1

2

3 Exploring Characteristics of Quality in Language Teaching & Learning:

4 The Mother Tongue Adapted Coding Scheme (MACS)

5
6

7

8

9 Beth O'Brien[ab], Nurul Taqiah Yussof[a], Poorani Vijayakumar[a],

10 Malikka Begum Habib Mohamed[a], Leong Xuan En Rachel[c]

11 [a]Centre for Research in Child Development, National Institute of Education
12 [b]Centre for Research and Development in Learning,
13 Nanyang Technological University, Singapore
14 [c]Singapore University of Social Sciences

15

16

17

18

19

20

21

22

23 Correspondence should be directed to Beth A. O'Brien, beth.obrien@nie.edu.sg, (ORCHID #0000-0002-
24 1187-5908), 1 Nanyang Walk, National Institute of Education, Nanyang Technological University,
25 Singapore, 637616
26
31
32 Declarations of interest: None.
33
34

35

**Abstract**

Teacher-child interactions are an integral factor influencing the quality of early childhood education, and multilingualism is increasingly prevalent in many contexts. In the multilingual society of Singapore this is particularly relevant because early childhood classrooms follow a bilingual policy and include the teaching and learning of Mother Tongue languages. To evaluate what constitutes quality for preschool language teaching in this multilingual context, the Mother Tongue Adapted Coding Scheme (MACS) was developed to examine teacher-child interactions during blocks of Mother Tongue language lessons. The comprehensive observational rating scheme for language learning included four domains: language input, language output, varied teaching strategies, and English use for teaching Mother Tongue language. This rating scheme was applied to fifty-one classrooms where teachers were videotaped while they conducted Chinese, Malay, or Tamil language lessons. Within the observed classrooms, four hundred and ninety-one children were assessed on their Mother Tongue language learning over three years for receptive vocabulary, reading, and morphological awareness. Overall interrater reliability on the MACS was high across language classrooms ($K > 0.72$). Correlation and exploratory factor analyses indicated a main factor for the language input and output domains, and separate factors for English language use. Teachers' factor scores were found unrelated to their scores on the CLASS. Higher factor scores on input/output quality were correlated with having more high progress learners (compared to peers) for receptive vocabulary, but less high progress learners for reading. Implications for future applications and adaptations of the MACS are discussed in light of the use of rating tools to examine and better understand the constituents of ECE quality within multilingual contexts.

**Introduction**

58

59        Early childhood is a critical period of development, and children are spending an

60    increasing amount of time in classrooms in their early years. Based on evidence that the quality

61    of early childhood education (ECE) impacts children's development (e.g., Burchinal, Magnuson,

62    Powell & Hong, 2015), there is heightening interest in understanding, assessing, and ensuring

63    quality in ECE. Quality in ECE is currently conceptualized in terms of process quality involving

64    caregiver-child interactions, and structural quality including teacher experience and

65    characteristics of the ECE program (Burchinal, 2018). In the quest to measure such components

66    of quality of the ECE environment there is also an awareness that either universal features or

67    socio-culturally contextualized ones may play a role (Schwartz, 2018). This paper focuses on

68    what constitutes quality of early bilingual education, and preschool classroom practices that best

69    support language outcomes (Palviainen et al., 2016). As this study is situated in Singapore - a

70    multilingual context with unique sociopolitical factors contributing to the convergence of diverse

71    language competencies - this question is key yet unresolved. Here, as in many multilingual

72    societies, the school and home language may differ for some students, creating a complex

73    classroom environment in which individuals may vary widely on proficiency levels. This

74    translates to a need for teachers to cater to different language competencies using a variety of

75    instructional strategies, to provide comprehensible input for all (Jiang, Garcia & Willis, 2014;

76    Probyn, 2015). Many of the currently established ECE quality rating systems were developed

77    within Western contexts and may not be sensitive to these unique teaching and learning needs.

78        The multilingual society of Singapore includes a 'bilingual education' model, in which

79    English is the main medium of instruction, but students must also demonstrate proficiency in a

80    designated Mother Tongue language (MTL) by the end of primary school. The official MTLs are

81    Mandarin (Chinese), Behasa Melayu (Malay) or Tamil, and are assigned based on children's

82    ethnicity, though these languages are currently used less frequently in homes. A proportion of

83    class time is allocated for MTL lessons (which we refer to as the "target" language), and this so-

84    called MTL at times does not match the actual home language (Dixon, 2009), which has become

85    increasingly English-dominant (Bolton & Ng, 2014). Consequently, at primary school entry

86    children are differentially prepared for the subject of MTL due to variation in home language

87    exposure (Sun, Yin, Amsah & O'Brien, 2018). This places the preschool classroom in a pivotal

88    role for children's early MTL learning, acquisition and school readiness.

89          Given this key period of preschool for child development, coupled with the specific

90    socio-cultural demands and educational objectives in the Singapore context, we sought to

91    examine characteristics of quality ECE regarding language learning within a bilingual context.

92    We focus on MTL learning and teaching within the ECE context given the concerns with

93    increasing English dominance. To this end, we constructed a measure of early education quality

94    specific to MTL learning called the Mother Tongue Adapted Coding Scheme (MACS). We then

95    applied the coding scheme to videotaped classroom language lessons. In sections of this paper

96    we describe the process of developing the coding scheme, the content and rating scales of the

97    coding scheme, and several approaches for examining the validity of the coding scheme.

98    **Conceptual Framework**

99          The MACS was developed as part of the Singapore Kindergarten Impact Project (SKIP,

100   Ng et al., 2014), which is a large-scale longitudinal study of children's school readiness. A key

101   research question of SKIP asks what characterizes quality preschool instruction, including

102   language instruction. Through several phases (see Figure 1), the rating scheme was developed

103   with a broad framework to be comprehensive, as an initial attempt to characterize the ECE

104 environment in Singapore. We began with a review of existing observational tools for preschool

105 and bilingual contexts (summarized in Appendix A) and extended these with consideration of

106 local language planning and language-specific factors, and with characteristics based on

107 language or second language learning theory. These sources contributed to the conceptualization

108 of the MACS (Figure 2) across 4 broad domains – teachers' language *input*, children's elicited

109 language *output*, *varied* pedagogical *strategies* for language learning, and *use of English* as

110 bilingual children's stronger language. They are described in more detail in the following

111 sections: (1) considering linguistic frameworks for the input and output strategies domains, (2)

112 considering pedagogical frameworks for the varied strategies domain, (3) considering context-

113 specific factors for the English use domain, and (4) overall adaptation from existing tools.

114 **1. Considering linguistic frameworks.** Language input and use are widely accepted as

115 the basis for language acquisition (Unsworth, Persson, Prins & de Bot, 2015). This occurs

116 ultimately through comprehensible input (i.e., the Input Hypothesis, Krashen, 1985), and so

117 teachers' provision of simplified input is key for learner comprehension. At the same time,

118 language use or 'output' is considered necessary for facilitating linguistic knowledge of the

119 learner (i.e., the Output hypothesis, Swain, 1985), and for leading to further input from teachers

120 (Pearson, 2007) (refer to Figure 2).

121 *Input* quantity and quality are consistently linked to bilingual children's language and

122 literacy outcomes (Hoff & Core, 2013; Huang & Kuo, 2020). Optimal input stems from highly

123 proficient speakers: input from native speakers in the home was found most beneficial for young

124 children's development (Place & Hoff, 2011), and input from highly proficient speakers was

125 reported as an integral factor in quality language instruction (Canh & Renandya, 2017; Richards,

126 2010). Such proficient input includes the qualities of correct usage (pronunciation, syntax and

127   grammar) as well as linguistic richness (variety of words and sentence types) (Bornstein, Haynes

128   & Painter, 1998). Quality input is especially important for languages with insufficient exposure

129   (De Houwer, 2005).

130          Acquiring both language form (knowledge about organizing sounds, words and

131   sentences) and content (meaning) knowledge are important for language learning (Lahey, 1988).

132   Hence, language teaching may focus on either or both of these, as the target language may be the

133   subject (teaching proper oral and written language forms) or the medium of instruction for

134   teaching other content (discussing meaning and developing concepts) (Ellis & Shintani, 2014).

135   For example, proficient teachers, who not only model oral and written language, but also

136   elaborate and explain words and concepts, may increase children's learning and continued

137   engagement (Pakarinen et al., 2010; La Paro, Pianta & Stuhlman, 2004).

138          *Output* is also considered key to language learning, including for bilingual vocabulary

139   development and proficiency (e.g., Bohman, Bedore, Pena, Mendez-Perez & Gillam, 2010;

140   Swain, 2005). Output allows learners to notice gaps in their knowledge (de Bot, 1996), to access

141   meaning and to generalize, which can serve as building blocks for language representations (e.g.,

142   Ibbotson, 2013). Encouraging children to use the target language and providing feedback on their

143   output enhances learning. This is particularly important because even with extended input,

144   children's target language output often lags behind their listening comprehension (e.g., DePalma,

145   2010).

146          Schwartz and Gorbatt (2016) noted that bilingual preschool teachers need to create a low-

147   anxiety, secure environment to encourage children's target language production. Thus,

148   encouraging children to use the target language within a positive classroom climate was

149    considered an important aspect of eliciting language output. Teachers' use of language modelling

150    and scaffolding learner output was also considered essential (e.g., Gibbons, 2006).

151         **2. Considering pedagogical frameworks.** Besides the linguistics perspective,

152    pedagogical factors were included in a set of *Varied Strategies* that teachers may use to further

153    assist children's understanding of the target language. A socio-cultural perspective holds social

154    interactions in the classroom as central to language acquisition and instruction (Vygotsky, 1978),

155    and is considered in terms of the input-output cycle and the role of feedback for second language

156    learning (Gass & Mackey, 2007). Since the current study focuses on preschoolers who may lack

157    exposure to MTL input, or who may be learning the MTL for the first time, we considered

158    different strategies recommended for beginning language learners (see Table 1), including

159    adapted (child-directed) speech (Wesche, 1994), and non-verbal cues like gestures and facial

160    expressions (Kersten, Steinlen, Tiefenthal, Wippermann & Matsson, 2010; Weitz, Pahl,

161    Mattsson, Buyl, & Kalbe, 2010).

162         Another set of strategies of particular relevance to multilingual contexts is the use of the

163    dominant or first language in target language teaching. Teachers may switch to children's

164    dominant or stronger language (often English) for explaining and translating to ensure target

165    language comprehension (Chimbutane, 2013; Lin, 2005), for managing classroom behaviors, and

166    facilitating task completion (Enama, 2016; Gort & Pontier, 2013). These strategies are detailed in

167    Table 1 and consider purposes of switching to English (instructional vs. managerial) (see Figure

168    2).

