

# Flat versus expressive storytelling: young children's learning and retention of a social robot's narrative

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

Jacqueline M. Kory Westlund<sup>1\*</sup>, Sooyeon Jeong<sup>1</sup>, Hae Won Park<sup>1</sup>, Samuel Ronfard<sup>2</sup>, Aradhana

Adhikari<sup>1</sup>, Paul L. Harris<sup>2</sup>, David DeSteno<sup>3</sup>, & Cynthia L. Breazeal<sup>1</sup>

<sup>1</sup> MIT Media Lab, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>2</sup> Harvard Graduate School of Education, Harvard University, Cambridge, MA, USA

<sup>3</sup> Department of Psychology, Northeastern University, Boston, MA, USA

\*Corresponding author:

Jacqueline M. Kory Westlund

[jakory@media.mit.edu](mailto:jakory@media.mit.edu)

## 20 **1 Abstract**

21 Prior research with preschool children has established that dialogic or active book reading  
22 is an effective method for expanding young children’s vocabulary. In this exploratory study,  
23 we asked whether similar benefits are observed when a robot engages in dialogic reading  
24 with preschoolers. Given the established effectiveness of active reading, we also asked  
25 whether this effectiveness was critically dependent on the expressive characteristics of the  
26 robot. For approximately half the children, the robot’s active reading was expressive; the  
27 robot’s voice included a wide range of intonation and emotion (*Expressive*). For the  
28 remaining children, the robot read and conversed with a flat voice, which sounded similar  
29 to a classic text-to-speech engine and had little dynamic range (*Flat*). The robot’s  
30 movements were kept constant across conditions. We performed a verification study using  
31 Amazon Mechanical Turk to confirm that the *Expressive* robot was viewed as significantly  
32 more expressive, more emotional, and less passive than the *Flat* robot.  
33 We invited 45 preschoolers with an average age of 5 years who were either English  
34 Language Learners (ELL), bilingual, or native English speakers to engage in the reading  
35 task with the robot. The robot narrated a story from a picture book, using active reading  
36 techniques and including a set of target vocabulary words in the narration. Children were  
37 post-tested on the vocabulary words and were also asked to retell the story to a puppet. A  
38 subset of 34 children performed a second story retelling 4-6 weeks later.  
39 Children reported liking and learning from the robot a similar amount in the *Expressive* and  
40 *Flat* conditions. However, as compared to children in the *Flat* condition, children in the  
41 *Expressive* condition were more concentrated and engaged as indexed by their facial  
42 expressions; they emulated the robot’s story more in their story retells; and they told  
43 longer stories during their delayed retelling. Furthermore, children who responded to the  
44 robot’s active reading questions were more likely to correctly identify the target  
45 vocabulary words in the *Expressive* condition than in the *Flat* condition. Taken together,  
46 these results suggest that children may benefit more from the expressive robot than from  
47 the flat robot.

### 48 **1.1 Keywords**

49 Preschool children; emotion; expressiveness; language development; peer modeling; social  
50 robotics; storytelling

## 51 **2 Introduction**

52 Prior research with preschool children has established that storytelling and story reading  
53 can promote oral language development and story comprehension (Cremin et al., 2016;  
54 Isbell et al., 2004; Speaker et al., 2004). Participating in storytelling can increase children’s  
55 verbal fluency, listening skills, and vocabulary. Book reading in particular can be an  
56 effective method for expanding young children’s vocabulary, especially when children are  
57 encouraged to actively process the story materials. For example, in an intervention study,  
58 middle class parents assigned to an experimental group were instructed to engage in  
59 ‘dialogic’ reading with their 2-year-olds, i.e., to ask more open-ended and  
60 function/attribute questions and to support the efforts of their children to answer these  
61 questions; parents in the control group were instructed to read in their usual fashion. In  
62 follow-up tests, children in the experimental group scored higher in assessments of

63 expressive vocabulary (Whitehurst et al., 1988). Subsequent studies have replicated and  
64 extended this result (e.g., Boteanu et al., 2016; Chang et al., 2012; Nuñez, 2015; Valdez-  
65 Menchaca & Whitehurst, 1992; Hargrave & Sénéchal, 2000). Taken together, these studies  
66 indicate that dialogic book reading is an effective method for boosting children's  
67 vocabulary. Indeed, the studies confirm that such an intervention is quite robust in its  
68 effects—it is effective for toddlers as well as preschoolers, for middle class and working  
69 class children and for typically developing as well as language-delayed children, when  
70 using print or digital storybooks.

71 In this exploratory study, we asked whether similar benefits could be observed when a  
72 social robot engages in dialogic story reading with preschoolers. Social robots share  
73 physical spaces with humans and leverage human means of communicating—such as  
74 speech, movement, and nonverbal cues, including gaze, gestures, and facial expressions—in  
75 order to interface with us in more natural ways (Breazeal, 2004; Breazeal et al., 2008; Feil-  
76 Seifer & Mataric, 2011). Given our expectation that children would learn from the robot, we  
77 also investigated how the emotional expressiveness of the robot's speech might modulate  
78 children's learning.

79 A growing body of research suggests that social robots have potential as learning  
80 companions and tutors for young children's early language education. For example, robots  
81 have played simple vocabulary games to help children learn new words in their own  
82 language or in a second language (Chang et al., 2010; Gordon et al., 2016; Kanda et al.,  
83 2004; Kennedy et al., 2016; Movellan et al., 2009; Tanaka & Matsuzoe, 2012). It is plausible  
84 that children's successful vocabulary learning in these experiments depended on their  
85 relating to the robots as interactive, social beings (Breazeal et al., 2016; Kahn et al., 2013;  
86 Kennedy et al., 2017). Social cues impact children's willingness to engage with and learn  
87 from interlocutors (Bloom, 2000; Corriveau et al., 2009; Harris, 2007; Meltzoff et al., 2009;  
88 Sage & Baldwin, 2010). Indeed, Kuhl (2007, 2011) has argued that a lack of social  
89 interaction with a partner can impair language learning. Thus, infants learn to differentiate  
90 new phonemes presented by a live person, but do not learn this information from a video of  
91 a person, or from mere audio. Because robots are seen by children as social agents—a peer,  
92 a tutor, or a companion—they seem to be providing the necessary social presence to  
93 engage children in a language learning task. Thus, social robots, unlike educational  
94 television programs (Naigles & Mayeux, 2001), may allow children to acquire more  
95 complex language skills and not just vocabulary. However, existing studies on robots as  
96 language learning companions have generally not assessed this possibility. Nearly all of the  
97 activities performed with social robots around language learning have been simple,  
98 vocabulary-learning tasks, with limited interactivity. For example, the robot might act out  
99 new verbs (Tanaka & Matsuzoe, 2012), show flashcard-style questions on a screen  
100 (Movellan et al., 2009), or play simple give-and-take games with physical objects (Movellan  
101 et al., 2009) (see also, Gordon et al., 2016).

102 A few studies have explored other kinds of activities for language learning. For example,  
103 Chang et al. (2010) had their robot read stories aloud, ask and answer simple questions,  
104 and lead students in reciting vocabulary and sentences. However, they primarily assessed  
105 children's engagement with the robot, rather than their language learning. One study used  
106 a story-based task in which the robot took turns telling stories with preschool and  
107 kindergarten children, for eight weeks (Kory 2014; Kory & Breazeal, 2014; Kory Westlund  
108 & Breazeal, 2015). In each session, the robot would tell two stories with key vocabulary

109 words embedded, and would ask children to make up their own stories for practice. For  
110 half the children, the robot personalized the level of the stories to the child's ability, telling  
111 more complex stories for children who had greater ability. This study found increases in  
112 vocabulary learning as well as in several metrics assessing the complexity of the stories  
113 that children produced, with greater increases when children heard appropriately leveled  
114 stories. These findings suggest that a social robot is especially likely to influence language  
115 learning if it conveys personal attunement to the child. Indeed, children were more trusting  
116 of novel information provided by a social robot whose nonverbal expressiveness was  
117 contingent on their behavior (Breazeal et al., 2016) and showed better recall of a story  
118 when the social robot teaching them produced high immediacy gestures in response to  
119 drops in children's attention (Szafir & Mutlu, 2012).

120 In the current study, we focus on a related but hitherto unexplored factor: the emotional  
121 expressiveness of the robot's speech. Nearly every study conducted so far on the use of  
122 social robots as learning companions for young children has used a computer-generated  
123 text-to-speech voice, rather than a more natural, human voice. We know very little about  
124 the effects of a more expressive, human-like voice as compared to a less expressive, flatter  
125 or synthetic voice on children's learning. Such expressive qualities may have an especially  
126 strong impact during storytelling activities. For example, if a potentially engaging story is  
127 read with a flat delivery, children might find it anomalous or even aversive. Using robots to  
128 study questions about expressivity is quite feasible, because we can carefully control the  
129 level of vocal expressiveness across conditions and between participants. Robots afford a  
130 level of control that it is difficult to achieve with human actors with the same consistency.  
131 A small number of human-robot interaction (HRI) studies have investigated the effects of a  
132 robot's voice on an interaction. However, these studies tested adults (e.g., Eyssel et al.,  
133 2012), compared different synthetic voices (e.g., Sandygulova et al., 2015; Tamagawa et al.,  
134 2011; Walters et al., 2008), or compared qualities of the same voice, such as pitch (e.g.,  
135 Lubold et al., 2016; Niculescu et al., 2013), rather than varying the expressiveness of a  
136 given voice. Eyssel et al. (2012) did compare human voices to synthetic voices, but the  
137 adult participants merely watched a short video clip of the robot speaking, and did not  
138 interact with it directly. These participants perceived the robot more positively when the  
139 voice shared their gender, and anthropomorphized the robot more when the voice was  
140 human.

141 Some related work in speech-language pathology and education has compared children's  
142 learning from speakers with normal human voices or voices with a vocal impairment,  
143 specifically, dysphonic voices. Children ages 8-11 years performed better on language  
144 comprehension measures after hearing passages read by a normal human voice than when  
145 the passages were read by a dysphonic voice (Lyberg-Ahlander et al., 2015; Morton et al.,  
146 2009; Rogerson et al., 2005). These studies suggest that vocal impairment can be  
147 detrimental to children's speech processing, and may force children to allocate processing  
148 to the voice signal at the expense of comprehension. However, it is unclear whether a lack  
149 of expressivity or the use of a synthetic voice would impair processing relative to a normal  
150 human voice.

151 Given the lack of research in this area, we compared the effect of an expressive as  
152 compared to a flat delivery by a social robot. We also focused on a more diverse population  
153 as compared to much prior work with regard to both age and language proficiency.  
154 Previous studies have tended to focus on just one population of children—either native

155 speakers of the language, or children learning a second language—whereas we included  
156 both. In addition, few previous studies have included preschool children (Kory, 2014;  
157 Movellan et al., 2009; Tanaka & Matsuzoe, 2012); the majority of studies has targeted older  
158 children. More generally, young children comprise an age group that is typically less  
159 studied in human-robot interaction (Baxter et al., 2016).  
160 In this study, we invited preschoolers with an average age of 5 years and considerable  
161 variation in language proficiency to engage in a dialogic reading task with a social robot.  
162 Thus, some children were English Language Learners (ELL), some were bilingual, and some  
163 were monolingual, native English speakers. All children were introduced to a robot who  
164 first engaged them in a brief conversation and then proceeded to narrate a story from a  
165 picture book using dialogic reading techniques. Two versions of the study were created;  
166 each version contained a unique set of three novel words. In post-story testing, children's  
167 comprehension of the novel words they had heard was compared to their comprehension  
168 of the novel words embedded in the story version they had not heard. We predicted that  
169 children would display superior comprehension of the novel words that they had heard.  
170 Given the established effectiveness of dialogic reading with young children, the robot  
171 always asked dialogic questions. We asked two related questions: First, we asked whether  
172 children would learn from a dialogic storytelling robot. Second, we asked whether its  
173 effectiveness was critically dependent on the expressiveness of the robot's voice—how  
174 might the robot's vocal expressivity impact children's engagement and learning? For  
175 approximately half the children, the robot's dialogic reading was expressive in the sense  
176 that the robot's voice included a wide range of intonation and emotion. For the remaining  
177 children, the robot read and conversed with a flat voice, which sounded similar to a classic  
178 text-to-speech engine and had little dynamic range. To control for the many differences  
179 that computer-generated voices have from human voices (e.g., pronunciation and quality),  
180 an actress recorded both voices, and we performed a manipulation check to ensure the  
181 expressive recording was perceived to be sufficiently more emotional and expressive than  
182 the flat recording. We anticipated that children would be more attentive, show greater gain  
183 in vocabulary, and use more of the target vocabulary words themselves if the dialogic  
184 reading was delivered by the expressive as compared to the flat robot. To further assess the  
185 potentially distinct impact of the two robots, children were also invited to retell the  
186 picture-book story that the robot had narrated. More specifically, they were invited to retell  
187 the story to a puppet who had allegedly fallen asleep during the robot's narration and was  
188 disappointed at having missed the story. Finally, a subset of children was given a second  
189 opportunity to retell the story approximately 4-6 weeks later.

