

# The development of children's ability to track and predict turn structure in conversation

Marisa Casillas<sup>a,\*</sup>, Michael C. Frank<sup>b</sup>

<sup>a</sup>*Max Planck Institute for Psycholinguistics, Nijmegen*

<sup>b</sup>*Department of Psychology, Stanford University*

---

## Abstract

Children begin developing turn-taking skills in infancy but take several years to assimilate their growing knowledge of language into their turn-taking behavior. In two eye-tracking experiments, we measured children's anticipatory gaze to upcoming responders while controlling linguistic cues to turn structure. In Experiment 1, we showed English and non-English conversations to English-speaking adults and children. In Experiment 2, we phonetically controlled lexicosyntactic and prosodic cues in English-only speech. Children spontaneously made anticipatory gaze switches by age two and continued improving through age six. In both experiments, children and adults made more anticipatory switches after hearing questions. Consistent with prior findings on adult turn prediction, prosodic information alone did not increase children's anticipatory gaze shifts. But, unlike prior work with adults, lexical information alone was not sufficient either—children's performance was best overall with lexicosyntax and prosody together. Our findings support an account in which turn tracking and prediction emerges in infancy, and then only gradually becomes integrated with linguistic processing.

*Keywords:* Turn taking, Conversation, Development, Questions, Eye-tracking, Anticipation

---

\*Corresponding author.

Address: Wundtlaan 1, 6525 XD, Nijmegen, The Netherlands

Email: marisa.casillas@mpi.nl

Telephone: +31 024 3521 566; Fax: +31 024 3521 213

**1** **Introduction**

**2** Spontaneous conversation is a universal context for using and learning  
**3** language. Like other types of human interaction, it is organized at its core  
**4** by the roles and goals of its participants. But, what sets conversation apart is  
**5** its structure: sequences of interconnected, communicative actions that take  
**6** place across alternating turns at talk. Sequential, turn-based structures in  
**7** conversation are strikingly uniform across language communities and linguis-  
**8** tic modalities. Turn-taking behaviors are also cross-culturally consistent in  
**9** their basic features and the details of their implementation (De Vos et al.,  
**10** 2015; Dingemanse et al., 2013; Stivers et al., 2009).

**11** Children participate in sequential coordination (proto-turn taking) with  
**12** their caregivers starting at three months of age—before they can rely on  
**13** any linguistic cues (see, among others, Bateson, 1975; Hilbrink et al., 2015;  
**14** Jaffe et al., 2001; Snow, 1977). However, infant turn taking is different from  
**15** adult turn taking in several ways: it is heavily scaffolded by caregivers, has  
**16** different timing from adult turn taking, and lacks semantic content (Hilbrink  
**17** et al., 2015; Jaffe et al., 2001). But children’s early, turn-structured social  
**18** interactions are presumably a critical precursor to their later conversational  
**19** turn taking. Early non-verbal interactions likely establish the protocol by  
**20** which children come to use language with others. How do children integrate  
**21** linguistic knowledge with these preverbal turn-taking abilities, and how does  
**22** this integration change over the course of childhood?

**23** In this study, we investigate when children begin to make predictions  
**24** about upcoming turn structure in conversation, and how they integrate lan-  
**25** guage into their predictions as they grow older. In the remainder of the  
**26** introduction, we first give a basic review of turn-taking research and the  
**27** state of current knowledge about adult turn prediction. We then discuss  
**28** recent work on the development of turn-taking skills before turning to the  
**29** details of our own study.

**30** *Adult turn taking*

**31** Turn taking itself is not unique to conversation. Many other human ac-  
**32** tivities are organized around sequential turns at action. Traffic intersections  
**33** and computer network communication both use turn-taking systems. Chil-  
**34** dren’s early games (e.g., give-and-take, peek-a-boo) have built-in, predictable  
**35** turn structure (Ratner & Bruner, 1978; Ross & Lollis, 1987). Even monkeys  
**36** take turns: non-human primates such as marmosets and Campbell’s monkeys

37 vocalize contingently with each other in both natural and lab-controlled environments (Lemasson et al., 2011; Takahashi et al., 2013). In all these cases, 38 turn taking serves as a protocol for interaction, allowing the participants to 39 coordinate with each other through sequences of contingent action.

40  
41 Conversational turn taking distinguishes itself from other turn-taking behaviors by the complexity of the sequencing involved. Conversational turns 42 come grouped into semantically-contingent sequences of action. The groups 43 can span turn-by-turn exchanges (e.g., simple question-response, “How are 44 you?”—“Fine.”) or sequence-by-sequence exchanges (e.g., reciprocals, “How 45 are you?”—“Fine, and you?”—“Great!”). Compared to other turn-taking behaviors, the possible sequence and action types in everyday talk are can be 46 diverse and unpredictable.

47  
48 Despite this complexity, conversational turn taking is precise in its timing. 49 Across a diverse sample of conversations in 10 languages, one study found 50 a consistent average turn transition time of 0–200 msec at points of speaker 51 switch (Stivers et al., 2009). Experimental results and current models of 52 speech production suggest that it takes approximately 600 msec to produce 53 a content word, and even longer to produce a simple utterance (Griffin & 54 Bock, 2000; Levelt, 1989). So in order to achieve 200 msec turn transitions, 55 speakers must begin formulating their response before the prior turn has 56 ended (Levinson, 2013, 2016). Moreover, to formulate their response early 57 on, speakers must track and anticipate what types of response might become 58 relevant next. They also need to predict the content and form of upcoming 59 speech so that they can launch their articulation at exactly the right moment. 60 Prediction thus plays a key role in timely turn taking.

61  
62 Adults have a lot of information at their disposal to help make accurate 63 predictions about upcoming turn content. Lexical, syntactic, and prosodic 64 information (e.g., *wh*- words, subject-auxiliary inversion, and list intonation) 65 can all inform addressees about upcoming linguistic structure (De Ruiter 66 et al., 2006; Duncan, 1972; Ford & Thompson, 1996; Torreira et al., 2015). 67 Non-verbal cues (e.g., gaze, posture, and pointing) often appear at turn- 68 boundaries and can sometimes act as late indicators of an upcoming speaker 69 switch (Rossano et al., 2009; Stivers & Rossano, 2010). Additionally, the 70 sequential context of a turn can make it clear what will come next: answers 71 after questions, thanks or denial after compliments, etc. (Schegloff, 2007).

72 Prior work suggests that adult listeners primarily use lexicosyntactic in- 73 formation to accurately predict upcoming turn structure. De Ruiter and 74 colleagues (2006) asked participants to listen to snippets of spontaneous con-

75 versation and to press a button whenever they anticipated that the current  
76 speaker was about to finish his or her turn. The speech snippets were con-  
77 trolled for the amount of linguistic information present; some were normal,  
78 but others had flattened pitch, low-pass filtered speech, or further manip-  
79 ulations. With pitch-flattened speech, the timing of participants' button  
80 responses was comparable to their timing with the full linguistic signal. But  
81 when no lexical information was available, participants' responded signifi-  
82 cantly earlier within the turn. The authors concluded that lexicosyntactic  
83 information<sup>1</sup> was necessary and possibly sufficient for turn-end projection,  
84 while intonation was neither necessary nor sufficient. Congruent evidence  
85 comes from studies varying the predictability of lexicosyntactic and prag-  
86 matic content: adults anticipate turn ends better when they can more accu-  
87 rately predict the exact words that will come next (Magyari & De Ruiter,  
88 2012; see also Magyari et al., 2014). They can also identify speech acts within  
89 the first word of an utterance (Gísladóttir et al., 2015), allowing them to start  
90 planning their response at the first moment possible (Bögels et al., 2015).

91 Despite this body of evidence, the role of prosody for adult turn pre-  
92 diction is still a matter of debate. De Ruiter and colleagues' (2006) exper-  
93 iment focused on the role of intonation, which is only a partial index of  
94 prosody. Prosody is tied closely to the syntax of an utterance, so the two  
95 linguistic signals are difficult to control independently (Ford & Thompson,  
96 1996). Torreira, Bögels & Levinson (2015) used a combination of button-  
97 press and verbal responses to investigate the relationship between lexicosyn-  
98 tactic and prosodic cues in turn-end prediction. Critically, their stimuli were  
99 cross-spliced so that each item had full prosodic cues to accompany the lex-  
100 icosyntax. Because of the splicing, they were able to create items that had  
101 syntactically-complete units with no intonational phrase boundary at the  
102 end. Participants never verbally responded or pressed the "turn-end" but-  
103 ton when hearing a syntactically-complete phrase without an intonational  
104 phrase boundary. And when intonational phrase boundaries were embedded  
105 in multi-utterance turns, participants were tricked into pressing the "turn-  
106 end" button 29% of the time. Their results suggest that listeners actually  
107 do rely on prosodic cues to execute a response (see also de De Ruiter et al.

---

<sup>1</sup>The "lexicosyntactic" condition only included flattened pitch and so was not exclusively lexicosyntactic—the speech would still have residual prosodic structure, including syllable duration and intensity.

<sup>108</sup> (2006):525). These experimental findings corroborate other corpus and ex-  
<sup>109</sup> perimental work promoting a combination of cues (lexicosyntactic, prosodic,  
<sup>110</sup> and pragmatic) as key for accurate turn-end prediction (Duncan, 1972; Ford  
<sup>111</sup> & Thompson, 1996; Hirvenkari et al., 2013).

<sup>112</sup> *Turn taking in development*

<sup>113</sup> The majority of work on children's early turn taking has focused on ob-  
<sup>114</sup> servations of spontaneous interaction. Children's first turn-like structures  
<sup>115</sup> appear as early as two to three months after birth, in proto-conversation with  
<sup>116</sup> their caregivers (Bruner, 1975, 1985). During proto-conversations, caregivers  
<sup>117</sup> treat their infants as capable of making meaningful contributions: they take  
<sup>118</sup> every look, vocalization, arm flail, and burp as "utterances" in the joint dis-  
<sup>119</sup> course (Bateson, 1975; Jaffe et al., 2001; Snow, 1977). Infants catch onto the  
<sup>120</sup> structure of proto-conversations quickly. By three to four months they notice  
<sup>121</sup> disturbances to the contingency of their caregivers' response and, in reaction,  
<sup>122</sup> change the rate and quality of their vocalizations (Bloom, 1988; Masataka,  
<sup>123</sup> 1993; Toda & Fogel, 1993).

<sup>124</sup> The timing of children's responses to their caregivers' speech shows a  
<sup>125</sup> non-linear pattern. Infants' contingent vocalizations in the first few months  
<sup>126</sup> of life show very fast timing (though with a lot of vocal overlap). But by  
<sup>127</sup> nine months, their timing slows down considerably, only to gradually speed  
<sup>128</sup> up again after 12 months (Hilbrink et al., 2015). For children, taking turns  
<sup>129</sup> with brief transitions between speakers is more difficult than avoiding speaker  
<sup>130</sup> overlap; children's incidence of overlap is nearly adult-like by nine months,  
<sup>131</sup> but the timing of their non-overlapped (i.e., gapped) responses remains longer  
<sup>132</sup> than the adult 200 msec standard for the next few years (Casillas et al., In  
<sup>133</sup> press; Garvey, 1984; Garvey & Berninger, 1981; Ervin-Tripp, 1979). This  
<sup>134</sup> puzzling pattern is likely due to children's linguistic development: taking  
<sup>135</sup> turns on time is easier when their response is a simple vocalization rather  
<sup>136</sup> than a linguistic utterance. Integrating language into the turn-taking system  
<sup>137</sup> may therefore be a major factor in children's delayed responses (Casillas  
<sup>138</sup> et al., In press).

<sup>139</sup> Before children manage to integrate linguistic cues into their turn-taking  
<sup>140</sup> behaviors (for both turn prediction and production), they can rely on non-  
<sup>141</sup> verbal interactional cues, including silence (e.g., between turns), eye gaze,  
<sup>142</sup> body orientation, and gesture, to identify the boundaries of social actions.  
<sup>143</sup> For example, with little to no linguistic knowledge, children are often able

144 to infer and anticipate desired responses to offers and requests by taking ac-  
145 count of their interlocutor's non-verbal communicative behavior, the struc-  
146 ture of routine events, and the affordances of the current interactional context  
147 (Reddy et al., 2013; Nomikou & Rohlfsing, 2011; Shatz, 1978). With respect  
148 to turn taking in particular, children's spontaneous vocalizations during in-  
149 teraction demonstrate a sensitivity to short inter-speaker gaps from infancy  
150 (Hilbrink et al., 2015), even when they are not yet using language to partic-  
151 ipate. Thus, before they can anticipate turn structure from linguistic cues,  
152 children might instead react to silence as a cue to upcoming speaker change.  
153 Interactional silence itself may then serve as one of children's first cues to  
154 turn structure, giving them information about when to respond before they  
155 can rely on language.

156 As children's language competence increases, they can use linguistic cues  
157 to make predictions about upcoming turn structure. Studies of early linguis-  
158 tic development point to a possible early advantage for prosody over lexi-  
159 cosyntaxis in children's turn-taking predictions. Infants can distinguish their  
160 native language's rhythm type from others soon after birth (Mehler et al.,  
161 1988; Nazzi & Ramus, 2003); they show preference for the typical stress pat-  
162 terns of their native language over others by 6–9 months (e.g., iambic vs.  
163 trochaic), and can use prosodic information to segment the speech stream  
164 into smaller chunks from 8 months onward (Johnson & Jusczyk, 2001; Mor-  
165 gan & Saffran, 1995). Four- to five-month-olds also prefer pauses in speech to  
166 be inserted at prosodic boundaries, and by 6 months infants can use prosodic  
167 markers to pick out sub-clausal syntactic units, both of which are useful for  
168 extracting turn structure from ongoing speech (Jusczyk et al., 1995; Soder-  
169 strom et al., 2003). In comparison, children show at best a very limited  
170 lexical inventory before their first birthday (Bergelson & Swingley, 2013; Shi  
171 & Melancon, 2010).

172 Keitel and colleagues (2013) were one of the first to explore how chil-  
173 dren use linguistic cues to predict upcoming turn structure. They asked 6-,  
174 12-, 24-, and 36-month-old infants, and adult participants to watch short  
175 videos of conversation and tracked their eye movements at points of speaker  
176 change. They showed their participants two types of conversation videos—  
177 one normal and one with flattened pitch—to test the role of intonation in  
178 participants' anticipatory predictions about upcoming speech. Comparing  
179 children's anticipatory gaze frequency to a random baseline, they found that  
180 only 36-month-olds and adults made anticipatory gaze switches more often  
181 than expected by chance, and only 36-month-olds were affected by a lack of

182 intonation contours. This finding led Keitel and colleagues to conclude that  
183 children's ability to predict upcoming turn structure relies on their ability to  
184 comprehend the stimuli lexicosemantically. They also suggest that intona-  
185 tion might play a secondary role in turn prediction, but only after children  
186 acquire more sophisticated, adult-like language comprehension abilities (also  
187 see Keitel & Daum, 2015).

188 Although the Keitel et al. (2013) study constitutes a substantial ad-  
189 vance over previous work in this domain, it has some limitations. Because  
190 these limitations directly inform our own study design, we review them in  
191 some detail. First, their estimates of baseline gaze frequency ("random" in  
192 their terminology) were not random. Instead, they used gaze switches dur-  
193 ing ongoing speech as a baseline. But ongoing speech is the period in which  
194 switching is least likely to occur (Hirvenkari et al., 2013)—their baseline thus  
195 maximizes the chance of finding a difference in gaze frequency at turn transi-  
196 tions compared to the baseline. A more conservative baseline would compare  
197 participants' looking behavior at turn transitions to their looking behavior  
198 during randomly selected windows of time throughout the stimulus, includ-  
199 ing turn transitions. We follow this conservative approach in the current  
200 study.

201 Second, the conversation stimuli Keitel et al. (2013) used were some-  
202 what unusual. The average gap between turns was 900 msec, a duration  
203 much longer than typical adult timing, which averages around 200 msec  
204 (Stivers et al., 2009). The speakers in the videos were also asked to mini-  
205 mize their movements while performing scripted, adult-directed conversation,  
206 which would have created a somewhat unnatural interaction. Additionally,  
207 to produce more naturalistic conversation, it would have been ideal to lo-  
208 calize the sound sources for the two voices in the video (i.e., to have the  
209 voices come out of separate left and right speakers). But both voices were  
210 recorded and played back on the same audio channel, which may have made  
211 it difficult to distinguish the two talkers. Again, we attempt to address these  
212 issues in our current study. Despite these minor methodological issues, the  
213 Keitel et al. (2013) study still demonstrates intriguing age-based differences  
214 in children's ability to predict upcoming turn structure. Our current work  
215 takes these findings as a starting point.<sup>2</sup>

---

<sup>2</sup>But also see Casillas & Frank (2012, 2013).

216 *The current study*

217 Our goal in the current study is to find out when children begin to make  
218 predictions about upcoming turn structure and to understand how their pre-  
219 dictions are affected by linguistic cues across development. We present two  
220 experiments in which we measured children’s anticipatory gaze to respon-  
221 ders while they watched conversation videos with natural (people speaking  
222 English vs. non-English; Experiment 1) and non-natural (puppets with pho-  
223 netically manipulated speech; Experiment 2) control over the presence of  
224 lexical and prosodic cues. We tested children across a wide range of ages  
225 (Experiment 1: 3–5 years; Experiment 2: 1–6 years), with adult control  
226 participants in each experiment. We additionally tested for the use of one  
227 non-verbal cue: inter-turn silence.

228 We highlight four primary findings: first, although children and adults  
229 use linguistic cues to make predictions about upcoming turn structure, they  
230 do so primarily to predict speaker transitions after questions (a speech act ef-  
231 fect). This intriguing effect, which has not been reported previously, suggests  
232 that participants track unfolding speech for cues to upcoming speaker change,  
233 which may affect how they use linguistic cues more generally for anticipa-  
234 tory processing in conversation. Second, we find that children make more  
235 predictions than expected by chance starting at age two, but that this effect  
236 is small at first, and continues to improve through age six, along with chil-  
237 dren’s use of linguistic cues to anticipate answers after question turns. Third,  
238 children and adults at all ages often used inter-turn silence (a non-verbal cue  
239 to turn structure) to make more predictive gaze switches to the responder,  
240 suggesting that non-verbal cues are useful for predicting turn structure early  
241 on and continue to be important in adulthood. Finally, we find no evidence  
242 for an early prosodic advantage in children’s anticipations and, further, no  
243 evidence that lexical cues alone are comparable to the full linguistic signal in  
244 aiding participants’ predictions (as is proposed for adults; De Ruiter et al.,  
245 2006). Anticipation is strongest for stimuli with the full range of linguistic  
246 cues. Our findings support an account in which turn prediction emerges in in-  
247 fancy, but becomes fully integrated with linguistic processing only gradually  
248 across development.

249 **Experiment 1**

250 We recorded participants’ eye movements as they watched six short videos  
251 of two-person (dyadic) conversation interspersed with attention-getting filler

252 videos. Each conversation video featured an improvised discourse in one  
253 of five languages (English, German, Hebrew, Japanese, and Korean). Par-  
254 ticipants saw two videos in English and one in every other language. The  
255 participants, all native English speakers, were only expected to understand  
256 the two videos in English. We showed participants non-English videos to  
257 limit their access to lexical information while maintaining their access to  
258 other cues to turn boundaries (e.g., non-English prosody, gaze, in-breaths,  
259 phrase final lengthening). Using this method, we compared children and  
260 adult's anticipatory looks from the current speaker to the upcoming speaker  
261 at points of turn transition in English and non-English videos.

262 *Methods*

263 *Participants*

264 We recruited 74 children ages 3;0–5;11 and 11 undergraduate adults to  
265 participate in the experiment. We recruited adult participants through the  
266 Stanford University Psychology participant database. Adult participants  
267 were either paid or received course credit for their time. Our child sample in-  
268 cluded 19 three-year-olds, 32 four-year-olds, and 23 five-year-olds, all enrolled  
269 in a local nursery school. All participants were native English speakers and  
270 volunteered their time. Approximately one-third ( $N=25$ ) of the children's  
271 parents and teachers reported that their child regularly heard a second (and  
272 sometimes third or further) language, but only one child frequently heard a  
273 language that was used in our non-English video stimuli, and we excluded  
274 his data from the analyses.<sup>3</sup> None of the adult participants reported fluency  
275 in a second language.

276 *Materials*

277 *Video recordings.* We recorded pairs of talkers while they conversed in  
278 a sound-attenuated booth (see a sample frame in Figure 1). Each talker  
279 was a native speaker of the language being recorded, and each talker pair  
280 was male-female. Using a Marantz PMD 660 solid state field recorder, we  
281 captured audio from two lapel microphones, one attached to each participant,

---

<sup>3</sup>Multilingual children may make predictions about upcoming turn structure differently from their monolingual peers due to their more varied experiences with language-relevant turn taking cues, but we are unable to test this hypothesis here due to the variability in multilingual language input and languages being learned in our sample. The same note applies to Experiment 2 below.



