1    Comparing different measures of bilingual input derived from naturalistic daylong recordings

2    Yufang Ruan*[1,2], Adriel John Orena[3], Linda Polka[1,2]

3    [1] School of Communication Sciences and Disorders, Faculty of Medicine and Health Sciences,

4    McGill University, Montreal, Canada

5    [2] Centre for Research on Brain, Language and Music, Montreal, Canada

6    [3] Department of Psychology, University of British Columbia, British Columbia, Canada

7

8    **Correspondence Author:** Yufang Ruan, 2001 McGill College Avenue 8[th] floor, Montréal,

9    Québec, H3A 1G1, Canada, yufang.ruan@mail.mcgill.ca.

17

18                                              **Abstract**

19   **Purpose:** Measuring language input, especially for infants growing up in bilingual

20   environments, is challenging. Although the ways to measure input have expanded rapidly in

21   recent years, there are many unresolved issues. In the current study, we compared different

22   measurement units and sampling methods used to estimate bilingual input in naturalistic daylong

23   recordings.

24   **Method:** We used the Language Environment Analysis (LENA) system to obtain and process

25   naturalistic daylong recordings from 21 French-English bilingual families with an infant at 10

26   and 18 months of age. We examined global and context-specific input estimates and their relation

27   with infant vocal activeness (i.e., volubility) when input was indexed by different units (Adult

28   Word Counts, speech duration, 30-second segment counts) and using different sampling methods

29   (every-other-segment, top-segment).

30   **Results:** Input measures indexed by different units were strongly and positively correlated with

31   each other and yielded similar results regarding their relation with infant volubility. As for

32   sampling methods, sampling every other 30-second segment was representative of the entire

33   corpus. However, sampling the top segments with the densest input was less representative and

34   yielded different results regarding their relation with infant volubility.

35   **Conclusions:** How well the input that a child receives throughout a day is portrayed by a

36   selected sample and correlates with the child's vocal activeness depends on the choice of input

37   units and sampling methods. Different input units appear to generate consistent results, while

38   caution should be taken when choosing sampling methods.

39

40      Comparing different measures of bilingual input derived from naturalistic daylong recordings

41          Measuring language input, especially for infants growing up in bilingual environments, is

42      challenging. Some researchers have used diaries and questionnaires completed by caregivers to

43      estimate the proportion of each language in a child's input (e.g., Carbajal & Peperkamp, 2020;

44      Place & Hoff, 2011, 2016). Other researchers have documented children's real-world input using

45      audio- or video-recordings. In recent years, a growing number of researchers have adopted the

46      Language Environment Analysis system (LENA, LENA Research Foundation, Boulder, CO) to

47      obtain and process daylong audio recordings of language input in bilingual households (e.g.,

48      Orena et al., 2020; Marchman et al., 2017; Ramírez-Esparza et al., 2017a, 2017b; VanDam et al.,

49      2016). The LENA system includes a recorder that children wear in a vest and algorithms that

50      automatically process and estimate language input by adult word counts (AWCs) or speech

51      duration. Researchers have tested LENA's accuracy in different languages and have found most

52      input estimates to be fairly reliable (Cristia et al., 2020, 2021; Orena et al., 2019).

53          Although the ways to measure input have expanded rapidly, there are still many

54      unresolved issues. First, various units have been used to measure bilingual input. Some

55      researchers, using diary or recording methods, have divided time in a day into equal-sized

56      segments and measured segment counts or durations where each language was used (e.g., Place

57      & Hoff, 2011, 2016; Ramírez-Esparza et al., 2017a, 2017b). This method of *counting segments* is

58      widely used for its simplicity and efficiency. However, bilingual caregivers do not always use the

59      same language within a given segment, especially in a relatively longer segment (e.g., 30

60      minutes). Thus, in some cases, researchers have asked caregivers to estimate the proportion of

61      time that each language was used within a segment (Carbajal & Peperkamp, 2020). This

62      approach improves the accuracy of bilingual exposure estimation, but still overlooks the fact that

63    caregivers might not continuously speak for the length of a segment. This limitation can be

64    addressed by using more fine-grained units, such as *speech duration* or *AWCs* extracted from

65    recordings via speech processing algorithms (e.g., Marchman et al., 2017; Ruan, Orena, & Polka,

66    2020; Ruan, Orena, Xu, et al., 2020). However, the accuracy of these algorithms is imperfect

67    (Cristia et al., 2020, 2021; Lehet et al., 2021), and these fine-grained units are not always

68    available (e.g., when using diaries and questionnaires). Therefore, it is important to examine

69    whether using different units impacts input estimation. Until we do so, it is difficult to compare

70    results across studies where input was indexed by different units.

71           Another unresolved issue is that algorithms are so far unable to automatically and reliably

72    classify bilingual input into two languages; this task requires manual annotation which is a

73    laborious endeavor. This challenge underscores the need for an effective and reliable sampling

74    method which allows researchers and clinicians to achieve their goals by processing only a

75    portion of their data. Orena and colleagues (2020) examined daylong recordings in the Montréal

76    Bilingual Infants corpus and found that the proportion of bilingual input in each language

77    estimated from a periodically-selected (e.g., *every-other-30s-segment sampling*, EOS) sample

78    was well-correlated with parental estimations. Taking the bilingual input proportions estimated

79    from an EOS sample as the gold standard, another study further showed that a smaller randomly-

80    selected sample (7% of the total recording or 11% of speech clips) yielded an estimation close to

81    the gold standard (Cychosz et al., 2021). This study suggests that a reliable estimation of the

82    bilingual input distribution can be achieved by annotating a modest amount of recording. Others

83    have tried to achieve the same goal using more selective sampling methods. For example, in one

84    study, researchers focused on the most input-dense portions of the recordings by selecting 40

85    temporally-scattered segments *with the highest AWCs* each day (i.e., the top segments), from

86     which they composed a sample of around 160 top segments across four days for each child

87     (Ramírez-Esparza et al., 2017b). However, there are concerns around this *top sampling method*:

88     the features of the input in the top segments may differ from what a child experiences throughout

89     a day (Bergelson et al., 2019; Tamis-LeMonda et al., 2017). It is also unknown whether the

90     distribution of a child's bilingual exposure remains the same in the top segments. Taken together,

91     examining the representativeness of different samples is a crucial step towards identifying

92     effective and reliable sampling methods.

93          Representativeness aside, researchers and clinicians may sample the top segments for

94     different objectives. For example, children growing up in bilingual families are simultaneously

95     exposed to two languages. Many bilingual children have more opportunities to receive input in

96     one language (dominant) than the other (non-dominant). Thus, it is informative to control for the

97     variance in the opportunity that a child has with each language when examining how language

98     dominancy shapes the relation between the input and child vocal behaviours. Selecting the same

99     number of segments and selecting the ones with the densest input for each language (e.g., top

100    sampling of input in a specific language) is a way to obtain two comparable samples with the

101    same duration and optimal density for language comparisons (Xu et al., 2019). The same

102    sampling approach can also be used to compare input received in different social contexts.

103    Because the implementation of this top sampling method requires annotating the input (i.e., in

104    what language and what social context) prior to sampling, we call it post-annotation top

105    sampling.

106         In our ongoing project, we are interested in using LENA recordings in the Montréal

107    Bilingual Infants corpus to investigate the relation between infants' vocal activeness (i.e.,

108    volubility) and their French-English bilingual input in different social (overhearing, one-on-one)

109    and language (dominant, non-dominant) contexts when infants were 10 and 18 months old.