## 190 **3 Methods**

### 191 **3.1 Design**

192 The experiment was designed to include two between-subjects conditions: Robot  
193 expressiveness (*Expressive* voice versus *Flat* voice) and Robot redirection behaviors  
194 (*Present* versus *Absent*). Regarding the robot's voice, the expressive voice included a wide  
195 range of intonation and emotion, whereas the flat voice sounded similar to a classic text-to-  
196 speech engine with little dynamic range. The robot redirection behaviors were a set of re-  
197 engagement phrases that the robot could employ to redirect a distracted child's attention

198 back to the task at hand. However, the conditions under which the robot would use  
199 redirection behaviors did not arise—i.e., all the children were attentive and the  
200 opportunity to redirect their attention did not occur. Thus, the experiment ultimately had a  
201 two-condition, between-subjects design (*Expressive* vs. *Flat*).

## 202 3.2 Participants

203 This study was carried out in accordance with the recommendations of the MIT Committee  
204 on the Use of Humans as Experimental Subjects. Children's parents gave written informed  
205 consent prior to the start of the study and all children assented to participate, in  
206 accordance with the Declaration of Helsinki. The protocol was approved by the MIT  
207 Committee on the Use of Humans as Experimental Subjects and by the Boston Public  
208 Schools Office of Data and Accountability.

209 We recruited 50 children aged 4-7 (23 female, 27 male) from a Boston-area school (36  
210 children) and the general Boston area (14 children) to participate in the study. Five  
211 children were removed from the analysis because they did not complete the study. The  
212 children in the final sample included 45 children (22 female, 23 male; 34 from the school  
213 and 11 from the general Boston area) with a mean age of 5.2 years ( $SD = 0.77$ ). 17 children  
214 were English Language Learners (ELL), 8 were bilingual, 18 were native English speakers,  
215 and 3 were not reported.

216 Children were randomly assigned to conditions. There were 23 children (13 male, 10  
217 female; 10 ELL, 5 native English, 5 bilingual, 2 unknown) in the *Expressive* condition and 22  
218 children (9 male, 13 female; 7 ELL, 12 native English, 3 bilingual, 1 unknown) in the *Flat*  
219 condition. The two conditions were not perfectly balanced due to the fact that we did not  
220 obtain information about children's language learning status until the completion of the  
221 study, and thus could not assign children evenly between conditions.

222 We created two versions of the story that the robot told (version A and version B); each  
223 version was identical except for the inclusion of a different set of target vocabulary words.  
224 Approximately half of the participants heard story version A (*Expressive*: 11, *Flat*: 10); the  
225 other half heard story version B (*Expressive*: 13, *Flat*: 11).

226 We used the Peabody Picture Vocabulary Test, 4<sup>th</sup> edition (PPVT; Dunn & Dunn, 2007), to  
227 verify that the children in the *Expressive* and *Flat* conditions did not have significantly  
228 different language abilities. The PPVT is commonly used to measure receptive language  
229 ability for standard American English. On each test item, the child is shown a page with four  
230 pictures, and is asked to point to the picture showing the target word. PPVT scores for  
231 three of the 45 children could not be computed due to missing data regarding their ages.

232 For the remaining 42 children, there were, as expected, no significant differences between  
233 the *Expressive* and *Flat* conditions in PPVT scores,  $t(40) = 0.64$ ,  $p = 0.53$ . A one-way analysis  
234 of variance (ANOVA) with age as a covariate revealed that children's PPVT scores were  
235 significantly related to their age,  $F(3,37) = 5.83$ ,  $p = 0.021$ ,  $\eta^2 = 0.114$ , as well as to their  
236 language status,  $F(3,37) = 2.72$ ,  $p = 0.058$ ,  $\eta^2 = 0.160$ . As expected, post-hoc pairwise  
237 comparisons indicated that children who were native English speakers had higher PPVT  
238 scores ( $M = 109.4$ ,  $SD = 18.2$ ) than ELL children ( $M = 92.0$ ,  $SD = 14.6$ ),  $p = 0.004$ . There  
239 were no differences between the bilingual children ( $M = 103.5$ ,  $SD = 15.7$ ) and either the  
240 native English-speaking children or the ELL children.

### 241 3.3 Hypotheses

242 The effects of the robot's expressivity might be transient or long-term, subtle or wide-  
243 ranging. Accordingly, we used a variety of measures, including immediate assessments as  
244 well as the delayed retelling task, to explore whether the effect of the robot's expressivity  
245 was immediate and stable and whether it impacted all measures, or selected measures  
246 only.

247 We tentatively expected the following results:

#### 248 3.3.1 Learning

- 249 • In both conditions, children would learn the target vocabulary words presented in  
250 the story version that they heard.
- 251 • Children who learned the target words would also use them in their story retells.
- 252 • The robot's expressivity would lead to differences in children's long-term retention  
253 of the story. Children in the *Expressive* condition would better retain the story, and  
254 thus tell longer stories, incorporating the phrases that appeared in the initial story  
255 into their delayed retells.

#### 256 3.3.2 Behavior

- 257 • In both conditions, children would typically respond to the robot's dialogic reading  
258 questions, but children who responded more often to the dialogic reading questions  
259 would show greater learning gains.
- 260 • The *Expressive* robot would promote greater modeling by the children of the robot's  
261 story. Children in the *Expressive* condition would produce more vocabulary and  
262 phrase mirroring.

#### 263 3.3.3 Engagement

- 264 • Although most children would express liking for the robot, indirect behavioral  
265 measures would show that children were more attentive and engaged with the  
266 *Expressive* robot than with the *Flat* robot.
- 267 • The surprising moments in the story would have greater impact on children in the  
268 *Expressive* condition, because suspense and surprise in the story were strongly  
269 reflected in the robot's voice.

### 270 3.4 Procedure

271 Each child was greeted by an experimenter and led into the study area. The experimenter  
272 wore a hand puppet, a purple Toucan, which she introduced to the child: "This is my friend,  
273 Toucan." Then the puppet spoke: "Hi, I'm Toucan!" The experimenter used the puppet to  
274 invite the child to do a standard vocabulary test, the PPVT, by saying "I love word games.  
275 Want to play a word game with me?" The experimenter then administered the PPVT.  
276 For the children who participated in the study at their school, the PPVT was administered  
277 during an initial session. The children were brought back on a different day for the robot  
278 interaction. This second session began with the puppet asking children if they remembered  
279 it: "Remember me? I'm Toucan!" Children who participated in the lab first completed the  
280 PPVT, and were then given a five-minute break before returning to interact with the robot.

281 For the robot interaction, the experimenter led the child into the robot area. The robot sat  
282 on a low table facing a chair, in which children were directed to sit. A tablet was positioned  
283 in an upright position in a tablet stand on the robot's right side. A smartphone sat in front  
284 of the robot; it ran software to track children's emotional expressions (See Figure 1). The  
285 experimenter sat to the side and slightly behind the children with the puppet. The  
286 interaction began with the puppet introducing the robot, Tega: "This is my friend, Tega!"  
287 The robot introduced itself, shared personal information, and prompted children to do the  
288 same, e.g., "Hi, I'm Tega! My favorite color is blue. What is your favorite color?" and "Do you  
289 like to dance? I like to dance!"

290 After this brief introductory conversation, the robot asked the children if they wanted to  
291 hear a story. At this point, the puppet interjected that it was sleepy, but would try to stay  
292 awake for the story. The experimenter made the puppet yawn and fall asleep; it stayed  
293 asleep for the duration of the story. The robot then told the story which consisted of a 22-  
294 page subset of the wordless picture book "Frog, Where Are you?" by Mercer Mayer. This  
295 book has been used before in numerous studies, especially in research on speech pathology  
296 (e.g., Boudreau & Hedberg, 1999; Diehl et al., 2006; Greenhalgh & Strong, 2001; Heilmann  
297 et al., 2010).

298 The pages of the book were shown one at a time on the tablet screen. Each page was  
299 accompanied by 1-2 sentences of text, which the robot read in either an expressive or a flat  
300 voice depending on the condition. For every other page, the robot asked a dialogic reading  
301 comprehension question about the events in the story, e.g., "What is the frog doing?", "Why  
302 did the boy and the dog fall?", and "How do you think the boy feels now?" (11 questions  
303 total). The robot responded to children's answers with encouraging, but non-committal,  
304 phrases such as "Mmhm," "Good thought," and "You may be right."

305 We embedded three target vocabulary words (all nouns) into the story. We did not test  
306 children on their knowledge of these words prior to the storytelling activity because we did  
307 not want to prime children to pay attention to these words, since that could bias our results  
308 regarding whether or not children would learn or use the words after hearing them in the  
309 context of the robot's story. Instead, in order to assess whether children were more likely  
310 to know or use the words after hearing the robot use them in the story, two versions of the  
311 story (version A and version B) were created with different sets of target words. The two  
312 versions of the story were otherwise identical. We identified six key nouns in the original  
313 story: animal, rock, log, hole, deer, and hill. Then, in each of our two story versions, we  
314 replaced three of the words with our target words, so that each story version included  
315 three target words and three original words. Version A included the target words "gopher"  
316 (original word: animal), "crag" (rock), and "lily pad" (log); version B included the words  
317 "hollow" (hole), "antlers" (deer), and "cliff" (hill). We anticipated that children would  
318 display selective learning and/or use of these words, depending on which story they heard.  
319 We looked both at children's later receptive knowledge of the words as well as expressive  
320 or productive abilities, since children who can recognize a word may or may not be able to  
321 produce it themselves.

322 At the end of the story, the Toucan woke up and exclaimed, "Oh no! Did I miss the story?"  
323 This presented an opportunity for children to retell the story to the puppet, thereby  
324 providing a measure of their story recall. Children were allowed to go through the story on  
325 the tablet during their retelling. Thus, the depictions on each page could serve as a  
326 reminder during retelling.



327 After the story-retelling task, the experimenter administered a PPVT-style vocabulary test  
328 for the six target words used across the two versions of the story. For each word, four  
329 pictures taken from the story's illustrations were shown to children and they were asked to  
330 point to the picture matching the target word. Finally, the experimenter asked children a  
331 set of questions regarding their perception of the robot and their enjoyment of the story.

332 These questions were as follows:

- 333 1. How much did you like the story the robot read? Really really liked it, liked it quite a  
334 lot, liked it a little bit, sort of liked it, didn't really like it
- 335 2. Why did you like or not like the story?
- 336 3. How much do you like Tega? Really really liked Tega, liked Tega quite a lot, liked  
337 Tega a little bit, sort of liked Tega, didn't really like Tega
- 338 4. Why do you like or not like Tega?
- 339 5. Would Tega help you feel better if you were feeling sad? Really really helpful,  
340 quite helpful, a little helpful, sort of helpful, not really helpful
- 341 6. Why would Tega help or not help?
- 342 7. How helpful was Tega in helping you learn the story? Really really helpful, quite  
343 helpful, a little helpful, sort of helpful, not really helpful
- 344 8. Why was Tega helpful or not helpful?
- 345 9. Would one of your friends would want to read stories with Tega? Really  
346 really want to, want to quite a lot, want to a little, sort of want to, won't really want  
347 to
- 348 10. Why your friend would or wouldn't want to read stories with Tega?
- 349 11. Can you describe Tega to your friend?
- 350 12. Who would you want to tell another story to: Toucan or Tega?
- 351 13. Why would you want to read another story to: Toucan or Tega?

352  
353 Where appropriate, we used a Smiley-o-meter to gather responses on a 1-5 scale (Read &  
354 McFarlane, 2006). Although Read and McFarlane (2006) suggest that this measure is not  
355 useful with children younger than 10 years, previous research has successfully used it, or  
356 similar measures, with modest pre-training (Harris et al., 1985; Leite et al., 2014). Thus, we  
357 did a practice item before the test questions so children could learn how the measure  
358 worked. Children were also asked to explain their answers, such as "Why do you like or not  
359 like Tega?" and "Why was Tega helpful or not helpful?" Children's parents or teachers  
360 provided demographic data regarding the children's age and language status (ELL,  
361 bilingual, or native English speaker).

362 A subset of 34 children from the school sample participated in a second, follow-up session  
363 approximately 4-6 weeks later at their school. Children who participated in the lab did not  
364 have a follow-up session due to logistical reasons. During this follow-up session, we  
365 administered the PPVT a second time, then asked children to retell the story to the puppet.  
366 The puppet prompted children by saying, "I tried to tell the story to my friend last week,  
367 but I forgot most of it! Can you tell it to me again?" This allowed us to observe children's  
368 long-term memory for the story.

369 Four different experimenters (3 female adults and 1 male adult) ran the study in pairs. One  
370 experimenter interacted with the child. The other experimenter acted as the robot  
371 teleoperator and equipment manager; she could be seen by the children, but she did not  
372 interact directly with them.