169         **3. Considering context-specific factors.** In multilingual societies with changing

170    demographics and differences between generational cohorts, the term 'Mother Tongue' language

171    (MTL) is not necessarily synonymous with the first and most proficient language (Wright, Boun

172  & Garcia, 2015). In Singapore schools, the MTL may not always correspond to children's home

173  or dominant language: for example, home languages could include Telugu, Urdu, Hindi,

174  Hokkien, or Teochew (Dixon, 2009); and younger generations tend to come from more English

175  dominant homes (Bolton & Ng, 2014; Singapore Department of Statistics, 2018). As such,

176  'MTLs' may appear more like a heritage language for many Singaporean school children, where

177  the MTL curriculum serves the education policy role of 'cultural ballasts', ensuring transmission

178  of cultural knowledge (Ministry of Education, 2015), while the teacher serves the role of  a good

179  model of MTL use. A policy for English use in the MTL classroom is not widely acknowledged,

180  and the common approach is to minimize switching to English and to maximize MTL input (e.g.,

181  Mukhlis & Pang, 2015). We therefore explored the prevalence of English use in preschool MTL

182  classrooms with the MACS.

183        Second, there are aspects of the MTLs themselves to be considered, in terms of the

184  impact on instruction. Two of Singapore's MTLs, Malay and Tamil, are characterized by

185  diglossia, with distinct spoken and written forms. An educator's guide for Malay language

186  teachers (Ministry of Education, 2015) expresses the need for using the standardized 'baku' form

187  of Malay during teaching, especially for phonological instruction and reading aloud. At the same

188  time, studies indicated that for spoken language a Johor-Riau variety of Malay is still more

189  prevalent than the standard form in Singapore (Subhan, 2013). For the Tamil language, a spoken

190  variety in Singapore referred to as Standard Spoken Tamil (SST) sees prevalent use in Tamil-

191  speaking households and the wider speech community, while a literary form of Tamil is used in

192  classrooms (Lakshmi & Saravanan, 2009; Schiffman, 2004). Teachers are recommended to use

193  SST in classrooms to bridge the gap with the home language and to facilitate children's

194  understanding and expression (Lakshmi & Saravanan, 2009). Thus, communicative language in

195    and out of classrooms differs from the literary forms taught in MTL classrooms for Malay and

196    Tamil. These sociolinguistic aspects influenced the conceptualization of the Teacher Input

197    indicator within MACS, as illustrated in Figure 3.

198        Varieties of spoken Chinese, by teachers who come from outside of Singapore (Teng,

199    2018), may also lead to gaps with children's community language. For example, Taiwanese

200    Mandarin has the tendency towards the neutralization of contrasts between past and present

201    (Cheng, 1985), while for mainland Chinese Mandarin, a past specific action is almost always

202    marked by *le* (completive or past). Another difference between these varieties of Mandarin is

203    with reduplication of the verb, which is generally used to express the tentative aspect in mainland

204    Chinese Mandarin, but not in Taiwanese Mandarin, which uses the term *yixia* (a while) instead.

205    This is seen in the phrase 'wait a minute', which is expressed as *deng-deng* (wait-wait), while in

206    Taiwanese Mandarin it is expressed as *deng yi-xia* (wait a while). As such, the Teacher Input

207    indicator for correctness was considered in relation to Chinese Mandarin (outlined in Table 1).

208        **4. Adapting from existing tools.** Existing observation schemes that informed the

209    development of the MACS include the IQOS (Weitz et al., 2010), the LISn (Atkins-Burnett,

210    Sprachman, & Caspe, 2010), and the CLASS (Pianta, La Paro, & Hamre, 2008) (see Appendix

211    A). The IQOS (Weitz et al., 2010) is geared towards classroom quality in bilingual preschool

212    contexts and shares the aim in the current paper to compare teachers' language use as it relates to

213    children's development, using quantitative observation methods. The LISn (Atkins-Burnett,

214    Sprachman, & Caspe, 2010) similarly focuses on the bilingual preschool context, but in diverse

215    settings where teachers cater to children from households with varying levels of target language

216    exposure. The CLASS (Pianta et al., 2008) was developed to assess different aspects of

217    children's preschool experience and teacher-child interactions holistically rather than focusing on

218    language strategies, and was adopted as a widely used scoring scheme within the larger SKIP

219    study (Ng et al., 2014).

220          With regard to the IQOS, we adapted many of the behavior indictors, including teachers'

221    comprehensible, simplified and contextualized input that is rated under ritualized language,

222    adapted speech, contextualization strategies under the domain of varied strategies, rich and

223    varied input indicators, and output indicators of positive climate, whereby children were

224    encouraged but not forced to produce and maintain target language output. After piloting and

225    discussions to reduce rating subjectivity, we use a reduced number of indicators and changed

226    some items from 'high' to 'low inference', through observing frequency rather than quality of

227    the indicators for ritualized language, adapted speech, and contextualization strategies.

228          Also, due to differences in sociolinguistic context, we did not adopt the IQOS language

229    separation stance where minimal translation or language switching reflected higher quality input.

230    With the MACS, both the potentially negative and positive instances of teacher language

231    switching to English were considered. Teachers were scored lower if they switched to English

232    without instructional or administrative purposes on the MACS Teacher Input indicator,

233    particularly in the case of colloquial English mixed language use. Within the MACS Use of

234    English domain, teachers' translations between English and Mother Tongue were observed for

235    providing learning opportunities or for classroom management purposes. Separate ratings of

236    these frequencies allow us to consider the benefits, or lack thereof, for each type.

237          With regard to the LISn (Atkins-Burnett, Sprachman, & Caspe, 2010), which captures the

238    varying ways that language switching occurs (repetition, confirmation, elaboration, and so on),

239    we were also interested in the ways teachers switched to English from Mother Tongue, but not to

240    the same level of detail. The LISn was used and scored differently, involving five-minute

241    'snapshots' per observed child, whereas MACS scoring was focused at the class level for the 10-

242    to 20-minute video observations. While more intensive, child-focused rating has its merits, the

243    approach was not intended to capture other elements of classroom language exchanges, such as

244    the quality or complexity of each type of talk, which was the focus for the MACS. The LISn goal

245    was to closely observe the type of individualized input and interactions a child would receive or

246    be engaged in, and thus provided a useful model for types of talk data and activity types that are

247    important for observation.

248            With the CLASS (Pianta et al., 2008), one of its three domain scores was most relevant to

249    the objectives of MACS: Instructional Support, which includes the indicators of Language

250    Modelling and Quality of Feedback. These were both incorporated within the Language

251    Modelling and Scaffolding indicator of the MACS. However, the MACS descriptors were

252    revised to include effectiveness of teacher strategies through the observation of the proportion of

253    children who spoke in class, and the extent to which they spoke (words, phrases or sentences).

254    The Varied Strategies domain of the MACS also adapted activities and directions from the

255    CLASS Instructional Learning Format dimension (La Paro, Pianta, & Stuhlman, 2004), but these

256    were further detailed in the MACS (see Table 1), while CLASS adopted a more general

257    evaluation. Contrary to the purposes of the other schemes above, CLASS was not intended for

258    rating bilingual practices in classrooms and did not include language switching.

259            The scoring scales of the CLASS were adapted into the MACS, in terms of the 'depth,

260    duration and frequency' of teacher input and strategies for each indicator coded. This was due to

261    the similar way the video data was captured and prepared for scoring, whereby 10- to 20-minute

262    video segments were scored, and the focus of the scoring was on teacher-child interactions.

263    However, for MACS' focus on specific aspects of mother tongue language learning, necessary

264    modifications were made to the rating scales after piloting (as described below).

265    **Current Study**

266        The contextual factors discussed above reflect diversity in preschoolers' MTL

267    proficiencies and learning needs that teachers have to address in order to provide quality

268    language instruction. Therefore, the MACS was developed somewhat differently from other

269    tools, in order to assess a range of teacher input and strategies. This was planned to better

270    understand aspects of quality MTL instruction and input, given the unique local context with

271    variation in children's MTL proficiency. This study explores the degree to which the MACS

272    achieves this aim with the following research questions (RQ):

273        **RQ 1.** What contributes to the conceptualization and development of a classroom quality

274    scoring scheme in preschool language classrooms for bilingual learners?

275        **RQ 2**. To what extent can different raters reliably code classrooms with the MACS?

276        **RQ 3**. Do the dimensions of the MACS fit with theoretical expectations?

277        **RQ 4**. Do ratings across the dimensions predict child learning outcomes?

278        Different methods were applied to address each research question. The first research

279    question involves steps in the construction of the MACS from initial development, to revisions,

280    and finally the application of the scheme, as seen in Figure 1. The literature review, summarized

281    above, was followed up with revising the scheme through piloting, consensus discussions and

282    consultation with an expert panel, resulting in a manual with descriptors for each indicator and

283    examples derived from the first set of classroom videos (Table 1). These phases are elaborated in

284    the results section. To address the second research question, a new group of raters was exposed

285    to a set of training videos, and their interrater reliability was established across the full sample of

286 the classrooms. To address the third research question, we use points outlined in Lahti, Elicker,

287 Zellman, and Fine (2015) and Zellman and Fiene (2012) to establish reliability, and construct

288 and concurrent validity. Finally, to examine research question four regarding the predictive

289 validity of the MACS, child language measures were examined across the rated classrooms.

290      The initial version of MACS consisted of the four domains noted above: *input strategies,*

291 *output strategies, varied strategies* and the *use of English*. These consisted of 12 behavioral

292 indicators initially, and 11 in the final version (Figure 2).

293                                                **Methods**

294      The samples, measures and procedures for the study are provided below for each research

295 question, based on the different phases of the study. Samples include the teachers in the observed

296 classrooms, the coders who rated the classroom teachers with the MACS, and students from the

297 classrooms. Overall there were 51 classroom teachers, with 27 observed for the initial phase (RQ

298 1) and 48 double scored for estimation of interrater reliability (RQ 2), and 491 students (RQ 4).

299 **Sampling**

300      Across the 51 classrooms selected from a larger longitudinal study (SKIP, Ng et al.,

301 2014), all teachers are female and acquired at least a diploma to teach in Singapore preschools

302 (Early Childhood Development Agency, 2018). Teachers' self-reported proficiency indicated

303 that they are fluent in MTL and conversational in English as rated on a 1-7 scale (English - $M =$

304 4.11, $SD = 1.85$; MTL - $M = 6.16$, $SD = 1.32$), speaking (English - $M = 4$, $SD = 1.78$; MTL – $M$

305 $= 6.11$, $SD = 1.35$), reading (English - $M = 4.03$, $SD = 1.91$; MTL – $M = 6.14$, $SD = 1.46$) and

306 writing (EL - $M = 3.57$, $SD = 2.08$; MTL – $M = 5.92$, $SD = 1.46$); where 3 = "can communicate

307 routine, basic information", 4 = "can independently understand/use language to carry on

308 conversation", and 5 = "can independently interact with fluency and spontaneity").

309    From the larger study, all Malay and Tamil classes and a random subset of Chinese

310    classes were included for coding and analysis. Differences in the number of classes per language

311    are due to population demographics (74.3% Chinese, 13.4% Malay, and 9.1% Indian; Singapore

312    Department of Statistics, 2018). Classrooms included in the sample of 51 videos were from

313    government-run (21.6%), public (17.6%), not-for-profit (45.1%), and commercial (15.7%)

314    preschool education providers recruited within the larger study.

315    During each observation session, one microphone and receiver was attached to a teacher-

316    focused camera to capture teacher-child interactions. Typical MTL lessons last 30 to 60 minutes,

317    depending on the school program, and at least one video session with the best angles and sound

318    quality was chosen per class for coding. If there was more than one video for the MTL lesson, all

319    videos (up to three) were coded and item scores were averaged. The maximum duration for each

320    video recording is 20 minutes, and the average duration across the 51 classrooms was 18min 20s

321    ($SD$ = 2min 31s).

