data from two independent Australian cohorts (6 child language and 3
 598 parenting behaviors). Although correlational findings were consistent with extant literature,
 599 Bland-Altman plots revealed substantial variation in agreement between parent-reported and
 600 directly measured child language and parenting behaviors across the distribution of scores.
 601 Agreement was generally stronger for children with poorer or exceptional language abilities,
 602 and weaker for children with average language abilities. Particularly for comparisons
 603 involving the ASQ, parents tended to underestimate their children's language abilities, when
 604 children's language was either poor or exceptional. Agreement between measures of
 605 parenting behaviors was slightly weaker than child language. Proportional bias between child
 606 language measures was minimal, but considerable bias was evident between parent-reported

607 and directly measured parenting behaviors. Differences between child language measures
608 were associated with child age and temperament, and differences between parenting behavior
609 measures were associated with parent age and speaking a language other than English.
610 Findings provide strong evidence that simple correlations are grossly insufficient for method
611 comparisons.

612 **4.1 Child language**

613 Findings suggest that parent-reported measures are most accurate for children who display
614 either language difficulties or exceptional language abilities. Overall, the strongest agreement
615 was observed for children with the poorest language. This may reflect parental concern and a
616 tendency to more closely observe and monitor child development. Children at either end of
617 the language spectrum may “stand out” from their peers. Reflecting the phenomenon
618 identified in Festiger’s (1954) Social Comparison Theory, parents may rely on social
619 comparisons to inform their decision about their children’s development. Children whose
620 abilities reflect the norm may not generate the same close attention from their parents as
621 children at either end of the spectrum. The variability in child language in the early years is
622 well-established (Ukoumunne et al., 2012), however it is possible that children at the extreme
623 ends of the spectrum are more stable in their language over time, supporting more accurate
624 measurement for these groups. Whereas parent-reported measures may be sufficient to
625 identify children with very poor or very strong language skills, multiple or gold standard
626 direct measures would be necessary to delineate the language skills of children across the mid
627 ranges of child language.

628

629 It should be noted that for comparisons involving the ASQ (Figures 1A, 1C, 2A, 2C) parents
630 tended to underestimate children’s language abilities for children with very poor or
631 exceptional language. This may reflect the limited variability captured by the ASQ, given that
632 it is a 6-item measure scored on a 3-point scale. For comparisons involving the CDI or the
633 UK version of the CDI (SSLM), a different pattern emerged, whereby agreement with direct
634 measures was stronger for children with poorer language ability and progressively worsened
635 as children’s language abilities strengthened. This may reflect a ceiling effect for this
636 commonly-used parent-reported measure of expressive vocabulary, where variation in
637 children with exceptional skills cannot be accurately captured. Indeed, the potential for
638 ceiling effects on the CDI for children aged 27 months and above has been reported
639 elsewhere, particularly for children with more advanced language (Fenson et al., 2000).
640 Together, these findings suggest that accurately capturing the full spectrum of language
641 abilities using parent-reported measures with a small number of items may be problematic.

642

643 The strongest agreement between child language measures was for the Language for
644 Learning cohort. This may reflect the study sample of slow-to-talk toddlers, as well as the use
645 of a standardized language assessment for this cohort, compared with the videotaped
646 observational measure used in the Early Home Learning Study cohort. Some disagreement
647 between measures may be attributable to differences in the constructs captured using each
648 measure. While the Sure Start Language Measure, Communicative Development Inventory,
649 and Preschool Language Scale specifically measure children’s expressive language, the Early
650 Communication Indicator and Ages & Stages Questionnaire include some aspects of non-
651 verbal communication. For example, the Early Communication Indicator includes the
652 frequency of a child’s communicative gestures, as well as vocalizations, single words and
653 multiple words. The Ages & Stages communication subscales include items which measure
654 both expressive and receptive language. The RMA plots provided a clear means of

655 identifying the presence of proportional bias; the six child language plots showed minimal
 656 proportional bias, suggesting that any bias between the measures was relatively consistent
 657 across the distribution of scores.

658

659 The strongest predictor of the difference between language measures was child age; however
 660 the direction of this association varied for each cohort. Parents of older children in the
 661 Language for Learning cohort tended to report higher scores on parent-reported measures,
 662 whereas parents of older children in the Early Home Learning Study cohort tended to report
 663 higher scores on the direct measure. Previous research has shown that parents' ability to
 664 accurately report on their child's language development may deteriorate as children grow
 665 older and their vocabulary expands and language use becomes more complex (Law & Roy,
 666 2008). Differences between these cohorts may also be attributable to the child age ranges (24
 667 months and 6-36 months, respectively), as well as the nature of the selected measures. For
 668 example, parents of children less than 18 months participating in the Early Home Learning
 669 Study were asked about receptive as well as expressive vocabulary. In addition, the Early
 670 Communication Indicator only assessed *observable* features, such as gestures, vocalizations,
 671 single and multiple words. Regardless, it is remarkable that child age was such a highly
 672 significant predictor for the Language for Learning cohort, given the narrow range of child
 673 ages ($M=24.4$ months; $SD= 1.1$ months). At this young age, language develops rapidly and a
 674 small amount of time can produce quite different data. This finding highlights the complexity
 675 of measuring language in young children, as well as the importance of selecting measures
 676 specific to child age in years and months.

677

678 Temperament also emerged as a predictor of child language difference scores, particularly for
 679 the Early Home Learning Study cohort. Perhaps surprisingly, more difficult child
 680 temperament was generally associated with less discrepancy between language measures.
 681 This may be due to parents of children with challenging behaviors having greater awareness
 682 of their child's behavior and development, permitting greater accuracy in parent-reported
 683 measures. Again, this could be more apparent through parents' use of social comparison with
 684 the child's peers. It is also possible that children with behavioral difficulties are the children
 685 with poorer language abilities, for whom the strongest agreement was evident. Indeed, there
 686 is evidence that language and behavioral difficulties can occur comorbidly (Carpenter &
 687 Drabick, 2011). The nature of the assessment – structured assessment or videotaped
 688 observation – as well as the presence of the researcher in the home, may also contribute to
 689 differences between measures of children's expressive language.