### 373 3.5 Materials

374 We used the Tega robot, a squash and stretch robot designed for educational activities with  
375 young children (Kory Westlund et al., 2016). The robot is shown in Figure 1. It uses an  
376 Android phone to run its control software as well as display an animated face. The face has  
377 two blue oval eyes and a white mouth, which can all morph into different shapes. This  
378 allows the face to show different facial expressions and to show appropriate visemes (i.e.,  
379 mouth shapes) when speech is played back. The robot can move up and down, tilt its head  
380 sideways or forward/backward, twist to the side, and lean forward or backward. Some  
381 animations played on the robot use only the face; others incorporate both facial  
382 expressions and physical movements of the body. The robot is covered in red fur with blue  
383 stripes, giving it a whimsical, friendly appearance. The robot was referred to in a non-  
384 gendered way by the experimenters throughout the study.

385 A female adult recorded the robot's speech. These utterances were shifted into a higher  
386 pitch to make them sound child-like. For the *Expressive* condition, the utterances were  
387 emotive with a larger dynamic range; the actress was instructed to speak in an expressive,  
388 human-like way. For the *Flat* condition, the actress imitated a computer-generated text-to-  
389 speech voice, keeping her intonation very flat. We did not use an actual computer-  
390 generated voice for the *Flat* voice because there would have been many differences in  
391 pronunciation and quality compared to the *Expressive* voice. Similarly, we did not use a  
392 computer-generated voice for the *Expressive* voice because no computer-generated voices  
393 can currently imitate the dynamic, expressive range that human voices are capable of.  
394 Many of the physical actions the robot can perform are expressive. We used the same  
395 physical movements in both conditions; however, in the *Expressive* condition, some  
396 movements were accompanied by expressive sounds (such as "Mm hm!"), whereas in the  
397 *Flat* condition, these movements were either accompanied by a flat sound ("Mm hm.") or,  
398 in cases where the sound was a short, non-linguistic expressive utterance, no sound.

399 We used a Google Nexus 9 8.9-inch tablet to display the storybook. Touchscreen tablets  
400 have been shown to effectively engage children and social robots in a shared task (Park et  
401 al., 2014). We used custom software to display the story pages that allowed a teleoperator  
402 to control when the pages were turned; this software is open-source and available online  
403 under the MIT License at <https://github.com/mitmedialab/SAR-opal-base/>.

404 We used a Samsung Galaxy S4 android smartphone to run Affdex, which is emotion  
405 measurement software from Affectiva, Inc<sup>1</sup>. Affdex performs automatic facial coding in four  
406 steps: face and facial landmark detection, face feature extraction, facial action, and emotion  
407 expression modeling based on the EMFACS emotional facial action coding system (McDuff  
408 et al., 2016; Ekman et al., 1978; Friesen & Ekman, 1983). Although no data has been  
409 published yet specifically comparing the performance of the software on adults versus  
410 children, FACS coding is generally the same for adults and for children and has been used  
411 with children as young as two years (e.g., LoBue & Thrasher, 2015; Camras et al., 2006; also  
412 see Ekman & Rosenberg, 1997). Furthermore, this software has been trained and tested on  
413 tens of thousands of manually coded images of faces from around the world (McDuff et al.,  
414 2015; McDuff et al., 2013; Senechal et al., 2015).

---

<sup>1</sup> Affectiva, Inc., <http://affectiva.com/>, retrieved September 19, 2016.

### 415 3.6 Teleoperation

416 We used a custom teleoperation interface to control the robot and the digital storybook.  
417 Using teleoperation allowed the robot to appear autonomous to participants while  
418 removing technical barriers such as natural language understanding, because the  
419 teleoperator could be in the loop as the language parser. The teleoperator used the  
420 interface to trigger when the robot should begin its next sequence of actions (a list of  
421 speech, physical motions, and gaze) and also when the storybook should proceed to the  
422 next page. Thus, the teleoperator needed to pay attention to timing in order to trigger the  
423 robot's next action sequence at the appropriate times relative to when the experimenter  
424 spoke (i.e., when introducing the robot to the child), or when the child responded to one of  
425 the robot's questions. Since the teleoperator did not manage the timing of actions within  
426 each sequence, the robot's behavior was highly consistent for all children.  
427 The four experimenters were all trained to control the robot by an expert teleoperator;  
428 they had all controlled robots before in multiple prior studies.

### 429 3.7 Manipulation Check

430 To check that the *Expressive* robot was, in fact, perceived to be more expressive than the  
431 *Flat* robot, we performed a verification study using Amazon Mechanical Turk (AMT). We  
432 recorded video of the robot performing all the speech and behavior used in the main study.  
433 We then selected samples of the robot's speech and behavior from the introductory  
434 conversation, the beginning, middle, and end of the story, and the closing of the interaction  
435 to create a video clip that was approximately two and a half minutes in length. We created  
436 one video of the *Flat* robot and one video of the *Expressive* robot. In the two videos, we used  
437 the same speech and behavior samples such that the only difference was the  
438 expressiveness of the robot's voice.

439 We recruited 40 AMT workers from the United States. Half the participants (11 male, 9  
440 female) viewed the video of the *Flat* robot and half (13 male, 7 female) viewed the video of  
441 the *Expressive* robot. After viewing the video, participants were asked to rate their  
442 impression of the robot and report demographic information. We used the following  
443 questions, each of which was measured on a 1-5 Likert-type scale anchored with "1: Not  
444 \_\_\_ at all" and "5: Extremely \_\_\_":

- 445 1. Overall, how expressive or not expressive was the robot in the video?
- 446 2. Overall, how emotional or not emotional was the robot in the video?
- 447 3. Overall, how passive or not passive was the robot in the video?
- 448 4. How expressive or not expressive was the robot's voice in the video?
- 449 5. How emotional or not emotional was the robot's voice in the video?
- 450 6. How passive or not passive was the robot's voice in the video?
- 451 7. How expressive or not expressive was the robot's movement in the video?
- 452 8. How emotional or not emotional was the robot's movement in the video?
- 453 9. How passive or not passive was the robot's movement in the video?

454

455 Table 1 shows a summary of participant responses. We found that participants who  
456 watched the *Expressive* robot video rated the robot as significantly more emotional overall

457 than participants who watched the *Flat* robot video,  $t(39) = 2.39, p = 0.022$ . Participants  
458 who watched the *Expressive* robot video rated the robot's voice as significantly more  
459 expressive,  $t(39) = 4.44, p < 0.001$ ; more emotional,  $t(39) = 5.15, p < 0.001$ ; and less  
460 passive,  $t(39) = 2.96, p = 0.005$ , than participants who watched the *Flat* robot video. There  
461 were no statistically significant differences in participants' ratings of the robot's  
462 movement.

463 The results demonstrate that the *Expressive* and *Flat* robot conditions were indeed  
464 sufficiently different from each other, with the voice of the *Expressive* robot being viewed as  
465 more expressive, more emotional, and less passive than the *Flat* robot.

### 466 **3.8 Data**

467 We recorded video and audio data for each session using two different cameras set up on  
468 tripods behind the robot, facing the child. We recorded children's facial expressions using  
469 Affdex, emotion measurement software from Affectiva, Inc. Children's responses to the  
470 PPVT, target word vocabulary test, and interview questions were recorded on paper during  
471 the experiment and later transferred to a spreadsheet.

### 472 **3.9 Data Analysis**

473 We coded whether or not children responded to each of the questions the robot asked  
474 during the initial conversation and during the story, and if they did respond, how many  
475 words their response consisted of. We also counted the number of questions that children  
476 asked the puppet when retelling the story.

477 To assess how children perceived Tega as a function of their assignment to the *Expressive*  
478 and *Flat* conditions, we coded children's responses to the open-ended question inviting  
479 them to describe Tega to a friend (i.e., "Can you describe Tega to your friend?") for positive  
480 traits (e.g., nice, helpful, smart, fun). All children provided a response to this question.

481 Children's responses to the Smiley-o-meter questions were coded on a 1-5 scale.

482 Children's transcribed story retells were analyzed in terms of their story length, overall  
483 word usage and target word usage, and phrase similarity compared to the robot's original  
484 story. Automatic tools were developed such that each word was converted into its original  
485 form for comparison (stemming), words with no significant information (i.e., stopwords)  
486 were removed, and an N-gram algorithm was implemented to match phrases between the  
487 child's and the robot's stories. N-gram refers to a contiguous sequence of N items from a  
488 given sequence of text. In our analysis, we used N=3 for matching and comparison. We  
489 chose N=3 because a smaller N (e.g., N=2) often retains too little information to constitute  
490 actual phrase matching, and a larger N may encompass more information than would  
491 constitute a single phrase. For example, the robot's story included the section, "The frog  
492 jumped out of an open window. When the boy and the dog woke up the next morning, they  
493 saw that the jar was empty." After stemming and stopword removal, this section would be  
494 converted to "frog jump open window boy dog wake next morning see jar empty." One  
495 child retold this section of the story by saying "Frog was going to jump out the window. So  
496 whe... then the boy and the dog woke up, the jar was empty." This was converted to "frog  
497 jump window boy dog wake jar empty." The N-gram phrase matching for this segment  
498 reveals multiple phrase matches, e.g., (robot) "*frog jump open window*" / (child) "*frog jump*

499 *window*”, and (robot) “*boy dog wake next morning see jar empty*” / (child) “*boy dog wake*  
500 *jar empty*”.

501 Children's affect data were collected using Affdex whenever a face was detected with the  
502 front-facing camera on the Samsung Galaxy S4 device (McDuff et al., 2016). Affdex is  
503 capable of measuring 15 expressions, which are used to calculate the likelihood that the  
504 detected face is displaying each of 9 different affective states. We analyzed the four  
505 affective states most relevant to our research questions: attention, concentration, surprise  
506 and engagement. *Attention* is a measure of focus based on head orientation—i.e., is the  
507 child attending to the task or not. The likelihood of *concentration* is increased by *brow*  
508 *furrow* and *smirk*, and decreased by *smile*. Thus, *concentration* reflects the effort and  
509 affective states associated with attending, rather than merely whether the child is looking  
510 in the correct direction or not. *Surprise* is increased by *inner brow raise*, *brow raise* and  
511 *mouth open*, and decreased by *brow furrow*. *Engagement* measures facial muscle activation  
512 reflective of the subject's expressiveness, and is calculated as a weighted sum of the *brow*  
513 *raise*, *brow lower*, *nose wrinkle*, *lip corner depressor*, *chin raise*, *lip pucker*, *lip press*, *lips part*,  
514 *lip suck* and *smile*. Thus, the *Engagement* score reflects total facial muscle activation during  
515 the task. On every video frame (up to 32 frames per second), each of these affective states  
516 was scored by Affdex in the range 0 (no expression present) to 100 (expression fully  
517 present). Values in the middle (e.g., 43 or 59) indicate that the expression is somewhat  
518 present; these values are relative and Affdex does not indicate what the exact difference is  
519 between each score. See Senechal et al. (2015) for more detail regarding the algorithms  
520 uses for classification.

521 For the story retelling, the audio quality of 40 out of 45 participants was sufficiently good  
522 enough for transcription (22 female, 18 male; age  $M = 5.2$ ,  $SD = 0.76$ ; 14 ELL, 7 bilingual, 16  
523 native English, 3 unknown). There were 21 children (10 female, 11 male; age  $M = 5.3$ ,  $SD =$   
524  $0.80$ ; 9 ELL, 4 bilingual, 6 native English, 2 unknown) in the *Expressive* condition and 19  
525 children (12 female, 7 male; age  $M = 5.1$ ,  $SD = 0.71$ ; 5 ELL, 3 bilingual, 10 native English, 1  
526 unknown) in the *Flat* condition. Half of the participants had heard story version A  
527 (*Expressive*: 10, *Flat*: 10); the other half had heard story version B (*Expressive*: 11, *Flat*: 9).  
528 To perform analyses across the two sessions, immediate and delayed retell pairs from 29  
529 children were used (14 female, 15 male; age  $M = 5.2$ ,  $SD = 0.68$ ; 14 ELL, 3 bilingual, 12  
530 native English). There were 15 children (6 female, 9 male; age  $M = 5.3$ ,  $SD = 0.70$ ; 9 ELL, 2  
531 bilingual, 4 native English) from the *Expressive* condition and 14 children (8 female, 6 male;  
532 age  $M = 5.1$ ,  $SD = 0.66$ ; 5 ELL, 1 bilingual, 8 native English) from the *Flat* condition. Half of  
533 the participants heard story version A (*Expressive*: 8, *Flat*: 8); the other half heard story  
534 version B (*Expressive*: 7, *Flat*: 6).

535 In the following analyses, we ran Shapiro-Wilk (S-W) tests to check for normality and  
536 Levene's test to check for equal variance, where applicable. Levene's null hypothesis was  
537 rejected for all data in our dataset ( $p > 0.05$ ) and constant variance was assumed across  
538 conditions and sessions. Parametric (paired/unpaired t-test) and non-parametric  
539 (Wilcoxon signed-rank and Mann-Whitney's U) tests were used based on the S-W result.