322    The average number of children per class was 14.8 ($SD$ = 3.99) for Malay, 15.18 ($SD$ =

323    3.45) for Tamil, and 15.08 ($SD$ = 3.32) for Chinese classes. Out of the four hundred and ninety-

324    one children with outcome data, 96.5% matched in their ethnic background and classroom MTL

325    assignment. Assignment to MTL classes is based on ethnicity, except in cases where the child's

326    MTL is not taught in their school due to low enrollment, they will join an available MTL class

327    (Wong, 2018), with the teacher sometimes codeswitching to the shared language of English

328    (Yussof & Sun, 2020).

329    All coders had at least a bachelor's degree in early childhood, education or psychology.

330    All were Singaporean, and were fluent bilinguals in English and their MTL, having completed

331    their formal education locally, and having acquired and practiced spoken and academic registers

332    of their MTL. Thus, all Malay coders used the Johor-Riau variant of spoken Malay (Subhan,

333    2013), and learned formal Malay in school. The Tamil coders used Standard Spoken Tamil

334    (Lakshmi & Saravanan, 2009), and literary Tamil for reading and writing. These language

335    experiences allowed them to distinguish and code teachers' use of different registers within the

336    classroom.

337          **RQ 1. Phase 1-2 scheme development and revision.**

338          *Teachers.* Twenty-seven Kindergarten 1 (K1) MTL teachers (6 Malay, 7 Tamil, 14

339    Chinese) were observed. Their mean age = 41.37 ($SD$ = 11.80), and average preschool teaching

340    experience = 11.05 years ($SD$ = 6.93).

341          *Coders.* Six coders (2 per MTL) rated classroom videos from Phase 1, while fifteen total

342    coders rated videos in Phase 2 (refer to Figure 1).

343          **RQ 2. Phase 3 scheme application to classrooms.**

344          *Teachers.* Forty-eight observed classrooms were double coded, including15 K1 teachers

345    and 33 K2 teachers (12 Malay, 4 Tamil, 32 Chinese). Their ages ranged from 21 – 69 years (K1 -

346    $M$ = 41.00, $SD$ = 11.22; K2 – $M$ = 41.88, $SD$ = 12.44), and preschool teaching experience ranged

347    from 1 – 30 years (K1 - $M$ = 12.45., $SD$ = 7.43; K2 – $M$ = 10.85, $SD$ = 7.09).

348          *Coders.* Eighteen coders (6 Malay, 3 Tamil, 9 Chinese) rated the 48 classroom videos

349    using the final version of the MACS. Nine coders (2 Malay, 2 Tamil, 5 Chinese) were also

350    trained in CLASS.

351          **RQ 3. Phase 3 scheme fit with theoretical expectations.**

352          *Teachers.* Ratings for all 51 MTL teachers (13 Malay; 6 Tamil; 32 Chinese) were

353    conducted by the eighteen coders described above.

354          **RQ 4. Phase 3 scheme prediction of child outcomes.**

355          *Teachers.* Ratings were used for the same 51 teachers from the 18 coders.

356          *Children.*

357          *Classrooms Sample.* There were 491 children within the 51 classrooms with teachers'

358    MACS ratings, including 325 Chinese, 75 Malay, and 91 Tamil learners, with an average of

359    10.13 children within each classroom cluster ($SD = 6.72$, Range of $1 - 41$). They were assessed

360    on the MTL vocabulary, reading, and morphological awareness measures at three time points:

361    Kindergarten 1 (K1), Kindergarten 2 (K2) and Primary 1 (P1).

362          *Total Comparison Sample.* To appraise the learning of these children, rather than using raw

363    test scores, we considered their learning trajectories relative to same-language peers from the

364    larger sample of SKIP. This was to account for differences in the assessments across MT

365    languages, and because of a lack of locally-normed assessments. Data from the children in the

366    MACS classrooms along with the larger SKIP sample ($N = 1538$) were included in growth

367    mixture models (GMM) to determine latent classes of learners (high versus low progress).

368    Separate GMMs were run per language group (CL-Chinese, ML-Malay, TL-Tamil) for each

369    outcome. The total samples for these models included: (1) 1075 CL, 195 ML, 238 TL for

370    vocabulary; (2) 323 CL, 137 ML, 161 TL for reading; and (3) 319 CL, 137 ML, 168 TL for

371    morphological awareness. Total sample size on each measure varies due to missing data (e.g.,

372    child absences) and to task administration (subsamples of SKIP). All children participating in

373    SKIP were typically developing learners.

374          The classroom sample did not differ from the total sample in terms of age, gender

375    proportion, non-verbal intelligence, or family income: age in years at K1 $M = 5.28$ ($SD = 0.32$)

376    vs. 5.26 ($SD = 0.32$); 46.7% vs. 49.7% female; Raven's Coloured Progressive Matrices (Raven,

377    1947) raw scores $M = 15.91$ ($SD = 5.10$) vs. $M = 15.76$ ($SD = 5.23$); and average monthly

378    household income rating $M = 13.47$ ($SD = 5.85$) vs. $M = 12.01$ ($SD = 5.97$) on a scale of 19,

379     corresponding to an average range of S\$7000 – S\$7499 vs. S\$6500 – S\$6999.

380    **Measures**

381         **MACS Ratings.** There are 12 behavioral indicators initially, and 11 indicators included

382    in the final rating scales of MACS (see Figure 2). The final indicators were: 3 Input (teacher

383    input, language forms, concept development); 2 Output (positive climate, language

384    modelling/scaffolding); 4 Varied Strategies (ritualized language, adapted speech, use of gestures,

385    use of pictures/objects); and 2 Use of English (learning opportunities, behavior and task

386    management) indicators. The Input and Output domains were scored on a 1 to 7 scale, similar to

387    CLASS (Pianta et al., 2008). This was done to observe the level of complexity in teacher input

388    and strategies (low – medium – high). Indicators under these domains are high inference items,

389    and coders were instructed to gauge consistency and depth of teacher input and strategies across

390    the whole video segment. Given this element of judgement, the MACS adopted the within-one

391    reliability scoring approach (allowing coders to be one point off from each other or from the

392    consensus score), a procedure used in other coding systems (e.g., CLASS). The Varied Strategies

393    and English Use domains were scored on a scale of 1 to 3, as these were low inference items.

394    These domains constitute behavior frequency ratings that quantify the occurrence of the

395    described behaviors indicated at different levels (absent – minimal – regular). Since items in

396    these domains have to be definitively observed, coder scores had to be exact to achieve inter-

397    rater reliability. Scores per indicator were entered into the analyses. The mean ratings across all

398    classrooms for the dimensions of MACS are indicated in Table 2.

399          For the final version of the MACS (used in Phase 3), a training manual included the

400     following operationalized definitions per indicator, with examples from sample classes to

401     illustrate the target behaviors. Table 1 provides a summarized version of all MACS descriptors.

402          ***Input Strategies.*** This domain includes 3 aspects of teacher linguistic behaviors and

403     language instruction strategies constituting the MTL input children are exposed to.

404          *Teacher input.* Given the importance of teacher proficiency for providing quality

405     instruction and facilitating children's language acquisition, this indicator firstly observed

406     'correctness' of teachers' input in terms of pronunciation, vocabulary use, sentence formation

407     and grammar (Canh & Renandya, 2017; Richards, 2010). Additionally, this indicator considered

408     if input included 'consistently varied vocabulary and sentence structures', since richness of input

409     is important (Bornstein et al., 1998). Further, teacher input was coded for 'appropriateness' of

410     formal and informal register use in relation to the context. For instance, the formal register was

411     expected during reading and phonological instruction, while the informal register was expected

412     in social interactions to build rapport with children (Lakshmi & Saravanan, 2009; Ministry of

413     Education, 2015; Schiffman, 2004; Subhan, 2013; see Figure 3 for descriptors of 'appropriate'

414     language register use). The Teacher Input indicator was thus scored using the 1 to 7 scale to

415     differentiate between teachers who provided consistently incorrect, inappropriate, colloquial, and

416     repetitive MTL input at the low end of the scale, and those who consistently provided correct,

417     appropriate and varied MTL input at the high end.

418          *Language forms.* Instructional methods for developing children's phonological,

419     morphological, syntactic, semantic, pragmatic, and vocabulary knowledge (Ishwaran et al., 2005;

420     Sze & Leung, 2010) were rated on a scale geared towards capturing the different levels of

421     teacher facilitation strategies. In the Language Forms section in Table 1, teachers were scored

422    low when they either did not engage children in learning or discussing language forms at all, or

423    any requests for children to read or spell that were not scaffolded. Teacher scaffolding strategies

424    included linking letters to letter sounds, and breaking words into syllables to facilitate decoding.

425    Teachers were scored higher when they were consistently observed to provide scaffolding for

426    decoding and making meaning.

427          *Concept development.* While the local curriculum framework on MTL teaching and

428    learning places heavy emphasis on the role of MTL in providing cultural education (Ministry of

429    Education, 2013), we considered application of the MTLs across a wider range of topics that

430    could facilitate and encourage broader language use and understanding. This indicator thus

431    focused on teachers' MTL use for developing concepts and modelling thinking processes in

432    other relevant topics such as science, math, and art (e.g. Pianta et al., 2008). With this indicator,

433    teachers were scored low if classroom talk is teacher-dominated, with little to no opportunity for

434    children to express their own ideas and opinions. Use of only close-ended questions and

435    decontextualized discussion of topics were considered as limiting to children's opportunities to

436    speak, think, and understand topics discussed in the MTL. Teachers who employed these

437    strategies were thus scored lower. Conversely, teachers who endeavored to elicit and encourage

438    children's speaking, thinking and understanding of different topics discussed in MTL through,

439    for instance, brainstorming together, asking open-ended questions and linking topics to

440    children's daily life, were scored higher.

441          ***Output Strategies.*** This domain includes 2 indicators related to how teachers elicit

442    children's MTL use, and how effective the strategies were, as evidenced by children's responses.

443          *Positive Climate.* For this indicator, we considered if teachers were patient, encouraging,

444    and if they scaffolded children's MTL output (Pianta et al., 2008). Coders observed if teachers

445   were effective in making children comfortable and confident enough to use MTL, and to what

446   degree. Coders also considered if teachers were permissive rather than punitive with children's

447   English use, since children's language mixing is indicative of emerging language skills (Byers-

448   Heinlein & Lew-Williams, 2013; Weitz et al, 2010). Teachers who limited children's expression,

449   especially through ignoring or scolding children when they switch to English and provided no

450   help as children struggled to use MTL, were scored low. When teachers were patient,

451   incorporated and affirmed children's talk regardless of the language used, and provided help and

452   wait time to support children's MTL use, they were scored higher.

453         *Language Modelling and Scaffolding.* As with positive climate, the effectiveness of

454   teacher strategies for modelling and scaffolding language acquisition and production were

455   assessed through child output. Coders considered if these strategies could effectively elicit more

456   child output: open and closed-ended questions, repetitions and extensions of children's talk,

457   mapping language to actions, introduction and contextualization of new words and concepts

458   (Pianta et al., 2008). Teachers were scored for consistency and variety of different language

459   modelling and scaffolding strategies used. When more children were observed to use more MTL

460   in response to these strategies, the teachers were scored higher.