110    Following Orena et al. (2019), we plan to estimate input using the AWCs from a sample of every

111    other 30-second segment containing adult speech across the entire corpus. Prior to that, the

112    reliability of this method was examined in the current study. Specifically, we examined whether

113    input measures indexed by AWC in the EOS sample were consistent with measures indexed

114    using other units (speech duration, segment counts) and sampling methods (the entire corpus,

115    simple top sampling, and post-annotation top sampling; see a summary of variables in Table 1).

116    We also examined whether the distribution of different social and language input in the top

117    segments resembled the distribution in the EOS sample. Lastly, we compared the relation

118    between infant volubility and input when the input was estimated using different units and

119    sampling methods.

120         Findings from this study can help aggregate results from studies using different input

121    units and/or sampling methods as well as provide methodological guidance for future studies

122    using daylong recordings. The findings can potentially be applied to all research regardless of

123    participants' language background (monolingual, bilingual, etc.), while some will be particularly

124    relevant to bilingualism research. For example, counting segments in each language is widely

125    used to determine the bilingual input distribution, thus knowing the reliability of this method will

126    have a wide implication. Moreover, the post-annotation sampling method derives comparable

127    samples with equal duration and maximal density, which enable us to compare input in two

128    languages while controlling for the inherent difference in quantity to some extent.

129                                          **Methods**

130      **Participants**

131            We analyzed data from the Montréal Bilingual Infants corpus (Orena et al., 2020).

132      Twenty-one families participated when the infant was 10 months old (13 males, 8 females; Age

133      *Mean* = 303 days, *Range* = 289 – 319 days) and 16 of them participated again when the child

134      was 18 months old (10 males, 6 females; Age *Mean* = 576 days, *Range* = 551 – 635 days). All

135      caregivers had knowledge of both French and English and most of them (27 out of 42) reported

136      speaking both languages to their child. According to parental estimates, their child was exposed

137      to each language for at least 20% of time. Four families reported a small amount of exposure to a

138      third language (< 5%). At 10 months, 12 infants were raised in a French-dominant language

139      environment and nine were English-dominant. At 18 months, eight were French-dominant and

140      eight were English-dominant. Parents provided consent to participate and declared no auditory

141      and neurocognitive disorders for their child.

142      **Procedure and Measures**

143            Measures used in this study are summarized in Table 1. Naturalistic audio recordings

144      were collected using a LENA digital language processor (DLP). Infants wore the DLP in a vest

145      for 16 hours per day. Three full-day recordings (2 weekdays and 1 weekend day) were made

146      when infants were 10 months old. For 16 families, a fourth recording was completed on a

147      weekend day when infants were 18 months old. In total, the families contributed 1,264 hours of

148      audio recordings ([21 families at 10 months × 3 days × 16 hours] + [16 families at 18 months ×

149      1 day × 16 hours]). Recordings were divided into 30-second segments. Estimates of child

150      vocalizations and language input were derived for each segment using the LENA algorithms. The

151      Child Vocalization Count (CVC) is the number of vocalizations produced by the key child. A

152    child vocalization is defined as a speech/speech-like sound produced by the key infant that is

153    preceded and followed by 300 milliseconds of silence or nonspeech. We summed CVCs across

154    the entire corpus for each child at each age to index *infant volubility*. Language input was

155    measured using different units and sampling methods. How each sampling method was

156    conducted while using different units is described in Figure 1.

157    ***Input Units (AWC, Duration, Segment Count)***

158            LENA algorithms estimate the number of words spoken near the key child (Adult Word

159    Counts, AWC). Previous research showed that LENA algorithms were reliable at estimating

160    AWCs in both English and French (Orena et al., 2019). Algorithms also estimate the duration of

161    these words and derive Adult Female Speech Duration and Adult Male Speech Duration. For

162    each infant, the sum of Adult Female and Male Speech Duration provided an approximation of

163    speech duration (Duration). The Segment Count referred to the number of 30-second segments.

164    ***LENA Sample***

165            The LENA sample consisted of all the recordings in the corpus for each age. There were

166    21 families × 3 days × 1920 segments per day = 120,960 segments in the 10-month LENA

167    sample and 16 families × 1 day × 1920 segments per day = 30,720 segments in the 18-month

168    LENA sample. One segment in the 10-month sample was excluded because of an evident

169    technical error (AWCs > 3000 in a 30-second segment, accounting for less than 0.5% of total

170    AWC). As the Segment Count was identical for all infants, we utilized AWC and Duration to

171    measure input in the LENA sample.

172    ***Every-Other-Segment (EOS) Sample***

173            As we were interested in caregivers' input, we first removed segments in the LENA

174    sample that did not contain any adult speech. From the remaining segments containing adult

175   speech, we selected *every other segment*. In total, 18,979 and 6,180 segments were included in

176   the 10-month and 18-month EOS samples respectively.

177        Segments in the EOS sample were manually annotated. Trained English-French bilingual

178   research assistants listened to each segment and coded for social contexts (i.e., how many

179   speakers and listeners, who was speaking to whom) and language contexts (i.e., what language

180   was being spoken). Seven research assistants completed this work after each of them

181   successfully completed a training file. Inter-coder reliability in the training file was high (on

182   average 94.2% agreement for speaker context and 92.4% agreement for language context, Orena

183   et al., 2020). Speech in which one caregiver spoke directly to the infant was tagged as one-on-

184   one input. Overheard input was tagged for speech spoken in the presence of the infant, but not

185   directly addressing the infant. Speech in which two or more caregivers spoke directly to the

186   infant was not included as another level of social contexts because this was rarely observed in

187   our corpus. Input tagged as "English" or "French" were recoded as "dominant" or "non-

188   dominant" with the dominance assigned according to the *parent-reported* relative exposure to

189   each language for each child at each age. Mixed-language input was not included as another

190   level of language contexts because it accounted for less than 10% of the total input on average at

191   each age.

192        For each child at each age, we summed the total input in the EOS sample (global) and

193   computed input measures by social contexts (one-on-one, overhearing) and language contexts

194   (dominant, non-dominant). As the number of segments containing adult speech varied across

195   infants, there was a considerable variation in the Segment Count in the EOS sample. Thus, we

196   utilized all three units (Segment Count, AWC, and Duration) to measure input in the EOS

197   sample.

198   ***Top150 Sample***

199        Following the work of Ramírez-Esparza and colleagues (Ramírez-Esparza et al., 2017a,

200   2017b), we selected the top 50 segments with the highest AWCs each day across three days in

201   the 10-month EOS sample for a total of 150 segments per child. For 18 months, despite having

202   only one daylong recording, we sampled the top 150 segments with the highest AWCs for each

203   child. This allowed us to examine whether the size of top samples relative to the original sample

204   affects input estimation.

205        Because the Top150 sample was selected from the EOS sample, social and language

206   context annotation was also accessible for the Top150 sample. Again, for each child at each age,

207   we summed the total input in the Top150 sample (global) and indexed it in AWC and Duration

208   (Segment Count was identical for all infants, n = 150). We also summed input by social and

209   language contexts and used all three measurement units (Segment Count, AWC, and Duration) to

210   index the input in each context.

211   ***Top40/20 Samples***

212        Segments in the EOS sample were initially categorized by social (one-on-one or

213   overhearing) and language (dominant or non-dominant) contexts according to the manual

214   annotation. For the two social contexts, top 40 segments with the highest AWCs in each context

215   were sampled for each child at each age. We chose 40 segments because most infants had at least

216   40 segments for each social context except one child at 18 months (one-on-one context analysis

217   was based on 23 segments for this child). For the two language contexts, top 20 segments with

218   the highest AWCs in each language were sampled for each child at each age. Again, we chose 20

219   segments because most infants had at least 20 segments for each language context with only a

220   few exceptions (non-dominant language analysis was based on less than 20 segment for one

221 child at 10 months (Segment Count = 7) and two children at 18 months (Segment Count = 4 and

222 9)). Together, at each age, we had one Top40 sample for each of the two social contexts (one-on-

223 one and overhearing) and one Top20 sample for each of the two language contexts (dominant

224 and non-dominant, see Figure 1). In total, we had four Top40/20 samples for 10 and 18 months

225 respectively. By definition, Segment Count was identical for each context; thus, input was

226 indexed only by AWC and Duration in the Top40/20 samples.