#### 540 **4 Results**

541 We present our results in three parts, with each part addressing one of our three main  
542 hypotheses: (1) *Learning*: our primary question was whether children would learn from a  
543 robot that led a dialogic storytelling activity, and specifically whether the expressiveness of

544 the robot's voice would impact children's learning; (2) *Behavior*: we asked whether  
545 children would learn more if they responded to the dialogic reading questions, and  
546 whether the robot's expressiveness would produce greater lexical and phrase modeling;  
547 (3) *Engagement*: we asked whether the robot's expressiveness would lead to greater  
548 attention or engagement. Finally, we also examined whether children's learning was  
549 impacted by their language status.

## 550 4.1 Learning

### 551 4.1.1 Target vocabulary word identification

552 Overall, children correctly identified a mean of 4.0 of the 6 target vocabulary words ( $SD =$   
553  $1.38$ ). A  $2 \times 2 \times 2$  mixed ANOVA with condition (*Expressive* vs. *Flat*), the story children  
554 heard (version A vs. version B), and the words correctly identified (number of version A  
555 words correct vs. number of version B words correct, where children were asked to  
556 identify both sets of words), with age as a covariate, revealed a trend toward age affecting  
557 how many words children identified correctly,  $F(1,81) = 3.40$ ,  $p = 0.069$ ,  $\eta^2 = 0.045$ . Post-  
558 hoc pairwise comparisons showed that older children identified more target words  
559 correctly, with four-year-olds identifying fewer words than five-year-olds ( $p = 0.016$ ), six-  
560 year-olds ( $p = 0.016$ ), and the seven-year-old ( $p = 0.077$ ) (Table 2). There was no difference  
561 between the total number of target vocabulary words that children identified correctly in  
562 the *Expressive* ( $M = 3.8$  correct of 6,  $SD = 1.48$ ) versus *Flat* ( $M = 4.23$ ,  $SD = 1.27$ ) conditions.  
563 We also found the expected interaction between story version heard and number of words  
564 correctly identified from each version (Figure 2). Children who heard story version A were  
565 likely to correctly identify more version A words ( $M = 2.00$  correct of 3,  $SD = 0.853$ ) than  
566 version B words ( $M = 1.62$ ,  $SD = 0.813$ ), whereas children who heard story version B were  
567 more likely to correctly identify more version B words ( $M = 2.21$  correct of 3,  $SD = 1.03$ )  
568 than version A words ( $M = 1.92$ ,  $SD = 0.626$ ),  $F(1,81) = 4.21$ ,  $p = 0.043$   
569 In summary, performance in the vocabulary test improved with age. Nevertheless, there  
570 was evidence of learning from the story in that children performed better on those items  
571 they had encountered in the story version they heard.

### 572 4.1.2 Target word use

573 First, because the two story versions (A and B) differed both in terms of the target words  
574 included and the original words (i.e., the lower level words that we replaced with the target  
575 words), we analyzed how often children used either type of word. This was to provide  
576 context in terms of children's overall word reuse rates after hearing the words in the  
577 robot's story. Thus, among 40 children, 35 children either used the target words or the  
578 original words in their story retelling ( $M = 2.15$  out of 6,  $SD = 1.48$ ). As in the target word  
579 identification, we also found significant differences in children's word usage behavior  
580 based on the story version they heard. A Wilcoxon signed-rank test revealed that children  
581 who heard story version A were more likely to use version A words in their story retelling  
582 ( $M = 1.75$ ,  $SD = 1.37$ ) than version B words ( $M = 1.00$ ,  $SD = 0.920$ );  $W = 12$ ,  $Z = -2.34$ ,  $p =$   
583  $0.019$ ,  $r = 0.52$ , whereas children who heard story version B were more likely to use  
584 version B words ( $M = 2.00$ ,  $SD = 1.69$ ) than version A words ( $M = 0.700$ ,  $SD = 0.660$ );  $W =$   
585  $12.5$ ,  $Z = -2.87$ ,  $p = 0.004$ ,  $r = 0.64$  (Figure 3).

586 Then, to analyze children's learning of new words from the robot, we focused on children's  
587 reuse of the target words. There was no significant difference in overall target word usage  
588 between the *Flat* and *Expressive* conditions. In the immediate retell, children used a mean  
589 of 0.45 target words (out of 3),  $SD = 0.69$ . However, out of the seventeen children who used  
590 at least one of the target words in their retell (*Expressive*: 10 children, *Flat*: 7), children in  
591 the *Expressive* condition used significantly more target words ( $M = 1.6$ ,  $SD = 0.70$ ) than  
592 children in the *Flat* condition ( $M = 1.00$ ,  $SD = 0.00$ ),  $t(15) = 2.248$ ,  $p = 0.040$ .

593 A trend toward older children using more target words than younger children was also  
594 observed; age 4 ( $M = 0.14$ ,  $SD = 0.38$ ), age 5 ( $M = 0.58$ ,  $SD = 0.77$ ), age 6 ( $M = 0.62$ ,  $SD =$   
595  $0.65$ ), age 7 ( $M = 3.0$ ,  $SD = 0.00$ ); Kendall's rank correlation  $\tau(38) = 0.274$ ,  $p = 0.059$ . In the  
596 delayed retell, time was significant ( $M = 0.21$ ,  $SD = 0.49$ ;  $W = 10$ ,  $Z = -2.77$ ,  $p = 0.05$ ). The  
597 correlation between the number of target words that children used in the immediate retell  
598 and their score on the target-word test was significant,  $\tau(38) = 0.348$ ,  $p = 0.011$ . This trend  
599 was significant in the *Expressive* condition ( $\tau(19) = 0.406$ ,  $p = 0.031$ ), but not in the *Flat*  
600 condition ( $\tau(17) = 0.246$ ,  $p = 0.251$ ) (Figure 4).

601 In summary, children tended to use more of the target words encountered in the story  
602 version they heard, and older children tended to use more of the target words.

### 603 4.1.3 Story length

604 The length of the story told by the robot was 365 words. In the immediate retell, the mean  
605 length of children's stories was 200.7 words ( $SD = 80.8$ ). No statistically significant  
606 difference in story length was observed between the two conditions (*Expressive*:  $M = 191.8$   
607 words,  $SD = 82.5$ , *Flat*:  $M = 210.6$ ,  $SD = 79.9$ ),  $t(38) = -0.73$ ,  $p = 0.47$ . Story length also did  
608 not vary with age, Pearson's  $r(7) = 0.06$ ,  $p = 0.71$ .

609 A 2x2 mixed ANOVA with time (within: Immediate vs. Delayed) and condition (between:  
610 *Expressive* vs. *Flat*) for the subset of children who produced both immediate and delayed  
611 retells revealed significant main effects of time,  $F(1,27) = 17.9$ ,  $p < 0.001$ ,  $\eta^2 = 0.398$ , as well  
612 as a significant interaction between time and condition,  $F(1,27) = 15.0$ ,  $p < 0.001$ ,  $\eta^2 =$   
613  $0.357$ . In the delayed retell, the overall length of children's story decreased to  $M = 147.9$   
614 ( $SD = 58.3$ ;  $t(13) = 5.35$ ,  $p < 0.001$ ). Children in the *Flat* condition showed a significant  
615 decrease (Immediate:  $M = 210.9$ ,  $SD = 85.4$ , Delayed:  $M = 125.4$ ,  $SD = 57.2$ ), while in the  
616 *Expressive* condition, the decrease was not statistically significant (Immediate:  $M = 173.3$ ,  
617  $SD = 79.33$ , Delayed:  $M = 168.9$ ,  $SD = 52.8$ ;  $t(14) = 0.33$ ,  $p = 0.75$ ). Furthermore, the length  
618 of stories in the two conditions were significantly different at the delayed retelling  
619 (*Expressive*:  $M = 168.9$ ,  $SD = 52.8$ , *Flat*:  $M = 125.4$ ,  $SD = 57.2$ ),  $t(27) = 2.13$ ,  $p = 0.043$ .

620 Thus, children in the *Flat* condition told shorter stories at the delayed retell as compared to  
621 the immediate retell whereas no such reduction was seen among children in the *Expressive*  
622 condition. Their stories were just as lengthy after 1-2 months (Figure 5). To further  
623 understand the impact of expressivity on retelling, we analyzed children's phrase  
624 production as reported in the following section.

## 625 4.2 Behavior

### 626 4.2.1 Responses to the robot's dialogic questions

627 Forty-two children had data regarding their responses to the robot-posed dialogic reading  
628 questions. Thirty-five (83.3%) responded to at least some of the questions; twenty-three

629 (54%) responded to all eleven questions; seven (16.7%) responded to none. There was no  
630 significant difference between the number of questions responded to by children in the  
631 *Expressive* and *Flat* conditions.

632 A simple linear regression model revealed that children who had responded to the robot's  
633 dialogic questions were likely to correctly identify more of the target vocabulary words,  
634  $F(1,38) = 5.84, p = 0.021, \eta^2 = 0.118$ . The interaction between the condition and the number  
635 of questions responded showed a trend,  $F(1,38) = 4.094, p = 0.0501, \eta^2 = 0.083$ , such that  
636 question answering in the *Expressive* condition was related to correct identification of  
637 target words, while question answering was not related to correct identification of words  
638 in the *Flat* condition. The correlation was driven primarily by the *Expressive* condition,  
639  $r(20) = 0.619, p = 0.002$ , i.e., children in the *Expressive* condition who answered the robot's  
640 questions were more likely to identify more of the target words; there was no significant  
641 correlation for the *Flat* condition,  $r(18) = 0.134, p = 0.57$  (Figure 6). Thus, answering the  
642 dialogic questions was linked to better vocabulary learning, but this link was only found in  
643 the *Expressive* condition.

644 Children who answered more dialogic questions also used significantly more target words  
645 in the immediate story retell as indicated by a Spearman's rank-order correlation  $r_s(38) =$   
646  $0.352, p = 0.026$ . These children also told longer stories,  $r_s(38) = 0.447, p = 0.003$  (Figure 7-  
647 A). They displayed greater emulation of the robot in terms of phrase usage,  $r_s(38) = 0.320,$   
648  $p = 0.044$ , but again this was driven primarily by the *Expressive* condition,  $r_s(19) = 0.437, p$   
649  $= 0.048$ , and not by the *Flat* condition,  $r_s(17) = 0.274, p = 0.257$ . Children in the *Expressive*  
650 condition also showed significant correlation to phrase usage in the delayed retell,  $r_s(13) =$   
651  $0.554, p = 0.032$  (Figure 7-B).

652 From the above observations, we can conclude that children were, in general, actively  
653 engaged in the robot's storytelling. When children were more engaged, as indexed by how  
654 often they responded to the robot's questions, their vocabulary learning was greater, and  
655 they were more likely to emulate the robot. However, these links between engagement and  
656 learning were evident in the *Expressive* rather than the *Flat* condition.

#### 657 **4.2.2 Emulating the robot's story**

658 An analysis of children's overall word usage reveals their word-level mirroring of the  
659 robot's story. In total, the robot used 96 unique words after stopword removal and the  
660 calculation of non-overlapping words. In the immediate retell, children used a mean of 58.7  
661 words ( $SD = 12.4$ ) emulating the robot. There was no significant difference between  
662 conditions. In the delayed retell, however, children in the *Expressive* condition used more  
663 words emulating the robot than children in the *Flat* condition (*Expressive*:  $M = 48.6, SD =$   
664  $13.5, Flat$ :  $M = 38.7, SD = 8.62; t(27) = 2.33, p = 0.028$ ).

665 We also analyzed the phrase-level similarity between the robot's story and the children's  
666 stories. In the immediate retell, a mean of 5.63 phrases ( $SD = 3.55$ ) were matched. A  
667 statistically significant difference was observed between conditions (*Expressive*:  $M = 6.67,$   
668  $SD = 3.98, Flat$ :  $M = 4.47, SD = 2.65; t(38) = 2.03, p = 0.049$ ) with robot's expressivity  
669 increasing children's phrase-level similarity. In the delayed retell, the overall usage of  
670 matched phrases decreased ( $M = 3.34, SD = 2.26$ ),  $t(28) = 5.87, p < 0.001$ . However, a Mann-  
671 Whitney U test showed that participants in the *Expressive* condition ( $M = 4.20, SD = 2.40$ )  
672 continued to use more similar phrases than participants in the *Flat* condition ( $M = 2.42, SD$



673 = 1.74),  $Z = 2.07$ ,  $p = 0.039$ ,  $r = 0.38$  (Figure 8). Thus, at both retellings, children were more  
674 likely to echo the expressive than the flat robot in terms of their phrasing.

675 The overall correlation between children's score on the target-word vocabulary test and  
676 the number of matched phrases they used in the retell was significant both for the  
677 immediate retell,  $r_s(38) = 0.375$ ,  $p = 0.017$ ; and for the delayed retell,  $r_s(27) = 0.397$ ,  $p =$   
678  $0.033$  (Figure 9). However, further analysis showed that this link was significant in the  
679 *Expressive* condition (Immediate:  $r_s(19) = 0.497$ ,  $p = 0.022$ ; Delayed:  $r_s(13) = 0.482$ ,  $p =$   
680  $0.031$ ), but not in the *Flat* condition (Immediate:  $r_s(19) = 0.317$ ,  $p = 0.186$ ; Delayed:  $r_s(12) =$   
681  $0.519$ ,  $p = 0.067$ ).