Figure 1: Example frame from a conversation video used in Experiment 1.

282 while simultaneously recording video from the built-in camera of a MacBook  
283 laptop computer. The talkers were volunteers and were acquainted with their  
284 recording partner ahead of time.

285 Each recording session began with a 20-minute warm-up period of spontaneous  
286 conversation during which the pair talked for five minutes on four topics (favorite foods, entertainment, hometown layout, and pets). Then we  
287 asked talkers to choose a new topic—one relevant to young children (e.g.,  
288 riding a bike, eating breakfast)—and to improvise a dialogue on that topic.  
289 We asked them to speak as if they were on a children’s television show in  
290 order to elicit child-directed speech toward each other. We recorded until the  
291 talkers achieved at least 30 seconds of uninterrupted discourse with enthusiastic,  
292 child-directed speech. Most talker pairs took less than five minutes  
293 to complete the task, usually by agreeing on a rough script at the start. We  
294 encouraged talkers to ask at least a few questions to each other during the  
295 improvisation. The resulting conversations were therefore not entirely spontaneous,  
296 but were as close as possible while still remaining child-oriented in  
297 topic, prosodic pattern, and lexicosyntactic construction.<sup>4</sup>

299 After recording, we combined the audio and video recordings by hand,  
300 and cropped each recording to the (approximate) 30-second interval with  
301 the most turn activity. Because we recorded the conversations in stereo, the

---

<sup>4</sup>All of the non-English talkers were fluent in English as a second language, and some fluently spoke three or more languages. We chose male-female pairs as a natural way of creating contrast between the two talker voices.

302 male and female voices came out of separate speakers during video play-  
303 back. This gave each voice in the videos a localized source (from the left or  
304 right loudspeaker). We coded each turn transition in the videos for language  
305 condition (English vs. non-English), inter-turn gap duration (in millisec-  
306 onds), and speech act (question vs. non-question). The non-English stimuli  
307 were coded for speech act from a monolingual English-speaker’s perspective,  
308 i.e., which turns “sound like” questions, and which do not: we asked five  
309 native American English speakers to listen to the audio recording for each  
310 non-English turn and judge whether it sounded like a question. We marked  
311 non-English turns as questions when at least 4 of the 5 listeners (80%) said  
312 that the turn “sounded like a question”. This procedure means that the cues  
313 for recognizing “questions” in the non-English condition only resembled na-  
314 tive English question cues, and therefore were likely harder to identify than  
315 cues to questionhood in the English condition. However, since participants  
316 did not speak the non-English languages and would only therefore ever treat  
317 “question-sounding” turns as questions, we proceeded with these analyses to  
318 see how pervasive question effects were—could they show up even without  
319 lexical access in a foreign language? If participants primarily rely on prosodic  
320 cues to question turns, it’s possible that even non-English prosody can elicit  
321 anticipatory gaze switches for question-like turns.

322 Because the conversational stimuli were recorded semi-spontaneously, the  
323 duration of turn transitions and the number of speaker transitions in each  
324 video was variable. We measured the duration of each turn transition from  
325 the audio recording associated with each video. We excluded turn transitions  
326 longer than 550 msec and shorter than 90 msec from analysis, additionally  
327 excluding overlapped transitions.<sup>5</sup> This left approximately equal numbers  
328 of turn transitions available for analysis in the English (N=20) and non-  
329 English (N=16) videos. On average, the inter-turn gaps for English videos  
330 (mean=318, median=302, stdev=112 msec) were slightly longer than for non-  
331 English videos (mean=286, median=251, stdev=122 msec).

332 Questions made up exactly half of the turn transitions in the English

---

<sup>5</sup>Overlap occurs when a responder begins a new turn before the current turn is finished. When overlap occurs, observers cannot switch their gaze in anticipation of the response because the response began earlier than expected. Participants expect conversations to proceed with “one speaker at a time” (Sacks et al., 1974). They would therefore still be fixated on the prior speaker when the overlap started, and would have to switch their gaze *reactively* to the responder.

333 (N=10) and non-English (N=8) videos. In the English videos, inter-turn  
334 gaps were slightly shorter for questions (mean=310, median=293, stdev=112  
335 msec) than non-questions (mean=325, median=315, stdev=118 msec). Non-  
336 English videos did not show a large difference in transition time for questions  
337 (mean=270, median=257, stdev=116 msec) and non-questions (mean=302,  
338 median=252, stdev=134 msec).

339 *Procedure*

340 Participants sat in front of an SMI 120Hz corneal reflection eye-tracker  
341 mounted beneath a large flatscreen display. The display and eye-tracker were  
342 secured to a table with an ergonomic arm that allowed the experimenter to  
343 position the whole apparatus at a comfortable height, approximately 60 cm  
344 from the viewer. We placed stereo speakers on the table, to the left and right  
345 of the display.

346 Before the experiment started, we warned adult participants that they  
347 would see videos in several languages and that, though they weren't expected  
348 to understand the content of non-English videos, we *would* ask them to an-  
349 swer general, non-language-based questions about the conversations. Then  
350 after each video we asked participants one of the following randomly-assigned  
351 questions: "Which speaker talked more?", "Which speaker asked the most  
352 questions?", "Which speaker seemed more friendly?", and "Did the speak-  
353 ers' level of enthusiasm shift during the conversation?" We also asked if the  
354 participants could understand any of what was said after each video. The  
355 participants responded verbally while an experimenter noted their responses.

356 Children were less inclined to simply sit and watch videos of conversation  
357 in languages they didn't speak, so we used a different procedure to keep them  
358 engaged: the experimenter started each session by asking the child about  
359 what languages he or she could speak, and about what other languages he  
360 or she had heard of. Then the experimenter expressed her own enthusiasm  
361 for learning about new languages, and invited the child to watch a video  
362 about "new and different languages" together. If the child agreed to watch,  
363 the experimenter and the child sat together in front of the display, with  
364 the child centered in front of the tracker and the experimenter off to the  
365 side. Each conversation video was preceded and followed by a 15–30 second  
366 attention-getting filler video (e.g., running puppies, singing muppets, flying  
367 bugs). If the child began to look bored, the experimenter would talk during  
368 the fillers, either commenting on the previous conversation ("That was a neat  
369 language!") or giving the language name for the next conversation ("This

370 next one is called Hebrew. Let's see what it's like.") The experimenter's  
371 comments reinforced the video-watching as a joint task.

372 All participants (child and adult) completed a five-point calibration rou-  
373 tine before the first video started. We used a dancing Elmo for the children's  
374 calibration image. During the experiment, participants watched all six 30-  
375 second conversation videos. The first and last conversations were in American  
376 English and the intervening conversations were Hebrew, Japanese, German,  
377 and Korean. The presentation order of the non-English videos was shuffled  
378 into four lists, which participants were assigned to randomly. The entire  
379 experiment, including instructions, took 10–15 minutes.

380 *Data preparation and coding*

381 To determine whether participants predicted upcoming turn transitions,  
382 we needed to define a set of criteria for what counted as an anticipatory gaze  
383 shift. Prior work using similar experimental procedures has found that adults  
384 and children make anticipatory gaze shifts to upcoming talkers within a wide  
385 time frame; the earliest shifts occur before the end of the prior turn, and the  
386 latest occur after the onset of the response turn, with most shifts occurring  
387 in the inter-turn gap (Keitel et al., 2013; Hirvenkari, 2013; Tice and Henetz,  
388 2011). Following prior work, we measured how often our participants shifted  
389 their gaze from the prior to the upcoming speaker *before* the shift in gaze  
390 could have been initiated in reaction to the onset of the speaker's response.  
391 In doing so, we assumed that it takes participants 200 msec to plan an eye  
392 movement, following standards from adult anticipatory processing studies  
393 (e.g., Kamide et al., 2003).

394 We checked each participant's gaze at each turn transition for three char-  
395 acteristics (Figure 2): (1) that the participant fixated on the prior speaker for  
396 at least 100 msec at the end of the prior turn, (2) that immediately thereafter  
397 the participant switched to fixate on the upcoming speaker for at least 100  
398 ms, and (3) that the switch in gaze was initiated within the first 200 msec of  
399 the response turn, or earlier. These criteria guarantee that we only counted  
400 gaze shifts when: (1) participants were tracking the previous speaker, (2)  
401 switched their gaze to track the upcoming speaker, and (3) did so before  
402 they could have simply reacted to the onset of speech in the response. Under  
403 this assumption, a gaze shift that was initiated within the first 200 msec of  
404 the response (or earlier) was planned *before* the child could react to the onset  
405 of speech itself.

406 As mentioned, most anticipatory switches happen in the inter-turn gap,

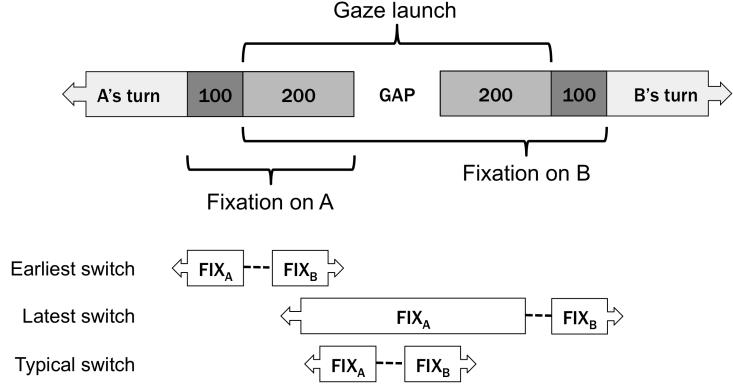


Figure 2: Schematic summary of criteria for anticipatory gaze shifts from speaker A to speaker B during a turn transition.

407 but we also allowed anticipatory gaze switches that occurred in the final  
 408 syllables of the prior turn. Early switches are consistent with the distribution  
 409 of responses in explicit turn-boundary prediction tasks. For example,  
 410 in a button press task, adult participants anticipated turn ends approxi-  
 411 mately 200 msec in advance of the turn’s end, and anticipatory responses  
 412 to pitch-flattened stimuli came even earlier (De Ruiter et al., 2006). We  
 413 therefore allowed switches to occur as early as 200 msec before the end of  
 414 the prior turn. Again, because it takes 200 msec to plan an eye movement,  
 415 we counted anticipatory switches, at the latest, 200 msec after the onset of  
 416 speech. Therefore, for very early and very late switches, our requirement of  
 417 100 msec of fixation on each speaker would sometimes extend outside of the  
 418 transition window boundaries (200 msec before and after the inter-turn gap).  
 419 The maximally available fixation window was therefore 100 msec before and  
 420 after the earliest and latest possible switch point (300 msec before and after  
 421 the inter-turn gap). We did not count switches made during the fixation win-  
 422 dow as anticipatory. We *did* count switches made during the inter-turn gap.  
 423 The period of time from the beginning of the possible fixation window on  
 424 the prior speaker to the end of the possible fixation window on the responder  
 425 was our total analysis window (300 msec + the inter-turn gap + 300 msec).

426 *Predictions.* We expected participants to show greater anticipation in the  
 427 English videos than in the non-English videos because of their increased  
 428 access to linguistic information in English. We also predicted that anticipa-

429 tion would be greater following questions compared to non-questions; ques-  
430 tions have early cues to upcoming turn transition (e.g., *wh*- words, subject-  
431 auxiliary inversion), and also make a next response immediately relevant.  
432 Our third prediction was that anticipatory looks would increase with devel-  
433 opment, along with children’s increased linguistic competence. Finally, we  
434 predicted that transitions with longer inter-turn gaps would show greater an-  
435 ticipation because longer gaps provide (a) more time to make a gaze switch  
436 and (b) are themselves a cue to possible upcoming speaker switch.

### 437 *Results*

438 Participants looked at the screen most of the time during video playback  
439 (81% and 91% on average for children and adults, respectively). They pri-  
440 marily kept their eyes on the person who was currently speaking in both  
441 English and non-English videos: they gazed at the current speaker between  
442 38% and 63% of the time, looking back at the addressee between 15% and  
443 20% of the time (Table 1). Even three-year-olds looked more at the current  
444 speaker than anything else, whether the videos were in a language they could  
445 understand or not. Children looked at the current speaker less than adults  
446 did during the non-English videos. Despite this, their looks to the addressee  
447 did not increase substantially in the non-English videos, indicating that their  
448 looks away were probably related to boredom rather than confusion about  
449 ongoing turn structure. Overall, participants’ pattern of gaze to current  
450 speakers demonstrated that they performed basic turn tracking during the  
451 videos, regardless of language. Figure 3 shows participants’ anticipatory gaze  
452 rates across age, language condition, and transition type.

### 453 *Statistical models*

454 We identified anticipatory gaze switches for all 36 usable turn transitions,  
455 based on the criteria outlined in Section b, and analyzed them for effects of  
456 language, transition type, and age with two mixed-effects logistic regressions  
457 (Bates et al., 2014; R Core Team, 2014). We built one model each for children  
458 and adults. We modeled children and adults separately because effects of age  
459 are only pertinent to the children’s data. The child model included condi-  
460 tion (English vs. non-English)<sup>6</sup>, transition type (question vs. non-question),  
461 age (3, 4, 5; numeric; intercept as age=0), and duration of the inter-turn

---

<sup>6</sup>Because each non-English language was represented by a single stimulus, we cannot treat individual languages as factors. Gaze behavior might be best for non-native languages

Age group	Condition	Speaker	Addressee	Other onscreen	Offscreen
3	English	0.61	0.16	0.14	0.08
4	English	0.60	0.15	0.11	0.13
5	English	0.57	0.15	0.16	0.12
Adult	English	0.63	0.16	0.16	0.05
3	Non-English	0.38	0.17	0.20	0.25
4	Non-English	0.43	0.19	0.21	0.18
5	Non-English	0.40	0.16	0.26	0.18
Adult	Non-English	0.58	0.20	0.16	0.07

Table 1: Average proportion of gaze to the current speaker and addressee during periods of talk.

462 gap (seconds, e.g., 0.441) as predictors, with two-way interactions between  
 463 gap duration and the other simple predictors (language condition, transi-  
 464 tion type, and age), also adding a three-way interaction between language  
 465 condition, transition type, and age. We included the two-way interactions  
 466 between gap duration and the other predictors in case the effect of inter-turn  
 467 silence changes with age or linguistic cueing (e.g., if children older children  
 468 rely less on silence as a cue)<sup>7</sup>. We also included random effects of item (turn  
 469 transition) and participant, with maximal random slopes of condition, tran-  
 470 sition type, and their interaction for participants (Barr et al., 2013).<sup>8</sup> The  
 471 adult model included fixed effects of condition, transition type, and their  
 472 interaction, plus two-way interactions between gap duration and language  
 473 condition and transition type (as in the child model). The adult model also

---

that have the most structural overlap with participants' native language: English speakers can make predictions about the strength of upcoming Swedish prosodic boundaries nearly as well as Swedish speakers do, but Chinese speakers are at a disadvantage in the same task (Carlson et al., 2005). We would need multiple items from each of the languages to check for similarity effects of specific linguistic features.

<sup>7</sup>We test these two-way interactions with gap duration in all of the models reported in this paper. Higher-order interactions with gap duration usually resulted in model non-convergence due to distributional sparsity when three or more predictor values were considered.

<sup>8</sup>The models we report in this paper are all qualitatively unchanged by the exclusion of their random slopes. We have left the random slopes in because of minor participant-level variation in the predictors modeled.

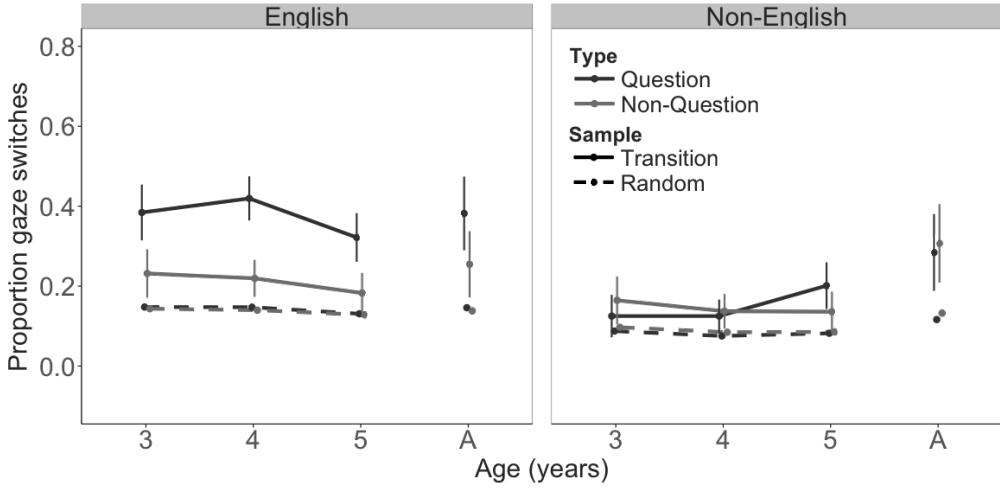


Figure 3: Anticipatory gaze rates across language condition and transition type for the real and randomly permuted datasets. Vertical bars represent 95% confidence intervals.

474 included random effects of item and participant with maximal random slopes  
 475 of condition, transition type, and their interaction for participant.

476 Children's anticipatory gaze switches showed effects of language condition  
 477 ( $\beta=-3.65$ ,  $SE=1.16$ ,  $z=-3.15$ ,  $p<.01$ ) and transition type ( $\beta=-2.95$ ,  
 478  $SE=1.13$ ,  $z=-2.61$ ,  $p<.01$ ) with additional effects of an age-by-language con-  
 479 dition interaction ( $\beta=0.5$ ,  $SE=0.212$ ,  $z=2.35$ ,  $p<.05$ ), a language condition-  
 480 by-transition type interaction ( $\beta=2.69$ ,  $SE=1.35$ ,  $z=1.99$ ,  $p<.05$ ), and an  
 481 interaction between transition type and gap duration ( $\beta=5.52$ ,  $SE=2.28$ ,  
 482  $z=2.42$ ,  $p<.05$ ). There were no significant effects of age or gap duration  
 483 alone ( $\beta=-0.002$ ,  $SE=0.26$ ,  $z=-0.009$ ,  $p=.99$  and  $\beta=2.25$ ,  $SE=3.19$ ,  $z=0.7$ ,  
 484  $p=.48$ , respectively).

485 Adults' anticipatory gaze switches showed an effect of transition type  
 486 ( $\beta=-3.3$ ,  $SE=0.93$ ,  $z=-3.54$ ,  $p<.001$ ) and significant interactions between  
 487 language condition and transition type ( $\beta=1.23$ ,  $SE=0.63$ ,  $z=1.96$ ,  $p<.05$ )  
 488 and transition type and gap duration ( $\beta=7.12$ ,  $SE=2.2$ ,  $z=3.24$ ,  $p<.01$ ).  
 489 There were no significant effects of language condition or gap duration alone  
 490 ( $\beta=-0.06$ ,  $SE=0.75$ ,  $z=-0.08$ ,  $p=.94$  and  $\beta=0.13$ ,  $SE=1.77$ ,  $z=0.08$ ,  $p=.94$ ,  
 491 respectively).

*Children*

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.604	1.242	-0.486	0.627
Age	-0.002	0.261	-0.009	0.993
LgCond= <i>non-English</i>	-3.65	1.16	-3.146	0.002 **
Type= <i>non-Question</i>	-2.95	1.13	-2.61	0.009 **
GapDuration	2.247	3.194	0.704	0.482
Age*LgCond= <i>non-English</i>	0.5	0.212	2.353	0.019 *
Age*Type= <i>non-Question</i>	0.009	0.196	0.044	0.965
LgCond= <i>non-English</i> *	2.692	1.347	1.999	0.046 *
Type= <i>non-Question</i>				
Age*GapDuration	-0.577	0.627	-0.921	0.357
LgCond= <i>non-English</i> *GapDuration	1.143	2.287	0.5	0.617
Type= <i>non-Question</i> *GapDuration	5.519	2.282	2.418	0.016 *
Age*LgCond= <i>non-English</i> *	-0.433	0.304	-1.426	0.154
Type= <i>non-Question</i>				

*Adults*

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.584	0.64	-0.913	0.361
LgCond= <i>non-English</i>	-0.059	0.751	-0.079	0.937
Type= <i>non-Question</i>	-3.298	0.933	-3.536	0.0004 ***
GapDuration	0.132	1.766	0.075	0.941
LgCond= <i>non-English</i> *	1.234	0.629	1.961	0.0498 *
Type= <i>non-Question</i>				
LgCond= <i>non-English</i> *GapDuration	-1.519	2.192	-0.693	0.488
Type= <i>non-Question</i> *GapDuration	7.116	2.195	3.241	0.001 **

Table 2: Model output for children and adults' anticipatory gaze switches in Experiment 1.