690 **4.2 Parenting behaviors**

691 Slightly poorer agreement was observed between measures of parenting behaviors compared
 692 to the language measures. We found relatively strong agreement between the parent-reported
 693 Home Activities with Child and the Indicator of Parent-Child Interactions Positive Caregiver
 694 Score, compared with the parent-reported Parental Verbal Responsivity and the IPCI. As a 4-
 695 item measure, the PVR performed more poorly as an indicator of parental responsiveness,
 696 whereby a ceiling effect led to restricted variation in scores. This measure also showed low
 697 internal consistency, making it a less reliable measure. Both the PVR and HAC showed a
 698 tendency to underestimate parental responsiveness at the lower and upper extremes. Overall,
 699 our findings suggest that a brief parent-reported measure of the frequency of engagement in
 700 parent-child activities in the home (HAC) may represent a reliable indicator of parental
 701 responsiveness and engagement, which shows relatively good agreement with a
 702 comprehensive observational measure. For studies with limited resources, the HAC could be

703 a feasible alternative to time-intensive and costly observation required for the IPCI. It should
704 be acknowledged that some disagreement between the measures of parenting behaviors could
705 be explained by differences in the construct being measured or coded. For example, the PVR
706 measures parents' *verbal* responsiveness specifically, whereas the HAC assesses parent
707 engagement and responsiveness more broadly, including both verbal and non-verbal
708 behaviors. Both the PVR and HAC ask parents about the frequency with which they engage
709 in everyday activities, such as reading books or talking about the day during mealtimes. The
710 Positive Caregiver Total score derived from the IPCI captured the frequency of both verbal
711 and non-verbal parenting behaviors, such as using descriptive language, and following the
712 child's lead (i.e. quantity and quality of parenting behaviors).

713
714 Language other than English was the strongest explanatory factor of the difference between
715 parent-reported and directly measured parenting behaviors. Families with a non-English
716 speaking background tended to report lower scores on the directly measured videotaped
717 observation, the IPCI, and higher scores on both the PVR and HAC. This may be attributable
718 to potential acquiescence (i.e. consistently indicating positive responses). Acquiescence has
719 been shown to vary cross-culturally, for example, strong cultural preferences to avoid
720 uncertainty can lead to a tendency to select more extreme values (Smith, 2004). Findings may
721 also reflect cultural differences in the frequency with which parents and children engage in
722 the activities being measured (e.g. HAC: "telling stories to your child" or PVR: "playing
723 peek-a-boo or hide-and-seek"). It is also possible that parents and children with a non-English
724 speaking background felt less comfortable than native English speakers during the videotaped
725 activities. Furthermore, these differences could be attributable to difficulties understanding
726 the verbal instructions of the videotaped activities, or difficulties in coding parent utterances
727 during these activities. Lastly, we acknowledge that parents' English proficiency may vary to
728 that of the child, particularly in early childhood when children have not yet been exposed to
729 English in the school environment.

730
731 The small proportion of variance explained by the socio-demographic factors for parenting
732 behaviors suggests that other unmeasured factors may be responsible for differences between
733 these measures. The current study was limited by the data collected in the two datasets
734 analyzed; it is possible that other factors may have greater explanatory power than variables
735 assessed in these studies. For example, the parent or child's unique and subjective experience
736 of the assessments, understanding of the task requirements or the questionnaire items, cultural
737 factors affecting parent-child interactions, discomfort during the assessment, rapport with the
738 assessor, experiences of fatigue or illness at the time of the assessment or external factors
739 causing stress or distraction may have been more relevant predictors of agreement. Quantile
740 regression analyses revealed that the associations between the socio-demographic factors and
741 the difference scores remained stable across the quantiles of agreement. The only exception
742 was the comparison between the parent-reported HAC and the directly measured videotaped
743 observation, the IPCI. Greater discrepancy between these measures was associated with
744 parents with a lower income. The five HAC items refer to everyday parent-child activities;
745 however, many of these activities require resources such as books and toys, which may be
746 less readily available for parents who have a very low income. Indeed, this link between
747 families of a lower socio-economic status and the provision of a less stimulating home
748 environment is well-established (e.g. Davis-Kean, 2005)

749 4.3 Implications

750 This study provides evidence to guide the selection of appropriate measures for parents and
751 their children aged 6-36 months. Method comparisons such as this are critical for supporting
752 the collection of high data quality and the appropriate allocation of limited resources. Our
753 data suggest that brief parent-reported measures of child language may be used with
754 reasonable confidence for children up to three years of age. Particularly for children who are
755 slow-to-talk, parent-reported measures may provide an accurate and cost-effective means of
756 monitoring development over time. Findings indicate that agreement between measures of
757 parenting behaviors is generally poorer than child language measures. Parenting behaviors
758 can be difficult to accurately measure, given that social desirability can cause parents to
759 consciously or unconsciously change the way they respond on parent-reported questionnaires
760 (Law & Roy, 2008; Zaslow et al., 2006), or the way they behave during observations (Arney,
761 2004). It is also conceivable that parents are more able to objectively report on their child's
762 language but are less objective when evaluating their own behaviors (e.g. parenting
763 responsiveness). Despite this, the parent-reported Home Activities with Child measure
764 showed relatively strong agreement with the direct videotaped observation, the Indicator of
765 Parent-Child Interactions, with minimal proportional bias. This suggests that measuring the
766 frequency of developmentally beneficial activities such as reading, story-telling, singing, or
767 involving the child in everyday tasks at home, provides a valid indication of parents' general
768 level of engagement and responsiveness.

769
770 When selecting measures, it is important to consider the purpose for which the data is being
771 generated; a brief parent-reported measure of children's expressive language or
772 communicative development such as the Sure Start Language Measure, Communicative
773 Development Inventory or Ages & Stages Questionnaire may be sufficient for large-scale
774 studies where time and resources are limited and a large pool of data is required. Whereas a
775 clinician making decisions about treatment options for a young child may be best to draw on
776 both direct and parent-reported measures to ensure a comprehensive assessment.

777
778 The study has significant implications for the analysis of method comparisons. We
779 demonstrate how the Bland-Altman Method and RMA regression permit a comprehensive
780 assessment of agreement across the distribution of scores. While correlational analyses
781 reported here were comparable to those reported elsewhere for similar constructs, analyses
782 using the Bland-Altman Method and RMA regression clearly show how correlations have the
783 potential to be misleading. Correlations represent a single figure which summarizes a linear
784 association across the spectrum of scores, whereas agreement may vary between higher and
785 lower scores. The level of detail generated by these more comprehensive techniques is crucial
786 for identifying groups of children or parents for whom one method may be sufficient (in the
787 case of strong agreement), or for whom multiple methods or an agreed "gold standard"
788 measure may be necessary (in the case of poor agreement).

789 4.4 Strengths and limitations

790 To our knowledge, this is the first study to apply the Bland-Altman Method to a comparison
791 of parent-reported and directly measured behaviors. This technique permitted the
792 identification of patterns of bias across the distribution of scores. As a result, we were able to
793 identify groups of children or parents for whom multi-method administration may be
794 necessary, or for whom one method of measurement may be permissible. Rarely used in non-
795 medical fields, the Bland-Altman method represents a relatively simple and visually
796 appealing technique. The approach lends itself to other comparisons such as parent-, teacher-

797 and child-report of the same questionnaire (e.g. Gabbe et al., 2010; Stolarova et al., 2014), or
798 comparisons of the same measure across time points (e.g. Eadie et al., 2014). Another
799 strength is the use of RMA regression to identify the magnitude of proportional bias between
800 methods. Together, Bland-Altman and RMA regression plots represent powerful visuals for
801 comparing measures which can be executed and interpreted with relative ease. The use of
802 quantile regression analyses also allowed us to determine whether associations between
803 socio-demographic factors and agreement varied across quantiles of agreement, which is not
804 possible using standard ordinary least squares regression.