682 In summary, children were more likely to use similar words and phrases as the robot in the  
683 *Expressive* than in the *Flat* condition during both retellings. Furthermore, given that scores  
684 on the target-word vocabulary test were not significantly different between the two  
685 conditions, the correlation results suggest that the robot's expressivity did not impact  
686 initial encoding, but did encourage children to emulate the robot in their subsequent  
687 retelling of the story.

### 688 4.3 Engagement

#### 689 4.3.1 Interview questions

690 We found no difference between conditions in children's responses to the interview  
691 questions. Children reported that they liked the story (*Median* = 5, *Mode* = 5, *Range* = 1-5,  
692 *Inter-Quartile Range (IQR)* = 1) and that they liked Tega (*Median* = 5, *Mode* = 5, *Range* = 3-5,  
693 *IQR* = 0). For example, one child said he liked the story because "in the end they found  
694 a new pet frog." Children's reasons for liking Tega included physical characteristics, such as  
695 "furry," "cute," and "red", as well as personality traits including "kind" and "nice."

696 Children thought Tega could help them feel better (*Median* = 5, *Mode* = 5, *Range* = 1-5, *IQR* =  
697 0), saying, for example, that "he's cute, funny, and makes me smile," and "would give a big  
698 hug." They thought Tega helped them learn the story (*Median* = 5, *Mode* = 5, *Range* = 2-5,  
699 *IQR* = 0). One child reported Tega was helpful because "the story was a little bit long", while  
700 another said "because he asked me what happened in the story." Another child also noted  
701 the questions, saying "stopped to ask questions and talked slowly so I could understand".  
702 Children thought their friends would like reading with Tega (*Median* = 5, *Mode* = 5, *Range* =  
703 1-5, *IQR* = 0), because "he's a nice robot and will be nice to them," and "Tega's got a lot of  
704 good stories, and is good at telling them."

705 When asked if they would prefer to play again with Tega or with the Toucan puppet, 26  
706 children picked Tega, 11 picked Toucan, and 8 either said "both", "not sure", or did not  
707 respond. They justified picking Toucan with reasons such as "Toucan didn't hear the story,"  
708 "because he fell asleep and is super, super soft," and "because she's very sleepy and never  
709 listens." They justified picking Tega with various reasons including "because Tega can  
710 listen and Toucan is just a puppet," "because she read the story to me," "because he's fun,"  
711 and "I like her." Thus, we see that children felt the desire to be fair in making sure Toucan  
712 got a chance to hear the story, and a desire to reciprocate Tega's sharing of a story with  
713 them, as well as expressing general liking for the robot.

714 When asked to describe Tega to a friend, 44% of children described the robot using  
715 positive traits (e.g., nice, helpful, smart, fun) in the *Expressive* condition and 48% in the *Flat*  
716 condition, *ns*. For example, one child said, "he told me about antlers. Tega is very helpful",

717 while another reported "that he read me a story and will be a nice robot to them". In sum,  
718 the expressiveness of the robot did not influence how children described the robot to a  
719 peer. Many of the other 56% of children in the *Expressive* condition and the 52% of children  
720 in the *Flat* condition focused on the robot's physical characteristics, for example, "red and  
721 blue, stripes, big eyes, tuft of blue hair, phone for face, fuzzy, cute smile". One child said  
722 Tega "looks like a rock star."

### 723 **4.3.2 Children's expressivity**

724 We analyzed affect data for 36 children (19 in the *Expressive* condition and 17 in the *Flat*  
725 condition). For the remaining nine children, no affect data were collected either because  
726 the children's faces were not detected by the system, or because of other system failures.  
727 As described earlier, we focused our analysis on the four affective states most relevant to  
728 our research questions: concentration, engagement, surprise, and attention. All other  
729 affective states were measured by Affdex very rarely (less than 5% of the time). We found  
730 that overall, children maintained attention throughout most of the session, were engaged  
731 by the robot, showed some concentration, and displayed surprise during the story (Table  
732 3).

733 To evaluate whether the robot's vocal expressiveness influenced children's facial  
734 expressiveness, we examined the mean levels of the four affective states across the entire  
735 session by condition. We conducted a one-way ANCOVA with condition (*Expressive* vs. *Flat*)  
736 for each Affdex score, with age as a covariate. The analysis revealed that children in the  
737 *Expressive* condition showed significantly higher mean levels of concentration,  $F(1, 32) =$   
738  $4.77, p = 0.036, \eta^2 = 0.127$ ; engagement,  $F(1, 32) = 4.15, p = 0.049, \eta^2 = 0.112$ ; and  
739 surprise,  $F(1, 32) = 5.21, p = 0.029, \eta^2 = 0.13$ , than children in the *Flat* condition, but that  
740 children's attention was not significantly different,  $F(1, 32) = 0.111, p = 0.741$ .

741 Furthermore, these differences were not affected by children's age (Table 3, Figure 10).  
742 The lack of difference in children's attention demonstrated that the differences in the  
743 concentration, engagement and surprise levels across the two conditions were not a result  
744 of children paying less attention to the *Flat* robot's story.

745 Next, we asked whether children's affect changed during the session. We split the affect  
746 data into two halves—the first half of the session and the second half of the session—using  
747 the data timestamps to determine the session halfway point. We ran a 2 x 2 mixed design  
748 ANOVA with time (within: first half vs. second half) x condition (between: *Expressive* vs.  
749 *Flat*) for each of the affect scores. These analyses revealed main effects of condition on  
750 children's concentration scores,  $F(1, 34) = 4.71, p = 0.037, \eta^2 = 0.067$ ; engagement scores,  
751  $F(1, 34) = 4.16, p = 0.049, \eta^2 = 0.075$ ; and surprise scores,  $F(1, 34) = 5.36, p = 0.027, \eta^2 =$   
752  $0.090$ . In all three cases, children displayed greater affect in the *Expressive* condition than  
753 the *Flat* condition (See Figure 11B-D). There were no main effects of time or any significant  
754 interactions for these affect measures. However, we did see a main effect of time for  
755 children's attention scores,  $F(1, 34) = 7.84, p = 0.008, \eta^2 = 0.044$ . In both conditions,  
756 children's attention scores declined over time (Figure 11A).

757 In summary, although all children were less attentive over time, they showed more facial  
758 expressiveness throughout the whole session with the expressive robot than with the flat  
759 robot.

## 760 **4.4 Language status**

761 We completed our analyses by checking whether the results were stronger or weaker  
762 based on children's language status (i.e., native English speakers, ELL, or bilingual). The  
763 differences were modest and are reported here.  
764 First, with regards to learning new vocabulary, a one-way ANOVA with age as a covariate  
765 revealed that children's language status affected how many target vocabulary words they  
766 identified correctly,  $F(3,37) = 4.10$ ,  $p = 0.012$ ,  $\eta^2 = 0.230$ , but vocabulary learning was not  
767 affected by age (Figure 12). Post-hoc pairwise comparisons showed that children who were  
768 native English speakers correctly identified more words ( $M = 4.53$  correct,  $SD = 1.23$ ) than  
769 ELL children ( $M = 3.13$  correct,  $SD = 1.30$ ),  $p = 0.002$ . Bilingual speakers also identified  
770 more words correctly ( $M = 4.86$ ,  $SD = 1.07$ ) than ELL children,  $p = 0.005$ , but were not  
771 significantly different from the native English speakers.  
772 Second, native English speakers used more target words ( $M = 0.94$ ,  $SD = 0.93$ ) than ELL  
773 students ( $M = 0.14$ ,  $SD = 0.36$ ) in the immediate retell,  $t(20) = -3.16$ ,  $p = 0.005$ . Bilingual  
774 students were in-between ( $M = 0.29$ ,  $SD = 0.49$ ). This trend was primarily driven by the  
775 *Expressive* condition,  $F(2,34) = 5.458$ ,  $p = 0.009$ , rather than the *Flat* condition. Post-hoc  
776 pairwise comparison within the *Expressive* condition showed that native speakers used  
777 more target words than bilingual speakers,  $t(8) = 3.00$ ,  $p = 0.017$ , and more than ELL  
778 speakers,  $t(13) = 7.45$ ,  $p < 0.001$ . Bilingual speakers also used more target words than the  
779 ELL group in the immediate retell,  $t(11) = 2.75$ ,  $p = 0.019$ .  
780 Lastly, native English speakers showed stronger phrase mirroring behavior ( $M = 6.56$ ,  $SD =$   
781  $4.49$ ) than ELL students ( $M = 4.43$ ,  $SD = 2.38$ ) in the *Expressive* condition in the immediate  
782 retell,  $t(13) = 3.41$ ,  $p = 0.005$ . The robot's expressivity had a significant effect on native  
783 English speakers' usage of similar phrases in both the immediate retell (*Expressive*  $M =$   
784  $10.17$ ,  $SD = 4.40$ , *Flat*  $M = 4.40$ ,  $SD = 2.99$ ),  $t(14) = 3.139$ ,  $p = 0.007$ ; and in the delayed retell  
785 (*Expressive*  $M = 5.50$ ,  $SD = 2.52$ , *Flat*  $M = 2.38$ ,  $SD = 1.30$ ),  $t(10) = 2.904$ ,  $p = 0.016$ . Though  
786 not significant, ELL children trended toward also using more similar phrases when they  
787 heard the story from the *Expressive* robot in both the immediate retell (*Expressive*  $M = 4.55$ ,  
788  $SD = 1.94$ , *Flat*  $M = 4.20$ ,  $SD = 3.27$ ) and the delayed retell (*Expressive*  $M = 3.78$ ,  $SD = 1.86$ ,  
789 *Flat*  $M = 2.20$ ,  $SD = 2.49$ ).  
790 In summary, as might be expected, native English speakers performed better on the  
791 vocabulary test, used more target words, and showed more phrase matching than either  
792 ELL or bilingual children.

## 793 **5 Discussion**

794 We asked whether children would learn from a dialogic, storytelling robot and whether the  
795 robot's effectiveness as a narrator and teacher would vary with the expressiveness of the  
796 robot's voice. We hypothesized that a more expressive voice would lead to greater  
797 engagement and greater learning. Below, we review the main findings pertinent to each of  
798 these questions and then turn to their implications.

799 Whether the robot spoke with a flat or expressive voice, children were highly attentive in  
800 listening to the robot—as indexed by their head orientation—when it was recounting the  
801 picture book story. Moreover, irrespective of the robot's voice, children were able to  
802 acquire new vocabulary items embedded in the story. Although some children may have  
803 already known some of the target words, as indicated by their above-zero recognition of  
804 the target words from the story version they did *not* hear, the interaction between story  
805 version heard and scores on each set of words (shown in Figure 2) shows that genuine

806 learning did occur. Children could also retell the story (with the help of the picture book)  
807 both immediately afterwards and some weeks later. At their initial retelling, children  
808 typically produced a story about half as long as the one they had heard, sometimes  
809 including a newly acquired vocabulary item. Finally, when they were invited to provide  
810 both an explicit evaluation and a free-form description, children were equally positive  
811 about the robot whether they had listened to the flat or the expressive robot.  
812 Despite this equivalence with respect to attentiveness, encoding and evaluation, there were  
813 several indications that children's mode of listening was different for the two robots. First,  
814 as they listened to the expressive rather than the flat robot, children's facial expressions  
815 betrayed more concentration (i.e., more brow furrowing and less smiling), more  
816 engagement (i.e., greater overall muscle activation) and more surprise (i.e., more brow  
817 raising with open mouth). Thus, children were not only attentive to what the robot was  
818 saying, they also displayed signs of greater emotional engagement.  
819 Furthermore, inclusion of the newly acquired vocabulary items in the initial retelling was  
820 more frequent among children who listened to the expressive rather than the flat robot.  
821 Note that children's score on the target-word test was not significantly different between  
822 the two conditions, suggesting that children who correctly identified the target words in  
823 the *Expressive* condition tended to also use them in their story retell whereas children who  
824 correctly identified the target words in the *Flat* condition were less likely to use them in  
825 their story recall. Thus, although children were able to acquire new vocabulary from either  
826 robot (*receptive* vocabulary knowledge), they were more likely to subsequently use that  
827 vocabulary in their stories if the expressive robot had been the narrator (i.e., *productive*  
828 vocabulary knowledge). This pattern of findings implies that children could encode and  
829 retain new input from either robot, but they were more likely to engage with the  
830 expressive robot during the narration and more likely to emulate the expressive robot's  
831 narrative vocabulary in their own recounting. That is, interacting with the expressive robot  
832 led to greater *behavioral* outcomes—producing new words rather than merely identifying  
833 them.  
834 Further signs of the differential impact of the two robots were found at the delayed  
835 retelling. Whereas there was a considerable decline in story length among children who  
836 had heard the flat robot, there was no such decline among children who had heard the  
837 expressive robot. Again, we cannot ascribe this difference to differences in encoding.  
838 Children in each condition had told equally long stories on their initial retelling. A more  
839 plausible interpretation is that children who had heard the expressive robot were more  
840 inclined to emulate its narrative than children who had heard the flat robot. More detailed  
841 support for this interpretation emerged in children's story phrasing. At both retellings,  
842 children were more likely to echo the expressive rather than the flat robot in terms of using  
843 parallel phrases. This may also indicate that children were engaging with the expressive  
844 robot as a more socially dynamic agent, since past research has shown that children are  
845 more likely to use particular syntactic forms when primed by an adult (e.g., Huttenlocher et  
846 al., 2004). In addition, recent work by Kennedy et al. (2017) showed that a robot that used  
847 more nonverbal immediacy behaviors (e.g., gestures, gaze, vocal prosody, facial  
848 expressions, proximity and body orientation, touch) led to greater short story recall by  
849 children. The difference in the expressive versus flat robot's vocal qualities (e.g., intonation  
850 and prosody) could have led to a difference in the perceived nonverbal immediacy of the

851 robot, which may have led to the differences in children’s engagement with the robot as a  
852 socially dynamic agent.