492 *Random baseline comparison*

493 Our primary analysis (above) makes the assumption that participants'  
 494 eye movements generally follow the turn structure of the stimulus, i.e., that  
 495 participants track the current speaker and switch their gaze to the upcom-  
 496 ing speaker near turn transitions. Based on this assumption, we used linear  
 497 mixed effects regressions to see how anticipatory looking is affected by as-  
 498 pects of participant group (age) and stimulus (e.g., transition type, language

499 condition). But what if the assumption that participants generally track  
500 turn structure were wrong? Could these results have emerged if partici-  
501 pants' eye movements were *not* linked to turn structure? For example, if  
502 participants were randomly looking back and forth between the two speak-  
503 ers, we might still find some anticipatory switching by chance—we would see  
504 anticipatory switching no matter where the real turn transitions occurred. To  
505 test whether our primary results (the regression output above) could have  
506 arisen from random switching, not linked to turn structure, we conducted  
507 a secondary analysis comparing participants' anticipatory gaze at *real* and  
508 *randomly shuffled* points of turn transition.

509 We conducted this analysis by running the same regression models on  
510 participants' eye-tracking data, only this time calculating their anticipatory  
511 gaze switches with respect to randomly permuted turn transition windows.  
512 This process involved: (1) randomizing the order and temporal placement  
513 of the analysis windows within each stimulus (see Figure 4; “analysis win-  
514 dow” is defined in Figure 2) to randomly redistribute the analysis windows  
515 across the eye-tracking signal, (2) re-running each participant's eye track-  
516 ing data through switch identification (described in Section b) on each of  
517 the randomly permuted analysis windows, and (3) modeling the anticipatory  
518 switches from the randomly permuted data (our random baseline dataset)  
519 with the same statistical models we used for the original dataset (Section  
520 b; Table 2). Importantly, although the onset time of each transition was  
521 shuffled within the eye-tracking signal, the other intrinsic properties of each  
522 turn transition (e.g., prior speaker identity, transition type, gap duration,  
523 language condition, etc.) stayed constant across each permutation.

524 The random shuffling procedure de-links participants' gaze data from the  
525 turn structure in the original stimulus, thereby allowing us to compare turn-  
526 related (original) and non-turn-related (randomly permuted) looking behav-  
527 ior using the same eye movement data. We created 5,000 permutations of the  
528 original turn transitions (step 1 above), thereby creating 5,000 anticipatory  
529 gaze datasets with randomly de-linked gaze data (step 2 above). Because the  
530 randomly shuffled turn transitions could occur anywhere in the stimulus (so  
531 long as they didn't overlap each other), the resulting turn-transition windows  
532 collectively covered the entire stimulus—during speech and silence, during  
533 speaker change and speaker continuation, and during all turn transitions in  
534 the stimulus, even those excluded in the original analyses (e.g., because they

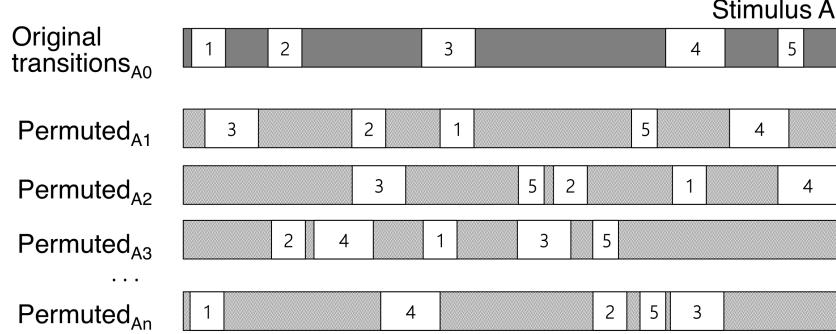


Figure 4: Example of analysis window permutations for a stimulus with five turn transitions. The windows included  $\pm 300$  msec around the inter-turn gap.

were overlapped).<sup>9</sup> Pooled together, the anticipatory gaze datasets yield an average anticipatory switch rate for each participant over all possible starting points in the stimuli: a random baseline.

Using this technique we can compare participants' anticipatory switches at turn transition windows to their anticipatory switches over the stimulus as a whole. If participants were looking randomly back and forth between the speakers, we should similar patterns in both cases. Rather than simply comparing participants' overall anticipatory switch rates with real and random transition windows, we estimated the likelihood that each of the predictor effects in the original data (e.g., the effect of language condition; Table 2) could have arisen with random gaze switching: we ran identical statistical models on the real and randomly permuted data sets. This tells us not only whether participants' switches were above chance, but whether the specific underlying effects of their anticipatory gaze patterns (e.g., the effect of linguistic condition) were above that expected by chance. But, because these analyses are complex and secondary to the main results, we report their full details in Appendix A.

Our baseline analyses revealed that none of the significant predictors from models of the original, turn-related data can be explained by random looking. For the children's data, the original  $z$ -values for language condition, transi-

---

<sup>9</sup>This technique crucially differs from that used by Keitel and colleagues (2013, 2015), which tests anticipatory gaze at turn transitions against anticipatory gaze during speech, thereby maximizing the possibility of finding a difference from the baseline measure.

555 tion type, the age-language condition interaction, the transition type-gap  
556 duration interaction, and the language condition-transition type interaction  
557 were all greater than 95% of  $z$ -values from models of the randomly permuted  
558 data (99.3%, 99.1%, 98.9%, 97%, and 96%, respectively, all  $p < .05$ ). Simi-  
559 larly, the adults' data showed significant differentiation from the randomly  
560 permuted data for all three significant predictors from the real transition  
561 dataset—transition type, the interaction between transition type and gap du-  
562 ration, and the interaction between language condition and transition type  
563 (greater than 100%, 99.8%, and 95% of random  $z$ -values, respectively, all  
564  $p \leq .05$ ). See Section Appendix A for more information on each predictor's  
565 random permutation distribution.<sup>10</sup>

566 *Developmental effects*

567 The models reported above revealed a significant interaction of age and  
568 language condition (Table 2) that was unlikely be due to random gaze switch-  
569 ing (Figure 3). To further explore this effect, we compared the effect of lan-  
570 guage condition across age groups: using the permuted datasets described  
571 above, we extracted the average difference score for the two language condi-  
572 tions (English minus non-English) for each participant, computing an overall  
573 average for each random permutation of the data. Then, within each per-  
574 mutation, we made pairwise comparisons of the average difference scores  
575 across participant age groups. This process yielded a distribution of ran-  
576 dom permutation-based difference scores that we could then compare to the  
577 difference score in the actual data. Details are given in Appendix B.

578 These analyses revealed that, while 3- and 4-year olds showed similarly-  
579 sized effects of language condition, 5-year-olds had a significantly smaller  
580 effect of language condition, compared to both younger age groups. The  
581 difference in the language condition effect between 5-year-olds and 3-year-  
582 olds was greater than would be expected by chance (99.52% of the randomly  
583 permuted data sets;  $p < .01$ ). Similarly, the difference in the language con-  
584 dition effect between 5-year-olds and 4-year-olds was greater than would be  
585 expected by chance (99.96% of the data sets;  $p < .001$ ). See Figure B.1 for

---

<sup>10</sup>This baseline analysis tests “random looking” against “turn-driven looking”, but it does not test subtypes of turn-driven looking. For example, children might switch their gaze from the current speaker to the addressee out of boredom with the ongoing speech rather than active anticipation of an upcoming response. We address this hypothesis about “turn-transition” gaze switchers vs. “boredom” gaze switchers in Appendix C

586 each difference score distribution.

587 Given these findings, when does spontaneous turn prediction emerge de-  
588 velopmentally? We tested whether the youngest age group (3-year-olds)  
589 already exceeded chance in their anticipatory gaze switches by comparing  
590 children's real gaze rates to the random baseline in the English condition  
591 with two-tailed *t*-tests. We used the English condition because we are most  
592 interested in finding out when children begin to make spontaneous turn pre-  
593 dictions for natural speech. We found that three-year-olds made anticipa-  
594 tory gaze switches significantly above chance, when all transitions were con-  
595 sidered ( $t(22.824)=-4.147$ ,  $p<.001$ ) as well as for question transitions alone  
596 ( $t(21.677)=-5.268$ ,  $p<.001$ ).

597 *Discussion*

598 Children and adults spontaneously tracked the turn structure of the con-  
599 versations, making anticipatory gaze switches at an above-chance rate across  
600 all ages and conditions. Children's anticipatory gaze rates were affected by  
601 language condition, transition type, age, and gap duration (Table 2), none of  
602 which could be explained by a baseline of random gaze switching (Appendix  
603 A; Figure A.1a). These data show a number of important features that bear  
604 on our questions of interest.

605 First, both adults' and children's anticipations were strongly affected by  
606 transition type. Both groups made more anticipatory switches after hear-  
607 ing questions, compared to non-questions, especially for the English stimuli  
608 compared to the non-English stimuli. Overall, participants made few antici-  
609 patory switches after non-questions, even in the English videos when they  
610 had full linguistic access. Prior work using online, metalinguistic tasks has  
611 shown that participants can use linguistic cues to accurately predict upcom-  
612 ing turn ends (Torreira et al., 2015; Magyari & De Ruiter, 2012; De Ruiter  
613 et al., 2006). The current results add a new dimension to our understanding  
614 of how listeners make predictions about turn ends: both children and adults  
615 spontaneously monitor the linguistic structure of unfolding turns for cues to  
616 imminent responses.

617 Second, children made more anticipatory switches overall in English videos,  
618 compared to non-English videos. This effect suggests that linguistic access is  
619 important for children's ability to anticipate upcoming turn structure, con-  
620 sistent with prior work on turn-end prediction in adults (De Ruiter et al.,  
621 2006; Magyari & De Ruiter, 2012) and children (Keitel et al., 2013).

622     Third, we saw that older children made anticipatory switches more re-  
623 liably than younger children, but only in the non-English videos. In the  
624 English videos, children anticipated well at all ages, especially after hear-  
625 ing questions. This interaction between age and language condition suggests  
626 that the 5-year-olds were able to leverage anticipatory cues in the non-English  
627 videos in a way that 3- and 4-year-olds could not, possibly by shifting more  
628 attention to the non-English prosodic or non-verbal cues. Prior work on chil-  
629 dren’s turn-structure anticipation has proposed that children’s turn-end pre-  
630 dictions rely primarily on lexicosyntactic structure (and not, e.g., prosody)  
631 as they get older (Keitel et al., 2013). The current results suggest more  
632 flexibility in children’s predictions; when they do not have access to lexical  
633 information, older children and adults are likely to find alternative cues to  
634 turn taking behavior.

635     Finally, children and adults made more anticipatory switches in tran-  
636 sitions with longer inter-turn gaps, though this effect was limited to non-  
637 question turns (Table 2). This finding suggests that gap duration indeed  
638 serves as a cue to upcoming turn structure; while short gaps may be per-  
639 ceived as within-turn pauses (Männel & Friederici, 2009), long gaps could  
640 instead be indicative of between-turn pauses (where speaker transition oc-  
641 curs). Participants might use these long silences to retroactively assign turn  
642 boundaries and anticipate speaker switches that were otherwise not antici-  
643 pated (in this case, because the preceding turn was not a question). An  
644 alternative explanation for gap duration effects is that longer inter-turn gaps  
645 result in longer analysis windows, which yields more time for participants to  
646 make an anticipatory gaze. However, if participants are generally more likely  
647 to make a switch at question transitions, as our results suggest, we would ex-  
648 pect that longer gaps would benefit questions more than non-questions—the  
649 opposite pattern from what the data show here. We take this as evidence  
650 that inter-turn silence may be most useful when participants have limited  
651 ability to make predictions about upcoming speaker transitions.

652     In Experiment 2, we followed up on these findings, improving on two  
653 aspects of the design: first, our language manipulation in this first experi-  
654 ment was too coarse to provide data regarding specific linguistic information  
655 channels (e.g., the effect of prosodic information alone). In Experiment 2, we  
656 compared lexicosyntactic and prosodic cues with phonetically altered speech  
657 and used puppets to eliminate non-verbal cues to turn taking. Second, we  
658 were not able to pinpoint the emergence of anticipatory switching because  
659 the youngest age group in our sample was already able to make anticipa-

660 tory switches at above chance rates. In Experiment 2, we explored a wider  
661 developmental range.

662 **Experiment 2**

663 Experiment 2 used English-only stimuli, controlled for lexical and prosodic  
664 information, eliminated non-verbal cues, and tested children from a wider age  
665 range. To tease apart the role of lexical and prosodic information, we phoneti-  
666 cally manipulated the speech signal for pitch, syllable duration, and lexical  
667 access. By testing 1- to 6-year-olds we hoped to find the developmental onset  
668 of turn-predictive gaze. We also hoped to measure changes in the relative  
669 roles of prosody and lexicosyntax across development.

670 Non-verbal gestural cues in Experiment 1 could have helped partici-  
671 pants make predictions about upcoming turn structure (Rossano et al., 2009;  
672 Stivers & Rossano, 2010). Since our focus here is on linguistic cues, we  
673 eliminated all gaze and gestural signals in Experiment 2 by replacing the  
674 videos of human actors with videos of puppets. Puppets are less realis-  
675 tic and expressive than human actors, but they create a natural context for  
676 having somewhat motionless talkers in the videos. Additionally, the prosody-  
677 controlled condition (described below) included small but global changes to  
678 syllable duration that would have required complex video manipulation or  
679 precise re-enactment with human talkers, neither of which was feasible. For  
680 these reasons, we decided to use puppet videos rather than human videos in  
681 the final stimuli. As in the first experiment, we recorded participants' eye  
682 movements as they watched six short videos of dyadic conversation, and then  
683 analyzed their anticipatory glances from the current speaker to the upcoming  
684 speaker at points of turn transition.

685 *Methods*

686 *Participants*

687 We recruited 27 undergraduate adults and 129 children ages 1;0–6;11 to  
688 participate in our experiment. We recruited adult participants via the Stan-  
689 ford University Psychology participant database. Adult participants were  
690 either paid or received course credit for their time. We recruited our child  
691 participants from the Children's Discovery Museum in San Jose, California<sup>11</sup>,

---

<sup>11</sup>We ran Experiment 2 at a local children's museum because it gave us access to children with a wider range of ages. Participants were volunteers.

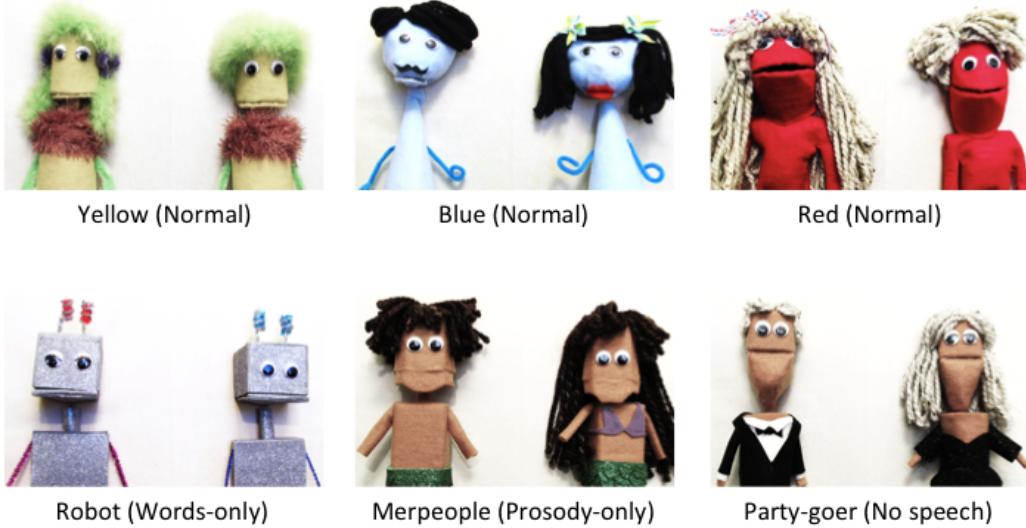


Figure 5: The six puppet pairs (and associated audio conditions). Each pair was linked to three distinct conversations from the same condition across the three experiment versions.

targeting approximately 20 children for each of the six one-year age groups (range: 20–23). All participants were native English speakers, though some parents ( $N=27$ ) reported that their child heard a second (and sometimes third) language at home. None of the adult participants reported fluency in a second language.

#### Materials

We created 18 short videos of improvised, child-friendly conversation (Figure 5). To eliminate non-verbal cues to turn transition and to control the types of linguistic information available in the stimuli we first audio-recorded improvised conversations, then phonetically manipulated those recordings to limit the availability of prosodic and lexical information, and finally recorded video to accompany the manipulated audio, featuring puppets as talkers.

*Audio recordings.* The recording session was set up in the same way as the first experiment, but with a shorter warm up period (5–10 minutes) and a pre-determined topic for the child-friendly improvisation ('riding bikes', 'pets', 'breakfast', 'birthday cake', 'rainy days', or 'the library'). All of the talkers were native English speakers, and were recorded in male-female pairs. As before, we asked talkers to speak "as if they were on a children's television

710 show” and to ask at least a few questions during the improvisation. We cut  
711 each audio recording down to the (approximate) 20-second interval with the  
712 most turn activity. The 20-second clips were then phonetically manipulated  
713 and used in the final video stimuli.

714 *Audio Manipulation.* We created four versions of each audio conversation:  
715 *normal*, *words only*, *prosody only*, and *no speech*. That is, one version  
716 with a full linguistic signal (*normal*), and three with incomplete linguistic  
717 information (hereafter “partial cue” conditions). The *normal* conversations  
718 were the unmanipulated, original audio clips.

719 The *words only* conversations were manipulated to have robot-like speech:  
720 we flattened the intonation contours to each talker’s average pitch ( $F_0$ ) and  
721 we reset the duration of every nucleus and coda to each talker’s average  
722 nucleus and coda duration.<sup>12</sup> We made duration and pitch manipulations  
723 using PSOLA resynthesis in Praat (Boersma & Weenink, 2012). Thus, the  
724 *words only* versions of the conversations had no pitch or durational cues  
725 to upcoming turn boundaries, but did have intact lexicosyntactic cues (and  
726 some residual phonetic correlates of prosody, e.g., intensity).

727 We created the *prosody only* conversations by low-pass filtering the orig-  
728 inal recording at 500 Hz with a 50 Hz Hanning window (following de Ruiter  
729 et al., 2006). This manipulation creates a “muffled speech” effect because  
730 low-pass filtering removes most of the phonetic information used to distin-  
731 guish between phonemes. The *prosody only* versions of the conversations  
732 lacked lexical information, but retained their intonational and rhythmic cues  
733 to upcoming turn boundaries.

734 The *no speech* condition served as a non-linguistic baseline. For this  
735 condition, we replaced the original audio clip for the conversation with multi-  
736 talker babble: we overlaid multiple child-oriented conversations (excluding  
737 the original one), and then cropped the result to the duration of the original  
738 conversation clip. Thus, the *no speech* conversation lacked any linguistic  
739 information to upcoming turn boundaries—the only cue to turn taking was  
740 the opening and closing of the puppets’ mouths.

741 Finally, because low-pass filtering removes significant acoustic energy, the  
742 *prosody only* conversations were much quieter than the other three conditions.  
743 Our last step was to downscale the intensity of the audio tracks in the three

---

<sup>12</sup>We excluded hyper-lengthened words like [wau:] ‘woooow!', but these were rare to begin with.

744 other conditions to match the volume of the *prosody only* clips. We referred  
745 to the conditions as “normal”, “robot”, “mermaid”, and “birthday party”  
746 speech when interacting with participants.

747 *Video recordings.* We created puppet video recordings to match the ma-  
748 nipulated 20-second audio clips. The puppets were minimally expressive;  
749 the puppeteer could only control the opening and closing of their mouths.  
750 The puppets’ heads, eyes, arms, and bodies stayed still. Puppets were posi-  
751 tioned side-by-side, looking in the same direction to eliminate shared gaze as  
752 a cue to turn structure (Thorgrímsson et al., 2015). We took care to match  
753 the puppets’ mouth movements to the syllable onsets as closely as possible,  
754 specifically avoiding mouth movement before the onset of a turn. We then  
755 added the manipulated audio clips to the puppet video recordings by hand  
756 with video editing software.

757 We used three pairs of puppets used for the *normal* condition—‘red’,  
758 ‘blue’ and ‘yellow’—and one pair of puppets for each partial cue condition:  
759 ‘robots’, ‘merpeople’, and ‘party-goers’ (Figure 8). We randomly assigned  
760 half of the conversation topics (‘birthday cake’, ‘pets’, and ‘breakfast’) to  
761 the *normal* condition, and half to the partial cue conditions (‘riding bikes’,  
762 ‘rainy days’, and ‘the library’). We then created three versions of the experi-  
763 ment, so that each of the six puppet pairs was associated with three different  
764 conversation topics across the different versions of the experiment (18 videos  
765 in total; 6 videos per experiment version). We ensured that the position of  
766 the talkers (left and right) was counterbalanced in each version by flipping  
767 the video and audio channels as needed.