805

806 We acknowledge that we were limited to the measures available within existing datasets, and
807 therefore cannot presume agreement findings are generalizable to other measures of child
808 language and parenting behaviors. Despite this, our measures are commonly used and well-
809 validated. It should be noted that the PLS-4 was only normed on US data at the time of data
810 collection; no Australian norms were available. Data also pertained to a sample of toddlers
811 identified as “slow-to-talk” at age 18 months, and another sample of families experiencing
812 social disadvantage; different populations may yield different results. We also recognize that
813 each of the measures used in this study will, naturally, capture slightly different aspects of
814 child language or parental responsiveness. As with any method comparison, total agreement
815 is not expected, nor is it feasible to strive for this; some degree of measurement error is
816 inevitable (Bland & Altman, 1999). Regardless, method comparisons are critical for
817 determining whether measures are potentially interchangeable, and may contribute to more
818 effective allocation of limited resources and strengthened data quality.

819 **4.5 Future research**

820 We suggest that researchers consider applying these techniques to method comparisons of
821 other commonly used early childhood language measures, such as Clinical Evaluation of
822 Language Fundamentals (CELF), and with larger sample sizes where possible to ensure
823 greater precision around the limits of agreement. Our future research will employ qualitative
824 methodologies to determine how parents’ unique and subjective experiences of assessments
825 may further explain and contextualize agreement findings. This is particularly important
826 given that a broad range of socio-demographic factors explained little variability in the
827 difference scores for a number of measures. It is possible that parents and children vary in
828 their level of comfort when behaviors are being measured directly (i.e. videotaped
829 observations or standardized assessments), especially for participants who are not native
830 English speakers. Exploring this qualitatively could go some way to understanding agreement
831 and supporting data collection methods which optimize the validity of parent and child data.

832 **4.6 Conclusions**

833 This study demonstrates how well-established statistical techniques from non-psychology
834 disciplines can be applied to method comparisons in the field of psychology. The Bland-
835 Altman Method is a useful visual technique for detecting bias and for determining potential
836 interchangeability between measurement methods, which can be used in combination with
837 RMA regression to identify the presence of both fixed and proportional bias. Although we
838 found correlations which were consistent with previous comparisons of child language and
839 parenting behaviors, agreement varied substantially across the distribution of scores,
840 demonstrating the need for these more comprehensive techniques. On the whole, poorer
841 agreement was observed for children with average expressive language abilities, and stronger
842 agreement was observed for children with very poor or more advanced language abilities.
843 Slightly poorer agreement was observed between measures of parenting behaviors, with the

844 weakest agreement seen for parents of average responsiveness. As would be expected,
 845 stronger agreement was observed between comparisons of two parent-reported measures.
 846 Further research is required to determine agreement between other commonly used measures
 847 and how the participant experience may explain agreement between parent-reported and
 848 directly measured behaviors. We recommend that journal editors encourage the use of the
 849 Bland-Altman Method and RMA regression techniques and discourage the use of correlations
 850 for method comparisons.

851 **5 Acknowledgements**

852 We thank the Language for Learning and Early Home Learning Study project teams and all
 853 participating families, and acknowledge the support of the NHMRC-funded Centre of
 854 Research Excellence in Child Language (#1023493). The Early Home Learning Study was
 855 funded by the Department of Early Education and Childhood Development (DEECD) and
 856 conducted by the Parenting Research Centre. The authors gratefully acknowledge the support
 857 of the Parenting Research Centre in providing access to this data. We also thank A/Prof
 858 Kathleen Baggett, A/Prof Jay Buzhardt and research staff at The University of Kansas for
 859 coding the videotaped data collected for the Early Home Learning Study.

860
 861 Shannon Bennetts is a PhD Candidate, supported by the NHMRC-funded Centre of Research
 862 Excellence in Child Language (#1023493), Dr Fiona Mensah is a Biostatistician supported by
 863 an NHMRC Early Career Fellowship (#1037449) and Career Development Fellowship
 864 (#1111160), Professor Sheena Reilly is supported by an NHMRC Practitioner Fellowship
 865 (#1041892) and Ms Shannon Bennetts, Dr Elizabeth Westrupp and Dr Naomi Hackworth are
 866 supported by Australian Communities Foundation through the Roberta Holmes Transition to
 867 Contemporary Parenthood Program (Coronella sub-fund) at La Trobe University. The
 868 contents of the published material herein are the sole responsibility of the authors and do not
 869 reflect the views of the NHMRC. Research at the Murdoch Childrens Research Institute is
 870 supported by the Victorian Government's Operational Infrastructure Support Program.

871 **6 Author Contributions**

872 SB was responsible for leading the preparation of this manuscript. SB, FM, EW, NH and SR
 873 all made substantial contributions to the study design, analysis, interpretation, writing and
 874 revision of the manuscript. All authors provided approval for publication.

875 **7 References**

- 876 Altman, D., & Bland, M. (1983). Measurement in medicine: The analysis of method
 877 comparison studies. *Journal of the Royal Statistical Society. Series D (The*
 878 *Statistician)*, 32(3), 307-317. doi: 10.2307/2987937
- 879 Arney, F. (2004). *A comparison of direct observation and self-report measures of parenting*
 880 *behaviour*. (PhD), Adelaide University, Adelaide.
- 881 Aspland, H., & Gardner, F. (2003). Observational measures of parent-child interaction: An
 882 introductory review. *Child and Adolescent Mental Health*, 8(3), 136-143. doi:
 883 10.1111/1475-3588.00061
- 884 Australian Bureau of Statistics. (2011). Socio-Economic Indexes for Areas (SEIFA).
 885 Retrieved May 5, 2016, from
 886 <http://www.abs.gov.au/ausstats/abs@.nsf/mf/2033.0.55.001>