461         ***Varied Strategies***. Varied strategies included 4 indicators frequently discussed in existing

462   literature, and included in other observation tools: *ritualized language* (phrases, songs or rhymes

463   used for routine activities), *adapted speech* (child-directed talk in terms of tone, pitch or

464   repetition), use of *gestures and expressions*, and use of *pictures, objects, and realia* (Kersten et

465   al., 2010; Weitz et al., 2010). As in Table 1, the scale used indicated whether the varied

466   strategies were absent, minimal (1 to 2 observed instances) or regular (more than 2 observed

467   instances).

468          *Use of English*. Two forms of English use for supporting MTL teaching were rated for

469    this domain. First is the presence and frequency of teachers' English use for *learning*

470    *opportunities*. These 'learning opportunities' were considered not only when teachers provided

471    explanations or translations of target vocabulary and concepts in English, but also when teachers

472    translated instructions between MTL and English to facilitate task completion and behavior

473    management. In these situations, the English use was still considered an opportunity for learners

474    to be exposed to MTL input. Second, English used for *behavior and task management* was coded

475    when only English was used for the sole purpose of disciplining children and facilitating task

476    completion (Gort and Pontier, 2013). Here it was considered that children did not have the

477    'opportunity' to make connections between English and MTL as part of learning MTL. Scores

478    indicated the presence and amount of English use.

479          **CLASS Ratings.** Videos were previously coded using the Classroom Assessment

480    Scoring System (CLASS) for the SKIP study (Ng et al., 2014). All coders who used CLASS

481    were Teachstone certified at the time of coding. Scores for these videos, on the three CLASS

482    domains of emotional support, classroom organization and instructional support, were derived

483    according to the training manual (Pianta et al., 2008). The mean CLASS ratings across all

484    classrooms were 3.46 ($SD = 0.726$) for the ES, 4.51 ($SD = 0.640$) for the CO, and 2.25 ($SD =$

485    0.824) for the IS domains respectively.

486          **Child outcome measures.** Assessments of children's MTL vocabulary, reading, and

487    morphological awareness were conducted across the noted three time points: K1, K2, P1. Details

488    of the assessment measures are included in Appendix B. These data were used to identify

489    learning trajectories over time. Five children (0.01%) completed tasks in the MTL assigned in

490    school rather than that used at home.

491         ***Receptive vocabulary.*** The Bilingual Language Assessment Battery (BLAB; Rickard-

492     Liow & Sze, 2008) is a locally developed test, which has been widely used for vocabulary

493     assessment in Singapore. For each language, the test was administered on a tablet, where the

494     participant chooses one of four pictures corresponding to a word presented aurally in MTL

495     through the tablet. Children completed three practice trials with corrective feedback, then 80

496     trials for a score of total correct responses. Spearman-Brown split-half reliabilities at each time

497     point were 0.85, 0.78 and 0.80, respectively.

498         ***Reading.*** Tasks were developed per MTL script that included grapheme recognition

499     (discrimination), naming and word naming components (see Appendix B for details). For

500     Chinese, the total score was calculated as the number of total correct responses across the three

501     tasks (character discrimination, stroke naming, character naming). Spearman-Brown split half-

502     reliabilities at each time point were 0.95, 0.94, and 0.96, respectively. For Malay, the total score

503     was calculated as the number of total correct responses across two tasks (letter naming and word

504     reading). Spearman-Brown split-half reliabilities at each time point were 0.98, 0.98, and 0.92,

505     respectively. For Tamil, the total score was calculated as the number of total correct responses

506     across the three tasks (letter discrimination, letter naming, and word reading). Spearman-Brown

507     split-half reliabilities at each time point were 0.98, 0.98, and 0.98, respectively.

508         ***Morphological awareness.*** Tasks to evaluate children's understanding on the morphemic

509     structure of words were developed per MTL, and included two measures adapted from previous

510     research: Compound word production, and compound structure judgement (Tong, McBride-

511     Chang, Shu, & Wong, 2009; Chen, Hao, Geva, Zhu, & Shu, 2008) (see Appendix B for details).

512     A total morphological awareness score was calculated as the sum of compound production and

513     compound structure task scores. Spearman-Brown split half-reliabilities at each time point were

514     0.75, 0.60, 0.73, respectively.

515     **Procedures**

516     **Classroom observations and coding.**

517     MTL lessons were video-recorded for offline coding. Figure 1 describes the phases of

518     development for the MACS. Two rounds of coding were conducted in Phase 1. The second

519     author introduced the initial scheme to 2 coders per MTL who then each coded one pilot video

520     with the MACS, and noted potential issues for coding. This feedback contributed to minor

521     elaborations and edits, then these original six coders plus one Chinese language coder reviewed

522     and applied the scheme to two videos each. They attained an average inter-rater reliability

523     agreement of 80% across these 2 videos. They next coded a total of 6 Malay, 7 Tamil and 14

524     Chinese language sessions.

525     In Phase 2, MACS revisions incorporated the feedback from the two coding rounds in

526     Phase 1, plus consultation with local language curriculum and linguistic experts, as well as

527     transcriptions and preliminary analyses of 27 language classrooms. A manual with detailed

528     descriptors and exemplars from the classroom videos was produced as part of finalizing the

529     MACS. This manual was used to train 15 coders (5 per MTL) for coding in Phase 2. Training

530     and reliability testing included two days of reviewing the scoring manual, discussion of all

531     indicators and examples, and some coding practice. Coders then scored 5 classroom videos and

532     were required to obtain on average 80% or better inter-rater agreement with the consensus

533     scores. Disagreements in ratings were resolved with consensus discussions, and these discussions

534     also informed the revised coding scheme.

535     In Phase 3, additional coders were trained to the reliability rate of 80% or better across 5

536     training videos. A total of 6 Malay, 3 Tamil, and 9 Chinese coders viewed 51 classroom videos.

537     Of these, 48 videos were double coded by randomly assigning pairs of coders to score them. The

538     dual ratings on these 48 classrooms were entered into an inter-rater reliability analysis. Across

539     the videos, final ratings were reached by consensus discussion if inter-rater agreement was below

540     80% (e.g., coders compared notes and provided justifications for their scores). For double-scored

541     videos with inter-rating reliability above 80%, scores from the first coder were used. Where

542     needed, transcript and video data were double-checked to ascertain the score to be given. These

543     consensus scores were then included in the data used for the exploratory factor analyses. Scores

544     from all 51 coded classrooms were used for the analyses of child outcomes.

545     **Analysis Plan**

546     We examined the 4 research questions using as a rough guide the components of

547     validation approaches that have been utilized in other quality measurement studies (e.g., Lahti et

548     al., 2015; Zellman & Fiene, 2012; see also Halle, Whittaker & Anderson, 2010) (keeping in mind

549     that our approach involves a quality rating *tool* rather than a *system*). RQ 1 involved the key

550     concepts of the face or content validity of the MACS, and these are reviewed in the results. RQ 2

551     involved reliability of scoring the MACS. RQ 3 involved a latent structural analysis and

552     exploration of convergent and divergent validity with another coding measure. RQ 4 involved

553     predictive validity with learning outcomes.

554     **RQ 2.** Inter-rater reliability was calculated using the Kappa statistic, or the weighted

555     Kappa for the low inference items such as ritualized language, adapted speech, use of gestures,

556     use of pictures/objects and English use for learning opportunities and behavior management. The

557     Kappa statistic is conservative in that it corrects for "chance" agreement, with more fine-grained

558    scales (such as our 7-point scales) scores that are close but fall within one point of each other

559    were considered to be in disagreement. Therefore, Kappa values were recalculated using scores

560    of the first 5 items (teacher input, language forms, concept development, positive climate,

561    language modelling/scaffolding) counted as in agreement if they were within one point of each

562    other. The other 3-point scale items were not included in the within-one re-calculation.

563         **RQ 3.** Inter-item correlations were run for all items on the MACS, using the consensus

564    scores from multiple raters of the 51 classrooms. To further examine the structure of the coding

565    scheme, and with the objective of reducing the number of items, we took a factor analytic

566    approach. While confirmatory factor analysis would be an ideal method to check whether items

567    align according to the four domains of the MACS, we opted to use exploratory factor analysis to

568    observe data-driven relations in the case that they would not fit according to our theoretical

569    domains, and due to our limited sample size (for samples smaller than 50 cases, EFA could yield

570    reliable results with 3-4 factors and moderate (0.7) factor loadings for about 12 variables; de

571    Winter, Dodou & Wieringa, 2009). Following these expectations, we proceeded to conduct the

572    EFA using all 11 items with a weighted least squares (WLSMV) estimator and oblique rotation

573    (Mplus version 8, Muthén & Muthén, 2017). Data were entered as ordered categorical variables,

574    having been rated on a scale of "low" to "high" dimensions of correctness/sophistication

575    (described above where the high dimension was considered to be more aligned with the goals of

576    the kindergarten curriculum framework from the Ministry of Education). To determine the

577    optimal number of factors, a scree plot was first examined, along with indices for model fit ($\chi^2$ p-

578    value > 0.05; RMSEA < 0.05, CFI > 95; SRMR < 0.08) (Clark & Bowles, 2018).

579         **RQ 4.** We first conducted growth mixture modelling (GMM) to identify latent classes of

580    learners based on their growth in several MTL skills: receptive vocabulary, reading, and

581     morphological awareness. Separate GMMs were conducted for each language group using

582     available data from the total SKIP (Ng et al., 2014) sample of children using Mplus v.8 Software

583     with an MLR estimator (Muthén & Muthén, 2017). This approach was taken because the

584     measures between MT languages differed in terms of the number of items and difficulty of items,

585     such that the forms were not parallel. The outcomes of these GMMs (provided in Appendix C)

586     supported a 2-class solution for each model, with one class having steeper slopes (and higher

587     intercepts in most cases, for vocabulary and reading), and another class with shallower slopes

588     (and lower intercepts in most cases). We refer to these classes, respectively, as the high progress

589     and low progress learners, and the latent class assignment was attributed to each individual child

590     for the outcome measures. These categorical data were then used to determine the proportion of

591     high progress learners within each classroom, and this was correlated with the MACS classroom

592     scores.

<div align="center">

**Results**

</div>

594      **RQ 1. Content for classroom quality**

595          Between Phases 1 and 2 (Figure 1), the following sets of revisions were made between

596     the initial and final versions of the MACS.

597          **Teacher input indicator and informal language.** The teacher input indicator initially

598     required coders to consider only correctness and variety in teachers' classroom talk, with

599     colloquial forms included at the low end of the scale. The informal language indicator was

600     initially a separate indicator under the varied strategies domain. Observations from the first

601     round of coding found prevalent use of non-standard and informal language forms such as the

602     Johor-Riau variant of Malay and Standard Spoken Tamil, which made the teacher input indicator

603     difficult to score. Following consultation with language and curriculum experts from the

604    Ministry of Education and National Institute of Education, the objectives of MTL instruction

605    (Ministry of Education, 2013) were highlighted: to develop communicative skills, it is essential

606    for teachers to also expose children to spoken language forms. Thus, the scheme was modified to

607    include informal language use within appropriate classroom settings, such as in social talk for

608    teacher-child rapport building (see Figure 3).