768 As before, the duration of turn transitions and the number of speaker  
769 changes across videos was variable because the conversations were recorded  
770 semi-spontaneously. We measured turn transitions from the audio signal of  
771 the *normal*, *words only*, and *prosody only* conditions. There was no audio  
772 from the original conversation in the *no speech* condition videos, so we mea-  
773 sured turn transitions from puppets’ mouth movements in the video signal,  
774 using ELAN video annotation software (Wittenburg et al., 2006).

775 There were 85 turn transitions for analysis after excluding transitions  
776 longer than 550 msec and shorter than 90 msec. The remaining turn transi-  
777 tions had more questions than non-questions (N=47 and N=38, respectively),  
778 with transitions distributed somewhat evenly across conditions (keeping in  
779 mind that there were three *normal* videos and only one partial cue video for  
780 each experiment version): *normal* (N=36), *words only* (N=13), *prosody only*  
781 (N=17), and *no speech* (N=19). Inter-turn gaps for questions (mean=366,

Age group	Speaker	Addressee	Other onscreen	Offscreen
1	0.44	0.14	0.23	0.19
2	0.50	0.13	0.24	0.14
3	0.47	0.12	0.25	0.16
4	0.48	0.11	0.29	0.12
5	0.54	0.11	0.20	0.14
6	0.60	0.12	0.18	0.10
Adult	0.69	0.12	0.09	0.10

Table 3: Average proportion of gaze to the current speaker and addressee during periods of talk across ages.

782 median=438, stdev=138 msec) were longer than those for non-questions  
 783 (mean=305, median=325, stdev=94 msec) on average, but gap duration  
 784 was overall comparable across conditions: *normal* (mean=334, median=321,  
 785 stdev=130 msec), *words only* (mean=347, median=369, stdev= 115 msec),  
 786 *prosody only* (mean=365, median=369, stdev=104 msec), and *no words*  
 787 (mean=319, median=329, stdev=136 msec).

#### 788 *Procedure*

789 We used the same experimental apparatus and procedure as in the first  
 790 experiment. Each participant watched six puppet videos in random order,  
 791 with five 15–30 second filler videos placed in-between (e.g., running puppies,  
 792 moving balls, flying bugs). Three of the puppet videos had *normal* audio  
 793 while the other three had *words only*, *prosody only*, and *no speech* audio. As  
 794 before, the experimenter immediately began each session with calibration and  
 795 then stimulus presentation. Participants were given no instruction about how  
 796 to watch the videos or what their purpose was, they were simply encouraged  
 797 to watch the “fun/nice puppet videos”. The entire experiment took less than  
 798 five minutes.

#### 799 *Data preparation and coding*

800 We coded each turn transition for its linguistic condition (*normal*, *words*  
 801 *only*, *prosody only*, and *no speech*) and transition type (question/non-question)<sup>13</sup>,

---

<sup>13</sup>We coded *wh*-questions as “non-questions” for the *prosody only* videos. Polar questions often have a final rising intonational contour, but *wh*-questions often do not (Hedberg et al.,

Condition	Speaker	Addressee	Other onscreen	Offscreen
Normal	0.58	0.12	0.17	0.13
Words only	0.54	0.11	0.24	0.10
Prosody only	0.48	0.12	0.26	0.15
No speech	0.44	0.13	0.26	0.18

Table 4: Average proportion of gaze to the current speaker and addressee during periods of talk across conditions.

and identified anticipatory gaze switches to the upcoming speaker using the methods from Experiment 1.

#### Results

Participants' pattern of gaze indicated that they performed basic turn tracking across all ages and in all conditions. Participants looked at the screen most of the time during video playback (82% and 86% average for children and adults, respectively), primarily looking at the person who was currently speaking (Table 2). They tracked the current speaker in every condition—even one-year-olds looked more at the current speaker than at anything else in the three partial cue conditions (40% for *words only*, 43% for *prosody only*, and 39% for *no speech*). There was a steady overall increase in looks to the current speaker with age and added linguistic information (Tables 3 and 4). Looks to the addressee also decreased with age, but the change was minimal. Figure 6 shows participants' anticipatory gaze rates across age, the four language conditions, and transition type.

#### Statistical models

We identified anticipatory gaze switches for all 85 usable turn transitions, and analyzed them for effects of language condition, transition type, and age with two mixed-effects logistic regressions. We again built separate models for children and adults because effects of age were only pertinent to the children's data. The child model included condition (*normal/prosody only/words only/no speech*; with *no speech* as the reference level), transition type (question vs. non-question), age (1, 2, 3, 4, 5, 6; numeric, intercept as age=0), and duration of the inter-turn gap (in seconds) as predictors, with

---

2010).

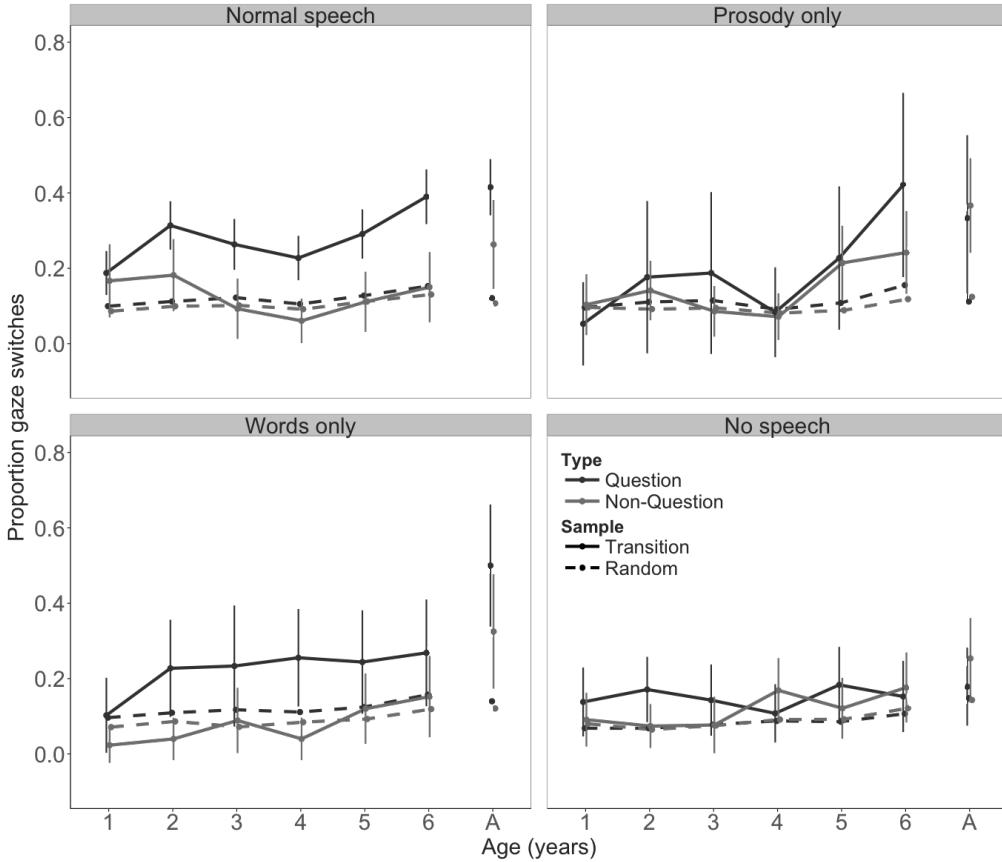


Figure 6: Anticipatory gaze rates across language condition and transition type for the real and randomly permuted datasets. Vertical bars represent 95% confidence intervals.

full interactions between language condition, transition type, and age and two-way interactions between gap duration and the other basic predictors (age, linguistic condition, and transition type). We also included random effects of participant and item (turn transition), with maximal random slopes of transition type for participant. The adult model included condition, transition type, their interactions, gap duration, and two-way interactions between gap duration and condition and transition type, with participant and item as random effects and maximal random slopes of condition and transition type for participant.

Children's anticipatory gaze switches showed an effect of gap duration

836 ( $\beta=3.85$ ,  $SE=1.73$ ,  $z=2.22$ ,  $p<.05$ ), a two-way interaction of age and lan-  
837 guage condition (for *prosody only* speech compared to the *no speech* reference  
838 level;  $\beta=0.38$ ,  $SE=0.19$ ,  $z=1.97$ ,  $p<.05$ ), a marginal two-way interaction of  
839 language condition and gap duration (for *prosody only* speech compared to  
840 the *no speech* reference level;  $\beta=-4.77$ ,  $SE=2.63$ ,  $z=-1.82$ ,  $p=.07$ ), and a  
841 three-way interaction of age, transition type, and language condition (for  
842 *normal* speech compared to the *no speech* reference level;  $\beta=-0.35$ ,  $SE=0.17$ ,  
843  $z=-2.05$ ,  $p<.05$ ). There were no significant effects of age or transition type  
844 alone (Table b) ( $\beta=-0.05$ ,  $SE=0.14$ ,  $z=-0.38$ ,  $p=.7$  and  $\beta=-1.22$ ,  $SE=0.96$ ,  
845  $z=-1.27$ ,  $p=.2$ , respectively)

846 Adults' anticipatory gaze switches showed a significant effect of language  
847 condition (for *words only* speech compared to the *no speech* reference level;  
848  $\beta=3.79$ ,  $SE=1.62$ ,  $z=2.34$ ,  $p<.05$ ) and a marginal two-way interaction be-  
849 tween language condition and transition type (for *words only* speech com-  
850 pared to the *no speech* reference level;  $\beta=-1.68$ ,  $SE=0.89$ ,  $z=-1.89$ ,  $p=.06$ ).  
851 There was no significant effect of transition type alone ( $\beta=-0.02$ ,  $SE=1.44$ ,  
852  $z=-0.02$ ,  $p=.99$ ).

853 *Random baseline comparison*

854 Using the same technique described in Experiment 1 (Section b), we cre-  
855 ated and modeled random permutations of participants' anticipatory gaze  
856 switches. These analyses revealed that the significant predictors from mod-  
857 els of the original, turn-related data were unlikely to be explained by random  
858 looking. In the children's data, the original model's  $z$ -values for gap dura-  
859 tion, the two-way interaction of age and language condition (*prosody only*)  
860 and the three-way interaction of age, transition type, and language condi-  
861 tion (*normal* speech) were all greater than 93% of the randomly permuted  
862  $z$ -values (95.6%, 94%, and 93.3%, respectively,  $p=.04$ , .06, and .07). Simi-  
863 larly, the adults' data showed significant differentiation from the randomly  
864 permuted data for the effect of language condition (*words only* speech; greater  
865 than 98.3% of random  $z$ -values,  $p<.02$ ). See Section Appendix A for more  
866 information on each predictor's random permutation distribution.

<i>Children</i>	Estimate	Std. Error	<i>z</i> value	Pr(>  <i>z</i>  )
(Intercept)	-3.452	0.76	-4.543	5.55e-06 ***
Age	-0.054	0.143	-0.379	0.705
Type= <i>non-Question</i>	-1.217	0.958	-1.27	0.204
GapDuration	3.852	1.735	2.221	0.026 *
Age*Type= <i>non-Question</i>	0.152	0.141	1.081	0.28
Age*GapDuration	0.214	0.266	0.805	0.421
Type= <i>non-Question</i> *	0.995	2.134	0.466	0.641
GapDuration				
Condition= <i>normal</i>	0.54	0.742	0.728	0.467
Age*Condition= <i>normal</i>	0.125	0.103	1.221	0.222
Condition= <i>normal</i> *	0.908	0.748	1.215	0.224
Type= <i>non-Question</i>				
Age*Condition= <i>normal</i> *	-0.355	0.173	-2.051	0.04 *
Type= <i>non-Question</i>				
Condition= <i>normal</i> *	-0.431	1.67	-0.258	0.797
GapDuration				
Condition= <i>prosody</i>	0.549	1.452	0.378	0.705
Age*Condition= <i>prosody</i>	0.375	0.191	1.967	0.049 *
Condition= <i>prosody</i> *	1.076	1.105	0.974	0.33
Type= <i>non-Question</i>				
Age*Condition= <i>prosody</i> *	-0.296	0.235	-1.257	0.209
Type= <i>non-Question</i>				
Condition= <i>prosody</i> *	-4.767	2.625	-1.816	0.069 .
GapDuration				
Condition= <i>words</i>	0.684	1.06	0.645	0.519
Age*Condition= <i>words</i>	0.127	0.136	0.934	0.350
Condition= <i>words</i> *	-1.244	1.031	-1.207	0.228
Type= <i>non-Question</i>				
Age*Condition= <i>words</i> *	0.111	0.225	0.495	0.621
Type= <i>non-Question</i>				
Condition= <i>words</i> *	-2.285	2.232	-1.024	0.306
GapDuration				

Table 5: Model output for children's anticipatory gaze switches in Experiment 2.

<b>Adults</b>	Estimate	Std. Error	<i>z</i> value	Pr(>  <i>z</i>  )
(Intercept)	-3.117	1.176	-2.649	0.008 **
Type= <i>non-Question</i>	-0.022	1.44	-0.015	0.988
GapDuration	4.073	2.947	1.382	0.167
Type= <i>non-Question</i> *	1.304	3.859	0.338	0.735
GapDuration				
Condition= <i>normal</i>	0.39	1.316	0.296	0.767
Condition= <i>normal</i> *	-0.709	0.754	-0.94	0.347
Type= <i>non-Question</i>				
Condition= <i>normal</i> *	2.1	3.336	0.629	0.529
GapDuration				
Condition= <i>prosody</i>	0.757	2.193	0.345	0.73
Condition= <i>prosody</i> *	0.386	1.065	0.362	0.717
Type= <i>non-Question</i>				
Condition= <i>prosody</i> *	-1.118	4.543	-0.246	0.805
GapDuration				
Condition= <i>words</i>	3.792	1.621	2.338	0.019 *
Condition= <i>words</i> *	-1.678	0.889	-1.888	0.059 .
Type= <i>non-Question</i>				
Condition= <i>words</i> *	-5.653	3.861	-1.464	0.143
GapDuration				

Table 6: Model output for adults' anticipatory gaze switches in Experiment 2.

869 *Developmental effects*

870 Our main goal in extending the age range to 1- and 2-year-olds in Experiment 871 was to find the age of emergence for spontaneous predictions about 872 upcoming turn structure. As in Experiment 1, we used two-tailed *t*-tests 873 to compare children's real gaze rates to the random baseline rates in the 874 *normal* speech condition, in which the speech stimulus is most like what 875 children hear every day. We tested real gaze rates against baseline rates 876 for three age groups: one-, two-, and three-year-olds. Two- and three-year- 877 old children made anticipatory gaze switches significantly above chance both 878 when all transitions were considered (2-year-olds:  $t(26.193)=-4.137$ ,  $p<.001$ ; 879 3-year-olds:  $t(22.757)=-2.662$ ,  $p<.05$ ) and for question transitions alone (2- 880 year-olds:  $t(25.345)=-4.269$ ,  $p<.001$ ; 3-year-olds:  $t(21.555)=-3.03$ ,  $p<.01$ ). 881 One-year-olds, however, only made anticipatory gaze shifts marginally above 882 chance for turn transitions overall and for question turns alone (overall: 883  $t(24.784)=-2.049$ ,  $p=.051$ ; questions:  $t(25.009)=-2.03$ ,  $p=.053$ ).

884 We also tested the two baseline linguistic conditions against each other—  
885 *no speech* and *normal speech*—to find out when linguistic information made  
886 a difference in children’s anticipations. Because, as we have seen, children  
887 primarily show linguistic effects in question-answer turn transitions, we in-  
888 vestigated the use of linguistic cues across age by testing anticipation sep-  
889 arately for question and non-question turns. Compared to the *no speech*  
890 condition, children made significantly more anticipatory switches in the *nor-*  
891 *mal* speech condition for questions at ages 6, 4, and 3, and also marginally at  
892 age 2 (6-year-olds:  $t(36.919)=3.8019$ ,  $p<.001$ ; 4-year-olds:  $t(41.449)=2.9777$ ,  
893  $p<.01$ ; 3-year-olds:  $t(35.724)=2.4286$ ,  $p<.05$ ; 2-year-olds:  $t(41.078)=1.8018$ ,  
894  $p=.079$ ). Children’s anticipatory switches for questions did not differ in  
895 the *no speech* and *normal* speech conditions at ages 5 or 1 (5-year-olds:  
896  $t(29.406)=1.2783$ ,  $p=.211$ ; 1-year-olds:  $t(35.907)=0.4961$ ,  $p=.623$ ). In con-  
897 trast, children’s anticipatory switch rates for non-question turns were not  
898 significantly different between the *no speech* and *normal* speech conditions  
899 at any age (all  $p>0.09$ ). Consistent with the regression output, children were  
900 more likely to show an effect of linguistic content as they got older, but only  
901 for question transitions.

902 The regression models for the children’s data also revealed two signifi-  
903 cant interactions with age. The first was a significant interaction of age and  
904 language condition (for *prosody only* compared to the *no speech* reference  
905 level), suggesting a different age effect between the two linguistic conditions.  
906 As in Experiment 1, we explored each age interaction by extracting an av-  
907 erage difference score over participants for the effect of language condition  
908 (*no speech* vs. *prosody only*) within each random permutation of the data,  
909 making pairwise comparisons between the six age groups. These tests re-  
910 vealed that children’s anticipation in the *prosody only* condition significantly  
911 improved at ages five and six (with difference scores greater than 95% of the  
912 random data scores;  $p<.05$ ). See Figure B.2 for these *prosody only* difference  
913 score distributions.

914 The second age-based interaction was a three-way interaction of age, tran-  
915 sition type, and language condition (for *normal* speech compared to the *no*  
916 *speech* baseline). We again created pairwise comparisons of the average dif-  
917 ference scores for the transition type-language condition interaction across  
918 age groups in each random permutation of the data, finding that the effect  
919 of transition type in the *normal* speech condition became larger with age,  
920 with significant improvements by age 4 over ages 1 and 2 (99.9% and 98.86%,  
921 respectively), by age 5 over age 4 (97.54%), and by age 6 over ages 1, 2, and 5

922 (99.5%, 97.36%, and 95.04%), all significantly different from chance ( $p < .05$ ).  
923 See Figure B.3 for these *normal* speech difference score distributions.

924 *Discussion*

925 The core aims of Experiment 2 were to gain better traction on the individual roles of prosody and lexicosyntax in children’s turn predictions, and  
926 to find the age of emergence for spontaneous turn anticipation. Many of our  
927 results replicate the findings from Experiment 1: participants often made  
928 more anticipatory switches when they had access to linguistic information  
929 and, when they did, tended to make more anticipatory switches for questions compared to non-questions.  
930

931 As in Experiment 1, children and adults spontaneously tracked the turn  
932 structure of the conversations. Participants made anticipatory gaze switches  
933 at above-chance rates starting at age one for question transitions and age two  
934 for both questions and non-questions. Longer gaps had a broader impact on  
935 participants’ anticipations in this experiment; we saw that, overall, longer  
936 inter-turn gaps resulted in more anticipatory switches, with the *no speech*  
937 condition showing equal or significantly stronger effects of gap duration than  
938 all other conditions.  
939

940 As before, participants made far more anticipations for questions than  
941 for non-question turns—at least for those two years old and older. But these  
942 effects were different for the conditions with partial linguistic information:  
943 *prosody only* and *words only*. In the *prosody only* condition, performance was  
944 initially low for young children and increased significantly with age. In the  
945 *words only* condition, children age two and older showed robust switching  
946 for questions (much like in *normal* speech), but never rose above chance for  
947 non-question turns (Figure 6), with no significant differences from the *no*  
948 *speech* baseline. These findings do not support an early role for prosody or  
949 lexical information alone in children’s spontaneous predictions about turn  
950 structure. They also give no evidence for the idea that lexical information  
951 is sufficient on its own to support children’s anticipatory switching. They  
952 do underscore the developing relationship between the use of linguistic cues  
953 and speech act (transition type) in spontaneous predictions about upcoming  
954 turn structure.

955    **General Discussion**

956    Children begin to develop conversational turn-taking skills long before  
957    their first words emerge (Bateson, 1975; Hilbrink et al., 2015; Jaffe et al.,  
958    2001; Snow, 1977). As they acquire language, they also acquire the infor-  
959    mation needed to make accurate predictions about upcoming turn structure.  
960    Until recently, we have had very little data on how children weave language  
961    into their already-existing turn-taking behaviors. In two experiments investi-  
962    gating children's anticipatory gaze to upcoming speakers, we found evidence  
963    that turn prediction develops early in childhood and that, when spontaneous  
964    predictions begin, they are primarily driven by participants' expectation of  
965    an immediate response in the next turn (e.g., after questions). In making  
966    predictions about upcoming turn structure, children used a combination of  
967    lexical and prosodic cues; neither lexical nor prosodic cues alone were suffi-  
968    cient to support increased anticipatory gaze. We also found no early advan-  
969    tage for prosody over lexicosyntax; children's anticipatory switch rates in the  
970    *prosody only* condition were initially low, but showed significant gains by age  
971    five. We discuss these findings with respect to: the role of linguistic cues and  
972    inter-turn silence for predicting upcoming turn structure, the importance of  
973    questions in predictions about conversation, and children's developing com-  
974    petence as conversationalists.