227 **Statistical Analysis**

228      Results and plots were generated using packages including languageR (Baayen &

229 Shafaei-Bajestan, 2019) and ggplot2 (Wickham, 2016) in R (R Core Team, 2021). The data and

230 code that support the findings of this study are available at https://osf.io/uqh35/.

231      To examine whether using different units provide similar input estimates, we correlated

232 input measures in different units. Next, to investigate whether each sampling method generates a

233 representative sample, we correlated input measures derived from a selected sample to the

234 measures derived from its original sample. Spearman's correlations were used because the input

235 distribution deviated from a normal distribution. Significance of these correlations was not tested

236 because we were interested in the degree of these correlations (i.e., the magnitude).

237      Next, to examine whether the proportions of input in different social and language

238 contexts remain the same in top segments as the ones observed throughout a day, we computed

239 these proportions in the Top150 and EOS samples. We used AWCs in a specific context divided

240 by the total AWCs in that sample. Then, for each context, we compared proportions estimated in

241 the two samples using Wilcoxon signed-rank tests. All $p$-values were adjusted using method of

242 Benjamini & Hochberg (1995).

243      Lastly, to test whether the relation between infant volubility and input changes depending

244      on how input was estimated, we compared Spearman's correlations between infant volubility and

245      language input when input was estimated using different units and sampling methods. We

246      repeated the analysis in different social and language contexts at each age. The original *p*-values

247      were reported because (1) The purpose of this set of analyses was not to test the hypothesis that

248      infant volubility was related to language input, but to examine the consistency across input units

249      and sampling methods; (2) We tried to mimic the reality where researchers and clinicians would

250      only select one measure of input and there would not be any *p*-value adjustment at the level of

251      input measurement.

<div align="center">

**Results**

</div>

**Does using different units and sampling methods provide similar estimations of language**

**input from daylong recordings?**

255      Spearman's correlations between different input measures are plotted in Figure 2, for

256      global input (a) as well as input in each social (b & c) and language (d & e) context. Results for

257      the 10-month dataset are plotted in the upper triangle and results for the 18-month dataset, in the

258      bottom triangle. The Spearman's correlation coefficient between each pair of input measures is

259      reported in each cell and the cell colour indicates the strength of the correlation, from weak

260      (yellow) to strong (red). Conventionally, a value of .80 or greater indicates a good consistency

261      across measures (Chiang et al., 2020, p. 98). A video-animated guide of Figure 2 is available in

262      the *Supplementary Material*.

263      First, we compared across different units (AWC, Duration, and Segment Count). Within

264      each sample, we correlated input indexed by different units. A stronger positive correlation

265      indicated a higher consistency between two units. We expected the correlation between AWC and

266    Duration to be positive and strong, while the correlation of each with Segment Count to be less

267    strong because Segment Count is a less fine-grained unit and less dependent on speech

268    processing algorithms compared to AWC and Duration. As expected, we observed strong

269    correlations across three units in all contexts, samples, and ages (shown in the cells close to the

270    diagonal line in Figure 2 and in Part 1 of the animation). For example, in Figure 2a, the cell

271    corresponding to the first column from the *left* (Column 1 or C1, LENA_AWC) and the second

272    row from the *bottom* (Row 2 or R2, LENA_Dur) shows the correlation between the global input

273    estimated in the entire 10-month corpus by AWC and speech duration, which is .99 (> .80)

274    suggesting a good consistency between these two unites. Correlations involving Segment Count

275    were relatively smaller, especially for overhearing context in Top150 samples (Figure 2b, [C4-5,

276    R6] and [C6, R4-5]).

277          Next, we compared across different sampling methods (LENA, EOS, Top150, and

278    Top40/20). The EOS sample was drawn from the LENA sample (i.e., the entire corpus) and the

279    annotation of context-specific input was not available for the LENA sample, thus we examined

280    the correlation between the global input estimated in the EOS and LENA samples (see Figure 2a,

281    [C1-2, R3-5] and [C3-5, R1-2] and Part 2 of the animation). We expected these correlations to be

282    positive and strong. Indeed, the correlation coefficients were beyond .80 (a few below this

283    threshold involving Segment Count). These results indicated that the EOS sample selected by

284    every-other-segment sampling was representative of the entire corpus.

285          The Top150 sample was drawn from the EOS sample, thus we examined the correlation

286    between the input estimated in these two samples. We expected weaker correlations here because

287    top sampling provides a narrow snapshot of the child's language exposure throughout a day

288    (Bergelson et al., 2019). Indeed, our results showed that compared to the correlations between

289    the EOS and LENA samples reported above, the correlations between the global input derived in

290    the Top150 and EOS samples were slightly smaller (Figure 2a, [C3-5, R6-7] and [C6-7, R3-5],

291    see also Part 3 of the animation). Although correlations were close to .80 for input in both

292    language contexts (Figure 2d & e, [C1-3, R4-6] and [C4-6, R1-3]) and in the 18-month dataset

293    (bottom triangles), correlations were below this threshold when input was indexed by Segment

294    Count and for both types of social input in the 10-month dataset (Figure 2b & c, [C1-3, R4-6]).

295    For instance, in 10-month dataset, the correlations between the overheard input estimated in the

296    EOS sample and in the Top150 sample ranged from .34 to .75 (Figure 2b, [C1-3, R4-6]),

297    evidently smaller than the threshold (.80), as well as smaller than their corresponding

298    correlations in the 18-month dataset that ranged from .56 to .98 (Fugure2b, [C4-6, R1-3]). These

299    results suggested that the most input-dense portions of the recordings might be less

300    representative of the daylong recordings.

301          The Top40/20 samples were drawn from the EOS sample with the goal of equating the

302    number and input density of the segments used to examine the variation across different types of

303    input. We did not have a clear hypothesis for the correlation between each type of social and

304    language input in Top40/20 samples and in the EOS sample as this analysis was exploratory.

305    Except for the dominant language input in 18-month dataset and measures involving Segment

306    Count, our results suggested a good representativeness of the Top40/20 samples by showing

307    correlations close to .80 (Figure 2 b-e, [C1-3, R7-8] and [C7-8, R1-3], see also Part 4 of the

308    animation). These results indicated that when we attempt to control for inherent differences in

309    infants' opportunity to receive each type of input, we still observe a similar pattern of individual

310    differences. Therefore, this post-annotation top sampling method can potentially provide a

311    representative sample of specific types of input.

312    **Do estimated proportions of input in different social and language contexts differ when**

313    **different sampling methods are used?**

314          Due to the discrepancies observed between context-specific input in the EOS and Top150

315    samples, we further compared the proportional estimates of social and language input across

316    these two samples. Given that a previous study found different input patterns in peak-hour versus

317    daylong samples (Bergelson et al., 2019), we expected the input proportions estimated in the

318    Top150 sample to differ from the ones in the EOS sample.

319          As shown in Table 2, differences were observed between the two samples with median

320    ranging from 1 to 7%, and they reached significance for overheard, one-on-one, and dominant

321    language input in 10-month samples (original $ps < .05$). However, none of these $p$-values was

322    significant after correcting for multiple comparisons. The results indicated that the input

323    distribution across different social and language contexts in the top segments differed from the

324    distribution observed throughout a day, but not substantially.