- 887 Baggett, K., & Carta, J. (2006). Using assessment to guide social-emotional intervention for
 888 very young children: An individual growth and development indicator (IGDI) of
 889 parent-child interaction. *Young Exceptional Children Monograph Series*, 8, 67-76.
- 890 Bennetts, S.K., Mensah, F.K., Westrupp, E.M., Hackworth, N., Nicholson, J.M., & Reilly, S.
 891 (2016). Establishing agreement between parent-reported and directly measured
 892 behaviours. *Australasian Journal of Early Childhood* [in review].
- 893 Bland, M. (2004). How can I decide the sample size for a study of agreement between two
 894 methods of measurement? Retrieved May 19, 2016, from [https://www-](https://www-users.york.ac.uk/~mb55/meas/sizemeth.htm)
 895 [users.york.ac.uk/~mb55/meas/sizemeth.htm](https://www-users.york.ac.uk/~mb55/meas/sizemeth.htm)
- 896 Bland, M., & Altman, D. (1986). Statistical methods for assessing agreement between two
 897 methods of clinical measurement. *The Lancet*, 327(8476), 307-310. doi:
 898 10.1016/S0140-6736(86)90837-8
- 899 Bland, M., & Altman, D. (1999). Measuring agreement in method comparison studies.
 900 *Statistical Methods in Medical Research*, 8(2), 135-160. doi:
 901 10.1177/096228029900800204
- 902 Bland, M., & Altman, D. (2003). Applying the right statistics: Analyses of measurement
 903 studies. *Ultrasound in Obstetrics and Gynecology*, 22(1), 85-93. doi: 10.1002/uog.122
- 904 Campbell, F., & Ramey, C. (1994). Effects of early intervention on intellectual and academic
 905 achievement: A follow-up study of children from low-income families. *Child*
 906 *Development*, 65(2), 684-698.
- 907 Carpenter, J.L., & Drabick, D.A.G. (2011). Co-occurrence of linguistic and behavioural
 908 difficulties in early childhood: a developmental psychopathology perspective. *Early*
 909 *Child Development and Care*, 181(8), 1021-1045. doi:
 910 10.1080/03004430.2010.509795
- 911 Carstensen, B. (2010). *Comparing Clinical Measurement Methods: A Practical Guide*.
 912 Chichester, UK: John Wiley & Sons Ltd.
- 913 Carta, J., Greenwood, C., Walker, D., & Buzhardt, J. (2010). *Using IGDIs: Monitoring*
 914 *Progress and Improving Intervention for Infants and Young Children*. Baltimore,
 915 MD.: Brookes Publishing Company.
- 916 Cartmill, E., Armstrong, B., Gleitman, L., Goldin-Meadow, S., Medina, T., & Trueswell, J.
 917 (2013). Quality of early parent input predicts child vocabulary 3 years later.
 918 *Proceedings of the National Academy of Sciences*, 110(28), 11278-11283. doi:
 919 10.1073/pnas.1309518110
- 920 Ching, T.Y., Leigh, G., & Dillon, H. (2013). Introduction to the longitudinal outcomes of
 921 children with hearing impairment (LOCHI) study: background, design, sample
 922 characteristics. *International Journal of Audiology*, 52 Suppl 2, S4-9. doi:
 923 10.3109/14992027.2013.866342
- 924 Cox, N., & Steichen, T.J. (2007). CONCORD: Stata module for concordance correlation.
 925 Retrieved from <https://ideas.repec.org/c/boc/bocode/s404501.html>
- 926 Davis-Kean, P.E. (2005). The influence of parent education and family income on child
 927 achievement: the indirect role of parental expectations and the home environment.
 928 *Journal of Family Psychology*, 19(2), 294. doi: 10.1037/0893-3200.19.2.294
- 929 Dreyer, B.P., Mendelsohn, A.L., & Tamis-LeMonda, C.S. (1996). Assessing the child's
 930 cognitive home environment through parental report; reliability and validity. *Early*
 931 *Development and Parenting*, 5(4), 271-287. doi: 10.1002/(SICI)1099-
 932 0917(199612)5:4<271::AID-EDP138>3.0.CO;2-D
- 933 Eadie, P., Nguyen, C., Carlin, J., Bavin, E., Bretherton, L., & Reilly, S. (2014). Stability of
 934 language performance at 4 and 5 years: measurement and participant variability.
 935 *International Journal of Language & Communication Disorders*, 49(2), 215-227. doi:
 936 10.1111/1460-6984.12065

- 937 Feldman, H., Dollaghan, C., Campbell, T., Kurs-Lasky, M., Janosky, J., & Paradise, J.
 938 (2000). Measurement properties of the MacArthur Communicative Development
 939 Inventories at ages one and two years. *Child Development*, 71(2), 310-322. doi:
 940 10.1111/1467-8624.00146
- 941 Fenson, L., Pethick, S., Renda, C., Cox, J.L., Dale, P.S., & Reznick, S.J. . (2000). Short-form
 942 versions of the MacArthur Communicative Development Inventories. *Applied*
 943 *Psycholinguistics*, 21(01), 95-116.
- 944 Festinger, L. (1954). A theory of social comparison processes. *Human Relations*, 7(2), 117-
 945 140. doi: 10.1177/001872675400700202
- 946 Gabbe, B., Simpson, P., Sutherland, A., Palmer, C., Butt, W., Bevan, C., & Cameron, P.
 947 (2010). Agreement between parent and child report of health-related quality of life:
 948 Impact of time postinjury. *Journal of Trauma and Acute Care Surgery*, 69(6), 1578-
 949 1582 doi: 10.1097/TA.0b013e3181f8fd5f
- 950 Gardner, F. (1997). Observational methods for recording parent-child interaction: How
 951 generalisable are the findings? *Child and Adolescent Mental Health*, 2(2), 70-74. doi:
 952 10.1111/j.1475-3588.1997.tb00049.x
- 953 Gardner, F. (2000). Methodological issues in the direct observation of parent-child
 954 interaction: Do observational findings reflect the natural behavior of participants?
 955 *Clinical Child and Family Psychology Review*, 3(3), 185-198. doi:
 956 10.1023/A:1009503409699
- 957 Gartstein, M., & Marmion, J. (2008). Fear and positive affectivity in infancy:
 958 Convergence/discrepancy between parent-report and laboratory-based indicators.
 959 *Infant Behavior and Development*, 31(2), 227-238. doi: 10.1016/j.infbeh.2007.10.012
- 960 Gliem, J.A., & Gliem, R.R. (2003). *Calculating, interpreting, and reporting Cronbach's*
 961 *alpha reliability coefficient for Likert-type scales*. Paper presented at the Midwest
 962 Research-to-Practice Conference in Adult, Continuing, and Community Education,
 963 Ohio, USA.
- 964 Greenwood, C.R., Walker, D., & Buzhardt, J. (2010). The Early Communication Indicator for
 965 infants and toddlers: Early Head Start growth norms from two states. *Journal of Early*
 966 *Intervention*, 32(5), 310-334. doi: 10.1177/1053815110392335
- 967 Hawes, D., & Dadds, M. (2006). Assessing parenting practices through parent-report and
 968 direct observation during parent-training. *Journal of Child and Family Studies*, 15(5),
 969 554-567. doi: 10.1007/s10826-006-9029-x
- 970 Hayden, E., Durbin, E., Klein, D., & Olino, T. (2010). Maternal personality influences the
 971 relationship between maternal reports and laboratory measures of child temperament.
 972 *Journal of Personality Assessment*, 92(6), 586-593. doi:
 973 10.1080/00223891.2010.513308
- 974 Jenkinson, C., & Layte, R. (1997). Development and testing of the UK SF-12. *Journal of*
 975 *Health Services Research*, 2(1), 14-18. doi: 10.1177/135581969700200105
- 976 Kessler, R.C., Andrews, G., Colpe, L.J., Hiripi, E., Mroczek, D.K., Normand, S.L., . . .
 977 Zaslavsky, A.M. (2002). Short screening scales to monitor population prevalences and
 978 trends in non-specific psychological distress. *Psychological Medicine*, 32(06), 959-
 979 976.
- 980 Law, J., & Roy, P. (2008). Parental report of infant language skills: A review of the
 981 development and application of the Communicative Development Inventories. *Child*
 982 *and Adolescent Mental Health*, 13(4), 198-206. doi: 10.1111/j.1475-
 983 3588.2008.00503.x
- 984 Lin, L.I. (1989). A concordance correlation coefficient to evaluate reproducibility.
 985 *Biometrics*, 45(1), 255-268. doi: 10.2307/2532051