609          **MTL specific coding descriptors.** The second major revision was to produce descriptors

610    specific to each MTL. Using the initial version, coders found the descriptors inadequate to

611    inform language specific factors. Hence, the teacher input indicator was adapted to incorporate

612    different registers and variants of the standard language for Tamil and Malay (Lakshmi &

613    Saravanan, 2009; Subhan, 2013), while Chinese grammar and vocabulary used by teachers hired

614    from outside Singapore was taken into account. Also, the language forms indicator added

615    descriptors that differentiated between modelling at the phoneme, alpha-syllable and syllable

616    level for Tamil and Malay language, versus discussion of strokes and stroke patterns for writing

617    characters and teaching meaning-based forms that aids character recognition for Chinese script.

618          **Reducing subjectivity.** Steps were taken to make descriptors more explicit, including the

619    MTL-specific descriptors. The first set of videos from K1 were transcribed, and a selection of

620    relevant examples were drawn out to illustrate the different descriptors and provide more detail

621    for coders. These transcriptions also served as a written account, reference check that facilitated

622    consensus-building for final ratings and for reliability training videos for new coders. To further

623    reduce subjectivity, quantitative guidelines were included where possible. For instance, under the

624    positive climate indicator, a medium score meant at least ¼ of children were comfortable

625    responding, while at least ⅓ were required for a high score. Similarly, the proportion of children

626    responding was used in the language modelling and scaffolding ratings. Behavior frequency

627    ratings were applied to the domains of varied strategies and English use, as 'absent' (0),

628    'minimal' (1-2), or 'regular' (> 2 occurrences).

629    **RQ 2. Rater Reliability**

630    Interrater reliability on the 48 double-coded classrooms showed overall acceptable Kappa

631    values using an exact estimate for Malay and Tamil ratings, but this was low for the Chinese

632    language group:  Chinese = 0.34, Malay = 0.60, Tamil = 0.60. Using within-one estimates

633    yielded generally higher Kappa values, which were within acceptable ranges for all languages:

634    Chinese = 0.75, Malay = 0.82, and Tamil = 0.72. Estimates per indicator are provided in Table 1.

635    **RQ 3. Construct Validity: Dimensional Structure of the MACS**

636    **Cross correlation between MACS indicators.** Relations between items within each of

637    the 4 MACS domains are shown in the shaded portions in Table 3. These inter-domain relations

638    were expected to be stronger than cross-domain relations, outside of the shaded areas in the

639    table. Moderate significant inter-domain correlations were found for indicators of Input (1-3),

640    showing that teacher's input quality covaried with both language forms and concept development

641    emphases. However, these indicators of input also showed strong correlations with those of other

642    domains, especially modelling/scaffolding. The two indicators for Output (4, 5) were

643    significantly correlated with each other, indicating that teachers who use one strategy tended to

644    use the other as well. As noted, these indicators also correlated with the input indicators. The

645    varied strategies domain indicators (6-9) were less coherent, especially for pictures/objects,

646    which was only correlated with concept development from Input. Significant moderate inter-

647    domain correlations here were found only for gesture with ritualized language, and with adapted

648    speech strategies. Adapted speech was also related to cross-domain indicators of input and

649    output. The final domain of English use showed no intra-domain relations between the two

650    indicators (10, 11), suggesting teachers selectively used English for different purposes.

651         **Factor Analysis of MACS ratings.** In the initial model with all items included, the

652    *Strategic use of gestures/expression* was highly correlated with and not statistically

653    distinguishable from two other variables, and the model was then re-run without this variable.

654    The EFA supported a three-factor structure for the MACS data, according to a scree plot and fit

655    indices:  $\chi^2$ (18) = 14.30, $p$ = 0.71, RMSEA = 0 (CI = 0.00-0.096), SRMR = 0.059, CFI = 1.00.

656    This model yielded a better fit than a 2-factor model, $\chi^2$ (26) = 23.51, $p$ = 0.60, RMSEA = 0 (CI

657    = 0.00-0.098), SRMR = 0.081, CFI=1.00, or a 1-factor model $\chi^2$ (35) = 40.29, $p$ = 0.25, RMSEA

658    = 0.054 (CI = 0.00-0.119), SRMR = 0.119, CFI= 0.988. Factor loadings from the final 3-factor

659    model are shown in Table 4. 'Input' and 'output' categories converged onto one factor along

660    with 'varied strategies' of ritualized language and adapted speech (F1), while the 'English use'

661    variables split into separate factors (F2, F3) in agreement with the intra-domain correlation noted

662    above (see Table 2). Instructional English use was coupled with teachers' lesser focus on

663    language forms, while task management English use was related to more of this focus plus

664    teachers' use of pictures/objects/realia (see Table 4).

665         **Relation between rating schemes.** Correlational analysis was run between the MACS

666    Factor Scores from the EFA and the CLASS domain scores. As shown in Table 5, there were no

667    significant relations between the three MACS factors and the three CLASS domains. This was in

668    spite of the similar indicators that were adapted from the CLASS; the instructional support

669    dimensions of language modelling and concept development, and the emotional support

670    dimension of positive climate.

671    **RQ 4. Predictive Validity: Link to Child Learning Outcomes**

672     Correlations were run between teachers' MACS factor scores and their proportion of high

673     progress learners within the classroom. In Table 5, it can be seen that MACS Factor 1 scores

674     were significantly and positively correlated with a greater proportion of high progress learners in

675     terms of MTL vocabulary, but significantly and negatively correlated with the proportion of high

676     progress learners for MTL reading. Scatterplots of these significant MACS factor score relations

677     to high progress learners are shown in Figure 4. Note that 1 case in vocabulary and 2 cases in

678     reading had significant Cook's distance and these classes were omitted. None of the other

679     relations were significant (all $p$'s > 0.1). For comparison, we also examined the correlations of

680     teachers' CLASS ratings with their proportion of high progress learners, and none of these

681     relations were significant (all $p$'s > 0.1), as shown in Table 5.

682                                                        **Discussion**

683     Given the pivotal role that early childhood education plays in multilingual learning, there

684     is a clear need to identify key aspects of quality language instruction. Available ECE quality

685     rating systems may not be sensitive to the specific issues related to bilingual learning across

686     different socio-cultural contexts such as Singapore. The current study describes the development

687     process for a new ECE rating tool to address this gap and to highlight characteristics of quality

688     ECE for MTL learning in an increasingly English-dominant multilingual society. The overriding

689     purpose for this tool is to inform research and practice in early bilingual development and

690     education. The MACS is not a teacher evaluation for high-stakes purposes. Rather, the MACS is

691     a first step toward understanding the features of ECE quality for multilingual outcomes, with

692     longer-term implications for bilingual education policy.

693     As a first step to identifying characteristics of MTL ECE quality, the MACS

694     observational tool is intended to be comprehensive. Multiple sources contributed to the content

695    of the coding scheme, which resulted in a four-domain conceptualization of quality that we felt

696    encompassed important aspects for MTL learning in the preschool setting: teachers' (1) input

697    strategies, (2) output strategies for eliciting child talk (3) other relevant varied strategies for early

698    language learning in bilingual contexts, and (4) use of English for classroom learning and

699    management purposes. Given the well-founded emphasis on comprehensible input for language

700    learning, in the current context we needed to account for teachers' use of different registers for

701    different purposes – as more instructional or more conversational talk (see Figure 3). Our input

702    domain also included indicators related to the purpose of the language lesson – as focused on

703    language forms, or concept development for topics such as science or arts (Pianta et al., 2008).

704    Emphasis in the literature on output for language learning also influenced the inclusion of

705    positive climate and language modelling and scaffolding strategies for increasing child output.

706    The recognition that children need to feel secure and supported in their language production is

707    established in preschool studies (Schwartz & Gorbatt, 2010; Weitz et al., 2010), and scaffolding

708    their efforts in language output is important to bilingual learning (Bohman, Bedore, Pena,

709    Mendez-Perez & Gillam, 2010). For each of these indicators, we arranged the examples of high

710    ratings to reflect the literature in terms of beneficial strategies for teaching and learning, while

711    also aligning them with the educational goals of the Ministry of Education's framework.

712           Rounding out the content of the MACS from these two key components of input and

713    output, we added domains with pedagogical indicators specific to the preschool and local

714    context. Strategies to increase comprehensible input for young children, who vary in MTL

715    proficiency were observed in teachers' simplified and child adapted speech, nonverbal gestures,

716    and English (Weitz et al., 2010; Atkins-Burnett et al., 2010). The literature is equivocal on which

717     strategies are more beneficial, such as the use of language switching, so we rated the frequency

718     with which such observed behaviors were used for different purposes.

719             To test and refine this conceptualization of MTL classroom quality, we then applied the

720     MACS to a set of video-recorded MTL kindergarten classroom sessions. After addressing initial

721     challenges with coding reliability, by detailing a coding manual with language-specific

722     illustrative examples and by modifying rating scales to include frequency ratings, overall

723     interrater agreement was within acceptable limits for within-one estimates (although there is

724     recent concern surrounding the "within one" approach, Mashburn, 2017). Reliability per

725     indicator varied, however, with better reliability on the frequency-based ratings (absent to

726     regular) (although some were still low, Kappa < 0.6), and improved reliability on the finer-

727     grained low-high scale items for within-one estimates.

728             Using consensus scores, we then examined the domain structure for the MACS with

729     correlation and factor analyses. The correlations between items did not follow the expected

730     pattern, whereby indicators within a domain would be more highly correlated than between

731     domains. Indicators across input and output domains were moderately related, while indicators

732     within the English use domain were not related to each other. Results from the exploratory factor

733     analysis coincide with this pattern of results, showing moderate to high factor loadings of the

734     input and output indicators on the first factor (F1), along with the adapted speech strategy. The

735     pattern of F1 factor loadings is suggestive that both comprehensible input (through teachers'

736     rich, accurate, and adapted MT speech and language modelling), and elicited child speech

737     (through a positive climate and scaffolding), contribute to quality language instruction. That both

738     of these domains load onto a single factor might also indicate the importance of an input-output

739     cycle for language learning and eliciting speech production (e.g., Gass & Mackey, 2007;

740    Pearson, 2007). The English use indicators loaded on separate factors (F2, F3), on the other

741    hand. English for learning opportunities loaded onto F2 with a negative contribution of language

742    forms strategies, which may be due to a tendency to emphasize oral language more than literary

743    forms and thereby using English translation to teach new vocabulary, as observed in some

744    classrooms. Thus, a re-organization of the MACS domain scores may be in order. These

745    preliminary results require replication with a larger sample of teachers, however.

746         Comparison of general classroom quality and our measure of MTL classroom quality

747    showed no overlap in these scores. The CLASS teacher ratings, which were completed on the

748    same classrooms as part of the larger study, did not correlate with the factor scores from the

749    MACS. This was in spite of the adapted indicators from the CLASS IS dimension (language

750    modelling and concept development) and ES (positive climate) into the MACS, and indicates

751    that our measure is unique from more broadly defined teacher-child interaction process quality.