975    *Predicting upcoming turn structure*

976    Prior work with adults has found a consistent role for lexicosyntax in pre-  
977    dicting upcoming turn structure (De Ruiter et al., 2006; Magyari & De Ruiter,  
978    2012), whereas the role of prosody still under debate (Duncan, 1972; Ford  
979    & Thompson, 1996; Torreira et al., 2015). Knowing that children compre-  
980    hend more about prosody than lexicosyntax early on (Section b; also see  
981    Speer & Ito, 2009 for a review), we thought it possible that young children  
982    would instead show an advantage for prosody in their predictions about turn  
983    structure in conversation. Our results suggest that, on the contrary, exclu-  
984    sively presenting prosodic information to children limits their spontaneous  
985    predictions about upcoming turn structure until age five.

986    Using prosody alone to accurately predict turn boundaries in conversa-  
987    tion appears to be difficult for adults and children. Prosodic information is  
988    continuous, multidimensional, and indexically complex—it encodes syntactic  
989    structure, speech act, and extralinguistic orientation, sometimes simultane-  
990    ously, though without clear one-to-one mappings between form and meaning

991 (Cutler et al., 1997; Shriberg et al., 1998; Lammertink et al., 2015). For  
992 these reasons, prosodic information alone may not be enough to both (a)  
993 make precise temporal predictions about turn structure, and (b) identify  
994 question turns, which otherwise appear to drive anticipatory gaze switching.  
995 Therefore, although children show early facility with prosodic discrimina-  
996 tion (Nazzi & Ramus, 2003; Soderstrom et al., 2003; Johnson & Jusczyk,  
997 2001; Jusczyk et al., 1995; Morgan & Saffran, 1995; Mehler et al., 1988),  
998 using prosodic knowledge for turn prediction may generally be too difficult  
999 without additional information from lexical or syntactic cues.

1000 Our findings suggest that there is one prosodic cue that serves as an excep-  
1001 tion to this rule: inter-turn silence. Generally speaking, participants showed  
1002 a greater anticipatory switches for longer inter-turn gaps, but the effect of  
1003 inter-turn gap duration is strongest in our data when upcoming responses are  
1004 less predictable, whether due to the question effect (Experiment 1) or the lack  
1005 of non-verbal cues and any linguistic information (Experiment 2). Notably,  
1006 there were no significant interactions of gap duration with participant age.  
1007 This pattern of results suggests that, when predictive information about up-  
1008 coming responses is absent, long silences may be used to retroactively assign  
1009 turn-end boundaries, thereby increasing the participant's expectation for a  
1010 speaker change and promoting more anticipatory gaze switches. The lack of  
1011 interactions between age and gap duration suggests that the use of inter-turn  
1012 silence remains important for older speakers under ambiguous turn structure  
1013 conditions and the interactions between transition type and gap duration  
1014 (Experiment 1) and condition and gap duration (Experiment 2; marginal),  
1015 suggest that this effect is not straightforwardly the result of having more  
1016 time to make a gaze switch.

1017 Because inter-turn silence works as a backwards-looking cue, it's pre-  
1018 dictive utility for speaker change is diminished compared to forward-looking  
1019 linguistic cues (e.g., syntactic constituent completion). However, it is a strat-  
1020 egy that can be usefully employed to predict upcoming responses some of the  
1021 time, even without access to language. Pauses are detected and related to  
1022 phrasal structure from early on; 5-month-old infants use pauses to parse in-  
1023 tonational phrases (Männel & Friederici, 2009). Our findings thus suggest  
1024 that silence is an early and lasting cue for identifying turn structure online  
1025 when other predictive information is not adequate.

1026 Perhaps surprisingly, we found no evidence that lexical information alone  
1027 is equivalent to the full linguistic signal in driving children's predictions, as  
1028 has been shown previously for adults (Magyari & De Ruiter, 2012; De Ruiter

1029 et al., 2006) and as is replicated with adult participants in the current study.  
1030 Unlike prosodic cues, lexicosyntactic cues are discreet and have much clearer  
1031 form-to-meaning mappings, with clear lexicosyntactic cues to questionhood  
1032 that occur early within turns (e.g., *wh*-words, *do*-insertion, and subject-  
1033 auxiliary inversion). That said, although children's lexical and syntactic  
1034 knowledge begins to develop early, it is limited for quite some time (Bergel-  
1035 son & Swingley, 2013; Shi & Melancon, 2010; Tomasello & Brooks, 1999).  
1036 Our stimuli were made in a child-friendly style, they are still other-directed  
1037 and fairly complex, with 20–30 seconds of continuous conversational speech.

1038 It is perhaps for this reason that children's performance was always best  
1039 with the full signal, where lexicosyntactic information was supported by  
1040 prosodic information and vice versa. There may be something extra infor-  
1041 mative about the combination of prosodic and lexical cues to questionhood  
1042 that helps to boost children's anticipations before they can use these cues  
1043 separately. Even in adults, Torreira and colleagues (2015) showed that the  
1044 trade-off in informativity between lexical and prosodic cues is more subtle in  
1045 semi-natural speech. The present findings are the first to show evidence of a  
1046 similar effect developmentally.

1047 Finally, because most anticipatory switches were made following question  
1048 turns, one could predict that the only lexical and prosodic cues that matter  
1049 for spontaneous turn prediction are those that cue questionhood. Our lin-  
1050 guistic manipulations here were focused on lexical and prosodic information  
1051 globally, which may have affected lexical and prosodic cues to questionhood  
1052 asymmetrically; while lexicosyntactic question cues were available on ev-  
1053 ery instance of *wh*- and *yes/no* questions in our stimuli, prosodic question  
1054 cues were only salient on *yes/no* questions. Relatedly, many non-verbal cues  
1055 also encode information about transition type, including gaze and gesture.  
1056 We did not systematically test those cues here but, like inter-turn silence,  
1057 they may play a critical role in parsing and making predictions about turn  
1058 structure when other linguistic information is not sufficient to make accurate  
1059 predictions.

#### 1060 *The question effect*

1061 In both experiments, anticipatory looking was primarily driven by ques-  
1062 tion transitions, a pattern that has not been previously reported in other an-  
1063 ticipatory gaze studies, on children or adults (Keitel et al., 2013; Hirvenkari,  
1064 2013; Tice and Henetz, 2011). Questions make an upcoming speaker switch  
1065 immediately relevant, helping the listener to predict with high certainty what

1066 will happen next (i.e., an answer from the addressee), and are often easily  
1067 identifiable by overt prosodic and lexicosyntactic cues.

1068 Prior work on children's acquisition of questions indicates that they may  
1069 already have some knowledge of question-answer sequences by the time they  
1070 begin to speak: questions make up approximately one third of the utter-  
1071 ances children hear, before and after the onset of speech, and even into  
1072 their preschool years, though the type and complexity of questions changes  
1073 throughout development (Casillas et al., In press; Fitneva, 2012; Henning  
1074 et al., 2005; Shatz, 1979).<sup>14</sup> For the first few years, many of the questions  
1075 directed to children are "test" questions—questions that the caregiver al-  
1076 ready has the answer to (e.g., "What does a cat say?"), but this changes as  
1077 children get older. Questions help caregivers to get their young children's  
1078 attention and to ensure that information is in common ground, even if the  
1079 responses are non-verbal or infelicitous (Bruner, 1985; Fitneva, 2012; Snow,  
1080 1977). Moreover, because of their high frequency and relatively limited num-  
1081 ber of formats, questions, especially *wh*-questions, may be more identifiable  
1082 and predictable compared to other types of speech acts. So, in addition to  
1083 having a special interactive status (for adults and children alike), questions  
1084 are a frequent, predictable and core characteristic of many caregiver-child  
1085 interactions, motivating a general benefit for questions in turn structure an-  
1086 ticipation.

1087 Two important questions for future work are then: (1) how does children's  
1088 ability to monitor for questions in conversation relate to their prior experience  
1089 with questions? and (2) what is it about questions that makes children and  
1090 adults more likely to anticipatorily switch their gaze to addressees? If this  
1091 effect is particular to turns that require an immediate response ("adjacency  
1092 pairs"; Schegloff, 2007), other turn types, such as imperatives, compliments,  
1093 and complaints should show similar patterns to questions. If the effect is  
1094 instead about overall predictability of the syntactic frame, children would in-  
1095 stead show similar patterns for other frequent frames in child-directed speech  
1096 (e.g., "Look at the X."); Mintz, 2003). Recognizability and predictability are  
1097 likely to play a role as children become more sophisticated language users,  
1098 even if the effect is truly about adjacency pairs; for example, rhetorical and  
1099 tag questions take a very similar form to prototypical polar questions, but

---

<sup>14</sup>There is substantial variation in question frequency by individual and socioeconomic class (Hart & Risley, 1992; Weisleder, 2012).

1100 usually do not require an answer. So, though it is clear that adults and chil-  
1101 dren anticipated responses more often for questions than non-questions, we  
1102 do not yet know whether their predictive action is limited to turns formatted  
1103 as questions, turns with high recognizability and predictability, or turns that  
1104 project an immediate response from the addressee.

1105 The question effect suggests that participants' spontaneous predictions,  
1106 at least while viewing third-party conversation, may be driven by what lies  
1107 *beyond* the end of the current turn—not just by the upcoming end of the turn  
1108 itself, as has been focused on in prior work (Torreira et al., 2015; Keitel et al.,  
1109 2013; Magyari & De Ruiter, 2012; De Ruiter et al., 2006). In future work, it  
1110 will be crucial to measure prediction from a first-person perspective to find  
1111 out what kinds of predictions are most relevant to addressees in conversation  
1112 (see also Holler & Kendrick, 2015).

1113 One possible scenario is that listeners in spontaneous, first-person con-  
1114 versation use multiple strategies to make predictions about upcoming turn  
1115 structure: they could use semi-passively attend to incoming speech for cues to  
1116 upcoming speaker transition (e.g., questions and other adjacency pairs) and,  
1117 when possible upcoming transition is detected, switch into a more precise  
1118 turn-end prediction mode (à la De Ruiter et al., 2006). A flexible prediction  
1119 system like this one allows listeners to continuously monitor ongoing conver-  
1120 sation for turn-related cues at a low cost while still managing to plan their  
1121 responses and come in quickly when needed.

1122 To test this hypothesis, we would need to look at prediction from a first-  
1123 person perspective, which very little work so far has accomplished (present  
1124 work included). Although third-party measures enable us to measure partic-  
1125 ipants' predictions without any interference from language production, they  
1126 also limit our knowledge about how the need to give a response might it-  
1127 self play an important role in addressees' prediction strategies. Recent work  
1128 has shown that shifts in addressee gaze similar to those measured here in-  
1129 deed occur in spontaneous conversation (Holler & Kendrick, 2015), but much  
1130 more work is needed to determine how participants make predictions about  
1131 turn structure in first-person contexts and whether those mechanisms shift  
1132 at points of imminent speaker change.

1133 *Early competence for turn taking?*

1134 One of the core aims of our study was to test whether children show an  
1135 early competence for turn taking, as is proposed by studies of spontaneous

1136 mother-infant proto-conversation and theories about the mechanisms under-  
1137 lying human interaction in general (Hilbrink et al., 2015; Levinson, 2006).  
1138 We found evidence that young children make spontaneous predictions about  
1139 upcoming turn structure: definitely at age two and marginally at age one.

1140 These results contrast with Keitel and colleagues' (2013) finding that chil-  
1141 dren cannot anticipate upcoming turn structure at above-chance rates until  
1142 age three. The current study used an appreciably more conservative ran-  
1143 dom baseline than the one used in Keitel and colleagues' study. Therefore,  
1144 this difference in age of emergence more likely stems from our use of a more  
1145 engaging speech style, stereo speech playback, and more typical turn tran-  
1146 sition durations. The child-friendly style of speech in particular may have  
1147 helped in two ways: keeping children more engaged with the stimuli in gen-  
1148 eral and using less syntactically complex and more prosodically exaggerated  
1149 cues (Fernald et al., 1989; Werker & McLeod, 1989; Snow, 1977) compared  
1150 to what they would get with adult-adult speech.

1151 To be clear, young children's "above chance" performance was often still  
1152 far from adult-like predictive behavior—children at ages one and two were  
1153 still very close to chance in their anticipations and, even at age six, chil-  
1154 dren were not fully adult-like in their predictions. This may indicate that  
1155 young children rely primarily on non-verbal cues (like inter-turn silence) in  
1156 anticipating turn transitions while adults use both verbal and non-verbal  
1157 cues to make predictions. Relatedly, adults may be more expert in flexibly  
1158 adapting to the turn-relevant cues present at any moment, e.g., responding  
1159 to non-English prosodic cues in Experiment 1.

1160 Taken together, our data suggest that turn-taking skills do begin to  
1161 emerge in infancy, but that children cannot consistently make effective pre-  
1162 dictions until they can identify and react to question turns. This finding  
1163 leads us to wonder how participant role (first- instead of third-person) and  
1164 differences in early interactional experience (e.g., frequent vs. infrequent  
1165 question-asking from caregivers) feed into this early predictive skill. It also  
1166 bridges prior work showing a predisposition for turn taking in infancy (e.g.,  
1167 Bateson, 1975; Hilbrink et al., 2015; Jaffe et al., 2001; Snow, 1977) with  
1168 children's apparently *late* acquisition of adult-like competence for turn tak-  
1169 ing in actual conversation (Casillas et al., In press; Garvey, 1984; Garvey  
1170 & Berninger, 1981; Ervin-Tripp, 1979) and reinforces the idea that it takes  
1171 children several years to fully integrate linguistic information into their turn-  
1172 taking systems (Casillas et al., In press; Garvey & Berninger, 1981).

1173 *Limitations and future work*

1174 There are at least two major limitations to our work: speech naturalness  
1175 and participant role. Following prior work (De Ruiter et al., 2006; Keitel  
1176 et al., 2013), we used phonetically manipulated speech in Experiment 2.  
1177 This decision resulted in speech sounds that children don't usually hear in  
1178 their natural environment. Many prior studies have used phonetically-altered  
1179 speech with infants and young children (cf. Jusczyk, 2000), but few of them  
1180 have done so in a conversational context. Future work could instead carefully  
1181 script speech or cross-splice sub-parts of turns to control for the presence of  
1182 linguistic cues for turn transition (see, e.g., Torreira et al., 2015).

1183 The prediction measure used in our studies is based on an observer's view  
1184 of third-party conversation but, because participants' role in the interaction  
1185 could affect their online predictions about turn taking, an ideal measure  
1186 would instead capture first-person predictions. If conversational participants'  
1187 predictions are partly shaped by their need to respond, first-person measures  
1188 of spontaneous turn prediction will be key to revealing how participants  
1189 distribute their attention over verbal and non-verbal cues while taking part  
1190 in everyday interaction, the implications of which relate to theories of online  
1191 language processing for both language learning and everyday talk.

1192 That said, the third-person paradigm used in the present study still has  
1193 much to tell us about turn prediction. The task is natural and intuitive  
1194 in that no instruction is required, which means that it captures spontaneous  
1195 predictive behavior and can be used with participants of all ages. Frequencies  
1196 of anticipatory gaze switching appear to be stable across language communi-  
1197 ties where similar tasks have been tested (Keitel et al., 2013; Keitel & Daum,  
1198 2015; Holler & Kendrick, 2015; Hirvenkari et al., 2013)—even from a first-  
1199 person perspective—so the task is one that measures robust predictive be-  
1200 havior relevant to conversational processing in many languages. It also lends  
1201 itself to many possibilities for controlling the presence of individual verbal  
1202 and non-verbal cues and has a clear method for assessing random switching  
1203 baselines across the entire stimulus (because the participant's task is con-  
1204 stant). Also, in the case that response preparation interferes with our ability  
1205 to see prediction at the ends of incoming turns, third-person paradigms are  
1206 one of the only ways to ensure that we are measuring prediction processes in  
1207 isolation.

1208 The current findings also make important predictions about what we will  
1209 see in first-person prediction paradigms. For example, the focus on upcom-  
1210 ing points of speaker transition is only more important when participants

1211 are embedded in a live interaction; we would thus expect question-like ef-  
1212 fects to occur in first-person paradigms too, and perhaps even be amplified.  
1213 If so, participants' use of linguistic information would still subserve this goal,  
1214 with prediction at a premium. So we would predict that linguistic cues to  
1215 upcoming speaker change would be most critical in a first-person measure.  
1216 Regarding development, the same facts about the complexity of prediction  
1217 with prosodic information and children's initial limited lexical inventories  
1218 would still hold, as would the use of silence and non-verbal cues to assess  
1219 and predict turn structure in the absence of clear predictive linguistic infor-  
1220 mation. We therefore think that the paradigm presented here has important  
1221 contributions to make in finding out how we attend to and make predictions  
1222 about conversational interaction.

1223 *Conclusions*

1224 Conversation plays a central role in children's language learning. It is  
1225 the driving force behind what children say and what they hear. Adults use  
1226 linguistic information to accurately predict turn structure in conversation,  
1227 which facilitates their online comprehension and allows them to respond rel-  
1228 evantly and on time. The present study offers new findings regarding the  
1229 role of speech acts and linguistic processing in online turn prediction, and  
1230 has given evidence that turn prediction emerges by age two, increases with  
1231 age, and is driven by question turns. However, turn prediction is not fully  
1232 integrated with linguistic cues until much later and, in the absence of pre-  
1233 dictive linguistic cues, children and adults alike rely on retroactive cues such  
1234 as inter-turn silence to predict upcoming speaker change. Using language to  
1235 make predictions about upcoming interactive content takes time to develop  
1236 and, for participants of all ages appears to be primarily driven by partici-  
1237 pants' orientation to what will happen next—beyond the end of the current  
1238 turn.

1239 **Acknowledgements**

1240 We gratefully acknowledge the parents and children at Bing Nursery  
1241 School and the Children's Discovery Museum of San Jose. This work was  
1242 supported by an ERC Advanced Grant to Stephen C. Levinson (269484-  
1243 INTERACT), an NSF Graduate Research Fellowship and NSF Dissertation  
1244 Improvement Grant to MC, and a Merck Foundation fellowship to MCF.

1245 Earlier versions of these data and analyses were presented to conference au-  
1246 diences (Casillas & Frank, 2012, 2013). We also thank Tania Henetz, Fran-  
1247 cisco Torreira, Stephen C. Levinson, and Eve V. Clark for their feedback on  
1248 earlier versions of this work. The analysis code for this project can be found  
1249 on GitHub at [https://github.com/langcog/turn\\_taking/](https://github.com/langcog/turn_taking/).

1250 **References**

- 1251 Allison, P. D. (2004). Convergence problems in logistic regression. In M. Alt-  
1252 man, J. Gill, & M. McDonald (Eds.), *Numerical Issues in Statistical Com-*  
1253 *puting for the Social Scientist* (pp. 247–262). Wiley-Interscience: New  
1254 York, NY.
- 1255 Allison, P. D. (2012). *Logistic Regression Using SAS: Theory and Applica-*  
1256 *tion*. SAS Institute.
- 1257 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects  
1258 structure for confirmatory hypothesis testing: Keep it maximal. *Journal*  
1259 *of Memory and Language*, 68, 255–278.
- 1260 Bates, D., Maechler, M., Bolker, B., & Walker, S. (2014).  
1261 *lme4: Linear mixed-effects models using Eigen and S4*. URL:  
1262 <https://github.com/lme4/lme4><http://lme4.r-forge.r-project.org/>  
1263 [Computer program] R package version 1.1-7.
- 1264 Bateson, M. C. (1975). Mother-infant exchanges: The epigenesis of conver-  
1265 sational interaction. *Annals of the New York Academy of Sciences*, 263,  
1266 101–113.
- 1267 Bergelson, E., & Sibley, D. (2013). The acquisition of abstract words by  
1268 young infants. *Cognition*, 127, 391–397.
- 1269 Bloom, K. (1988). Quality of adult vocalizations affects the quality of infant  
1270 vocalizations. *Journal of Child Language*, 15, 469–480.
- 1271 Boersma, P., & Weenink, D. (2012). *Praat: doing phonetics by computer*.  
1272 URL: <http://www.praat.org> [Computer program] Version 5.3.16.
- 1273 Bögels, S., Magyari, L., & Levinson, S. C. (2015). Neural signatures of  
1274 response planning occur midway through an incoming question in conver-  
1275 sation. *Scientific Reports*, 5.