- 986 Ludbrook, J. (1997). Comparing methods of measurements. *Clinical and Experimental*
 987 *Pharmacology and Physiology*, 24(2), 193-203.
- 988 Ludbrook, J. (2010). Linear regression analysis for comparing two measurers or methods of
 989 measurement: But which regression? *Clinical and Experimental Pharmacology and*
 990 *Physiology*, 37(7), 692-699. doi: 10.1111/j.1440-1681.2010.05376.x
- 991 Meisenberg, G., & Williams, A. (2008). Are acquiescent and extreme response styles related
 992 to low intelligence and education? *Personality and Individual Differences*, 44(7),
 993 1539-1550. doi: 10.1016/j.paid.2008.01.010
- 994 Miles, J., & Banyard, P. (2007). *Understanding and Using Statistics in Psychology: A*
 995 *Practical Introduction*. London: SAGE Publications Ltd.
- 996 Nicholson, J.M., Berthelsen, D., Abad, V., Williams, K., & Bradley, J. (2008). Impact of
 997 music therapy to promote positive parenting and child development. *Journal of*
 998 *Health Psychology*, 13(2), 226-238. doi: 10.1177/1359105307086705
- 999 Olino, T., Durbin, E., Klein, D., Hayden, E., & Dyson, M. (2013). Gender differences in
 1000 young children's temperament traits: Comparisons across observational and parent-
 1001 report methods. *Journal of Personality*, 81(2), 119-129. doi: 10.1111/jopy.12000
- 1002 Reese, E., & Read, S. (2000). Predictive validity of the New Zealand MacArthur
 1003 Communicative Development Inventory: Words and Sentences. *Journal of Child*
 1004 *Language*, 27(02), 255-266.
- 1005 Reilly, S., Bavin, E.L., Bretherton, L., Conway, L., Eadie, P., Cini, E., . . . Wake, M. (2009).
 1006 The Early Language in Victoria Study (ELVS): A prospective, longitudinal study of
 1007 communication skills and expressive vocabulary development at 8, 12 and 24 months.
 1008 *International Journal of Speech-Language Pathology*, 11(5), 344-357. doi:
 1009 10.1080/17549500903147560
- 1010 Reilly, S., Wake, M., Ukoumunne, O., Bavin, E., Prior, M., Cini, E., . . . Bretherton, L.
 1011 (2010). Predicting language outcomes at 4 years of age: Findings from the Early
 1012 Language in Victoria Study. *Pediatrics* 126(6), e1530-e1537 doi: 10.1542/peds.2010-
 1013 0254
- 1014 Ring, E., & Fenson, L. (2000). The correspondence between parent report and child
 1015 performance for receptive and expressive vocabulary beyond infancy. *First*
 1016 *Language*, 20(59), 141-159. doi: 10.1177/014272370002005902
- 1017 Roberts, J., Burchinal, M., & Durham, M. (1999). Parents' report of vocabulary and
 1018 grammatical development of African American preschoolers: Child and
 1019 environmental associations. *Child Development*, 70(1), 92-106. doi: 10.1111/1467-
 1020 8624.00008
- 1021 Roy, P., Kersley, H., & Law, J. (2005). The Sure Start Language Measure Standardisation
 1022 Study.
 1023 <http://tna.europarchive.org/20070101101348/http://www.dfes.gov.uk/research/progra>
 1024 mmeofresearch/projectinformation.cfm?projectid=14628&resultspage=1
- 1025 Sachse, S., & Von Suchodoletz, W. (2008). Early identification of language delay by direct
 1026 language assessment or parent report? *Journal of Developmental & Behavioral*
 1027 *Pediatrics*, 29(1), 34-41 doi: 10.1097/DBP.0b013e318146902a
- 1028 Sim, S. (2012). *Supporting children's language and literacy skills: The effectiveness of*
 1029 *shared book reading intervention strategies with parents*. (Doctor of Philosophy),
 1030 Queensland University of Technology, Queensland, Australia.
- 1031 Skeat, J., Eadie, P., Ukoumunne, O., & Reilly, S. (2010). Predictors of parents seeking help
 1032 or advice about children's communication development in the early years. *Child:*
 1033 *Care, Health and Development*, 36(6), 878-887. doi: 10.1111/j.1365-
 1034 2214.2010.01093.x

- 1035 Smith, P.B. (2004). Acquiescent response bias as an aspect of cultural communication style.
 1036 *Journal of Cross-Cultural Psychology*, 35(1), 50-61. doi:
 1037 10.1177/0022022103260380
- 1038 Squires, J., Twombly, E., Bricker, D., & Potter, L. . (2009). *ASQ-3 User's Guide*. Oregon:
 1039 Brookes Publishing Co. Inc.
- 1040 StataCorp. (2013). *Stata Statistical Software: Release 13*. College Station, TX: StataCorp LP.
- 1041 Stolarova, M., Wolf, C., Rinker, T., & Brielmann, A. (2014). How to assess and compare
 1042 inter-rater reliability, agreement and correlation of ratings: An exemplary analysis of
 1043 mother-father and parent-teacher expressive vocabulary rating pairs. *Frontiers in*
 1044 *Psychology*, 5(509). doi: 10.3389/fpsyg.2014.00509
- 1045 Stothard, S., Snowling, M., Bishop, D., Chipchase, B., & Kaplan, C. (1998). Language-
 1046 impaired preschoolers: A follow-up into adolescence. *Journal of Speech, Language,*
 1047 *and Hearing Research*, 41, 407–418.
- 1048 Ukoumunne, O., Wake, M., Carlin, J., Bavin, E., Lum, J., Skeat, J., . . . Reilly, S. (2012).
 1049 Profiles of language development in pre-school children: a longitudinal latent class
 1050 analysis of data from the Early Language in Victoria Study. *Child: Care, Health and*
 1051 *Development*, 38(3), 341-349. doi: 10.1111/j.1365-2214.2011.01234.x
- 1052 Wysocki, T. (2015). Introduction to the special issue: Direct observation in pediatric
 1053 psychology research. *Journal of Pediatric Psychology*, 40(1), 1-7. doi:
 1054 10.1093/jpepsy/jsu104
- 1055 Zaslow, M., Weinfield, N., Gallagher, M., Hair, E., Ogawa, J., Egeland, B., . . . De Temple, J.
 1056 (2006). Longitudinal prediction of child outcomes from differing measures of
 1057 parenting in a low-income sample. *Developmental Psychology*, 42(1), 27-37. doi:
 1058 10.1037/0012-1649.42.1.27
- 1059 Zimmerman, I.L., Steiner, V.G., & Pond, R.E. . (2002a). *Preschool Language Scale* (4th
 1060 ed.). San Antonio, TX: Psychcorp.
- 1061 Zimmerman, I.L., Steiner, V.G., & Pond, R.E. . (2002b). *Preschool Language Scale Fourth*
 1062 *Edition: Examiner's Manual*. San Antonia, TX: Harcourt Assessment.
- 1063 Zimmerman, I.L., Steiner, V.G., & Pond, R.E. . (2009). Technical Report: Preschool
 1064 Language Scale Fourth Edition.
- 1065 Zubrick, S.R., Taylor, C.L., & Rice, M.L. (2007). Late language emergence at 24 months: An
 1066 epidemiological study of prevalence, predictors, and covariates. *Journal of speech,*
 1067 *language, and hearing research : JSLHR*, 50(6), 1562-1592. doi: 10.1044/1092-
 1068 4388(2007/106)
- 1069
- 1070
- 1071
- 1072

