

# 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. These early non-verbal interactions likely establish the protocol  
**20** by which children come to use language with others. How then do children  
**21** integrate linguistic knowledge with these preverbal turn-taking abilities?

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

**28** *Adult turn taking*

**29** Turn taking itself is not unique to conversation. Many other human ac-  
**30** tivities are organized around sequential turns at action. Traffic intersections  
**31** and computer network communication both use turn-taking systems. Chil-  
**32** dren’s early games (e.g., give-and-take, peek-a-boo) have built-in, predictable  
**33** turn structure (Ratner & Bruner, 1978; Ross & Lollis, 1987). Even monkeys  
**34** take turns: non-human primates such as marmosets and Campbell’s monkeys  
**35** vocalize contingently with each other in both natural and lab-controlled en-  
**36** vironments (Lemasson et al., 2011; Takahashi et al., 2013). In all these cases,

37 turn taking serves as a protocol for interaction, allowing the participants to  
38 coordinate with each other through sequences of contingent action.

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

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

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

70 Prior work suggests that adult listeners primarily use lexicosyntactic in-  
71 formation to accurately predict upcoming turn structure. De Ruiter and  
72 colleagues (2006) asked participants to listen to snippets of spontaneous con-  
73 versation and to press a button whenever they anticipated that the current  
74 speaker was about to finish his or her turn. The speech snippets were con-

trolled for the amount of linguistic information present; some were normal, but others had flattened pitch, low-pass filtered speech, or further manipulations. With pitch-flattened speech, the timing of participants' button responses was comparable to their timing with the full linguistic signal. But when no lexical information was available, participants' responded significantly earlier within the turn. The authors concluded that lexicosyntactic information<sup>1</sup> was necessary and possibly sufficient for turn-end projection, while intonation was neither necessary nor sufficient. Congruent evidence comes from studies varying the predictability of lexicosyntactic and pragmatic content: adults anticipate turn ends better when they can more accurately predict the exact words that will come next (Magyari & De Ruiter, 2012; see also Magyari et al., 2014). They can also identify speech acts within the first word of an utterance (Gísladóttir et al., 2015), allowing them to start planning their response at the first moment possible (Bögels et al., 2015).

Despite this body of evidence, the role of prosody for adult turn prediction is still a matter of debate. De Ruiter and colleagues' (2006) experiment focused on the role of intonation, which is only a partial index of prosody. Prosody is tied closely to the syntax of an utterance, so the two linguistic signals are difficult to control independently (Ford & Thompson, 1996). Torreira, Bögels & Levinson (2015) used a combination of button-press and verbal responses to investigate the relationship between lexicosyntactic and prosodic cues in turn-end prediction. Critically, their stimuli were cross-spliced so that each item had full prosodic cues to accompany the lexicosyntax. Because of the splicing, they were able to create items that had syntactically-complete units with no intonational phrase boundary at the end. Participants never verbally responded or pressed the "turn-end" button when hearing a syntactically-complete phrase without an intonational phrase boundary. And when intonational phrase boundaries were embedded within multi-utterance turns, participants were tricked into pressing the "turn-end" button 29% of the time. These findings suggest that listeners actually do rely on prosodic cues to execute a response in interaction with their predictions about the unfolding syntactic structure (see also de De Ruiter et al. (2006):525). These experimental findings corroborate other corpus and ex-

---

<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.

108 experimental work promoting a combination of cues (lexicosyntactic, prosodic,  
109 and pragmatic) as key for accurate turn-end prediction (Duncan, 1972; Ford  
110 & Thompson, 1996; Hirvenkari et al., 2013).