853 Both the expressive and the flat robot asked dialogic questions about the story as they  
854 narrated it. The more often children answered these dialogic questions the more  
855 vocabulary items they learned. Here too, however, the robot’s voice made a difference. The  
856 link between question answering and vocabulary acquisition was only significant for the  
857 *Expressive* condition. Children who answered more dialogic questions also displayed  
858 greater fidelity to the robot’s story in terms of phrase usage when they retold it, but again  
859 this link was only significant for the *Expressive* condition. Thus, answering more of the  
860 robot’s questions was associated with the acquisition of more vocabulary and greater  
861 phrase emulation but only for the expressive robot. Finally children’s score on the target-  
862 word vocabulary test correlated with the number of matched phrases they used at both  
863 retellings. However, this correlation emerged only for children in the expressive condition,  
864 again consistent with the idea that robot expressivity enhanced emulation but not initial  
865 encoding.

866 In sum, we obtained two broad patterns of results. On the one hand, both robots were  
867 equally successful in capturing children’s attention, telling a story that children were  
868 subsequently able to narrate, and teaching the children new vocabulary items. On the other  
869 hand, as compared to the flat robot, the expressive robot provoked stronger emotional  
870 engagement in the story as it was being narrated, greater inclusion of the newly learned  
871 vocabulary into the retelling of the story and greater fidelity to the original story during the  
872 retelling. A plausible interpretation of these two patterns is that story narration per se was  
873 sufficient to capture children’s attention and sufficient to ensure encoding both of the story  
874 itself and of the new vocabulary. By contrast the mode in which the story was narrated—  
875 expressive or flat—impacted the extent to which the child eventually cast him or herself  
876 into the role enacted by the narrator. More specifically, it is plausible that children who  
877 were emotionally engaged by the expressive robot were more prone to re-enact the story-  
878 telling mode of the robot when it was their turn to tell the story to the puppet: they were  
879 more likely to reproduce some of the unfamiliar nouns that they had heard the robot use  
880 and more likely to mimic the specific phrases included in the robot’s narrative.

881 It is tempting to conclude that children identified more with the expressive robot and  
882 found it more appealing. It is important to emphasize, however, that no signs of that  
883 differentiation were apparent either in children’s explicit verbal judgments about the two  
884 robots or in the open-ended descriptions. In either case, children were quite positive about  
885 both of the robots. An important implication of these findings, therefore, is that children’s  
886 verbal ratings of the robots are not a completely accurate guide to the effectiveness of the  
887 robots as role models. Future research on social robots as companions and pedagogues  
888 should pay heed to such findings. More generally, the results indicate that it is important to  
889 assess the influence and impact of a robot via a multiplicity of measures rather than via  
890 questionnaires or self-report.

## 891 **5.1 Language Status**

892 When analyzing children’s learning and performance based on their language status, we  
893 saw only modest differences. These differences—in which native English speakers and  
894 bilingual children correctly identified more target vocabulary words with both robots, and

895 showed stronger phrase mirroring and use more target words than ELL students with the  
896 expressive robot—were not unexpected, given that bilingual and native English speaking  
897 children have greater familiarity with the language. Nevertheless, it is important to note  
898 that both native English speakers and ELL children who heard the story from the *Expressive*  
899 robot reused and retained more information from the robot’s story. Thus, despite the  
900 limitations listed in the following section, these results suggest that the storytelling activity  
901 was an effective intervention for all the native English speakers, the bilingual and the ELL  
902 children, leading to learning and engagement by all groups. This is an important finding  
903 given that ELL children arguably need the most additional support for their language  
904 development (Paez et al., 2007). Effective and engaging language learning interventions  
905 like this one that can benefit the entire classroom—native English speakers, bilingual, and  
906 ELL children alike—will be important educational tools in years to come.

## 907 **5.2 Limitations**

908 We should note several limitations of this study. First, some potentially important  
909 individual differences among children, such as their learning ability, socio-economic status,  
910 and sociability were not controlled for. Second, although 45 children participated in the  
911 study ranging from 4 to 7 years, we did not have an equal number of children at each age.  
912 We also did not have an equal number of children with each language status. In future  
913 work, it will be important to assess a more homogenous sample, as well as the degree to  
914 which our results remain stable across these individual differences and across the  
915 preschool and elementary school years.

916 In addition, we did not have complete story retelling data for all children. As reported  
917 earlier, the audio quality of some of the recordings of children’s retells prevented analysis,  
918 and not all children performed delayed retellings. As a result of this and the  
919 aforementioned imbalances in age and language, the analyses we report here are under-  
920 powered. This is exploratory work, and the result should be interpreted in light of this fact.  
921 Future work should take greater effort to collect quality audio recordings and to see all  
922 children at the delayed test.

923 Finally, while the target vocabulary words used were uncommon, some children may still  
924 have known them—particularly older children, given the correlation between age and  
925 target words identified. The rarity of the words may have also increased their saliency,  
926 being a cue for children to pay attention to the words. Follow-up studies should either  
927 consider using nonce words or include a vocabulary pretest for the target words.

## 928 **5.3 Future Directions**

929 The technology landscape continues to rapidly evolve from passively consumed content  
930 such as television and radio to interactive and social experiences enabled through digital  
931 technology and the internet. Each new technology transforms the ways we interact with  
932 one another, how we communicate and share, how we learn, tell stories, and experience  
933 imaginary worlds.

934 Today, the linguistic and interpersonal environment of children is comprised of other  
935 people, yet children are increasingly growing up talking *with* AI-based technologies, too.  
936 Despite the proliferation of such technologies, very little is understood about children’s

937 language acquisition in this emerging social-technological landscape. While it has been  
938 argued in the past that children cannot learn language from impersonal media because  
939 language acquisition is socially gated (e.g., Kuhl, 2007; Kuhl, 2011), the reality of social  
940 robots forces us to revisit our past assumptions. These assumptions need revisiting  
941 because numerous studies have now shown that children and adults interact with social  
942 robots as social others (Breazeal et al., 2008; Breazeal et al., 2016; DeSteno et al., 2012).  
943 Social robots represent a new and provocative psychological category betwixt and between  
944 inanimate things and socially animate beings. They bridge the digital world of content and  
945 information to the physically co-present and interpersonal world of people. Because of this,  
946 we are likely to interact with social robots differently than prior technologies. As such,  
947 social robots open new opportunities for how educational content and experiences can be  
948 brought to the general public, just as their technological predecessors have.  
949 Therefore, how *should* social robots be designed to best foster the learning, development,  
950 and benefit to children? This is very new territory, indeed. This work explores three key  
951 avenues, although there are many others to explore, and to explore deeply.  
952 In the context of language learning for preschool age children, we begin by applying  
953 knowledge and taking inspiration from how children learn language through storytelling  
954 with a peer-like companion. Children learn quite a lot from interacting with and socially  
955 modeling the behavior and attitudes of their peers, and in prior work, we have seen  
956 behaviors suggesting that children also socially model or emulate the behavior of social  
957 robots. For instance, we have found that children become more emotively expressive when  
958 a robot is more expressive (Spaulding et al., 2016). We see this effect again in this work. We  
959 have also observed that when children play with a “curious” robot that exhibits pro-curious  
960 behaviors and attitudes, children express and engage in more curious behaviors (Gordon et  
961 al., 2015) and are more willing to teach new tasks to a robot peer (Park et al., 2015). In the  
962 present study, we found evidence of this social modeling effect in terms of children  
963 emulating the linguistic phrases and vocabulary a robot uses. This peer-learning dynamic is  
964 quite different from how children learn with other technologies.  
965 Emotional expressivity is another characteristic that social robots bring to interaction.  
966 Understanding the impact of emotion and expressive behavior on learning with young  
967 children is an area worth further systematic investigation. It is generally accepted that  
968 telling a story more expressively will make it more engaging. Social robots enable us to  
969 study the impact of expressivity on children’s behavior and learning in a more systematic  
970 and carefully controlled way. Because of these attributes, social robots could serve as a  
971 compelling tool to gain insights into children’s social development and learning.  
972 In this work, we observed a greater tendency for children to emulate a storytelling robot’s  
973 phrasing when the robot was more vocally expressive; children reproduced this pattern  
974 after a month-long delay. Further research is warranted to understand whether children  
975 are encoding the information differently when the delivery is more expressive, or whether  
976 they are simply more apt to emulate the robot when it is more expressive.  
977 We see growing evidence that the more socially expressive and interactive a robot is, the  
978 more it “opens the spigot” to children’s social engagement and learning. This suggests a  
979 new paradigm for educational technology and how it promotes children’s learning and  
980 development. It is increasingly clear that it is not just the introduction of a social robot into  
981 an educational context that matters, *but how socially designed* the robot is that impacts  
982 children's behavior and learning.

983 Finally, for social robots to have a large-scale impact in the educational realm, research  
984 should extend beyond the context of 1:1 interaction of a social robot with a child. We need  
985 to also understand how to design social robots to support and foster peer learning among  
986 groups of children. We need to understand how social robots can best support and include  
987 the participation of adults, such as teachers and parents, or facilitate classroom  
988 orchestration (e.g., Dillenbourg et al., 2010). And we need to understand how to effectively  
989 integrate robots into the broader educational context of the classroom and continued  
990 learning at home. Much work remains to be done in order to understand how to best  
991 design social robots that can successfully engage and support learning over longitudinal  
992 time scales, where the opportunity to deeply attune to the individual child exists—not only  
993 in terms of curricular goals, but in order to foster positive attitudes toward learning and  
994 challenge, and to build trust and rapport as well. Finally, as research matures and social  
995 robots become an affordable mass consumer technology, there exists many opportunities  
996 for social robots to help support and augment learning experiences for children who are  
997 underserved, at-risk, or have other learning challenges.

## 998 **6 Author Contribution Statement**

999 The study was conceived and designed by JK, SJ, HP, SR, PH, DD, CB. Data analysis was  
1000 performed by JK, SJ, HP, SR, AA. The paper was drafted, written, revised, and approved by  
1001 JK, SJ, HP, SR, AA, PH, DD, CB.

## 1002 **7 Acknowledgements**

1003 We thank Mirko Gelsomini for his help in data collection. This research was supported by  
1004 the National Science Foundation (NSF) under Grant IIS-1122886, IIS-1122845, and IIS-  
1005 1123085. Any opinions, findings and conclusions, or recommendations expressed in this  
1006 paper are those of the authors and do not represent the views of the NSF.

## 1007 **8 References**

- 1008 Bartneck, C., Kulić, D., Croft, E., and Zoghbi, S. (2008). Measurement Instruments for the  
1009 Anthropomorphism, Animacy, Likeability, Perceived Intelligence, and Perceived Safety  
1010 of Robots. *Int J of Soc Robotics* 1, 71–81. doi:10.1007/s12369-008-0001-3.
- 1011 Baxter, P., Kennedy, J., Senft, E., Lemaignan, S., and Belpaeme, T. (2016). From  
1012 Characterising Three Years of HRI to Methodology and Reporting Recommendations.  
1013 in *The Eleventh ACM/IEEE International Conference on Human Robot Interaction HRI*  
1014 '16. (Piscataway, NJ, USA: IEEE Press), 391–398. Available at:  
1015 <http://dl.acm.org/citation.cfm?id=2906831.2906897> [Accessed August 23, 2016].
- 1016 Bloom, P. (2000). *How children learn the meanings of words*. Cambridge, MA: MIT press.
- 1017 Boteanu, A., Chernova, S., Nunez, D., and Breazeal, C. (2016). Fostering parent–child dialog  
1018 through automated discussion suggestions. *User Model User-Adap Inter*, 1–31.  
1019 doi:10.1007/s11257-016-9176-8.
- 1020 Boudreau, D. M., and Hedberg, N. L. (1999). A Comparison of Early Literacy Skills in  
1021 Children With Specific Language Impairment and Their Typically Developing Peers.  
1022 *American Journal of Speech-Language Pathology* 8, 249. doi:10.1044/1058-  
1023 0360.0803.249.