325    **Does the relation between input and infant volubility change when the input is estimated**

326    **using different input units or sampling methods?**

327          We compared the correlation between infant volubility (derived from the entire corpus)

328    and the input when the input was estimated by different units and sampling methods. When

329    comparing across different units (AWC, Duration, Segment Count), similar correlations indicated

330    consistency. We found that within the EOS sample, input-volubility correlations were generally

331    same in direction and significance, and similar in magnitude across different units (see Table 3

332    columns 2-4). For example, the correlation between volubility and overheard input at 10 months

333    was .49, .49, and .52 when input was indexed by Segment Count, AWC, and Duration

334    respectively. These correlations were uniformly positive, significant at 0.05 level, and

335     numerically close to each other. These results suggested that using different units to measure

336     input led to similar conclusions regarding the relation between input and infant volubility.

337         When comparing across sampling methods (EOS, Top150), we viewed the input-

338     volubility correlations where the input was estimated in the EOS sample as the gold standard.

339     Thus, deviations from this gold standard suggested potential problems with the top sampling

340     method. Based on the results from our previous research questions, we expected a discrepancy

341     between the correlations for the EOS and Top150 samples. Indeed, compared to the EOS

342     correlations (Table 3, columns 2-4), the Top150 correlations (columns 5-7) were consistently

343     smaller and sometimes in the opposite direction (e.g., input in overhearing contexts at both ages

344     and non-dominant language contexts at 18 months, all indexed by Segment Count). For example,

345     compared to the correlation between volubility and overheard input estimated using AWC and

346     EOS sampling at 10 months (*rho* = .49), the corresponding correlation was much smaller when

347     input was estimated in the Top150 sample (*rho* = .28). Therefore, using a simple top sampling

348     method to estimate input might lead to a different conclusion regarding the relation between

349     input and infant volubility.

350         Additionally, we examined the correlation between infant volubility and the context-

351     specific input estimated in Top40/20 samples. This post-annotation sampling method was used to

352     achieve a different goal, namely, to assess whether the input-volubility correlation would change

353     when comparable samples were used for each type of input. Therefore, if Top40/20 correlations

354     deviates from the EOS gold standard, it would not suggest problem with this post-annotation

355     sampling method but reveal that the quality of the input (language dominancy or social

356     interaction) played a role independent from any inherent differences in the opportunity that a

357     child has with a specific type of input. This comparison was exploratory hence we did not have a

358   clear expectation. As shown in the last two columns in Table 3, the Top40/20 correlation

359   coefficients were generally smaller than the corresponding EOS correlations (hence less likely to

360   be significant). Meanwhile, the Top40/20 correlation coefficients were numerically closer to the

361   EOS ones. For example, although the correlation between volubility and overheard input at 10

362   months when input was estimated using AWC and Top40 sampling was smaller and still not

363   significant (*rho* = .39), this value was closer to the one observed for the EOS sampling (*rho*

364   = .49). The pattern across two social contexts (overhearing versus one-on-one) as well as across

365   two language contexts (dominant versus non-dominant) also generally reassembled the pattern

366   observed in the EOS sample. These results advanced our findings from the first research question

367   by showing that in addition to input's variation, its covariation with infant volubility also

368   persisted in the Top40/20 samples.

369                                                   **Discussion**

370          In summary, our analyses yielded the following findings: (1) Input measures indexed by

371   different units (AWC, Duration, and Segment Count) were positively and strongly correlated

372   with each other and generated similar results regarding their relation with infant volubility; (2)

373   Input estimates derived using the every-other-segment sampling method were representative of

374   input in the entire corpus; (3) Sampling the top segments with the densest input might derive a

375   less representative sample for estimating input and its relation with infant volubility; and (4)

376   Context-specific input's variation and its covariation with infant volubility persisted in segments

377   selected by a post-annotation top sampling method.

378          Measures of language input using different units (AWC, Duration, and Segment Count)

379   and their relation with infant volubility were highly consistent. Hence, we validated the method

380   of using AWC to estimate the input in LENA daylong recording for our ongoing project.

381    Meanwhile, correlations involving Segment Count were slightly smaller and this deviation was

382    amplified when we compared units across samples. For instance, in the 10-month dataset, the

383    correlation between the Segment Count of overheard input in the Top150 sample and the AWC

384    of the same type of input in the EOS sample was only .34, much smaller than .80 (Figure 2b,

385    [C1, R6]). In addition, the only correlations that showed a negative relation between input and

386    infant volubility, were based on input measures indexed by Segment Count (Table 3, column 5).

387    These findings have important implications on how we assess bilingual exposure, given that

388    counting segments or segment duration is a common practice in previous bilingualism research

389    (e.g., Place & Hoff, 2011, 2016; Ramírez-Esparza et al., 2017a, 2017b). When counting

390    segments, we lose information such as how verbally active the speaker is and whether the

391    speaker consistently uses the same language for the entire segment. Some researchers have tried

392    to address the latter by asking caregivers to estimate the time that each language was used within

393    a segment (Carbajal & Peperkamp, 2020). In future studies, researchers could also estimate the

394    time that caregivers are actively speaking within a segment to quantify the input more precisely,

395    when fine-grained units like AWC and speech duration are not available. On the other hand, one

396    of the advantages of using Segment Count is that counting segments relies less on speech

397    processing algorithms, which spares it from concerns regarding the accuracy of these algorithms

398    (Cristia et al., 2020, 2021; Lehet et al., 2021).

399         Our results also showed that sampling every-other-segment achieved a good

400    representativeness of the entire corpus, which replicated our previous pilot study (Orena et al.,

401    2019). Meanwhile, a sample of the top segments with the densest input was less representative of

402    a child's language exposure throughout a day, shown by the correlations between the Top150 and

403    EOS samples (Figure 2a – e). In addition, the correlations between infant volubility and the input

404     estimated in the Top150 sample diverged from the ones where the input was estimated in the

405     EOS sample (Table 3). This deviation may arise for two reasons. One is biased sampling. The

406     Top150 sample consisted of segments containing the highest AWCs, essentially the moments

407     when caregivers were the most verbally active around the child (talking to the child or others).

408     The distribution of different types of input might differ in top segments. Indeed, we observed that

409     the proportion of overheard and one-on-one input differed between the Top150 and EOS samples

410     in the 10-month dataset (although no longer significant after $p$-adjustment, Table 2). These

411     differences corresponded to the weaker correlations observed in Figure 2 (b) and (c). Other

412     aspects of the input in the top segments might also differ from infants' language experience

413     throughout a typical day, as suggested by previous research where researchers found a denser

414     usage of nouns in peak-hour recordings (Bergelson et al., 2019). These differences observed for

415     the input in the top segments might help explain the deviant relation between infant volubility

416     and input found when the input was estimated using the Top150 sampling method (Table 3). If it

417     is true, periodic or random sampling which selects input without reference to input features,

418     might yield a less biased sample (Cychosz et al., 2021; Orena et al., 2019).