- 1276 Bruner, J. (1985). Child's talk: Learning to use language. *Child Language*  
1277 *Teaching and Therapy*, 1, 111–114.
- 1278 Bruner, J. S. (1975). The ontogenesis of speech acts. *Journal of Child Lan-*  
1279 *guage*, 2, 1–19.
- 1280 Carlson, R., Hirschberg, J., & Swerts, M. (2005). Cues to upcoming swedish  
1281 prosodic boundaries: Subjective judgment studies and acoustic correlates.  
1282 *Speech Communication*, 46, 326–333.
- 1283 Casillas, M., Bobb, S. C., & Clark, E. V. (In press). Turn taking, timing,  
1284 and planning in early language acquisition. *Journal of Child Language*, .
- 1285 Casillas, M., & Frank, M. C. (2012). Cues to turn boundary prediction in  
1286 adults and preschoolers. In *Proceedings of SemDial* (pp. 61–69).
- 1287 Casillas, M., & Frank, M. C. (2013). The development of predictive processes  
1288 in children's discourse understanding. In *Proceedings of the 35th Annual*  
1289 *Meeting of the Cognitive Science Society* (pp. 299–304).
- 1290 Cutler, A., Dahan, D., & Van Donselaar, W. (1997). Prosody in the com-  
1291 prehension of spoken language: A literature review. *Language and speech*,  
1292 40, 141–201.
- 1293 De Ruiter, J. P., Mitterer, H., & Enfield, N. J. (2006). Projecting the end of  
1294 a speaker's turn: A cognitive cornerstone of conversation. *Language*, 82,  
1295 515–535.
- 1296 De Vos, C., Torreira, F., & Levinson, S. C. (2015). Turn-timing in signed  
1297 conversations: coordinating stroke-to-stroke turn boundaries. *Frontiers in*  
1298 *Psychology*, 6.
- 1299 Dingemanse, M., Torreira, F., & Enfield, N. (2013). Is "Huh?" a univer-  
1300 sal word? Conversational infrastructure and the convergent evolution of  
1301 linguistic items. *PloS one*, 8, e78273.
- 1302 Duncan, S. (1972). Some signals and rules for taking speaking turns in  
1303 conversations. *Journal of Personality and Social Psychology*, 23, 283.
- 1304 Ervin-Tripp, S. (1979). Children's verbal turn-taking. In E. Ochs, & B. B.  
1305 Schieffelin (Eds.), *Developmental Pragmatics* (pp. 391–414). Academic  
1306 Press, New York.

- 1307 Fernald, A., Taeschner, T., Dunn, J., Papousek, M., de Boysson-Bardies,  
1308 B., & Fukui, I. (1989). A cross-language study of prosodic modifications  
1309 in mothers' and fathers' speech to preverbal infants. *Journal of Child  
1310 Language*, 16, 477–501.
- 1311 Fitneva, S. (2012). Beyond answers: questions and children's learning. In  
1312 J.-P. De Ruiter (Ed.), *Questions: Formal, Functional, and Interactional  
1313 Perspectives* (pp. 165–178). Cambridge University Press, Cambridge, UK.
- 1314 Ford, C. E., & Thompson, S. A. (1996). Interactional units in conversation:  
1315 Syntactic, intonational, and pragmatic resources for the management of  
1316 turns. *Studies in Interactional Sociolinguistics*, 13, 134–184.
- 1317 Garvey, C. (1984). *Children's Talk* volume 21. Harvard University Press.
- 1318 Garvey, C., & Berninger, G. (1981). Timing and turn taking in children's  
1319 conversations 1. *Discourse Processes*, 4, 27–57.
- 1320 Gísladóttir, R., Chwilla, D., & Levinson, S. C. (2015). Conversation electri-  
1321 fied: ERP correlates of speech act recognition in underspecified utterances.  
1322 *PloS one*, 10, e0120068.
- 1323 Griffin, Z. M., & Bock, K. (2000). What the eyes say about speaking. *Psy-  
1324 chological science*, 11, 274–279.
- 1325 Hart, B., & Risley, T. R. (1992). American parenting of language-learning  
1326 children: Persisting differences in family-child interactions observed in nat-  
1327 ural home environments. *Developmental Psychology*, 28, 1096.
- 1328 Hedberg, N., Sosa, J. M., Görgülü, E., & Mameni, M. (2010). The prosody  
1329 and meaning of Wh-questions in American English. In *Speech Prosody  
1330 2010* (pp. 100045:1–4).
- 1331 Henning, A., Striano, T., & Lieven, E. V. (2005). Maternal speech to infants  
1332 at 1 and 3 months of age. *Infant Behavior and Development*, 28, 519–536.
- 1333 Hilbrink, E., Gattis, M., & Levinson, S. C. (2015). Early developmental  
1334 changes in the timing of turn-taking: A longitudinal study of mother-  
1335 infant interaction. *Frontiers in Psychology*, 6.

- 1336 Hirvenkari, L., Ruusuvuori, J., Saarinen, V.-M., Kivioja, M., Peräkylä, A.,  
1337 & Hari, R. (2013). Influence of turn-taking in a two-person conversation  
1338 on the gaze of a viewer. *PloS one*, 8, e71569.
- 1339 Holler, J., & Kendrick, K. H. (2015). Unaddressed participants' gaze in  
1340 multi-person interaction. *Frontiers in Psychology*, 6.
- 1341 Jaffé, J., Beebe, B., Feldstein, S., Crown, C. L., Jasnow, M. D., Rochat,  
1342 P., & Stern, D. N. (2001). *Rhythms of dialogue in infancy: Coordinated  
1343 timing in development*. Monographs of the Society for Research in Child  
1344 Development. JSTOR.
- 1345 Johnson, E. K., & Jusczyk, P. W. (2001). Word segmentation by 8-month-  
1346 olds: When speech cues count more than statistics. *Journal of Memory  
1347 and Language*, 44, 548–567.
- 1348 Jusczyk, P. W. (2000). *The Discovery of Spoken Language*. MIT press.
- 1349 Jusczyk, P. W., Hohne, E., Mandel, D., & Strange, W. (1995). Picking up  
1350 regularities in the sound structure of the native language. *Speech perception  
1351 and linguistic experience: Theoretical and methodological issues in cross-  
1352 language speech research*, (pp. 91–119).
- 1353 Kamide, Y., Altmann, G., & Haywood, S. L. (2003). The time-course of  
1354 prediction in incremental sentence processing: Evidence from anticipatory  
1355 eye movements. *Journal of Memory and Language*, 49, 133–156.
- 1356 Keitel, A., & Daum, M. M. (2015). The use of intonation for turn anticipation  
1357 in observed conversations without visual signals as source of information.  
1358 *Frontiers in Psychology*, 6.
- 1359 Keitel, A., Prinz, W., Friederici, A. D., Hofsten, C. v., & Daum, M. M.  
1360 (2013). Perception of conversations: The importance of semantics and  
1361 intonation in childrens development. *Journal of Experimental Child Psy-  
1362 chology*, 116, 264–277.
- 1363 Lammertink, I., Casillas, M., Benders, T., Post, B., & Fikkert, P. (2015).  
1364 Dutch and english toddlers' use of linguistic cues in predicting upcoming  
1365 turn transitions. *Frontiers in Psychology*, 6.

- 1366 Lemasson, A., Glas, L., Barbu, S., Lacroix, A., Guilloux, M., Remeuf, K., &  
1367 Koda, H. (2011). Youngsters do not pay attention to conversational rules:  
1368 is this so for nonhuman primates? *Nature Scientific Reports*, 1.
- 1369 Levelt, W. J. (1989). *Speaking: From intention to articulation*. MIT press.
- 1370 Levinson, S. C. (2006). On the human “interaction engine”. In N. Enfield,  
1371 & S. Levinson (Eds.), *Roots of Human Sociality: Culture, Cognition and*  
1372 *Interaction* (pp. 39–69). Oxford: Berg.
- 1373 Levinson, S. C. (2013). Action formation and ascriptions. In T. Stivers, &  
1374 J. Sidnell (Eds.), *The Handbook of Conversation Analysis* (pp. 103–130).  
1375 Wiley-Blackwell, Malden, MA.
- 1376 Levinson, S. C. (2016). Turn-taking in Human Communication – Origins  
1377 and Implications for Language Processing. *Trends in Cognitive Sciences*,  
1378 20, 6–14.
- 1379 Magyari, L., Bastiaansen, M. C. M., De Ruiter, J. P., & Levinson, S. C.  
1380 (2014). Early anticipation lies behind the speed of response in conversation.  
1381 *Journal of Cognitive Neuroscience*, 26, 2530–2539.
- 1382 Magyari, L., & De Ruiter, J. P. (2012). Prediction of turn-ends based on  
1383 anticipation of upcoming words. *Frontiers in Psychology*, 3:376, 1–9.
- 1384 Männel, C., & Friederici, A. D. (2009). Pauses and intonational phrasing:  
1385 ERP studies in 5-month-old German infants and adults. *Journal of Cog-*  
1386 *nitive Neuroscience*, 21, 1988–2006.
- 1387 Masataka, N. (1993). Effects of contingent and noncontingent maternal stim-  
1388 ulation on the vocal behaviour of three-to four-month-old Japanese infants.  
1389 *Journal of Child Language*, 20, 303–312.
- 1390 Mehler, J., Jusczyk, P., Lambertz, G., Halsted, N., Bertoni, J., & Amiel-  
1391 Tison, C. (1988). A precursor of language acquisition in young infants.  
1392 *Cognition*, 29, 143–178.
- 1393 Mintz, T. H. (2003). Frequent frames as a cue for grammatical categories in  
1394 child directed speech. *Cognition*, 90, 91–117.

- 1395 Morgan, J. L., & Saffran, J. R. (1995). Emerging integration of sequential  
1396 and suprasegmental information in preverbal speech segmentation. *Child  
1397 Development*, *66*, 911–936.
- 1398 Nazzi, T., & Ramus, F. (2003). Perception and acquisition of linguistic  
1399 rhythm by infants. *Speech Communication*, *41*, 233–243.
- 1400 Nomikou, I., & Rohlfing, K. J. (2011). Language does something: Body  
1401 action and language in maternal input to three-month-olds. *IEEE Trans-  
1402 actions on Autonomous Mental Development*, *3*, 113–128.
- 1403 R Core Team (2014). *R: A Language and Environment for Statistical Com-  
1404 puting*. R Foundation for Statistical Computing Vienna, Austria. URL:  
1405 <http://www.R-project.org> [Computer program] Version 3.1.1.
- 1406 Ratner, N., & Bruner, J. (1978). Games, social exchange and the acquisition  
1407 of language. *Journal of Child Language*, *5*, 391–401.
- 1408 Reddy, V., Markova, G., & Wallot, S. (2013). Anticipatory adjustments to  
1409 being picked up in infancy. *PloS one*, *8*, e65289.
- 1410 Ross, H. S., & Lollis, S. P. (1987). Communication within infant social games.  
1411 *Developmental Psychology*, *23*, 241.
- 1412 Rossano, F., Brown, P., & Levinson, S. C. (2009). Gaze, questioning and cul-  
1413 ture. In J. Sidnell (Ed.), *Conversation Analysis: Comparative Perspectives*  
1414 (pp. 187–249). Cambridge University Press, Cambridge.
- 1415 Sacks, H., Schegloff, E. A., & Jefferson, G. (1974). A simplest systematics for  
1416 the organization of turn-taking for conversation. *Language*, *50*, 696–735.
- 1417 Schegloff, E. A. (2007). *Sequence organization in interaction: Volume 1: A  
1418 primer in conversation analysis*. Cambridge University Press.
- 1419 Shatz, M. (1978). On the development of communicative understandings:  
1420 An early strategy for interpreting and responding to messages. *Cognitive  
1421 Psychology*, *10*, 271–301.
- 1422 Shatz, M. (1979). How to do things by asking: Form-function pairings in  
1423 mothers' questions and their relation to children's responses. *Child Devel-  
1424 opment*, *50*, 1093–1099.

- 1425 Shi, R., & Melancon, A. (2010). Syntactic categorization in French-learning  
1426 infants. *Infancy*, 15, 517–533.
- 1427 Shriberg, E., Stolcke, A., Jurafsky, D., Coccaro, N., Meteer, M., Bates, R.,  
1428 Taylor, P., Ries, K., Martin, R., & Van Ess-Dykema, C. (1998). Can  
1429 prosody aid the automatic classification of dialog acts in conversational  
1430 speech? *Language and Speech*, 41, 443–492.
- 1431 Snow, C. E. (1977). The development of conversation between mothers and  
1432 babies. *Journal of Child Language*, 4, 1–22.
- 1433 Soderstrom, M., Seidl, A., Kemler Nelson, D. G., & Jusczyk, P. W. (2003).  
1434 The prosodic bootstrapping of phrases: Evidence from prelinguistic in-  
1435 fants. *Journal of Memory and Language*, 49, 249–267.
- 1436 Speer, S. R., & Ito, K. (2009). Prosody in first language acquisition—  
1437 Acquiring intonation as a tool to organize information in conversation.  
1438 *Language and Linguistics Compass*, 3, 90–110.
- 1439 Stivers, T., Enfield, N. J., Brown, P., Englert, C., Hayashi, M., Heinemann,  
1440 T., Hoymann, G., Rossano, F., De Ruiter, J. P., Yoon, K.-E. et al. (2009).  
1441 Universals and cultural variation in turn-taking in conversation. *Proceed-  
1442 ings of the National Academy of Sciences*, 106, 10587–10592.
- 1443 Stivers, T., & Rossano, F. (2010). Mobilizing response. *Research on Language  
1444 and Social Interaction*, 43, 3–31.
- 1445 Takahashi, D. Y., Narayanan, D. Z., & Ghazanfar, A. A. (2013). Coupled  
1446 oscillator dynamics of vocal turn-taking in monkeys. *Current Biology*, 23,  
1447 2162–2168.
- 1448 Thorgrímsson, G., Fawcett, C., & Liszkowski, U. (2015). 1- and 2-year-olds'  
1449 expectations about third-party communicative actions. *Infant Behavior  
1450 and Development*, 39, 53–66.
- 1451 Tice (Casillas), M., & Henetz, T. (2011). Turn-boundary projection: Looking  
1452 ahead. In *Proceedings of the 33rd Annual Meeting of the Cognitive Science  
1453 Society* (pp. 838–843).
- 1454 Toda, S., & Fogel, A. (1993). Infant response to the still-face situation at 3  
1455 and 6 months. *Developmental Psychology*, 29, 532.

- 1456 Tomasello, M., & Brooks, P. J. (1999). Early syntactic development: A  
1457 construction grammar approach. In M. Barrett (Ed.), *The development of*  
1458 *language* (pp. 161–190). Psychology Press.
- 1459 Torreira, F., Bögels, S., & Levinson, S. C. (2015). Intonational phrasing is  
1460 necessary for turn-taking in spoken interaction. *Journal of Phonetics*, 52,  
1461 46–57.
- 1462 Weisleder, A. (2012). *Richer language experience leads to faster understand-  
1463 ing: Links between language input, processing efficiency, and vocabulary  
1464 growth*. Ph.D. thesis Stanford University.
- 1465 Werker, J. F., & McLeod, P. J. (1989). Infant preference for both male and  
1466 female infant-directed talk: A developmental study of attentional and af-  
1467 fective responsiveness. *Canadian Journal of Psychology/Revue Canadienne  
1468 de Psychologie*, 43, 230.
- 1469 Wittenburg, P., Brugman, H., Russel, A., Klassmann, A., & Sloetjes, H.  
1470 (2006). Elan: a professional framework for multimodality research. In  
1471 *Proceedings of LREC*.

1472 **Appendix A. Permutation Analyses**

1473 How can we be sure that our primary dependent measure (anticipatory  
1474 gaze switching) actually relates to turn transitions? Even if children were  
1475 gazing back and forth randomly during the experiment, we would have still  
1476 captured some false hits—switches that ended up in the turn-transition win-  
1477 dows by chance.

1478 We estimated the baseline probability of making an anticipatory switch  
1479 by randomly permuting the placement of the transition windows within each  
1480 stimulus (Figure 4). We then used the switch identification procedure from  
1481 Experiments 1 and 2 (Section b) to find out how often participants made  
1482 “anticipatory” switches within these randomly permuted windows. This pro-  
1483 cedure de-links participants’ gaze data from turn structure by randomly re-  
1484 assigning the onset time of each turn-transition in each permutation. We  
1485 created 5,000 of these permutations for each experiment to get an anticipa-  
1486 tory switch baselines over all possible starting points.

1487 Importantly, the randomized windows were not allowed to overlap with  
1488 each other, keeping true to the original stimuli. We also made sure that the  
1489 properties of each turn transition stayed constant across permutations. So,  
1490 while “transition window A” might start 2 seconds into Random Permu-  
1491 tation 1 and 17 seconds into Random Permutation 2, it maintained the same  
1492 prior speaker identity, transition type, gap duration, language condition, etc.,  
1493 across both permutations.

1494 We then re-ran the statistical models from the original data on each of the  
1495 random permutations, e.g., using Experiment 1’s original model to analyze  
1496 the anticipatory switches from each random permutation of the Experiment  
1497 1 looking data. We could then calculate the proportion of random data  
1498  $z$ -values exceeded by the original  $z$ -value for each predictor. We used the  
1499 absolute value of all  $z$ -values to conduct a two-tailed test. If the original  
1500 effect of a predictor exceeded 95% of the random model effects for that same  
1501 predictor, we deemed that predictor’s effect to be significantly different from  
1502 the random baseline (i.e.,  $p < .05$ ).

1503 For example, children’s “language condition” effect from Experiment 1  
1504 had a  $z$ -value of  $|3.429|$ , which is greater than 99.9% of all  $|z\text{-value}|$  esti-  
1505 mates from Experiment 1’s random permutation models (i.e.,  $p = .001$ ). It is  
1506 therefore highly unlikely that the effect of language condition in the original  
1507 model derived from random gaze shifting.

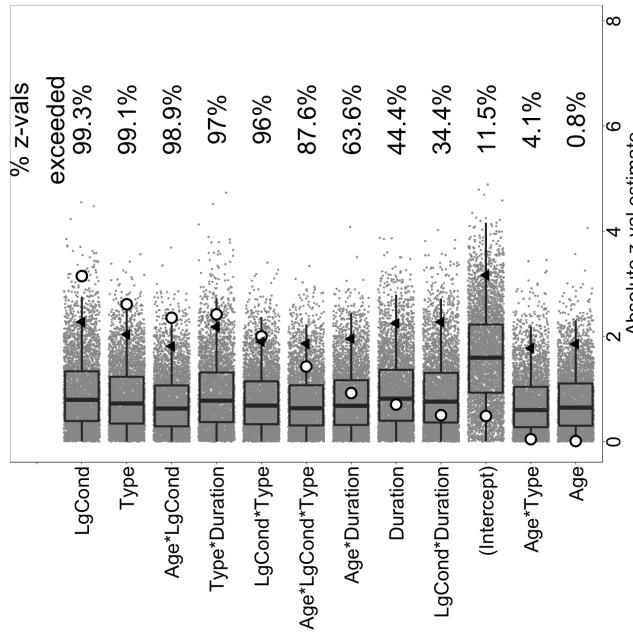
1508 We used this procedure to derive the random-baseline comparison values

1509 in the main text (above). However, we ran into two issues along the way:  
1510 first, we had to report  $z$ -values rather than beta estimates. Second, we had  
1511 to exclude a substantial portion of the models, especially in Experiment 2  
1512 because of model non-convergence. We address each of these issues below.

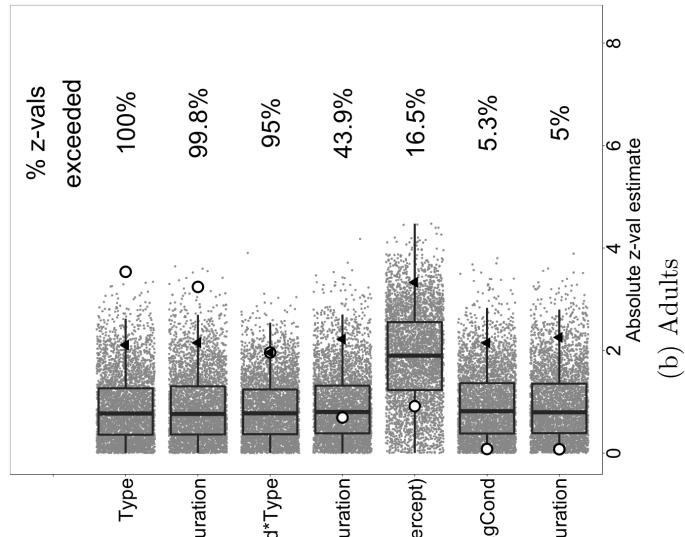
1513 *Appendix A.1. Beta, standard error, and  $z$  estimates*

1514 We reported  $z$ -values in the main text rather than beta estimates because  
1515 the standard errors in the randomly permuted data models were much higher  
1516 than for the original data. The distributions for each predictor's beta esti-  
1517 mate, standard error, and  $z$ -value for adults and children in each experiment  
1518 are shown in the graphs below (Figures A.1a–A.6b). In each plot, the gray  
1519 dots represent the absolute value of the 5,000 randomly permuted model es-  
1520 timates for the estimate type plotted (beta, standard error, or  $z$ ), the white  
1521 circles represent the model estimates from the original data, and the black  
1522 triangles represent the 95th percentile for each random distribution.