1073 Table 5. Adjusted analysis for Language for Learning difference scores and socio-demographic factors.

	ASQ vs. PLS-E			SSLM vs. PLS-E			ASQ vs. SSLM		
	Coeff	p	95% CI	Coeff.	p	95% CI	Coeff.	p	95% CI
Parent age (years)	*	*	*	*	*	*	.03	.05	.00, .07
Child age (months)	.23	<.001	.12, .34	.35	<.001	.23, .46	-.13	.02	-.24, -.02
Child gender (female)	*	*	*	.18	.14	-.06, .42	*	*	*1079
Single parent	*	*	*	.37	.24	-.24, .98	*	*	*1080
Household unemployment	*	*	*	*	*	*	-.20	.50	-.79, .38
No higher education	*	*	*	*	*	*	*	*	*1082
Income									1083
low vs. mid	*	*	*	-.22	.25	-.59, .16	.27	.12	-.07, .61
low vs. high	*	*	*	-.33	.08	-.70, .04	.38	.02	.06, .71
SEIFA/100 (less disadvantage)	-.23	.05	-.46, .00	-.15	.22	-.39, .09	*	*	*1086
LOTE	*	*	*	*	*	*	.29	.15	-.11, .69
Difficult child temperament	*	*	*	*	*	*	-.18	.02	-.33, -.03
	$R^2 = .09$			$R^2 = .19$			$R^2 = .10$		

*excluded due to $p > .1$ at univariate level; LOTE=Language other than English.

1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107

1108 Table 6. Adjusted analysis for the Early Home Learning Study difference scores and socio-demographic factors (child language measures).

	ASQ vs. ECI			CDI vs ECI			ASQ vs CDI		
	Coeff.	<i>p</i>	95% CI	Coeff.	<i>p</i>	95% CI	Coeff.	<i>p</i>	95% CI
Parent age (years)	*	*	*	*	*	*	*	*	*1111
Child age (months)	-.09	<.001	-.12, -.06	-.07	<.001	-.10, -.04	*	*	*1112
Child gender (female)	.24	.24	-.16, .64	*	*	*	*	*	*1113
Single parent	*	*	*	*	*	*	*	*	*1114
Household unemployment	*	*	*	*	*	*	*	*	*1115
No higher education	*	*	*	*	*	*	*	*	*1116
Income									1117
low vs mid	*	*	*	*	*	*	*	*	*1118
low vs high	*	*	*	*	*	*	*	*	*1119
SEIFA/100 (Less disadvantage)	*	*	*	*	*	*	*	*	*1120
LOTE	*	*	*	*	*	*	.31	.08	-.04, .66 *1121
Difficult child temperament	-.50	.05	-.99, -.01	-.51	.02	-.95, -.07	*	*	*1122
High parenting self-efficacy	.06	.63	-.18, .30	*	*	*	*	*	*1123
Poor health-related quality of life	-.07	.54	-.29, .15	-.10	.34	-.31, .11	*	*	*1124
Greater psychological distress	*	*	*	*	*	*	*	*	*1125
	R ² =.40			R ² =.33			R ² =.01		
									1126 1127

1128 *excluded due to $p > .1$ at univariate level; LOTE=Language other than English.

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142 Table 7. Adjusted analysis for the Early Home Learning Study difference scores and socio-demographic factors (parenting behavior measures).

	PVR vs IPCI			HAC vs IPCI			PVR vs HAC		
	Coeff.	<i>p</i>	95% CI	Coeff.	<i>p</i>	95% CI	Coeff.	<i>p</i>	95% CI
Parent age (years)	-.04	.07	-.08, .00	-.04	.02	-.08, .00	*	*	*
Child age (months)	-.01	.32	-.04, .01	*	*	*	-.02	<.01	-.04, -.01
Child gender (female)	*	*	*	-.47	.02	-.85, -.09	*	*	*
Single parent	*	*	*	*	*	*	*	*	*
Household unemployment	*	*	*	*	*	*	*	*	*
No higher education	.53	.01	.12, .93	*	*	*	*	*	*
Income									
low vs. mid	*	*	*	*	*	*	*	*	*
low vs. high	*	*	*	*	*	*	*	*	*
SEIFA/100	*	*	*	*	*	*	*	*	*
LOTE	1.23	<.001	.66, 1.79	1.09	<.001	.56, 1.62	*	*	*
Difficult child temperament	-.43	.11	-.96, .10	*	*	*	-.10	.59	-.46, .26
Low parenting self-efficacy	.17	.16	-.07, .40	*	*	*	*	*	*
Poor health-related quality of life	*	*	*	*	*	*	*	*	*
Greater psychological distress	*	*	*	*	*	*	*	*	*
	$R^2 = .18$			$R^2 = .14$			$R^2 = .05$		

1143 *excluded due to $p > .1$ at univariate level; LOTE=Language other than English.