419              The other possible reason is related to the size of the Top150 sample relative to its

420     original sample (i.e., the EOS sample). Recall that albeit 10-month dataset (3-day recordings)

421     being larger than 18-month dataset (1-day recording), we selected the same number of top

422     segments (n = 150) for each child from these two datasets. Therefore, the Top150 sample

423     accounted for a smaller proportion of the EOS sample for 10-month dataset (17%) than for the

424     18-month dataset (37%). This might explain why we observed relatively weaker correlations

425     between input measures estimated in the Top150 and EOS sample in the 10-month dataset

426     (Figure 2 upper triangles) than the 18-month dataset (Figure 2 bottom triangles). Therefore,

427    sampling a fixed number of top segments might be disadvantageous for larger samples: For a

428    given number of top segments, the larger the original sample is, the smaller proportion of

429    segments are selected, and thus less likely to be representative. To tackle this problem with larger

430    samples, it might be helpful to sample a fixed proportion, instead of a fixed number, of segments,

431    so that the size of the selected sample would changes with the size of the original sample. Future

432    studies should consider finding an optimal proportion for selecting a representative sample from

433    daylong recordings with the least segments. For example, previous research suggested that

434    around 7% randomly-selected segments from overall recordings was representative in terms of

435    the proportion of bilingual exposure and child-directed speech observed in the EOS sample

436    (Cychosz et al., 2021). However, the correlation between these estimations based on 7% of data

437    and parental reports were still not optimal. We should also keep in mind that no matter how

438    many segments are sampled from the recordings, samples only provide a snapshot of a child's

439    everyday life hence carry some extent of sampling bias.

440            We also examined a post-annotation top sampling for a different purpose which was to

441    assess whether the pattern of different types of social or language input and their relation with

442    infant volubility would change when comparable samples with equal duration and maximal

443    density (Top40/20 samples) were used for different types of input. Our results showed that when

444    being provided with the same opportunity and in its optimal condition, each type of social and

445    language input seems to preserve its variation and its covariation with infant volubility as

446    observed throughout a day. These results also suggested that this post-annotation top sampling

447    might be used for other purposes, such as more detailed annotations. For example, for

448    researchers and clinicians who are interested in comparing child-caregiver interaction when

449    caregivers speak one language or the other, they might consider to initially annotate language

450   contexts (i.e., which language(s) was used) for every other segment containing adult speech, and

451   then select a fix proportion of top segments from each language context to conduct detailed

452   annotations on child-caregiver interactions.

453          Infant age might also contribute to the different patterns observed in the correlation

454   between infant volubility and input received in different contexts. Although we presume that the

455   differences can be primarily attributed to the fact that language input and infant volubility

456   increase and the relation between them changes as infants grow, we cannot rule out potential

457   impacts from our methodological choices. For example, the recordings were collected on a

458   weekend day at 18 months while they were collected on three days including both weekdays and

459   weekends at 10 months. Many factors including children's daily routine and primary caregiver(s)

460   may vary across weekday and weekend. Additionally, when children were at 18 months, some of

461   them went to a daycare but language input at the daycare was not captured in our recordings.

462   Activities at settings outside home may alter children's language exposure (Larson et al., 2020;

463   Soderstrom et al., 2018), which might potentially alter our results.

464          There are limitations to the current study. First, although previous reviews and

465   evaluations have suggested the LENA-derived measures, especially Child Vocalization Counts

466   (CVCs, indexed infant volubility) and AWCs, to be reasonably accurate (Cristia et al., 2020,

467   2021), there are still some gaps between the accuracy of manual annotation and automatic

468   processing algorithms (Lehet et al., 2021). Second, our sample size is relatively small ($N = 21$

469   and 16 for 10- and 18-month respectively) due to the laborious work involved in manual

470   annotation, but the corpus consists of 1,264 hours of daylong recordings. Third, the families

471   contributed to the corpus lived in a French-English bilingual community in which both languages

472   have high social status. This characteristic may restrict the generalization of findings to bilingual

473     communities where language status is uneven. For example, we suspect that the top sampling

474     methods may yield less-representative results in bilingual contexts where languages have

475     different levels of social status and/or involve more complex patterns of language use.

476          There are several directions for future studies. First, future studies should try to replicate

477     the findings of this study using a larger sample, in different bilingual contexts, and recordings

478     made outside of the home. Second, when Ramírez-Esparza and colleague composed their sample

479     using the simple top sampling method, the authors made the effort to ensure selected segments

480     were 3-minute apart (Ramírez-Esparza et al., 2017a, 2017b). Whether this effort would improve

481     the representativeness of top sampling remains to be tested. Third, although manually coding

482     adult speech in every other segment reduced the work needed to code the full corpus by half, it

483     was still laborious. Future studies could investigate the reliability of other periodic, but less

484     dense sampling methods, such as sampling 1 minute every hour (Scaff et al., 2022).

485          In conclusion, while the methods to estimate children's language input have been

486     expanding rapidly in recent years, it is important to know that our research conclusions are not

487     built on methodological biases. Our results suggested a high consistency across different units

488     (AWC, speech duration, segment count). However, caution should be taken when choosing

489     sampling methods. While sampling every other 30-second segment might generate a unbiased

490     sample, there is more work needed to be done to improve the representativeness of top sampling

491     methods. That said, top sampling methods can still be used for different research purposes. Taken

492     together, findings from this study highlight the need for our field to direct more attention to the

493     exact measures used to estimate language input and to be thoughtful when selecting sampling

494     methods.

495     **Acknowledgement**

506    **Ethics Approval Statement**

507    We received ethics approval from the Institutional Review Board at McGill University (IRB #

508    A05-B20-16A).

509    **Data Availability Statement**

510    The data and code that support the findings of this study are available at https://osf.io/uqh35/.

511    **Supplemental Material**: https://doi.org/10.23641/asha.22335688

512                               References

513    Baayen, R. H., & Shafaei-Bajestan, E. (2019). *languageR: Analyzing Linguistic Data: A*

514         *Practical Introduction to Statistics* (1.5.0).

515         https://cran.rstudio.com/web/packages/languageR/index.html

516    Benjamini, Y., & Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and

517         Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series B*

518         *(Methodological)*, *57*(1), 289–300. https://doi.org/10.1111/j.2517-6161.1995.tb02031.x

519    Bergelson, E., Amatuni, A., Dailey, S., Koorathota, S., & Tor, S. (2019). Day by day, hour by

520         hour: Naturalistic language input to infants. *Developmental Science*, *22*(1), e12715.

521         https://doi.org/10.1111/desc.12715

522    Carbajal, M. J., & Peperkamp, S. (2020). Dual language input and the impact of language

523         separation on early lexical development. *Infancy*, *25*(1), 22–45.

524         https://doi.org/10.1111/infa.12315

525    Chiang, I.-C. A., Jhangiani, R. S., & Price, P. C. (2020). *Research Methods in Psychology – 2nd*

526         *Canadian Edition*. BCcampus. https://opentextbc.ca/researchmethods/

527    Cristia, A., Bulgarelli, F., & Bergelson, E. (2020). Accuracy of the Language Environment

528         Analysis System Segmentation and Metrics: A Systematic Review. *Journal of Speech,*

529         *Language, and Hearing Research*, *63*(4), 1093–1105.

530         https://doi.org/10.1044/2020_JSLHR-19-00017

531    Cristia, A., Lavechin, M., Scaff, C., Soderstrom, M., Rowland, C., Räsänen, O., Bunce, J., &

532         Bergelson, E. (2021). A thorough evaluation of the Language Environment Analysis

533         (LENA) system. *Behavior Research Methods*, *53*(2), 467–486.