752         Finally, we examined how the MACS ratings of quality for MTL learning in the

753    classroom may be related to the learning outcomes of students in those classrooms. We

754    examined several measures of MTL learning over time, including high progress (relative to

755    peers) in receptive vocabulary, reading, and morphological awareness. Teacher ratings on F1

756    (input and output) were correlated with having more high progress learners of MT vocabulary,

757    but negatively correlated with these types of learners for MT reading. This may follow from the

758    emphasis on oral language development in the kindergarten framework for MTL (Ministry of

759    Education, 2013) and the stronger teacher input and child elicited speech scores in some

760    classrooms may disproportionally affect oral vocabulary as compared with reading of scripts.

761    The majority of teachers (85%) were rated in the medium range for teacher input a (Table 1),

762    whereas there was a broader range of ratings across the low to medium ranges for language

763    forms (with 43% in the low range). Extending reading and writing lessons to a context of making

764    meaning with print (high rating) rarely occurred, in only 2% of teachers. Thus, explicit linkages

765    between oral and written language may not be made at this age – meaning that a focus on oral

766    language vs. written language may be more exclusive. On the other hand, children's scores on

767    language and literacy were positively related, consistent with the extant literature, so the

768    converse effects might instead be related to the MACS scales, which centered on spoken

769    language indicators to a greater degree. Teacher scores on the other factors, and on the CLASS

770    domains were unrelated to the proportion of high progress learners. Though preliminary, these

771    results suggest that the MACS ratings may capture elements of language classroom quality that

772    are important for children's MTL outcomes.

773    **Implications**

774           We identified teachers' language and elicitation of child speech as key components for

775    ECE quality related to MTL learning. Strategies of language modelling, scaffolding and adapted

776    speech, as well as using the MTL for concept development and creating a positive classroom

777    climate, appear to have positive effects on children's vocabulary progress. Scaffolding may be

778    especially important in the context of diverse language learner competencies and in the face of

779    shifting home language trends. On the other hand, a range of teachers' English use was observed

780    in the present classrooms, and for different purposes. These results do not inform the question

781    about language switching benefits, but may recommend training teachers to be aware of their

782    language quality and purpose for English use.

783    **Limitations**

784           The findings need to be considered in light of several limitations. First, the sample of

785    classrooms was relatively small and at the limits for conducting the analyses we wished to

786    enlighten about the MACS' latent structure. Also, Tamil, and Malay teachers were less

787    represented, but showed that some quality indicators needed adjustment to account for language-

788    specific issues such as diglossia use. Second, the data sampling method followed standard

789    protocols for observation schemes, with roughly 20-minute segments from a "typical day" in the

790    classroom. Other approaches may yield behavior stability over observations or at finer time

791    scales. Third, while interrater reliability improved over iterations of the coding scheme, many of

792    the indicators yielded poor agreement unless a with-one estimate was used for the 7-point items,

793    suggesting more master-code feedback may be required during training. Fourth, the evaluation

794    regarding student outcomes was limited by a lack of standardized assessments, so learners were

795    categorized as high- or low-progress relative to peers. However, in some cases the groups

796    differed in incoming scores (e.g., reading) and we cannot rule out extra-classroom variables on

797    their outcomes, which would require further study.

798                                                          **Conclusion**

799            The MACS was developed to apply generally across bilingual preschool contexts – as

800    compelled by the local Singapore context, wherein multiple Mother Tongue languages are

801    incorporated into ECE. Although we aimed for a generalizable coding scheme, this was balanced

802    by the need for language-specific aspects of the rating tool. By providing a detailed account of

803    the conceptualization and development phases of the MACS, we hope to provide a

804    demonstration of how rating tools can be adapted to assess and describe quality language

805    teaching in bilingual ECE contexts. The current intention was using this tool to understand the

806    constituents of ECE quality as related to child outcomes within multilingual contexts. Further

807    research applying the rating tool to larger, representative samples of teachers in classrooms may

808    contribute to bilingual education policy and teacher professional development.

**References**

809

810   Atkins-Burnett, S., Sprachman, S., & Caspe, M. (2010). *Language Interaction Snapshot + End*

811        *of Visit Ratings (LISn + EVR)*. Princeton, NJ: Mathematica Policy Research.

812   Bohman, T. M., Bedore, L. M., Pena, E. D., Mendez-Perez, A., & Gillam, R. B. (2010). What

813        you hear and what you say: Language performance in Spanish-English

814        bilinguals. *International Journal of Bilingual Education and Bilingualism, 13*(3), 325-

815        344.

816   Bolton, K., & Ng, B. C. (2014). The dynamics of multilingualism in contemporary Singapore.

817        *World Englishes, 33*(3), 307-318.

818   Bornstein, M. H., Haynes, M. O., & Painter, K. M. (1998). Sources of child vocabulary

819        competence: A multivariate model. *Journal of Child Language*, *25*(2), 367-393.

820   Burchinal, M. (2018). Measuring Early Care and Education Quality. *Child Development*

821        *Perspectives, 12*(1), 3-9.

822   Burchinal, M., Magnuson, K., Powell, D., & Hong, S. S. (2015). Early child care and education

823        and child development. In M. Bornstein, R. Lerner & T. Leventhal (Eds.), *Handbook of*

824        *child psychology and developmental science* (Vol. 4, 7th ed., pp. 223–267). Hoboken, NJ:

825        Wiley.

826   Byers-Heinlein, K., & Lew-Williams, C. (2013). Bilingualism in the early years: What the

827        science says. *Learning Landscapes, 7*(1)*,* 95-112.

828   Canh, L.V., & Renandya, W. (2017). Teachers' English proficiency and classroom language use:

829        A conversation analysis study. *RELC Journal*, *48*(1), 67-81.

830     Chen, X., Hao, M., Geva, E., Zhu, J., & Shu, H. (2008). The role of compound awareness in

831          Chinese children's vocabulary acquisition and character reading. *Reading and Writing: An*

832          *Interdisciplinary Journal, 21*, 559–586.

833     Cheng, R. L. (1985). A comparison of Taiwanese, Taiwan Mandarin, and Peking Mandarin.

834          *Language*, *61*(21), 352-377.

835     Chimbutane, F. (2013). Codeswitching in L1 and L2 learning contexts: Insights from a study of

836          teacher beliefs and practices in Mozambican bilingual education programmes. *Language*

837          *and Education*, *27*(4), 314-328. doi: 10.1080/09500782.2013.788022

838     Clark, D.A. & Bowles, R.P. (2018). Model fit and item factor analysis: Overfactoring,

839          underfactoring, and a program to guide interpretation. *Multivariate Behavioral Research*,

840          *53*(4), 544–558.

841     de Bot, K. (1996). The psycholinguistics of the output hypothesis. *Language Learning*, *46*(3),

842          529-555. doi: 10.1111/j.1467-1770.1996.tb01246.x.

843     de Houwer, A. (2005). Early bilingual acquisition. In J. F. Kroll & A. M. B. De Groot (Eds.),

844          *Handbook of bilingualism: Psycholinguistic approaches* (pp. 30-48). New York: Oxford

845          University Press.

846     DePalma, R. (2010). Language Use in the Two-way Classroom: Lessons from a Spanish-English

847          Bilingual Kindergarten. Clevedon, UK: Multilingual Matters.

848     de Winter, J. C. F., Dodou, D., & Wieringa, P. A. (2009). Exploratory factor analysis with small

849          sample sizes. *Multivariate Behavioral Research*, *44*(2), 147-181.

850     Dixon, L.Q. (2009). Assumptions behind Singapore's language-in-education policy: Implications

851          for language planning and second language acquisition. *Language Policy, 8*(2)*,* 117-137.

852    Early Childhood Development Agency. "Requirements for teacher certification". ECDA.gov.sg.

853         Published November 2018. Accessed 12 December 2018.

854         https://www.ecda.gov.sg/Documents/Requirements%20for%20Teacher%20Certification.

855         pdf

856    Ellis, R., & Shintani, N. (2014). *Exploring language pedagogy through second language*

857         *acquisition research*. London: Routledge.

858    Enama, P. (2016). The impact of English-only and bilingual approaches to EFL Instruction on

859         low-achieving bilinguals in Cameroon: An empirical study. *Journal of Language*

860         *Teaching and Research*, *7*(1), 19-30. doi: 10.17507/jltr.0701.03.

861    Gass, S. M., & Mackey, A. (2007). Input, interaction, and output in second language acquisition.

862         In B. Van Patten & J. Williams (Eds.), *Theories in second language acquisition: An*

863         *introduction* (pp. 175–199). New York/London: Routledge.

864    Gibbons, P. (2006). *Bridging discourses in the ESL classroom: Students, teachers and*

865         *researchers*. London: Continuum Books.

866    Gort, M., & Pontier, R. (2013). Exploring bilingual pedagogies in dual language preschool

867         classrooms. *Language and Education*, *27*(3), 223-245.

868    Halle, T., Whittaker, J. E. V., & Anderson, R. (2010). *Quality in Early Childhood Care and*

869         *Education Settings: A Compendium of Measures, Second Edition.* Washington, DC: Child

870         Trends. Prepared by Child Trends for the Office of Planning, Research and Evaluation,

871         Administration for Children and Families, U.S. Department of Health and Human

872         Services.

873    Hoff, E., & Core, C. (2013). Input and language development in bilingually developing children.

874         *Seminars in speech and language*, *34*(4), 215–226.

875    Huang, B. H., & Kuo, L. (2020). The role of input in bilingual children's language and literacy

876          development: Introduction to the special issue. *International Journal of Bilingualism,*

877          *24*(1), 3-7.

878    Ibbotson, P. (2013). The scope of usage-based theory. *Frontiers in Psychology*, *4*, 1-15.

879    Ishwaran, S., Shanmugam, K., Varaprasad, N., Sankaran, C., Lakshimi, S., Saravanan, V., &

880          Peng, H. (2005). *Report of the Tamil Language Curriculum and Pedagogy Review*

881          *Committee* (Rep.) Republic of Singapore: Ministry of Education.

882    Jiang, Y., García, G., & Willis, A. (2014). Code-Mixing as a bilingual instructional

883           strategy. *Bilingual Research Journal*, *37*(3), 311-326. doi:

884            10.1080/15235882.2014.963738.

885    Kersten, K., Steinlen, A. K., Tiefentahl, C., Wipperman, I., & Mattsson, A. F. (2010). Guidelines

886          for language use in bilingual preschools. In K. Kersten, A. Rohde, C. Schelletter & A.

887          K. Steinlen (Eds.), *Bilingual preschools: Best practices* (pp. 103-116). Trier:

888          WVT Wissenschaftlicher Verlag Trier.

889    Krashen, S. D. (1985). *The input hypothesis: Issues and implications*. London: Longman.

890    Lahey, M. (1988). *Language disorders and language development*. London, UK: Macmillan.

891    Lahti, M., Elicker, J., Zellman, G., & Fiene, R. (2015). Approaches to validating child care

892          quality rating and improvement systems (QRIS): Results from two states with similar

893          QRIS type designs. *Early Childhood Research Quarterly*, *30*, 280-290.

894    Lakshmi, S., & Saravanan, V. (2009). *An examination of the use of Standard Spoken Tamil in*

895          *Singapore – in the school and media domains in Tamil classrooms in order to establish*

896          *SST as and additional resource for the teaching and learning of Tamil.* Final Research

897  Report for Project no. CRP 6/04 SL & CRP 10/06 SL. Centre for Research in Pedagogy &

898  Practice, National Institute of Education, Nanyang Technological University, Singapore.

899  La Paro, K. M., Pianta, R. C., & Stuhlman, M. (2004). The Classroom Assessment Scoring

900  System: findings from the prekindergarten year. *The Elementary School Journal*, *104*(5),

901  409-426.