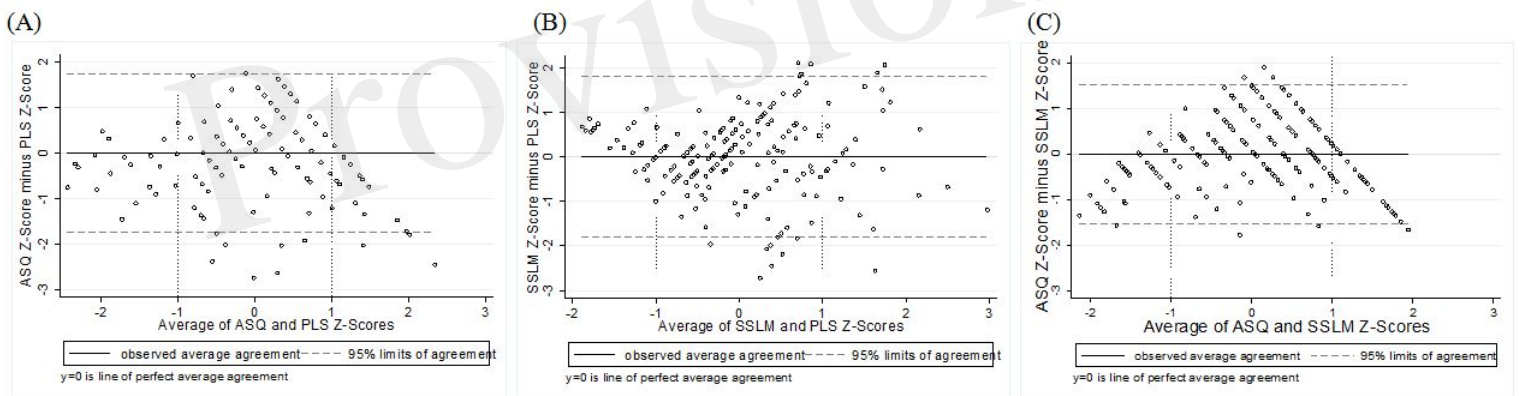


Figure 1. Bland-Altman Plots: (A) Ages & Stages Questionnaire and Preschool Language Scale [95% limits of agreement: -1.75 to 1.73]; (B) Sure Start Language Measure and Preschool Language Scale [95% limits of agreement: -1.80 to 1.80]; (C) Ages & Stages Questionnaire and Sure Start Language Measure [95% limits of agreement: -1.52 to 1.53].

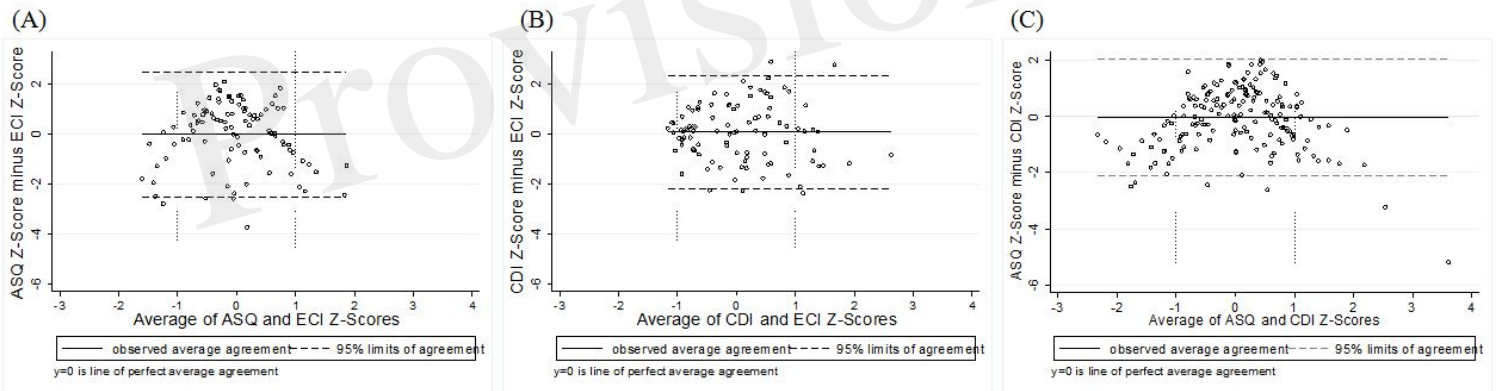


Figure 2. Bland-Altman Plots: (A) Ages & Stages Questionnaire and Early Communication Indicator [95% limits of agreement: -2.55 to 2.48]; (B) Communicative Development Inventory and Early Communication Indicator [95% limits of agreement: -2.21 to 2.33]; (C) Ages & Stages Questionnaire and Communicative Development Inventory [95% limits of agreement: -2.10 to 2.06].

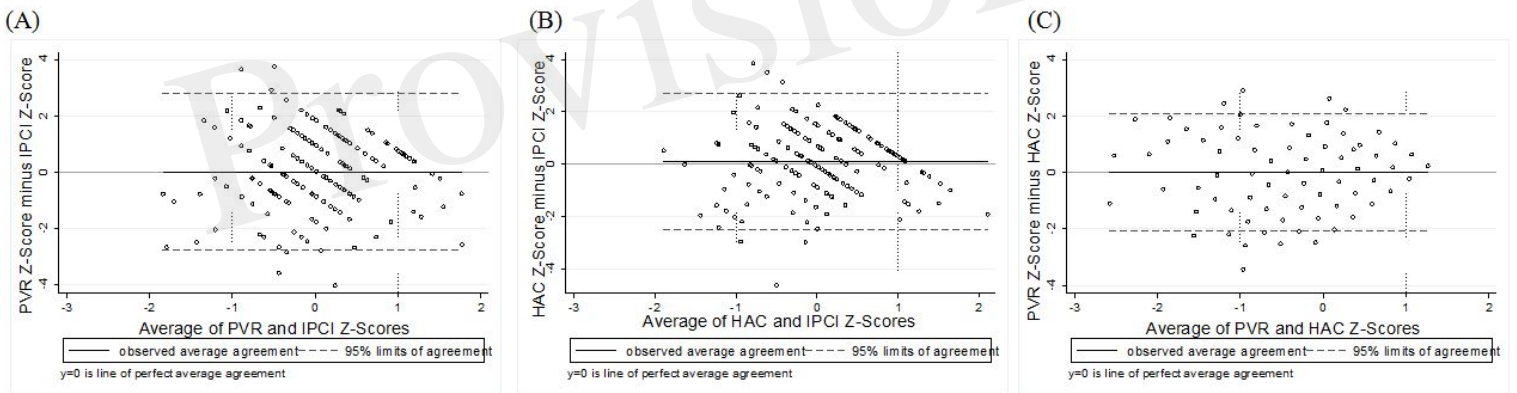


Figure 3. Bland-Altman Plots: (A) Parental Verbal Responsivity and Indicator of Parent-Child Interactions [95% limits of agreement: -2.78 to 2.80]; (B) Home Activities with Child and Indicator of Parent-Child Interactions [95% limits of agreement: -2.52 to 2.70]; (C) Parental Verbal Responsivity and Home Activities with Child [95% limits of agreement: - 2.06 to 2.06].