534         https://doi.org/10.3758/s13428-020-01393-5

535   Cychosz, M., Villanueva, A., & Weisleder, A. (2021). Efficient Estimation of Children's

536          Language Exposure in Two Bilingual Communities. *Journal of Speech, Language, and*

537          *Hearing Research*, *64*(10), 3843–3866. https://doi.org/10.1044/2021_JSLHR-20-00755

538   Larson, A. L., Barrett, T. S., & McConnell, S. R. (2020). Exploring Early Childhood Language

539          Environments: A Comparison of Language Use, Exposure, and Interactions in the Home

540          and Childcare Settings. *Language, Speech, and Hearing Services in Schools*, *51*(3), 706–

541          719. https://doi.org/10.1044/2019_LSHSS-19-00066

542   Lehet, M., Arjmandi, M. K., Houston, D., & Dilley, L. (2021). Circumspection in using

543          automated measures: Talker gender and addressee affect error rates for adult speech

544          detection in the Language ENvironment Analysis (LENA) system. *Behavior Research*

545          *Methods*, *53*(1), 113–138. https://doi.org/10.3758/s13428-020-01419-y

546   Marchman, V. A., Martínez, L. Z., Hurtado, N., Grüter, T., & Fernald, A. (2017). Caregiver talk

547          to young Spanish-English bilinguals: Comparing direct observation and parent-report

548          measures of dual-language exposure. *Developmental Science*, *20*(1), e12425.

549          https://doi.org/10.1111/desc.12425

550   Orena, A. J., Byers-Heinlein, K., & Polka, L. (2019). Reliability of the Language Environment

551          Analysis Recording System in Analyzing French–English Bilingual Speech. *Journal of*

552          *Speech, Language, and Hearing Research*, *62*(7), 2491–2500.

553          https://doi.org/10.1044/2019_JSLHR-L-18-0342

554   Orena, A. J., Byers-Heinlein, K., & Polka, L. (2020). What do bilingual infants actually hear?

555          Evaluating measures of language input to bilingual-learning 10-month-olds.

556          *Developmental Science*, *23*(2), e12901. https://doi.org/10.1111/desc.12901

557 Place, S., & Hoff, E. (2011). Properties of Dual Language Exposure That Influence 2-Year-Olds'

558      Bilingual Proficiency: Dual Language Exposure and Bilingual Proficiency. *Child*

559      *Development*, *82*(6), 1834–1849. https://doi.org/10.1111/j.1467-8624.2011.01660.x

560 Place, S., & Hoff, E. (2016). Effects and noneffects of input in bilingual environments on dual

561      language skills in 2 ½-year-olds. *Bilingualism: Language and Cognition*, *19*(5), 1023–

562      1041. https://doi.org/10.1017/S1366728915000322

563 R Core Team. (2021). *R: A language and environment for statistical computing.* (4.1.2). R

564      Foundation for Statistical Computing. https://www.R-project.org/

565 Ramírez-Esparza, N., García-Sierra, A., & Kuhl, P. K. (2017a). Look Who's Talking NOW!

566      Parentese Speech, Social Context, and Language Development Across Time. *Frontiers in*

567      *Psychology*, *8*, 1008. https://doi.org/10.3389/fpsyg.2017.01008

568 Ramírez-Esparza, N., García-Sierra, A., & Kuhl, P. K. (2017b). The Impact of Early Social

569      Interactions on Later Language Development in Spanish-English Bilingual Infants. *Child*

570      *Development*, *88*(4), 1216–1234. https://doi.org/10.1111/cdev.12648

571 Ruan, Y., Orena, A. J., & Polka, L. (2020, July). *The Relation of Language Input and Infant*

572      *Volubility in English-French Bilingual Families*. International Congress of Infant Studies

573      2020, Virtual.

574      https://www.researchgate.net/publication/342715654_The_Relation_of_Language_Input_

575      and_Infant_Volubility_in_English-French_Bilingual_Families

576 Ruan, Y., Orena, A. J., Xu, K., & Polka, L. (2020). Language input and volubility in French-

577      English bilingual infants. *The Journal of the Acoustical Society of America*, *148*, 2501.

578      https://doi.org/10.1121/1.5146940

579    Scaff, C., Casillas, M., Stieglitz, J., & Cristia, A. (2022). Characterizing children's verbal input in

580            a forager-farmer population using long-form audio recordings and diverse input

581            definitions. PsyArXiv. https://doi.org/10.31234/osf.io/mt6nz.

582    Soderstrom, M., Grauer, E., Dufault, B., & McDivitt, K. (2018). Influences of number of adults

583            and adult: Child ratios on the quantity of adult language input across childcare settings.

584            *First Language*, *38*(6), 563–581. https://doi.org/10.1177/0142723718785013

585    Tamis-LeMonda, C. S., Kuchirko, Y., Luo, R., Escobar, K., & Bornstein, M. H. (2017). Power in

586            methods: Language to infants in structured and naturalistic contexts. *Developmental*

587            *Science*, *20*(6), e12456. https://onlinelibrary.wiley.com/doi/10.1111/desc.12456

588    VanDam, M., Warlaumont, A. S., Bergelson, E., Cristia, A., Soderstrom, M., De Palma, P., &

589            MacWhinney, B. (2016). HomeBank: An Online Repository of Daylong Child-Centered

590            Audio Recordings. *SEMINARS IN SPEECH AND LANGUAGE*, *37*(2), 128–141.

591            https://doi.org/10.1055/s-0036-1580745

592    Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis.* Springer-Verlag.

593            https://ggplot2.tidyverse.org.

594    Xu, K., Orena, A. J., Ruan, Y., & Polka, L. (2019). The effects of parental interaction on French-

595            English bilingual infants' vocalization and turn-taking rates. *Journal of the Acoustical*

596            *Society of America*, *145*, 1765. https://doi.org/10.1121/1.5101466

597

598

Tables

Table 1 *Variables and descriptions.*

| Variables | Descriptions |
|---|---|
| **Infant Language Development** | |
| Infant volubility | LENA-derived child vocalization counts (CVC) in the entire corpus. |
| **Language Input** | |
| *Measuring Units* | |
| AWC | The LENA-derived estimate of the number of words spoken near the child. |
| Duration | The sum of LENA-derived Adult Female Speech Duration and Adult Male Speech Duration. |
| Segment Count | The number of 30-second segments. |
| *Sampling Methods* | |
| Every-other-segment sampling | A periodic sampling method selecting every other 30-second segment containing adult speech. |
| Top sampling | A sampling method selecting a certain number of segments with the highest AWCs. Two top sampling methods were examined in this study: a simple top sampling (see Top150) and a post-annotation top sampling (see Top40/20). |
| *Samples* | |
| LENA | The entire Montréal Bilingual Infants corpus, consisted of 1,264 hours of audio recordings. |
| EOS | The sample selected by every-other-segment (EOS) sampling method, consisted of every other 30-second segment containing adult speech from the entire corpus. Every segment was annotated for speaker(s), listeners(s), and language usage. |
| Top150 | The sample selected from the EOS sample by the simple top sampling method, consisted of top 150 segments with the highest AWCs. |
| Top40/20 | Samples selected from the EOS sample by the post-annotation top sampling method, consisted of top 40 segments with the highest AWCs in a specific social context (overhearing, one-on-one), or top 20 segments with the highest AWCs in a specific language context (dominant, non-dominant). |
| *Social and Language Contexts* | |
| Global | All input in the sample. |
| Overhearing | Caregivers spoke in the presence of the infant but not exclusively addressing the infant. |
| One-on-one | One caregiver (mother, father, nanny, older sibling, and other) talked to the infant. |
| Dominant language | Parent-reported language (French or English) that the infant has more exposure to at each age. |
| Non-dominant language | The language other than the dominant language (English or French). |

Table 2 *Comparison of proportions of language input (indexed by AWC) in different social and language contexts across Every-other-segment (EOS) and Top150 samples* [1].