### Experiment 1: $z$ -value estimates



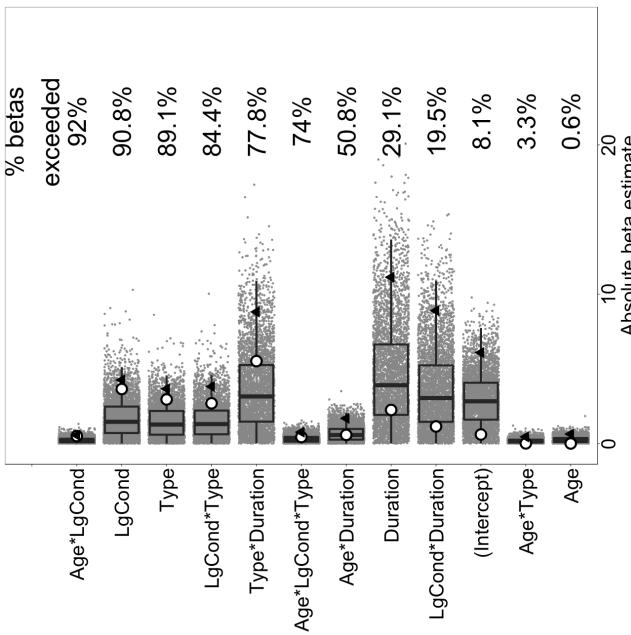
(a) Children



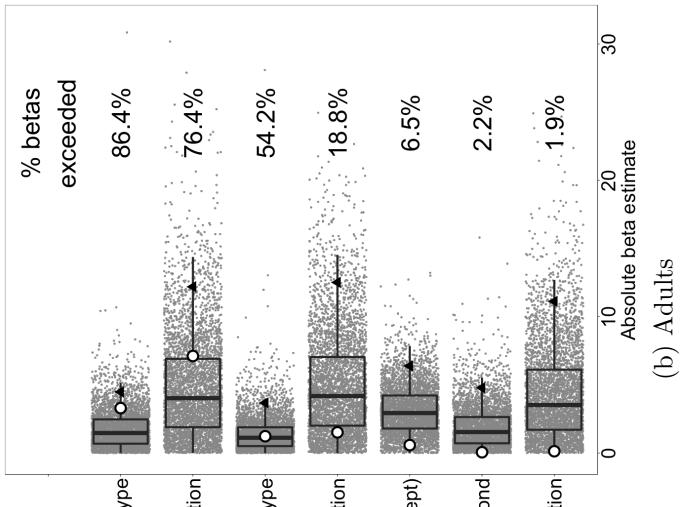
(b) Adults

Figure A.1: Random-permutation and original  $|z\text{-values}|$  for predictors of anticipatory gaze rates in Experiment 1.

### Experiment 1: $\beta$ estimates



55



(a) Children

(b) Adults

Figure A.2: Random-permutation and original  $|\beta\text{-values}|$  for predictors of gaze rates in Experiment 1.

### Experiment 1: *SE* estimates

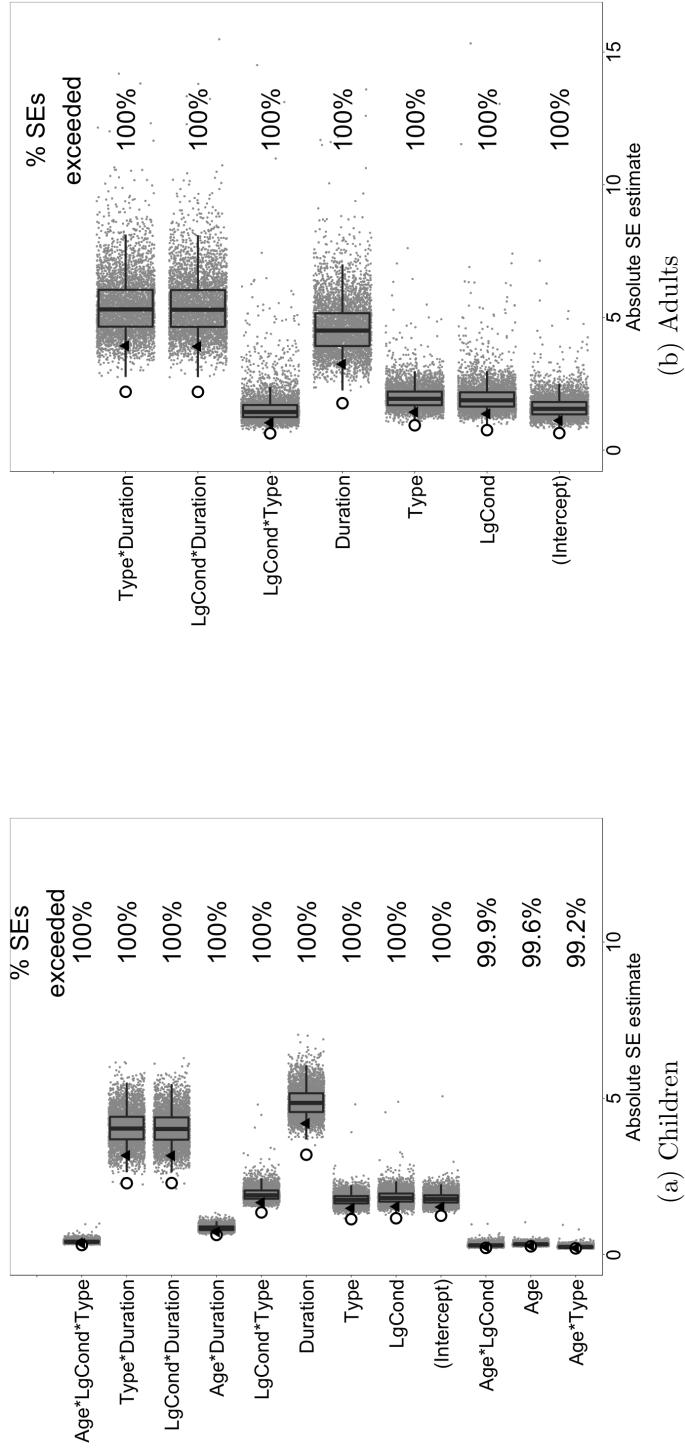


Figure A.3: Random-permutation and original *SE*-values for predictors of anticipatory gaze rates in Experiment 1.

## Experiment 2: $z$ estimates

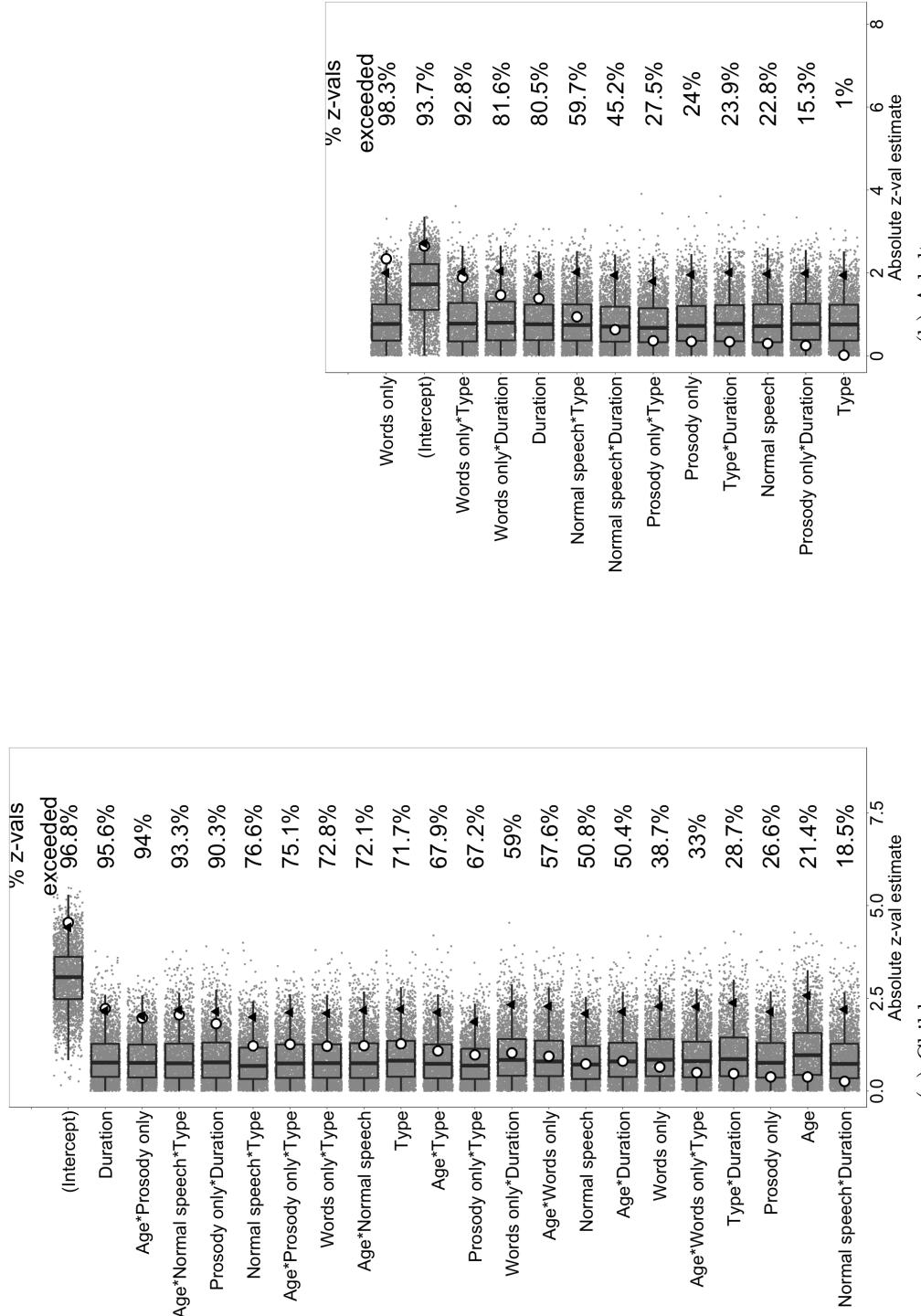


Figure A.4: Random-permutation and original  $|z$ -values for predictors of anticipatory gaze rates in Experiment 2.

## Experiment 2: $\beta$ estimates

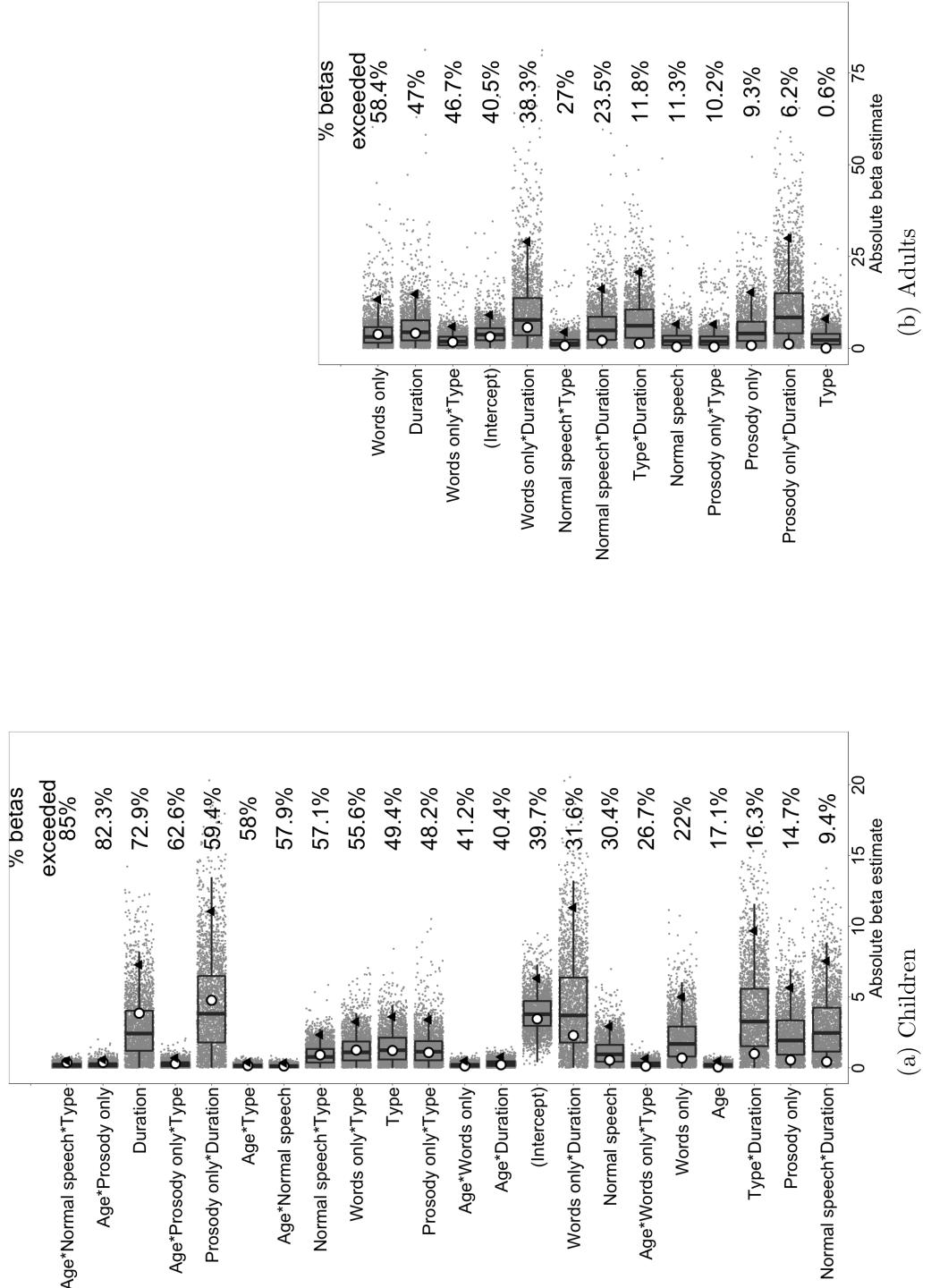


Figure A.5: Random-permutation and original  $|\beta\text{-values}|$  for predictors of anticipatory gaze rates in Experiment 2.

## Experiment 2: SE estimates

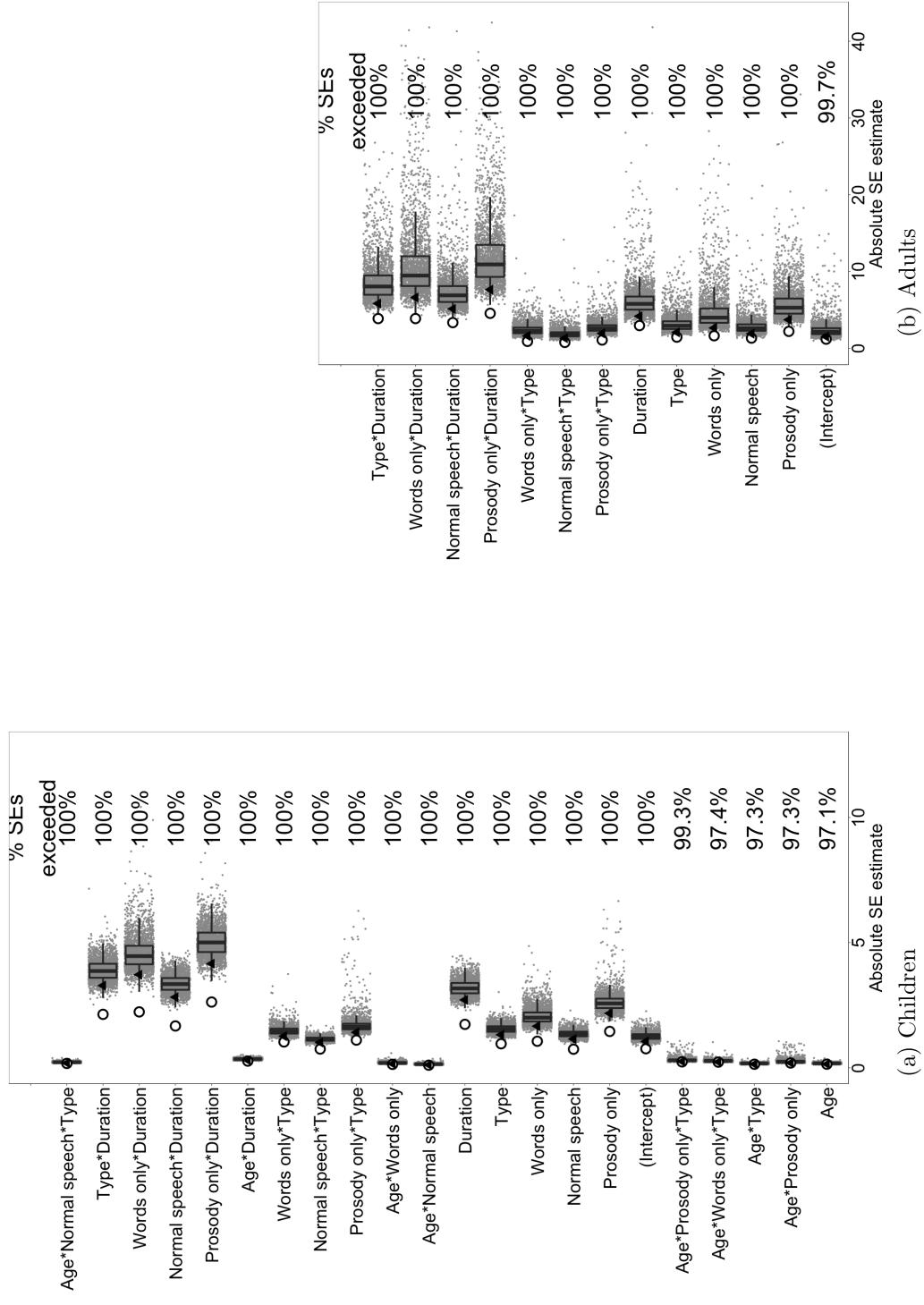


Figure A.6: Random-permutation and original SE-values for predictors of anticipatory gaze rates in Experiment 2.

1523 *Appendix A.2. Non-convergent models*

1524 In comparing the real and randomly permuted datasets, we excluded the  
1525 output of random-permutation models that gave convergence warnings to  
1526 remove erratic model estimates from our analyses. Non-convergent models  
1527 made up 22.4–24.4% of the random permutation models in Experiment 1  
1528 and 69–70% of the random permutation models in Experiment 2. The  $z$ -  
1529 values for each predictor in the converging and non-converging models from  
1530 Experiment 1 are shown in Table A.1.

1531 Although many of the non-converging models show estimates within range  
1532 of the converging models (e.g., with a mean difference of only 0.09 in median  
1533  $z$ -value across predictors), they also show many radically outlying estimates  
1534 (e.g., showing a mean difference of 237.3 in mean  $z$ -value across predictors).  
1535 Similar patterns were obtained in the non-converging models for Experiment  
1536 2 and persisted even when we tried other optimizers.

1537 We suspect that the issue derives from data sparsity in some of the ran-  
1538 dom permutations. This problem is known to occur when there are limited  
1539 numbers of binary observations in each of a design matrix’s bins (Allison,  
1540 2004). We could instead use zero-inflated poisson or negative binomial re-  
1541 gression models to allow for overdispersion in our data (Allison, 2012). How-  
1542 ever, these would give us baselines for the normal, convergent model, which  
1543 is not the aim of this analysis.