1024 Breazeal, C., Dautenhahn, K., and Kanda, T. (2016a). "Social Robotics," in *Springer Handbook*  
1025 *of Robotics*, eds. B. Siciliano and O. Khatib (Springer International Publishing), 1935–  
1026 1972. doi:10.1007/978-3-319-32552-1\_72.

1027 Breazeal, C., Harris, P. L., DeSteno, D., Kory Westlund, J. M., Dickens, L., and Jeong, S.  
1028 (2016b). Young Children Treat Robots as Informants. *Top Cogn Sci*, 1–11.  
1029 doi:10.1111/tops.12192.

1030 Breazeal, C. L. (2004). *Designing Sociable Robots*. MIT Press.

1031 Breazeal, C., Takanishi, A., and Kobayashi, T. (2008). "Social Robots that Interact with  
1032 People," in *Springer Handbook of Robotics*, eds. Bruno Siciliano and Oussama Khatib  
1033 (Springer Berlin Heidelberg), 1349–1369. doi:10.1007/978-3-540-30301-5\_59.

1034 Camras, L. A., Bakeman, R., Chen, Y., Norris, K., and Cain, T. R. (2006). Culture, ethnicity, and  
1035 children's facial expressions: A study of European American, mainland Chinese, Chinese  
1036 American, and adopted Chinese girls. *Emotion* 6, 103–114. doi:10.1037/1528-3542.6.1.103.

1037 Chang, A., Breazeal, C., Faridi, F., Roberts, T., Davenport, G., Lieberman, H., et al. (2012).  
1038 Textual Tinkerability: Encouraging Storytelling Behaviors to Foster Emergent Literacy.  
1039 in *CHI '12 Extended Abstracts on Human Factors in Computing Systems CHI EA '12*.  
1040 (New York, NY, USA: ACM), 505–520. doi:10.1145/2212776.2212826.

1041 Chang, C.-W., Lee, J.-H., Chao, P.-Y., Wang, C.-Y., and Chen, G.-D. (2010). Exploring the  
1042 possibility of using humanoid robots as instructional tools for teaching a second  
1043 language in primary school. *Educational Technology & Society* 13, 13–24.

1044 Corriveau, K. H., Harris, P. L., Meins, E., Fernyhough, C., Arnott, B., Elliott, L., et al. (2009).  
1045 Young children's trust in their mothers' claims: Longitudinal links with attachment  
1046 security in infancy. *Child development* 80, 750–761.

1047 Cremin, T., Flewitt, R., Mardell, B., and Swann, J. (2016). *Storytelling in Early Childhood:*  
1048 *Enriching language, literacy and classroom culture*. Routledge.

1049 DeSteno, D., Breazeal, C., Frank, R. H., Pizarro, D., Baumann, J., Dickens, L., et al. (2012).  
1050 Detecting the Trustworthiness of Novel Partners in Economic Exchange. *Psychological*  
1051 *science* 23, 1549–1556.

1052 Diehl, J. J., Bennetto, L., and Young, E. C. (2006). Story Recall and Narrative Coherence of  
1053 High-Functioning Children with Autism Spectrum Disorders. *J Abnorm Child Psychol*  
1054 34, 83–98. doi:10.1007/s10802-005-9003-x.

1055 Dillenbourg, P., and Jermann, P. (2010). "Technology for Classroom Orchestration," in *New*  
1056 *Science of Learning*, eds. M. S. Khine and I. M. Saleh (Springer New York), 525–552.  
1057 doi:10.1007/978-1-4419-5716-0\_26.

1058 Dunn, L. M., and Dunn, L. M. (2007). *Peabody Picture Vocabulary Test. 4th ed.* Pearson  
1059 Assessments.

1060 Ekman, P., Friesen, W. V., and Hager, J. C. (1978). Facial action coding system (FACS). *A*  
1061 *technique for the measurement of facial action. Consulting, Palo Alto* 22.

1062 Ekman, P., and Rosenberg, E. L. (1997). *What the face reveals: Basic and applied studies of*  
1063 *spontaneous expression using the Facial Action Coding System (FACS)*. Oxford University  
1064 Press, USA.

1065 Eyssel, F., Kuchenbrandt, D., Bobinger, S., de Ruitter, L., and Hegel, F. (2012). "If You Sound  
1066 Like Me, You Must Be More Human": On the Interplay of Robot and User Features on  
1067 Human-robot Acceptance and Anthropomorphism. in *Proceedings of the Seventh*  
1068 *Annual ACM/IEEE International Conference on Human-Robot Interaction HRI '12*. (New  
1069 York, NY, USA: ACM), 125–126. doi:10.1145/2157689.2157717.

1070 Feil-Seifer, D., and Mataric, M. J. (2011). Socially Assistive Robotics. *IEEE Robotics*  
1071 *Automation Magazine* 18, 24–31. doi:10.1109/MRA.2010.940150.

1072 Friesen, W. V., and Ekman, P. (1983). EMFACS-7: Emotional facial action coding system.  
1073 *Unpublished manuscript, University of California at San Francisco* 2, 1.

1074 Gordon, G., Breazeal, C., and Engel, S. (2015). Can Children Catch Curiosity from a Social  
1075 Robot? in *Proceedings of the Tenth Annual ACM/IEEE International Conference on*  
1076 *Human-Robot Interaction HRI '15*. (New York, NY, USA: ACM), 91–98.  
1077 doi:10.1145/2696454.2696469.

1078 Gordon, G., Spaulding, S., Kory Westlund, J., Lee, J. J., Plummer, L., Martinez, M., et al. (2016).  
1079 Affective personalization of a social robot tutor for children’s second language skill. in  
1080 *Proceedings of the 30th AAAI Conference on Artificial Intelligence* (Palo Alto, CA).

1081 Greenhalgh, K. S., and Strong, C. J. (2001). Literate Language Features in Spoken Narratives  
1082 of Children with Typical Language and Children with Language Impairments.  
1083 *Language, Speech, and Hearing Services in Schools* 32, 114–25.

1084 Hargrave, A. C., and Sénéchal, M. (2000). A book reading intervention with preschool  
1085 children who have limited vocabularies: the benefits of regular reading and dialogic  
1086 reading. *Early Childhood Research Quarterly* 15, 75–90. doi:10.1016/S0885-  
1087 2006(99)00038-1.

1088 Harris, P. L. (2007). Trust. *Developmental science* 10, 135–138.

1089 Harris, P. L., Guz, G. R., Lipian, M. S., and Man-Shu, Z. (1985). Insight into the Time Course of  
1090 Emotion among Western and Chinese Children. *Child Development* 56, 972–988.  
1091 doi:10.2307/1130109.

1092 Hart, B., and Risley, T. R. (1995). *Meaningful differences in the everyday experience of young*  
1093 *American children*. ERIC.

1094 Heilmann, J., Miller, J. F., Nockerts, A., and Dunaway, C. (2010). Properties of the Narrative  
1095 Scoring Scheme Using Narrative Retells in Young School-Age Children. *American*  
1096 *Journal of Speech-Language Pathology* 19, 154–166. doi:10.1044/1058-  
1097 0360(2009/08-0024).

1098 Hood, D., Lemaignan, S., and Dillenbourg, P. (2015). When Children Teach a Robot to Write:  
1099 An Autonomous Teachable Humanoid Which Uses Simulated Handwriting. in  
1100 *Proceedings of the Tenth Annual ACM/IEEE International Conference on Human-Robot*  
1101 *Interaction HRI '15*. (New York, NY, USA: ACM), 83–90.  
1102 doi:10.1145/2696454.2696479.

1103 Huttenlocher, J., Vasilyeva, M., and Shimpi, P. (2004). Syntactic priming in young children.  
1104 *Journal of Memory and Language* 50, 182–195. doi:10.1016/j.jml.2003.09.003.

1105 Isbell, R., Sobol, J., Lindauer, L., and Lowrance, A. (2004). The Effects of Storytelling and Story  
1106 Reading on the Oral Language Complexity and Story Comprehension of Young Children.  
1107 *Early Childhood Education Journal* 32, 157–163.  
1108 doi:10.1023/B:ECEJ.0000048967.94189.a3.

1109 Kahn, P. H., Gary, H. E., and Shen, S. (2013). Children’s Social Relationships With Current  
1110 and Near & Future Robots. *Child Development Perspectives* 7, 32–37.

1111 Kennedy, J., Baxter, P., and Belpaeme, T. (2017). Nonverbal Immediacy as a Characterisation of  
1112 Social Behaviour for Human–Robot Interaction. *Int J of Soc Robotics* 9, 109–128.  
1113 doi:10.1007/s12369-016-0378-3.

1114 Kennedy, J., Baxter, P., Senft, E., and Belpaeme, T. (2016). Social robot tutoring for child  
1115 second language learning. in *2016 11th ACM/IEEE International Conference on Human-*

1116 *Robot Interaction (HRI)* (IEEE), 231–238. Available at:  
1117 [http://ieeexplore.ieee.org/xpls/abs\\_all.jsp?arnumber=7451757](http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=7451757) [Accessed May 20,  
1118 2016].

1119 Kory, J. (2014). Storytelling with robots: Effects of robot language level on children’s  
1120 language learning.

1121 Kory, J., and Breazeal, C. (2014). Storytelling with robots: Learning companions for  
1122 preschool children’s language development. in *2014 RO-MAN: The 23rd IEEE*  
1123 *International Symposium on Robot and Human Interactive Communication*, 643–648.  
1124 doi:10.1109/ROMAN.2014.6926325.

1125 Kory Westlund, J., and Breazeal, C. (2015). The Interplay of Robot Language Level with  
1126 Children’s Language Learning During Storytelling. in *Proceedings of the Tenth Annual*  
1127 *ACM/IEEE International Conference on Human-Robot Interaction Extended Abstracts*  
1128 *HRI’15 Extended Abstracts*. (New York, NY, USA: ACM), 65–66.  
1129 doi:10.1145/2701973.2701989.

1130 Kory Westlund, J. M., Lee, J. J., Plummer, L., Faridi, F., Gray, J., Berlin, M., et al. (2016). Tega: A  
1131 Social Robot. in *2016 11th ACM/IEEE International Conference on Human-Robot*  
1132 *Interaction (HRI)* doi:10.1109/HRI.2016.7451856.

1133 Kuhl, P. K. (2007). Is speech learning “gated” by the social brain? *Developmental science* 10,  
1134 110–120.

1135 Kuhl, P. K. (2011). Social Mechanisms in Early Language Acquisition: Understanding  
1136 Integrated Brain Systems Supporting Language. Available at:  
1137 <http://www.oxfordhandbooks.com/view/10.1093/oxfordhb/9780195342161.001.0001/oxfordhb-9780195342161-e-043> [Accessed September 12, 2016].

1138 01/oxfordhb-9780195342161-e-043 [Accessed September 12, 2016].

1139 Lee, J. J., Knox, W. B., Wormwood, J. B., Breazeal, C., and DeSteno, D. (2013).  
1140 Computationally modeling interpersonal trust. *Front Psychol* 4.

1141 Leite, I., Castellano, G., Pereira, A., Martinho, C., and Paiva, A. (2014). Empathic Robots for  
1142 Long-term Interaction. *Int J of Soc Robotics* 6, 329–341. doi:10.1007/s12369-014-0227-1.  
1143 doi:10.3389/fpsyg.2013.00893.

1144 LoBue, V., and Thrasher, C. (2015). The Child Affective Facial Expression (CAFE) set: validity  
1145 and reliability from untrained adults. *Front. Psychol.* 5. doi:10.3389/fpsyg.2014.01532.

1146 Lubold, N., Walker, E., and Pon-Barry, H. (2016). Effects of Voice-Adaptation and Social  
1147 Dialogue on Perceptions of a Robotic Learning Companion. in *The Eleventh ACM/IEEE*  
1148 *International Conference on Human Robot Interaction HRI ’16*. (Piscataway, NJ, USA:  
1149 IEEE Press), 255–262. Available at:  
1150 <http://dl.acm.org/citation.cfm?id=2906831.2906876> [Accessed September 14, 2016].

1151 Lyberg-Åhländer, V., Haake, M., Brännström, J., Schötz, S., and Sahlén, B. (2015). Does the  
1152 speaker’s voice quality influence children’s performance on a language comprehension  
1153 test? *International Journal of Speech-Language Pathology* 17, 63–73.  
1154 doi:10.3109/17549507.2014.898098.

1155 McDuff, D., Kaliouby, R. el, and Picard, R. W. (2015). Crowdsourcing facial responses to online  
1156 videos: Extended abstract. in *2015 International Conference on Affective Computing and*  
1157 *Intelligent Interaction (ACII)*, 512–518. doi:10.1109/ACII.2015.7344618.

1158 McDuff, D., Kaliouby, R. el, Senechal, T., Amr, M., Cohn, J. F., and Picard, R. (2013).  
1159 Affectiva-MIT Facial Expression Dataset (AM-FED): Naturalistic and Spontaneous Facial  
1160 Expressions Collected #x0022;In-the-Wild #x0022; in *2013 IEEE Conference on Computer*  
1161 *Vision and Pattern Recognition Workshops*, 881–888. doi:10.1109/CVPRW.2013.130.