902  Lin, A. M. Y. (2005). Critical, transdisciplinary perspectives on language-in-education policy

903  and practice in postcolonial contexts: The case of Hong Kong. In A. M. Lin & P. W.

904  Martin (Eds.), *Decolonisation, globalisation: Language-in-education policy*

905  *and practice* (pp. 38–54). Clevedon, UK: Multilingual Matters.

906  Mashburn, A. J. (2017). Evaluating the validity of classroom observations in the Head Start

907  Designation Renewal System. *Educational Psychologist, 52*(1), 38-49.

908  Ministry of Education. (2013). *Nurturing early learners – A curriculum for kindergartens in*

909  *Singapore. Framework for Mother Tongue Languages.* Republic of Singapore.

910  Ministry of Education. (2015). Memupuk pelajar pada peringkat awal – Kurikulum untuk tadika

911  di Singapura. Panduan pendidik Bahasa Melayu untuk prasekolah. Republic of Singapore.

912  Mukhlis, A. B., & Pang, E. (2015). Learning to be biliterate in English and Malay using dual-

913  language books. (NIE Research Brief Series No. 15-013). Singapore: National Institute of

914  Education.

915  Muthén, L. K., & Muthén, B. O. (1998-2017). *Mplus user's guide.* Eighth edition. Los Angeles,

916  CA: Muthén & Muthén.

917  Ng, E. L., O'Brien, B. A., Khng, K. H., Poon, K. L. K., Karuppiah, N., Bull, R., Pang, E., Lee,

918  K., Hwee, L. M., Tan, C.T., & Tan, G. H. (2014). *Singapore Kindergarten Impact Project*

919 *(SKIP).* OER 09/14RB, Office of Education Research, National Institute of Education,

920 Singapore.

921 Pakarinen, E., Lerkkanen, M., Pikkeus, A., Kiuru, N., Siekkinen, M., Rasku-Puttonen, H., &

922 Nurmi, J. (2010). A validation of the class assessment scoring system in Finnish

923 Kindergartens. *Early Education and Development, 21*(1), 95-124. Doi:

924 10.1080/10409280902858764

925 Palviainen, A., Protassova, E., Mård-Miettinen, K., Schwartz, M. (2016). Two languages in the

926 air: A cross-cultural comparison of preschool teachers' reflections on their flexible

927 bilingual practices. *International Journal of Bilingual Education and Bilingualism, 19*(6),

928 614-630.

929 Pianta, R., La Paro, K., & Hamre, B. (2008). *Classroom Assessment Scoring System manual:*

930 *Pre-K.* Baltimore, MD: Paul H. Brookes Pub. Co.

931 Pearson, B. Z. (2007). Social factors in childhood bilingualism in the United States. *Applied*

932 *Psycholinguistics, 28*, 399-410.

933 Place, S., & Hoff, E. (2011). Properties of dual language exposure that influence 2-year-olds'

934 bilingual proficiency. *Child Development, 82*(6), 1834-1849.

935 Probyn, M. (2015). Pedagogical translanguaging: Bridging discourses in South African science

936 classrooms. *Language and Education*, *29*(3), 218-234. doi:

937 10.1080/09500782.2014.994525.

938 Raven, J. C. (1947). *Progressive matrices. Set A, Ab, B, book form*. London: H.K. Lewis.

939 Rickard-Liow, S. J., & Sze, W. P. (2008). *Bilingual Language Assessment Battery (BLAB).* In D.

940 o. P. Singapore: Psycholinguistics Lab (Ed.): National University of Singapore.

941     Richards, J. C. (2010). Competence and performance in language teaching. *RELC*

942          *Journal*, *41*(2), 101-122.

943     Schiffman, H. F. (2004). The Tamil case system. *South Indian horizons: felicitation volume*

944          *for  Francois Gros on the occasion of his 70th birthday*, 293-322.

945     Schwartz, M. (2018). Preschool bilingual education: Agency in interactions between children,

946          teachers and parents. In M. Schwartz (Ed.), *Preschool Bilingual Education* (pp. 1-24).

947          Cham, Switzerland: Springer.

948     Schwartz, M., & Gorbatt, N. (2016). 'Why do we know Hebrew and they do not know Arabic?'

949          Children's meta-linguistic talk in bilingual preschool. *International Journal of Bilingual*

950          *Education and Bilingualism*, *19*, 1–21.

951     Singapore Department of Statistics, Ministry of Trade and Industry. (2018). Population Trends,

952          2018. Retrieved from: https://www.singstat.gov.sg/-

953          /media/files/publications/population/population2018.pdf

954     Subhan, M. A. (2013). *Bilingualism and its effects on Malay language planning. (Unpublished*

955          *PhD thesis).* National Institute of Education Nanyang Technological University,

956          Singapore.

957     Sun, H., Yin, B., Amsah, N. F. B. B., & O'Brien, B. A. (2018). Differential effects of internal

958          and external factors in early bilingual vocabulary learning: The case of

959          Singapore. *Applied Psycholinguistics*, *39*(2), 383-411.

960     Swain, M. (1985). Communicative competence: Some roles of comprehensible input and

961          comprehensible output in its development. In S. M. Gass & C. M. Madden (Eds.), *Input in*

962          *Second Language Acquisition* (pp. 235-253). Rowley, MA: Newbury House.

963    Swain, M. (2005). The output hypothesis: Theory and research. In E. Hinkel (Ed.), *Handbook of*

964         *research in second language teaching and learning* (pp. 471-483). Mahwah, NJ:

965         Lawrence Erlbaum Associates.

966    Sze, P., & Leung, F. F. (2010). Enhancing learners' metalinguistic awareness of language form:

967         The use of eTutor resources. *Assessment and Learning, 3,* 79-96.

968    Teng, A. (2018). Pre-schools turn to native speakers to meet demand. *The Straits*

969         *Times.* Retrieved from: https://www.straitstimes.com/singapore/education/pre-schools-

970         turn-to-native-speakers-to-meet-demand .

971    Tong, X., McBride-Chang, C., Shu, H., & Wong, A. M. Y. (2009). Morphological awareness,

972         orthographic knowledge, and spelling errors: Keys to understanding early Chinese literacy

973         acquisition. *Scientific Studies of Reading, 13*(5), 426-452.

974    Unsworth, S., Persson, L., Prins, T., & de Bot, K. (2015). An investigation of factors affecting

975         early foreign language learning in the Netherlands. *Applied Linguistics, 36*(5), 527-548.

976    Vygotsky, L. S. (1978). *Mind in society: The development of higher psychological processes.*

977         Cambridge, MA: Harvard University Press.

978    Weitz, M., Pahl, S., Mattsson, A. F., Buyl, A., & Kalbe, E. (2010). The input quality observation

979         scheme (IQOS): The nature of L2 input and its influence on L2 development in bilingual

980         preschools. In K. Kersten, A. Rohde, C. Schelleter & A. K. Steinlen, (Eds.) *Bilingual*

981         *preschools volume 1: Learning and development* (pp.5-44). Trier:

982         WVT Wissenschaftlicher Verlag Trier.

983    Wesche, M. (1994). Input and interaction in second language acquisition. In C. Gallaway & B.

984         Richards (Eds.), *Input and interaction in language acquisition* (pp. 219-249). Cambridge:

985         Cambridge University Press.

986    Wong, D. (2018). Scaling up supply of Mother Tongue pre-school teachers. *The Straits*

987        *Times.* Retrieved from: https://www.channelnewsasia.com/news/singapore/scaling-up-

988        supply-of-mother-tongue-pre-school-teachers-9880148.

989    Wright, W. E., Boun, S., & Garcia, O. (2015). Introduction: Key Concepts and issues in bilingual

990        and multilingual education. In W. E. Wright, S. Boun & O. Garcia (Eds.), *Handbook of*

991        *Bilingual and Multilingual Education*. Malden, MA: Wiley-Blackwell.

992    Yussof, N. T., & Sun, H. (2020). Mismatches between teacher perceptions, practices and reasons

993        for English use in preschool Malay language classrooms. *Language & Education*. doi:

994        10.1080/09500782.2020.1720230.

995    Zellman, G. L., & Fiene, R. (2012). *Validation of quality rating and improvement systems for*

996        *early care and education and school-age care* (Research-to-Policy, Research-to-Practice

997        Brief OPRE 2012-29). Washington, DC: Office of Planning, Research and Evaluation,

998        Administration for Children and Families, U.S. Department of Health and Human

999        Services.

1000

Appendix A

*Comparison of a set of selected observational rating tools and their contributions to the development of the Mother Tongue Adapted*
*Coding Scheme (MACS)*

| Observation Tool | Aligned features that were adopted | Adaptations | Limitations |
|---|---|---|---|
| Classroom Assessment Scoring System (CLASS; Pianta, La Paro & Hamre, 2008) | • Application to preschool context<br>• Relevant dimensions incorporated into MACS:<br>-Concept development, positive climate, instructional learning formats, language modelling, quality of feedback<br>• Scoring system:<br>- 1-7 scale<br>- Including depth, duration, and frequency of observed behavior<br>-'Within one' reliability scoring | • Relevant dimensions included in different MACS domains (input, output, strategies)<br>• Positive climate made specific to language learning<br>• Modelling and scaffolding strategies modified to be appropriate for second or MT language learners<br>• Modelling/scaffolding strategies extended to include effectiveness in terms of elicited child talk<br>• Overall productive learning activities replaced with specific, independently scored varied teaching strategies<br>• Low inference items scored on 1-3 scale of frequency | • A more holistic observation tool for preschool classroom experience, not specific to language teaching and learning<br>• Behaviors related to bilingual contexts, such as language switching, not included |
| Input Quality Observation Scheme (IQOS; Weitz et al., 2010) | • Application to bilingual preschool context<br>• Subset of appropriate domains incorporated into MACS (input, output, strategies)<br>• Specific strategies for second language teaching included in MACS: | • Total number of domains reduced<br>• Different scoring scales used for low and high inference items<br>• Different observation and scoring procedures | • Emphasizes a target language-only input<br>• Consideration of language switching as negative instance<br>• Does not consider both negative and positive instances for use of children's stronger language |

| | | | |
|---|---|---|---|
| | -Ritualized language; facial expression, gestures; pictures, objects, realia; adapted speech; variety & richness of input | | |
| Language Interaction Snapshot (LISn; Atkins-Burnett, Sprachman, & Caspe, 2010) | • Application to bilingual preschool context with children from differing home language proficiency levels<br>• Subset of domains incorporated into MACS | • Number of domains reduced<br>• Types of talk and activity types used as examples to observe for input and output | • Observation focus on individual child 'snapshots' rather than classroom level<br>• More detailed set of types of talk by children and which language used for each<br>• Language switching scored at a more detailed scale of types of talk |

| Additions unique to MACS | |
|---|---|
| Additional aspects | Purpose for the additions |
| Language correctness<br><br>Gauging quality through correctness of pronunciation, language choice, syntax etc. | Correct and fluent input considered key for language teaching and learning. Considered impact of colloquial language use (Singlish) that may impact teacher input. |
| Language appropriateness<br><br>Addressing concerns regarding diglossia in the classroom<br>• Formal language for literacy instruction<br>• Informal language for developing communicative skills | Considered impact of diglossia in local languages that may impact children's language resources, and what kinds of language children need to be exposed to for school preparation, but also to facilitate their communicative skills within the relevant language communities |
| Teachers' English use<br><br>Attempts at optimizing English use in MTL classrooms | Given the dominance of English as the societal language, and studies that reflect how target-language only classrooms are rare, instructional purposes for English use were observed |
| Two scoring scales<br><br>Two different scales were used to facilitate scoring of low and high inference items | Different evaluation approaches different types of variables:<br>• 1-3 point scale for observing behavioural frequency<br>• 1-7 point scale for observing the level of teacher facilitation |

Appendix B

*Details and sensitivity of child mother tongue language (MTL) outcome measures.*

**Vocabulary.** The locally developed measure was sensitive to change over time,

according to the mean and range of scores for the whole sample: Means for Chinese at K1 =

32.3, K2 = 38.6, and P1 = 43.6; for Malay at K1 = 30.8, K2 = 37.7, and P1 = 41.9; and for Tamil

at K1 = 25.6, K2= 31.1, and P1 = 35.9. There was also no indication of floor or ceiling effects.