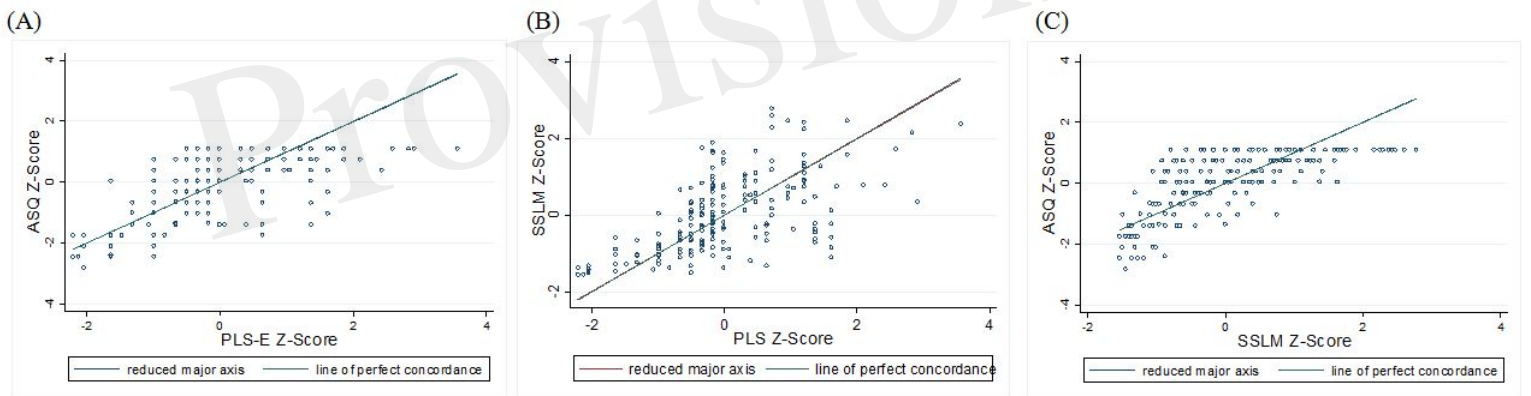


Figure 4. RMA plots: (A) Ages & Stages Questionnaire and Preschool Language Scale, Slope = 0.996; Intercept = -0.009; (B) Sure Start Language Measure and Preschool Language Scale, Slope = 0.990; Intercept = -0.001; (C) Ages & Stages Questionnaire and Sure Start Language Measure, Slope = 1.005; Intercept = 0.002.

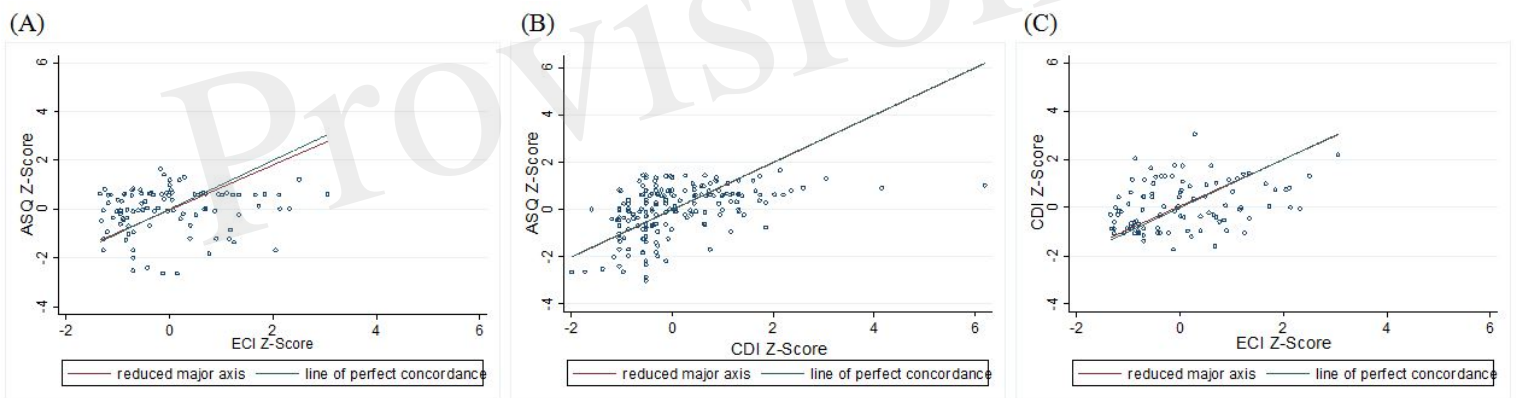


Figure 5. RMA Plots: (A) Ages & Stages Questionnaire and Early Communication Indicator, Slope = 0.924; Intercept = -0.036; (B) Ages & Stages Questionnaire and Communicative Development Inventory, Slope = 1.000; Intercept = -0.020; (C) Communicative Development Inventory and Early Communication Indicator, Slope = 0.967; Intercept = 0.064.

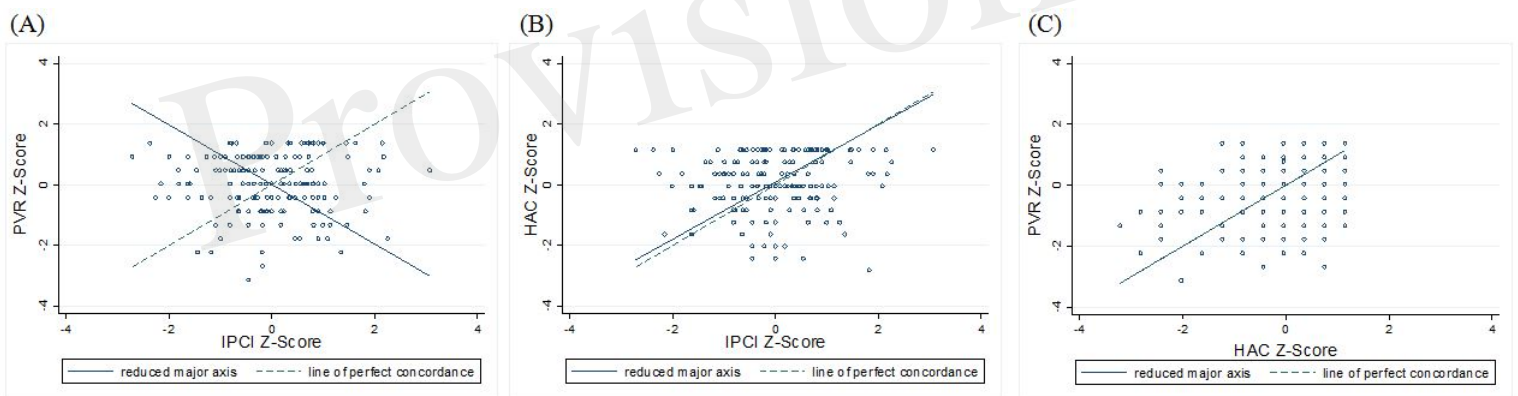


Figure 6. RMA Plots: (A) Parental Verbal Responsivity and Indicator of Parent-Child Interactions, Slope = -0.98; Intercept = 0.012; (B) Home Activities with Child and Indicator of Parent-Child Interactions, Slope = 0.94; Intercept = 0.09; (C) Parental Verbal Responsivity and Home Activities with Child, Slope = 1.00; Intercept = 0.00