1162 McDuff, D., Mahmoud, A., Mavadati, M., Amr, M., Turcot, J., and Kaliouby, R. el (2016).  
1163 AFFDEX SDK: A Cross-Platform Real-Time Multi-Face Expression Recognition Toolkit.  
1164 in *Proceedings of the 2016 CHI Conference Extended Abstracts on Human Factors in*  
1165 *Computing Systems CHI EA '16*. (New York, NY, USA: ACM), 3723–3726.  
1166 doi:10.1145/2851581.2890247.

1167 Meltzoff, A. N., Kuhl, P. K., Movellan, J., and Sejnowski, T. J. (2009). Foundations for a new  
1168 science of learning. *Science* 325, 284–288.

1169 Movellan, J., Eckhardt, M., Virnes, M., and Rodriguez, A. (2009). Sociable robot improves  
1170 toddler vocabulary skills. in *Proceedings of the 4th ACM/IEEE international conference*  
1171 *on Human Robot Interaction* (ACM), 307–308.

1172 Naigles, L. R., and Mayeux, L. (2001). Television as incidental language teacher. *Handbook of*  
1173 *children and the media*, 135–152.

1174 Nass, C., and Brave, S. (2005). *Wired for speech: How voice activates and advances the*  
1175 *human-computer relationship*. Cambridge, MA, US: MIT Press.

1176 Niculescu, A., Dijk, B. van, Nijholt, A., Li, H., and See, S. L. (2013). Making Social Robots More  
1177 Attractive: The Effects of Voice Pitch, Humor and Empathy. *Int J of Soc Robotics* 5, 171–  
1178 191. doi:10.1007/s12369-012-0171-x.

1179 Nuñez, D. S. (David S. (2015). GlobalLit : a platform for collecting, analyzing, and reacting to  
1180 children’s usage data on tablet computers. Available at:  
1181 <http://dspace.mit.edu/handle/1721.1/98622> [Accessed September 12, 2016].

1182 Park, H. W., Coogler, R. A., and Howard, A. (2014). Using a shared tablet workspace for  
1183 interactive demonstrations during human-robot learning scenarios. in *2014 IEEE*  
1184 *International Conference on Robotics and Automation (ICRA)*, 2713–2719.  
1185 doi:10.1109/ICRA.2014.6907248.

1186 Park, H. W., and Howard, A. M. (2015). Retrieving experience: Interactive instance-based  
1187 learning methods for building robot companions. in *2015 IEEE International*  
1188 *Conference on Robotics and Automation (ICRA)*, 6140–6145.  
1189 doi:10.1109/ICRA.2015.7140061.

1190 Read, J. C., and MacFarlane, S. (2006). Using the Fun Toolkit and Other Survey Methods to  
1191 Gather Opinions in Child Computer Interaction. in *Proceedings of the 2006 Conference*  
1192 *on Interaction Design and Children IDC '06*. (New York, NY, USA: ACM), 81–88.  
1193 doi:10.1145/1139073.1139096.

1194 Rogerson, J., and Dodd, B. (2005). Is There an Effect of Dysphonic Teachers’ Voices on  
1195 Children’s Processing of Spoken Language? *Journal of Voice* 19, 47–60.  
1196 doi:10.1016/j.jvoice.2004.02.007.

1197 Sage, K. D., and Baldwin, D. (2010). Social gating and pedagogy: Mechanisms for learning  
1198 and implications for robotics. *Neural Networks* 23, 1091–1098.

1199 Sandygulova, A., and O’Hare, G. M. P. (2015). “Children’s Perception of Synthesized Voice:  
1200 Robot’s Gender, Age and Accent,” in *Social Robotics Lecture Notes in Computer*  
1201 *Science*, eds. A. Tapus, E. André, J.-C. Martin, F. Ferland, and M. Ammi (Springer  
1202 International Publishing), 594–602. doi:10.1007/978-3-319-25554-5\_59.

1203 Scarborough, H. S. (1990). Index of Productive Syntax. *Applied Psycholinguistics* 11, 1–22.  
1204 doi:10.1017/S0142716400008262.

1205 Senechal, T., McDuff, D., and Kaliouby, R. el (2015). Facial Action Unit Detection Using  
1206 Active Learning and an Efficient Non-linear Kernel Approximation. in *Proceedings of the*

1207           2015 *IEEE International Conference on Computer Vision Workshop (ICCVW) ICCVW '15.*  
1208           (Washington, DC, USA: IEEE Computer Society), 10–18. doi:10.1109/ICCVW.2015.11.  
1209 Shiomi, M., Kanda, T., Howley, I., Hayashi, K., and Hagita, N. (2015). Can a Social Robot  
1210 Stimulate Science Curiosity in Classrooms? *Int J of Soc Robotics*, 1–12.  
1211 doi:10.1007/s12369-015-0303-1.  
1212 Spaulding, S., Gordon, G., and Breazeal, C. (2016). Affect-Aware Student Models for Robot  
1213 Tutors. in *Proceedings of the 2016 International Conference on Autonomous Agents &*  
1214 *Multiagent Systems AAMAS '16.* (Richland, SC: International Foundation for  
1215 Autonomous Agents and Multiagent Systems), 864–872. Available at:  
1216 <http://dl.acm.org/citation.cfm?id=2937029.2937050>.  
1217 Speaker, K. M., Taylor, D., and Kamen, R. (2004). Storytelling: Enhancing Language  
1218 Acquisition in Young Children. *Education* 125, 3–14.  
1219 Szafir, D., and Mutlu, B. (2012). Pay Attention!: Designing Adaptive Agents That Monitor  
1220 and Improve User Engagement. in *Proceedings of the SIGCHI Conference on Human*  
1221 *Factors in Computing Systems CHI '12.* (New York, NY, USA: ACM), 11–20.  
1222 doi:10.1145/2207676.2207679.  
1223 Tamagawa, R., Watson, C. I., Kuo, I. H., MacDonald, B. A., and Broadbent, E. (2011). The  
1224 Effects of Synthesized Voice Accents on User Perceptions of Robots. *Int J of Soc*  
1225 *Robotics* 3, 253–262. doi:10.1007/s12369-011-0100-4.  
1226 Tanaka, F., and Matsuzoe, S. (2012). Children teach a care-receiving robot to promote their  
1227 learning: Field experiments in a classroom for vocabulary learning. *Journal of Human-*  
1228 *Robot Interaction* 1, 78–95.  
1229 Valdez-Menchaca, M. C., and Whitehurst, G. J. (1992). Accelerating language development  
1230 through picture book reading: A systematic extension to Mexican day care.  
1231 *Developmental Psychology* 28, 1106–1114. doi:10.1037/0012-1649.28.6.1106.  
1232 Valerie Morton, D. R. W., and Watson, D. R. (2001). The impact of impaired vocal quality on  
1233 children's ability to process spoken language. *Logopedics Phoniatrics Vocology* 26, 17–  
1234 25. doi:10.1080/14015430118232.  
1235 Walters, M. L., Syrdal, D. S., Koay, K. L., Dautenhahn, K., and Boekhorst, R. te (2008). Human  
1236 approach distances to a mechanical-looking robot with different robot voice styles. in  
1237 *RO-MAN 2008 - The 17th IEEE International Symposium on Robot and Human*  
1238 *Interactive Communication*, 707–712. doi:10.1109/ROMAN.2008.4600750.  
1239 Whitehurst, G. J., Falco, F. L., Lonigan, C. J., Fischel, J. E., B. D. DeBaryshe, M. C. Valdez-  
1240 Menchaca, et al. (1988). Accelerating language development through picture book  
1241 reading. *Developmental Psychology* 24, 552–559. doi:10.1037/0012-1649.24.4.552.  
1242

1243 **Tables**

1244 **Table 1: Summary of participant responses for the *Expressive* versus *Flat* robot**  
1245 **verification study. The *Expressive* robot was viewed as more expressive, more**  
1246 **emotional, and less passive than the *Flat* robot.**

Question	Condition	Mean	Median	Mode	Range	Inter-quartile Range
Overall expressive	<i>Flat</i>	3.75	4	4	1-5	1
	<i>Expressive</i>	3.60	4	4	2-5	1
Overall emotional	<i>Flat</i>	2.90	3	2	1-5	2
	<i>Expressive</i>	3.60	4	4	2-5	1
Overall passive	<i>Flat</i>	3.15	3	3	1-5	2
	<i>Expressive</i>	2.56	3	3	1-5	1
Expressive voice	<i>Flat</i>	2.65	3	1	1-5	2.25
	<i>Expressive</i>	4.05	4	4	3-5	0
Emotional voice	<i>Flat</i>	2.15	2	1	1-5	2
	<i>Expressive</i>	3.85	4	4	3-5	1
Passive voice	<i>Flat</i>	3.45	3.5	3	1-5	2
	<i>Expressive</i>	2.30	2	2	1-5	1

1247

1248

1249

1250 **Table 2: Older children to correctly identified more of the target vocabulary words.**

Age	Number of children	Mean target words correct ( <i>Stdev.</i> )
4 years	9	3.22 ( <i>1.30</i> )
5 years	21	4.14 ( <i>1.28</i> )
6 years	14	4.21 ( <i>1.53</i> )
7 years	1	5.00 ( <i>N/A</i> )

1251

1252

1253

1254 **Table 3. Analysis of four facial expressions during the interaction by condition. Values can**  
1255 **range from 0 (no expression present) to 100 (expression fully present).**

<b>Expression</b>	<b>Overall Mean (SD)</b>	<b>Expressive Mean (SD)</b>	<b>Flat Mean (SD)</b>
Concentration	11.7 (7.63)	14.1 (8.33)	8.93 (5.83)
Engagement	20.8 (12.1)	24.5 (12.6)	16.6 (10.3)
Surprise	6.71 (4.57)	8.28 (4.99)	4.95 (3.39)
Attention	82.6 (7.45)	82.4 (7.47)	83.0 (7.63)

1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289

1290 **9 Figure Legends**

1291 Figure 1: (A, top) The Tega robot sat on a table facing the child. The tablet that displayed the  
1292 storybook was positioned to the right of the robot. Video cameras recorded the interaction from  
1293 behind the robot, and the phone in the front used Affdex to record children's emotional states. (B,  
1294 bottom) A child looks up at the experimenter at the end of the robot interaction. Tega and the  
1295 Toucan puppet have just said goodbye.

1296  
1297 Figure 2: Children who heard story version A correctly identified more version A words  
1298 than version B words, whereas children who heard story version B correctly identified  
1299 more version B words than version A words.

1300  
1301 Figure 3: Children who heard story version A used more version A target and original  
1302 words than version B words, whereas children who heard story version B used more  
1303 version B target and original words than version A words in their immediate story retell.

1304  
1305 Figure 4: Children who correctly identified more target words also used them more in their  
1306 story retell. The trend was primarily driven by the Expressive condition.

1307  
1308 Figure 5: Children's story-retell length significantly reduced after 1-2 months in the Flat  
1309 condition, but not in the Expressive condition.

1310  
1311 Figure 6: (A, left) Children who responded to the robot's dialogic questions were also more  
1312 likely to correctly identify more of the target vocabulary words. The correlation was  
1313 primarily driven by children in the Expressive condition. (B, right) The majority of children  
1314 responded to most or all of the robot's dialogic questions.

1315  
1316 Figure 7: (A, top) Children who responded to the robot's dialogic questions were more  
1317 likely to use the target words in their retells and tell longer stories. (B, bottom) These  
1318 children were also more likely to emulate robot's story in terms of phrase similarity, in  
1319 both immediate and delayed retell. The trend was primarily driven by the *Expressive*  
1320 condition.

1321  
1322  
1323 Figure 8: Children in the Expressive condition showed stronger emulation of the robot's  
1324 story in terms of phrase similarity, in both immediate and delayed retell.

1325  
1326 Figure 9: Children who correctly identified more target words were also more likely to  
1327 emulate the robot's story in terms of phrase similarity, in both immediate and delayed  
1328 retell. These trends were primarily driven by the *Expressive* condition.

1329  
1330 Figure 10. Children in the Expressive condition showed more concentration, engagement  
1331 and surprise during the session than children in the Flat condition. Attention levels were  
1332 not statistically different between the two conditions.

1333  
1334 Figure 11. (A, top left) Children's level of attention decreased in the course of the session  
1335 but showed no difference between conditions. (B, top right) The concentration level of



1336 children in the Expressive condition was consistently higher than that of children in the  
1337 Flat condition. (C, bottom left) The engagement level of children in the Expressive  
1338 condition was consistently higher than that of children in the Flat condition. (D, bottom  
1339 right) The surprise level of children in the Expressive condition was consistently higher  
1340 than that of children in the Flat condition.

1341

1342 Figure 12: Children who were native English speakers or bilingual correctly identified  
1343 more of the target vocabulary words than did ELL children.