| | EOS | | Top150 | | Difference[1] | | Wilcoxon *V* |
|---|---|---|---|---|---|---|---|
| | Median | Interquartile Range | Median | Interquartile Range | Median | Interquartile Range | |
| | | | | Social Contexts | | | |
| 10M: Overhearing | 73% | 68 – 76% | 81% | 70 – 87% | 6% | 4 – 10% | 43# |
| 10M: One-on-one | 25% | 20 – 32% | 19% | 11 – 28% | 7% | 3 – 9% | 183# |
| 18M: Overhearing[2] | 66% | 44 – 82% | 68% | 44 – 82% | 2% | 0.8 – 4% | 57 |
| 18M: One-on-one | 34% | 18 – 56% | 32% | 18 – 56% | -[2] | - | - |
| | | | | Language Contexts | | | |
| 10M: Dominant | 51% | 38 – 54% | 46% | 33 – 54% | 5% | 2 – 7% | 174# |
| 10M: Non-dominant | 21% | 14 – 30% | 21% | 15 – 25% | 2% | 1 – 4% | 80 |
| 18M: Dominant | 35% | 25 – 51% | 33% | 23 – 53% | 2% | 0.5 – 3% | 98 |
| 18M: Non-dominant | 15% | 12 – 21% | 15% | 11 – 22% | 1% | 0.8 – 1% | 95 |

*Note*: # $p < .05$, adjusted $p > .05$. The *p*-values were adjusted using method of Benjamini & Hochberg (1995).

[1] EOS: the sample of every other 30-second segment containing AWCs. Top150: top 150 segments with the highest AWCs.

Difference: the difference between proportions of the same type of input in the two samples.

[2] Statistical analysis was not performed for one-on-one input in 18-month dataset to avoid redundancy as the proportion of overheard and 1:1 input added up to 100% in the 18-month dataset. It was not the case for the 10-month dataset because there was a third type of social input, that is group input which we observed none in the 18-month dataset.

10M: 10-month sample; 18M: 18-month sample.

Table 3 *Comparison among Spearman's correlations between infant volubility and bilingual input estimated by different units and sampling methods[1].*

| Input Measures[1] | Every-other-segment sampling | | | Top sampling | | | | |
| | | | | Top150 | | | Top40/20 | |
| Correlations | Segment | AWC | Duration | Segment | AWC | Duration | AWC | Duration |
|---|---|---|---|---|---|---|---|---|
| 10M: Global | .68*** | .60** | .62** | - | .49* | .48* | - | - |
| 10M: Overhearing Contexts | .49* | .49* | .52* | −.02 | .28 | .25 | .39 | .39 |
| 10M: One-on-one Contexts | .50* | .50* | .48* | .05 | .07 | .11 | .32 | .32 |
| 10M: Dominant Language | .23 | .37 | .34 | .06 | .30 | .29 | .32 | .34 |
| 10M: Non-dominant Language | .39 | .45* | .43 | .03 | .30 | .31 | .40 | .39 |
| 18M: Global | .58* | .49 | .49 | - | .35 | .33 | - | - |
| 18M: Overhearing Contexts | .17 | .15 | .17 | −.28 | .05 | .02 | .05 | .04 |
| 18M: One-on-one Contexts | .44 | .57* | . 56* | .28 | .45 | .46 | .58* | .57* |
| 18M: Dominant Language | .50 | .61** | .60** | .33 | .51* | .53* | .37 | .47 |
| 18M: Non-dominant Language | .04 | .10 | .09 | −.26 | <.01 | .01 | .09 | .07 |

*Note*:  * *p* < .05; ** *p* < .01; *** *p* < .001.

[1] Segment: Segment Count, the number of 30-second segments. AWC: LENA-derived adult word counts. Duration: the sum of LENA-derived Female and Male Speech Duration. Every-other-segment sampling: a sample of every other 30-second segment containing AWCs. Top150: top 150 segments with the highest AWCs. Top 40: top 40 segments with the highest AWCs in one-on-one or overhearing social contexts. Top 20: top 20 segments with the highest AWCs in the dominant or non-dominant language.

10M: 10-month sample; 18M: 18-month sample.

Figures

```
┌─────────────────────────┐
│      Entire corpus      │
│         (LENA)          │
│  10M: 120959; 18M:30720 │
│      AWC, Duration      │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│ Every-other-segment Sample# │        Social Contexts
│         (EOS)           │      "Who is talking to whom"
│  10M: 18979; 18M: 6180  │  Human
│ Segment count, AWC, Duration │ Annotation   Language Contexts*
└─────────────────────────┘              "In which language(s)"
            │
            ▼
┌─────────────────────────┐        Social Contexts
│ Top-150-segment Sample  │
│        (Top150)         │
│ 10M: 3150, 17%; 18M: 2401, 37% │  Language Contexts
│      AWC, Duration      │
└─────────────────────────┘
```

**Overhearing**
10M: 12225; 18M: 3906
*Segment count, AWC, Duration*
→
**Top40**
10M:840, 7%; 18M:640, 18%
*AWC, Duration*

**One-on-one**
10M: 6351; 18M: 2274
*Segment count, AWC, Duration*
→
**Top40**
10M:840, 14%; 18M:623, 31%
*AWC, Duration*

**Dominant Language**
10M: 9527; 18M: 2606
*Segment count, AWC, Duration*
→
**Top20**
10M:420, 4%; 18M:320, 12%
*AWC, Duration*

**Non-dominant Language**
10M: 3922; 18M: 1177
*Segment count, AWC, Duration*
→
**Top20**
10M:407, 10%; 18M:293, 30%
*AWC, Duration*