	Mean <sub>C</sub>	Mean <sub>NC</sub>	Median <sub>C</sub>	Median <sub>NC</sub>	SD <sub>C</sub>	SD <sub>NC</sub>	Min <sub>C</sub>	Min <sub>NC</sub>	Max <sub>C</sub>	Max <sub>NC</sub>
<i>Children</i>										
(Intercept)	-2.52	-458.42	-2.54	-2.86	0.87	1319.22	-5.53	-8185.36	0.41	0.97
Age	-0.51	-17.83	-0.49	-0.53	0.79	83.78	-3.71	-672.2	2.3	342.8
LgCond	-0.53	-109.91	-0.55	-0.63	0.93	564.42	-3.93	-4418.74	3.23	2296.19
Type	-0.1	-29.66	-0.09	-0.1	0.98	515.12	-4.06	-4383.92	3.36	3416.68
Duration	0.99	345.53	0.98	1.15	1.07	1323.13	-2.44	-5048.24	5.78	9985.16
Age*LgCond	0.19	10.64	0.2	0.18	0.9	109.6	-3.31	-581.61	3.59	946.81
Age*Type	0.02	-1.8	0.001	-0.04	0.9	98.27	-3.36	-884.36	3.45	640.43
LgCond*Type	0.2	45.32	0.2	0.27	0.96	691.3	-3.12	-4160.06	3.39	5107.64
Age*LgCond*Type	-0.12	-14.23	-0.12	-0.15	0.93	156.72	-2.98	-1318.26	2.90	927.69
<i>Adults</i>										
(Intercept)	-1.63	-126.14	-1.71	-1.73	0.97	713.39	-4.08	-12111.22	2.15	649.55
LgCond	-0.26	-679.6	-0.3	-0.53	1.02	15894.33	-3.45	-494979.7	3.35	88581.58
Type	-0.11	6.29	-0.13	-0.04	1.11	501.5	-3.85	-6420.75	3.28	8177.88
Duration	0.25	84.09	0.27	0.26	1.1	1152.94	-3.25	-10864.51	3.46	18540.62
LgCond*Type	0.12	-242.27	0.1	0.34	1.07	26836.7	-3.41	-622642.7	3.81	509198.4
LgCond*Duration	0.15	780.03	0.16	0.39	1.04	44105.02	-3.84	-798498.6	3.55	1145951
Type*Duration	0.05	-6.56	0.05	0.02	1.13	1389.9	-3.54	-15979.22	3.87	16419.46
LgCond*Type*Duration	-0.06	1083.63	-0.08	-0.21	1.1	63116.54	-4.21	-1201895	4.02	1284965

Table A.1: Estimated  $z$ -values for each predictor in converging ( $C$ ) and non-converging ( $NC$ ) child and adult models from Experiment 1.

1544 **Appendix B. Pairwise developmental tests**

1545 Experiments 1 and 2 both showed effects of age in interaction with lin-  
1546 guistic condition and transition type (e.g., English vs. non-English). To  
1547 explore these effects in more depth, we recorded the average difference score  
1548 for the predictor that interacted with age for each participant (e.g., English  
1549 minus non-English anticipatory switches), using these values to compute an  
1550 average difference score over participants in each age group (e.g., age 3, 4,  
1551 and 5) within each random permutation. That averaging process produces  
1552 5,000 baseline-derived difference scores for each age group.

1553 We then made pairwise age comparisons of these difference scores (e.g.,  
1554 the linguistic condition effect in 3-year-olds vs. 4-year-olds), computing the  
1555 percent of random-permutation difference scores exceeded by the real-data  
1556 difference score. If the real-data difference score exceeded 95% of the random-  
1557 data age difference scores, we deemed it to be an age effect significantly  
1558 different from chance—e.g., a significant difference between ages three and  
1559 four in the effect of linguistic condition. This procedure is essentially a two-  
1560 tailed  $t$ -test, adapted for use with the randomly permuted baseline data.

1561 In each of the plots below, the black dot represents the real data value  
1562 for the effect being shown. The effect sizes from the 5,000 randomly per-  
1563 mitted data sets are shown in the distribution. The percentage displayed  
1564 is the percentage of random permutation values exceeded by the original  
1565 data value (taking the absolute value of all data points for a two-tailed test).  
1566 Comparisons marked with 95% or higher are significant at the  $p < 0.05$  level.

Experiment 1: Age and linguistic condition

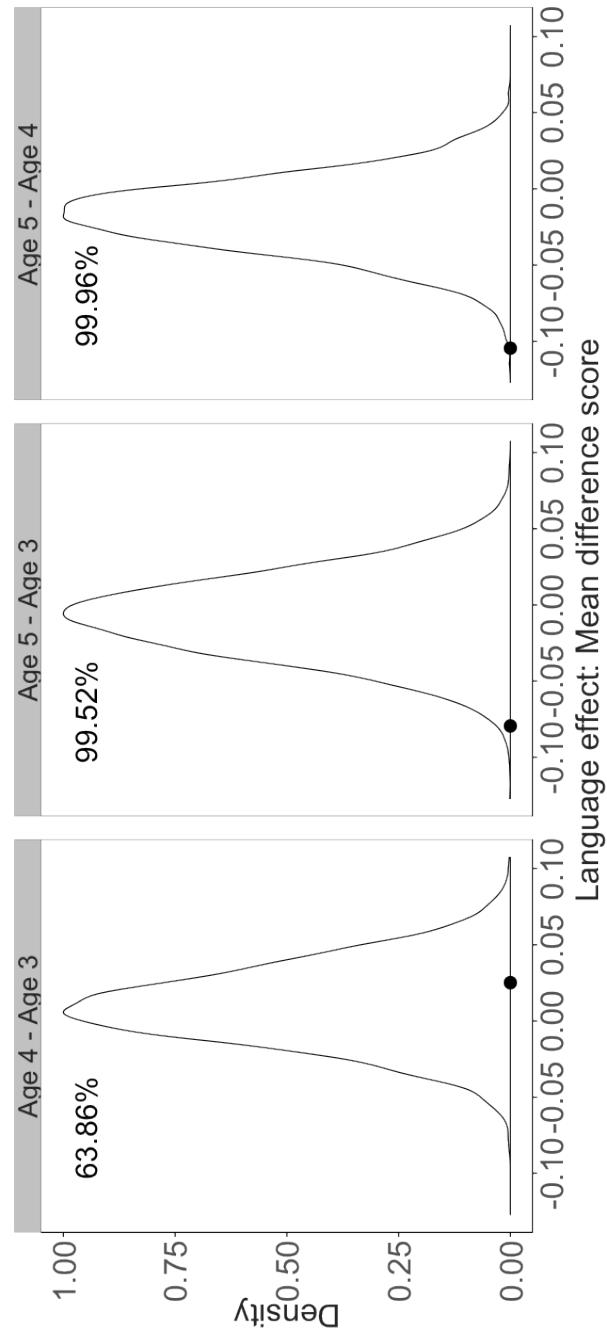


Figure B.1: Pairwise comparisons of the language condition effect across ages in Experiment 1.

### Experiment 2: Age and the *prosody only* condition

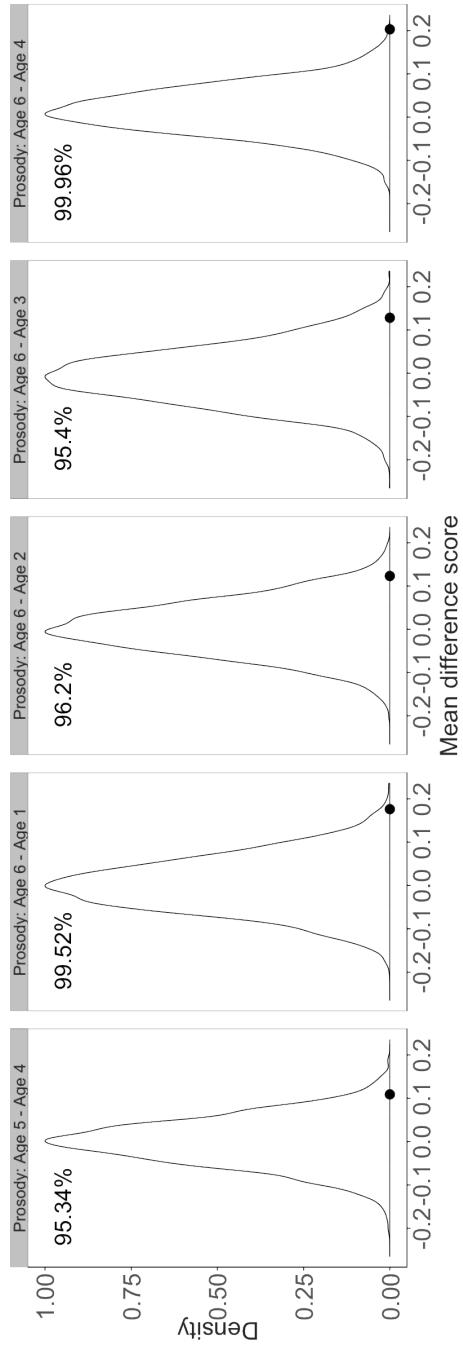


Figure B.2: Significant pairwise comparisons of the *prosody only-no speech* linguistic condition effect, across ages in Experiment 2

### Experiment 2: Age, transition type, and *normal* speech

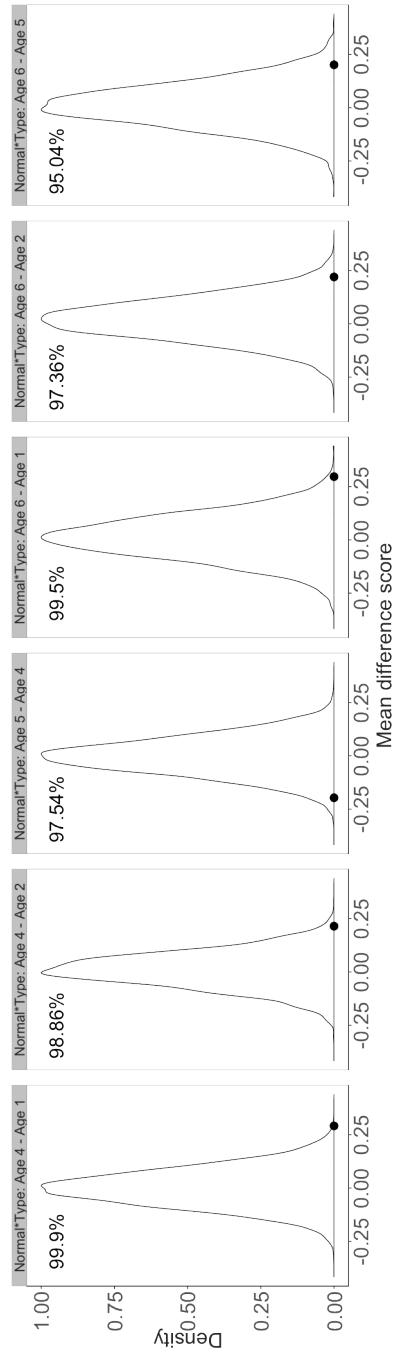


Figure B.3: Significant pairwise comparisons of the *normal speech-no speech* language condition effect for transition type, across ages, in Experiment 2.

1567 **Appendix C. Boredom-driven anticipatory looking**

1568 One alternative hypothesis for children’s anticipatory gazes is that they  
1569 are looking at the current speaker at the start of each turn, but then grow  
1570 bored and start looking away at a constant rate. Even though this alternative  
1571 hypothesis does not predict the primary effects in our data (e.g., the difference  
1572 between questions and non-questions), we cannot rule out the possibility that  
1573 a portion of participants’ saccades come from boredom.

1574 The data plotted here show a hypothetical group of boredom-driven par-  
1575 ticipants (gray dots) compared to participants from the actual data in Ex-  
1576 periment 2 (black dots). The hypothetical boredom-driven participants look  
1577 away from the current speaker at a linear rate, beginning one second after  
1578 the start of a turn.

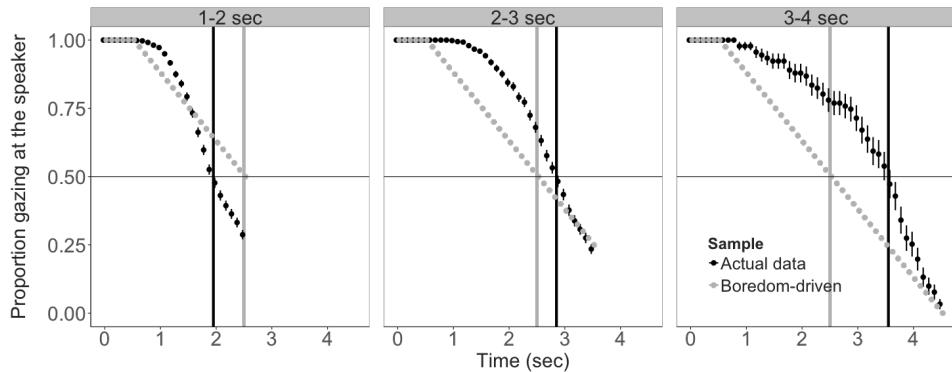


Figure C.1: Proportion of participants (hypothetical boredom-driven=gray; actual Experiment 2=black) looking at the current speaker, split by turn duration. Vertical bars indicate standard error in the experimental data.

1579 If children’s switches away from the current speaker were purely driven  
1580 by boredom, they would switch away equally quickly on long and short turns.  
1581 Therefore, their crossover point—the point in time at which 50% of the chil-  
1582 dren have switched away from the current speaker—would be the same for  
1583 all turns, no matter the length of the turn. This pattern is demonstrated  
1584 in the hypothetical boredom-driven crossover points, which always occur 2.5  
1585 seconds after the start of speech (gray vertical line; Figure C.1).

1586 In children’s *actual* looking data we see that crossover points increase with  
1587 turn duration: 2.0, 2.9, and 3.6 seconds after the start of speech for turns

1588 with durations of 1–2, 2–3, and 3–4 seconds, respectively (black vertical line;  
1589 Figure C.1). This pattern suggests that, though children do look away as  
1590 the turn is unfolding, their looks away are not simply driven by boredom.

1591 Are the looks away in Figure C.1 still too early to count as “turn-transition”  
1592 anticipation? It is true that children start looking away after one second has  
1593 passed (but then only gradually). It is possible that some of these early looks  
1594 away are boredom-driven, but it is equally plausible that some of them are  
1595 turn-driven. Early predictive behavior is common in turn-taking studies with  
1596 adults, in both constrained turn-taking tasks (e.g., De Ruiter et al., 2006;  
1597 Gísladóttir et al., 2015; Bögels et al., 2015) and in spontaneous conversation  
1598 (e.g., Holler & Kendrick, 2015; Torreira et al., 2015). Although this same  
1599 pattern has yet to be established for children’s turn predictions, the looking  
1600 behavior here is at least consistent with adult response patterns in previous  
1601 work. Additionally, because our analysis windows in the main study only  
1602 overlapped with the pre-gap utterance by 300 msec (Figure 2), our primary  
1603 results are unlikely to capture any very early or early boredom-driven gaze  
1604 switches.

1605 We conclude that the boredom-driven effects in our data are unlikely  
1606 to change our primary results, though we acknowledge that characterizing  
1607 different gaze switching strategies in this kind of data is an important avenue  
1608 for future work.

1609 **Appendix D. Puppet pair and linguistic condition**

1610     The design for Experiment 2 does not fully cross puppet pair (e.g., robots,  
1611     blue puppets) with linguistic condition (e.g., *words only* and *no speech*). Even  
1612     though each puppet pair is associated with different conversation clips across  
1613     children (e.g., robots talking about kitties, birthday parties, and pancakes),  
1614     the robot puppets themselves were exclusively associated with the *words only*  
1615     condition. Similarly, merpeople were exclusively associated with *prosody only*  
1616     speech, and the puppets wearing dressy clothes were exclusively associated  
1617     with the *no speech* condition. We designed the experiment this way to in-  
1618     crease its pragmatic felicity for older children (i.e., robots make robot sounds,  
1619     merpeople’s voices are muffled under the water, the party-going puppets are  
1620     in a ‘party’ room with many other voices). There is therefore a confound  
1621     between linguistic condition and puppet pair; for example, children could  
1622     have made fewer anticipatory switches in the *prosody only* condition because  
1623     the puppets were less interesting. To test whether puppet pair drove the  
1624     condition-based differences found in Experiment 2, we ran a short follow-up  
1625     study.

1626 **Methods**

1627 We recruited 30 children between ages 3;0 and 5;11 from the Children’s Dis-  
1628 covery Museum of San Jose, California to participate in our experiment. All  
1629 participants were native English speakers. Children were randomly assigned  
1630 to one of six videos (five children per video).

1631 *Materials.* We created 6 short videos from the stimulus recordings made for  
1632 Experiment 2. Each video featured a puppet pair (red/blue/yellow/robot/  
1633 merpeople/party-goer; Figure 5). Puppets in all six videos performed the  
1634 exact same conversation recording (‘birthday party’; Experiment 2) with  
1635 normal, unmanipulated speech; this experiment therefore holds all things  
1636 constant across stimuli except for the appearance of the puppets.

1637 *Procedure.* We used the same experimental apparatus and procedure as in  
1638 Experiments 1 and 2. Each participant was randomly assigned to watch only  
1639 one of the six puppet videos. Five children watched each video. As in Experi-  
1640 ment 2, the experimenter immediately began each session with calibration  
1641 and then stimulus presentation because no special instructions were required.  
1642 The entire experiment took less than three minutes.

1643 *Data preparation.* We identified anticipatory gaze switches to the upcoming  
1644 speaker using the same method as in Experiments 1 and 2.

1645 **Results and discussion**

1646 We modeled children’s anticipatory switches (yes or no at each transition)  
1647 with mixed effects logistic regression, including puppet pair (robots/mer-  
1648 people/party-goers/other-3) as a fixed effect and participant and turn transi-  
1649 tion as random effects. We grouped the red, blue, and yellow puppets  
1650 together because they collectively represented the puppets used in the *nor-*  
1651 *mal* speech condition—this follow-up experiment is meant to test whether  
1652 the condition-based differences from Experiment 2 arose from the puppets  
1653 used in each condition.

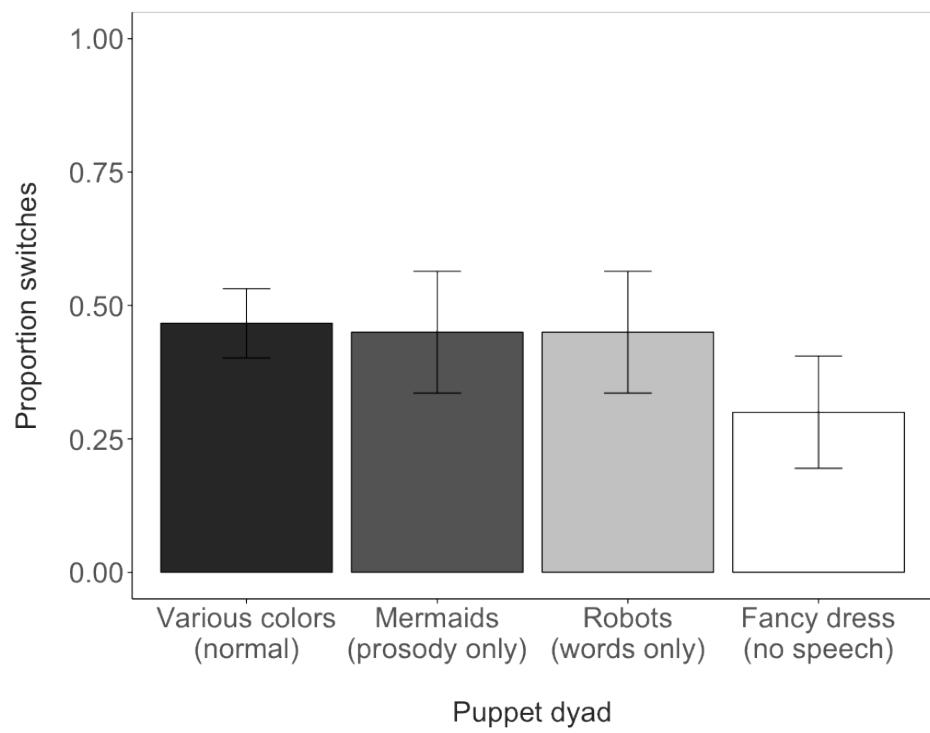


Figure D.1: Proportion gaze switches across puppet pairs when linguistic condition and conversation are held constant.

	Estimate	Std. Error	<i>z</i> value	Pr(>  <i>z</i>  )
<i>Reference level: normal-condition puppets</i>				
(Intercept)	-0.14790	0.32796	-0.451	0.652
Puppets= <i>mermaid</i>	-0.07581	0.65532	-0.116	0.908
Puppets= <i>robot</i>	-0.07104	0.65321	-0.109	0.913
Puppets= <i>party</i>	-0.78206	0.68699	-1.138	0.255
<i>Reference level: mer-puppets</i>				
(Intercept)	-0.22371	0.56832	-0.394	0.694
Puppets= <i>robot</i>	0.004763	0.80096	0.006	0.995
Puppets= <i>party</i>	-0.70626	0.82742	-0.854	0.393
<i>Reference level: robot puppets</i>				
(Intercept)	-0.21895	0.56565	-0.387	0.699
Puppets= <i>party</i>	-0.71102	0.82657	-0.860	0.390
<i>Reference level: party-goer puppets</i>				
(Intercept)	-0.9300	0.6067	-1.533	0.125

Table D.2: Model output for children’s anticipatory gaze switches with reference levels varied to show all possible pairwise differences between puppet pairs.

1654 In four versions of this model, we systematically varied the reference level  
 1655 of the puppet pair to check for any cross-condition differences. We found no  
 1656 significant effects of puppet pair on switching rate (all  $p > 0.25$ ; Table D.2).

1657 We take this finding as evidence that our decision to not fully cross puppet  
 1658 pairs and linguistic conditions in Experiment 2 was unlikely to have strongly  
 1659 affected children’s anticipatory gaze rates above and beyond the intended  
 1660 effects of linguistic condition.