Minimum-maximum scores at each time point = 13-60, 17-66, 17-71 for Chinese; 13-54, 19-55,

17-57 for Malay; and 12-47, 12-47,18-61 for Tamil. The $25^{th}$-$75^{th}$ percentile scores at each time

point = 26-38, 32-45, 37-51 for Chinese;  25-37, 32-43, 37-48 for Malay; and 21-30, 25-37, 29-

43 for Tamil.

**Reading**. For Chinese, grapheme recognition was based on the character discrimination

test of the Preschool and Primary Chinese Literacy Scale (Li, H. (2015). *Teaching Chinese*

*literacy in the early years*. New York, NY: Routledge) and included 19 items where the child

had to match one of four characters to the spoken word. Children then completed a naming task

where they gave the names of 15 strokes and radicals, and a character reading task, where they

read up to 100 characters. The word list was administered in blocks of 25 characters, with words

decreasing in frequency over the blocks (according to the corpus of Loo, 1989). The task was

discontinued when the child could not read any words in a block. Total correct responses were

summed over the three tasks.

For Malay, children completed a letter naming task with 11 items from the Latin script,

and a word reading task where they read up to 100 words. The word list was administered in

blocks of 25, with words decreasing in frequency over the blocks according to the corpus of Lee

and Low (Lee, L. W., & Low, H. M. (2011). Developing an online Malay language word corpus

for primary schools. *International Journal of Education and Development using Information and*

1029    *Communication Technology, 7*(3), 96-101), and with increasing word complexity within each

1030    block (in terms of word length and syllable structure). The task was discontinued when the child

1031    was not able to read any words in a block. Total correct responses were summed over the two

1032    tasks.

1033         For Tamil, grapheme recognition included a letter discrimination task with 18 items,

1034    where the child was shown four glyphs and had to match one to the spoken letter for each item.

1035    Children then completed a naming task where they gave the names of 12 Tamil letters

1036    (aksharas), and a word reading task, where they read up to 100 words. The word list was

1037    administered in blocks of 20, with words increasing in difficulty (according to grade level

1038    literacy resource materials, from the National Institute of Education and public libraries). The

1039    task was discontinued when the child was not able to read any words in a block. Total correct

1040    responses were totalled across the three tasks.

1041         These locally developed reading measures showed sensitive to change over time. Mean

1042    total raw scores over time for the whole sample are as follows: For Chinese at K1 = 28.3, K2 =

1043    35.5, and P1 = 40.7; for Malay at K1 = 11.6, K2 = 19.1, and P1 = 57.9; and for Tamil at K1 =

1044    20.4, K2 = 32.9, and P1 = 59.5. The scores did not yield floor or ceiling effects. The range

1045    (minimum to maximum) at each time point was 0-96, 0-99, 0-111 for Chinese; 0-80, 0-108, 0-

1046    108 for Malay; and 1-88, 2-118, 0-122 for Tamil, respectively. The 25th-75th percentile scores at

1047    each time point = 14.3-40, 20-49, 23.5-57 for Chinese;  7-11,10-14, 28-90 for Malay; and 6.5-29,

1048    9-50.8, 39-85 for Tamil.

1049         **Morphological awareness**. For each language, compound production requires the child

1050    to to generate as many compound words as they could from a base word (e.g. the word 'book'

1051    could be used to make the words 'bookstore' and 'bookshelf'). There were 15 items in Chinese,

1052    10 items in Malay, and 8 items in the Tamil tasks, and total correct responses across all items

1053    were scored. For the compound structure task, the experimenter would describe a scenario in a

1054    sentence, followed by multiple choices for the child to choose to best describe the scenario. For

1055    example, "A house that is built in a tree: It is a treehouse or a housetree?" There were 13 items

1056    for Chinese, 11 items for Malay, and 14 items for Tamil. Total correct items were summed and

1057    the total score was calculated with a correction for guessing. The sum of compound production

1058    and compound structure task scores was taken as the total raw score.

1059            Mean total scores for the whole sample indicate that the task was sensitive to change over

1060    time: for Chinese at K1 = 10.2, K2 = 14.0, and P1 = 18.1; for Malay at K1= 5.4, K2 = 7.4, and

1061    P1 = 9.3; and for Tamil at K1 = 3.4, K2 = 7.2, and P1 = 13.5. The range of scores suggested no

1062    floor or ceiling effects: Minimum-maximum scores at each time point = 0-36, 0-35, 0-43 for

1063    Chinese; 0-15, 0-22, 2-23 for Malay; and 0-14, 0-15, 4-32 for Tamil. The percent of 0 scores =

1064    3.1, 0.6, 0.3 % for Chinese at K1, K2, P1; 2.2, 4.4, 0.7% for Malay; 25, 0.6, 0% for Tamil. The

1065    25th-75th percentile scores at each time point = 5-14, 10-18, 12-23 for Chinese;  4-7, 5-9, 6.3-11

1066    for Malay; and 0-5, 6-9, 9-16 for Tamil.

1067                                     Appendix C

1068     *Growth mixture models to identify latent classes of MTL learners*

1069          Growth mixture models (GMMs) with 1 to 3 classes were fit to data over 3 waves (K1,

1070     K2, P1) using Mplus v.8 Software with an MLR estimator (Muthén & Muthén, 2017). Students

1071     were clustered in classrooms for all models (although intraclass correlation coefficients indicated

1072     a design effect for classroom level only for Chinese reading and Chinese and Tamil

1073     morphological awareness). For each Mother Tongue Language (MTL) group, model fit was

1074     compared for a 2-class model compared to 1-class or 3-class models. The best fitting model in

1075     each case was considered across several criteria, including model fit, likelihood ratio tests

1076     (Vuong-Lo-Mendell-Rubin Likelihood Ratio Test, Lo-Mendell-Rubin Adjusted Likelihood Ratio

1077     Test, p-values < 0.05), entropy (> 0.8), and the proportion of students within the smallest latent

1078     class (> 1%). In addition, plots were inspected to determine interpretability of the pattern of

1079     results (e.g., Fu, R., Chen, X., Wang, L., & Yang, F. (2016). Developmental trajectories of

1080     academic achievement in Chinese children: Contributions of early social-behavioral functioning.

1081     *Journal of Educational Psychology*, *108*(7), 1001-1012). Fit indices are shown in Table C1

1082     below. Further, the number of children with incoming scores below the low progress group's

1083     intercept, above the high progress group's intercept, and in between intercepts is also shown in

1084     Table C2.

1085

1086

Table C1.

*Growth mixture model fit indices for 2-class versus 1-class model, per variable for each language group*

| | *n* | Likelihood Ratio Tests (p-values) | | | | Classification Accuracy | | | | | | Comparison to 3-class model | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | VLMR-LRT | LMR-LRT | Smallest Class (%) | Entropy | Min | Max | Class 1 intercept | Class 1 slope | Class 2 intercept | Class 2 slope | 2- vs. 3-class model | 3-class entropy |
| **Vocabulary** | | | | | | | | | | | | | |
| CL | *1074* | 0.000 | 0.000 | 37.9% | 0.747 | 0.910 | 0.935 | 27.6 | 5.34 | 40.4 | 5.89 | ** | .690 |
| ML | *194* | 0.179 | 0.193 | 39.8% | 0.570 | 0.821 | 0.887 | 27.1 | 3.87 | 33.7 | 6.70 | ns | |
| TL | *236* | 0.000 | 0.000 | 40.2% | 0.733 | 0.905 | 0.928 | 29.6 | 7.67 | 22.6 | 3.67 | ns | |
| **Reading** | | | | | | | | | | | | | |
| CL | *323* | 0.003 | 0.003 | 38.4% | 0.822 | 0.933 | 0.955 | 17.8 | 5.13 | 45.4 | 8.41 | ns | |
| ML | *137* | 0.017 | 0.021 | 12.4% | 0.995 | 0.996 | 1.000 | 6.25 | 7.84 | 35.7 | 38.19 | ns | |
| TL | *160* | 0.267 | 0.280 | 19.3% | 0.928 | 0.949 | 0.981 | 53.1 | 23.1 | 11.5 | 15.02 | ns | |
| **Morphological Awareness** | | | | | | | | | | | | | |
| CL | *319* | 0.000 | 0.000 | 19.9% | 0.777 | 0.867 | 0.956 | 19.5 | 1.16 | 7.8 | 4.64 | ns | |
| ML | *137* | 0.016 | 0.020 | 13.6% | 0.798 | 0.885 | 0.949 | 5.6 | 1.29 | 4.5 | 6.00 | ns | |
| TL | *161* | 0.178 | 0.194 | 15.9% | 0.744 | 0.857 | 0.946 | 2.3 | 4.69 | 7.3 | 3.16 | ns | |

*Note.* Likelihood ratio tests compared 1-class to 2-class models (VLMR-LRT = Vuong-Lo-Mendell-Rubin likelihood ratio test; LMR-LRT = Lo-Mendell-Rubin likelihood ratio test). Smallest class percent, entropy, classification accuracy, and class intercepts and slopes are from the 2-class models. Comparing 2-class model vs. 3-class model fits with LRT: *ns* = 3-class was not a better fit; *$p <$ 0.05, ** $p < 0.01$, meaning 3-class model yielded a better fit, but had lower entropy than the 2-class model. CL= Chinese language learners; ML= Malay language learners; TL= Tamil language learners.

Table C2.

*Number of children scoring above, below or between the intercepts per latent class for Mother Tongue Language (MTL) measures*

|  | Below latent "low progress" class intercept | Between latent class intercepts | Above latent "high progress" class intercept |
|---|---|---|---|
| *Vocabulary* | 157 | 183 | 107 |
| *Reading* | 92 | 175 | 74 |
| *Morphological Awareness* | 48 | 143 | 146 |

*Note.* CL= Chinese language learners; ML= Malay language learners; TL= Tamil language learners.