**Overhearing**
10M: 2322; 18M: 1556
*Segment count, AWC, Duration*

**One-on-one**
10M: 781; 18M: 845
*Segment count, AWC, Duration*

**Dominant Language**
10M: 1393; 18M: 936
*Segment count, AWC, Duration*

**Non-dominant Language**
10M: 700; 18M: 434
*Segment count, AWC, Duration*

*Notes:*
*# All recording intervals with no LENA-detected speech were first removed from the entire corpus.*
*\*Mixed-language input was infrequent (<10% in every family) and was not included in the language context analysis.*
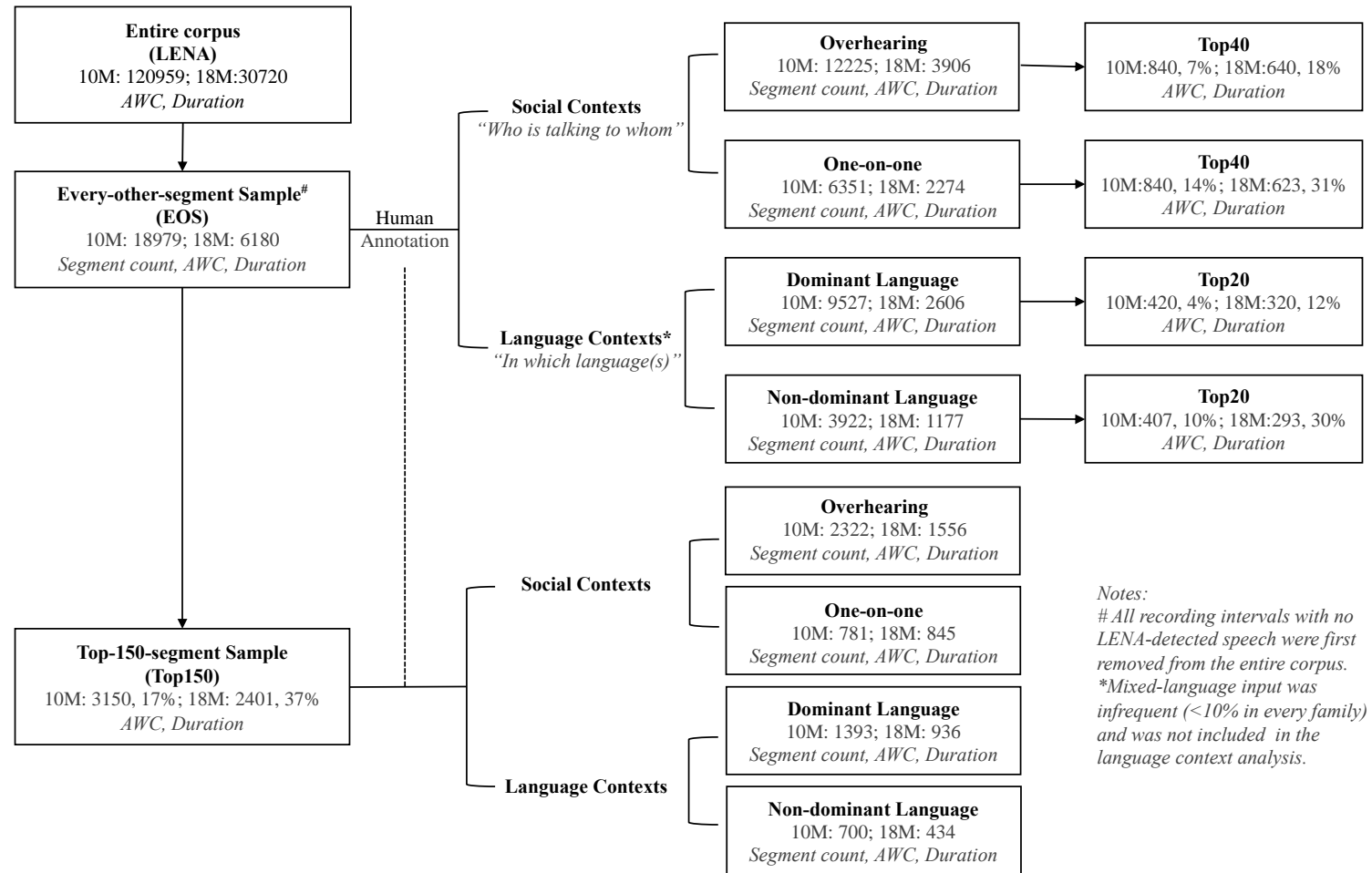
Figure 1. Flowchart describing how each sample was derived from the corpus, the number of segments included in each sample, the

proportion of segments selected from the original sample (median), and units used to index input in each sample (Italic). 10M: 10-

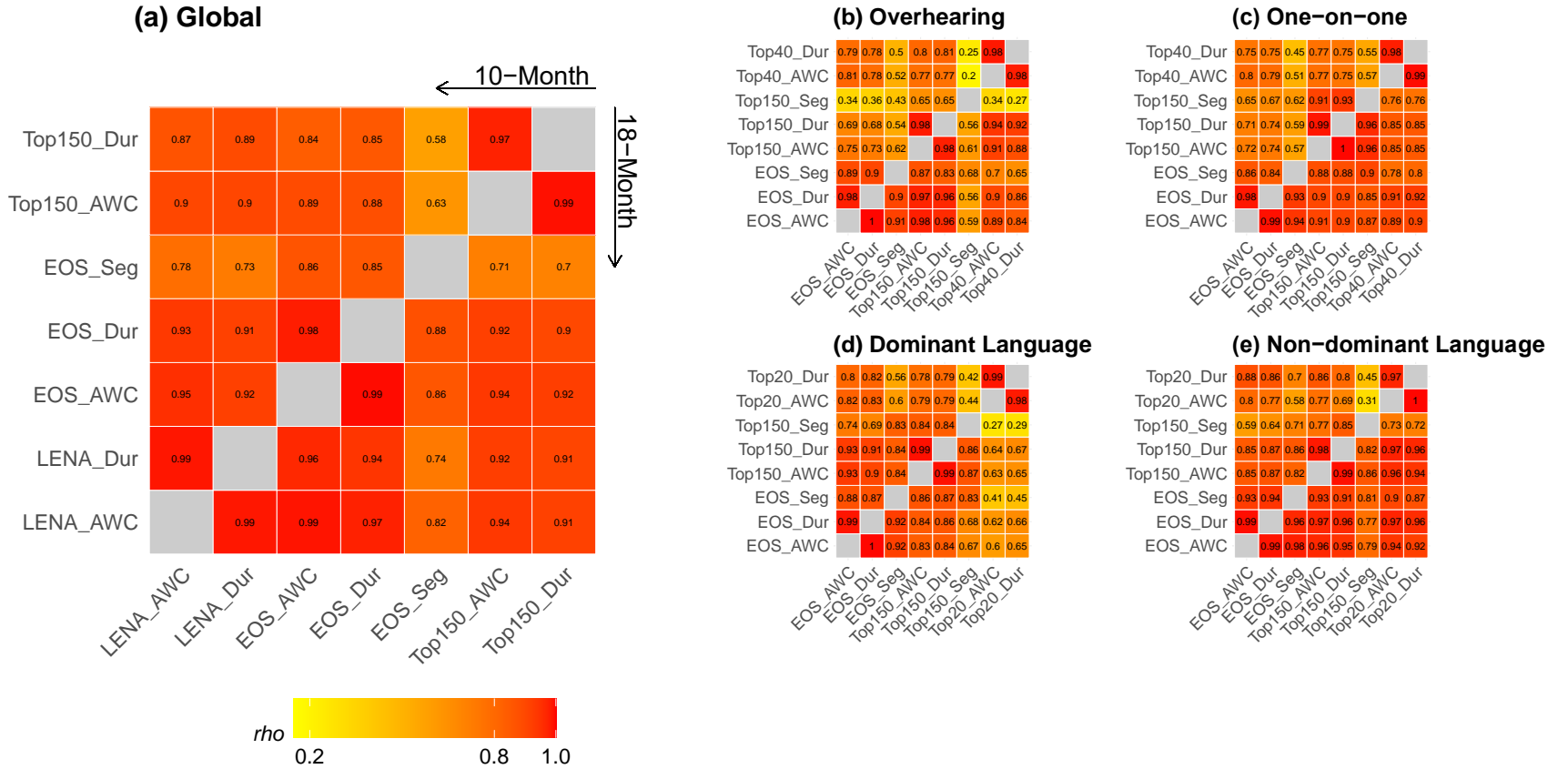month sample; 18M: 18-month sample.

Figure 2. Spearman's correlations between different language input measures in (a) global, (b) overhearing and (c) one-on-one social contexts, as well as (d) dominant and (c) non-dominant language contexts. Upper triangle: 10-month sample; Bottom triangle: 18-month sample. Each cell indicates the correlation between a pair of input measures. The Spearman's *rho* value is reported in each cell. The cell colour indicates the strength of the correlation, from weak (yellow) to strong (red). LENA: the entire corpus. EOS: every-other-segment sample. Top150: top 150 segments with the highest adult word counts (AWCs). Top 40: top 40 segments with the

highest AWCs in one-on-one or overhearing social context. Top 20: top 20 segments with the highest AWCs in the dominant or non-dominant language. AWC: LENA-derived adult word counts. Dur: Duration, the sum of LENA-derived female and male speech duration. Seg: Segment Count, the number of 30-second segments. The columns (C) and rows (R) are referred numerically from left (1) to right, and from bottom (1) to top. For example, [C1, R2] in (a) refers to the cell corresponding to the first column from the *left* (LENA_AWC) and the second row from the *bottom* (LENA_Dur), which shows the correlation between the global input estimated in the entire 10-month corpus by AWC and speech duration. A video-animated guide is available in the Supplementary Material.