111 *Turn taking in development*

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

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

138 Before children manage to integrate linguistic cues into their turn-taking  
139 behaviors (for both turn prediction and production), they can rely on non-  
140 verbal interactional cues, including silence, eye gaze, body orientation, and  
141 gesture, to identify the boundaries of social actions. For example, with little  
142 to no linguistic knowledge, children are often able to infer desired responses  
143 to offers and requests by taking account of their interlocutor's non-verbal  
144 communicative behavior, the structure of routine events, and the affordances

145 of the current interactional context (Reddy et al., 2013; Nomikou & Rohlfing,  
146 2011; Shatz, 1978). With respect to turn taking in particular, children's spontane-  
147 aneous vocalizations during interaction demonstrate a sensitivity to short  
148 inter-speaker gaps from infancy (Hilbrink et al., 2015). Thus, before they can  
149 anticipate turn structure from linguistic cues, children might react to silence  
150 as a cue to upcoming speaker change. Interactional silence itself may then  
151 serve as one of children's first cues to turn structure, giving them information  
152 about when to respond before they can rely on language.

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

169 Keitel and colleagues (2013) were one of the first to explore how children  
170 use linguistic cues to predict upcoming turn structure. They asked 6-, 12-,  
171 24-, and 36-month-old infants, and adult participants to watch short videos  
172 of conversation and tracked their eye movements at points of speaker change.  
173 They showed their participants two types of videos—one normal and one with  
174 flattened pitch—to test the role of intonation in participants' anticipatory  
175 predictions about upcoming speech. Comparing children's anticipatory gaze  
176 frequency to a random baseline, they found that only 36-month-olds and  
177 adults made anticipatory gaze switches more often than expected by chance,  
178 and only 36-month-olds were affected by flattened intonation contours. This  
179 finding led Keitel and colleagues to conclude that children's ability to predict  
180 upcoming turn structure relies on their ability to comprehend the stimuli  
181 lexicosemantically. They also suggest that intonation might play a secondary  
182 role in turn prediction, but only after children acquire more sophisticated,

183 adult-like language comprehension abilities (also see Keitel & Daum, 2015).

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

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

### 212 *The current study*

213 Our goal in the current study is to find out when children begin to make  
214 predictions about upcoming turn structure and to understand how their pre-  
215 dictions are affected by linguistic cues across development. We present two  
216 experiments in which we measured children’s anticipatory gaze to respon-

---

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

ders while they watched conversation videos with natural (people speaking English vs. non-English; Experiment 1) and non-natural (puppets with phonetically manipulated speech; Experiment 2) control over the presence of lexical and prosodic cues. We tested children across a wide range of ages (Experiment 1: 3–5 years; Experiment 2: 1–6 years), with adult control participants in each experiment. We additionally tested for the use of one non-verbal cue: inter-turn silence.

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

## Experiment 1

We recorded participants’ eye movements as they watched six short videos of two-person (dyadic) conversation interspersed with attention-getting filler videos. Each conversation video featured an improvised discourse in one of five languages (English, German, Hebrew, Japanese, and Korean). Participants saw two videos in English and one in every other language. The participants, all native English speakers, were only expected to understand the two videos in English. We showed participants non-English videos to limit

253 their access to lexical information while maintaining their access to other  
254 cues to turn boundaries (e.g., non-English prosody, gaze, in-breaths, phrase  
255 final lengthening). Using this method, we analyzed children and adult's an-  
256 ticipatory looks from the current speaker to the upcoming speaker at points  
257 of turn transition in English and non-English videos.

258 *Methods*

259 *Participants*

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

272 *Materials*

273 *Video recordings.* We recorded pairs of talkers while they conversed in  
274 a sound-attenuated booth (Figure 1). Each talker was a native speaker of  
275 the language being recorded, and each talker pair was male-female. Using  
276 a Marantz PMD 660 solid state field recorder, we captured audio from two  
277 lapel microphones, one attached to each participant, while simultaneously  
278 recording video from the built-in camera of a MacBook laptop computer.  
279 The talkers were volunteers and were acquainted with their recording partner  
280 ahead of time.

281 Each recording session began with a 20-minute warm-up period of spon-  
282 taneous conversation during which the pair talked for five minutes on four

---

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



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

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

295 After recording, we combined the audio and video recordings by hand,  
296 and cropped each one to the (approximate) 30-second interval with the most  
297 turn activity. Because we recorded the conversations in stereo, the male and  
298 female voices came out of separate speakers during video playback. This gave  
299 each voice in the videos a localized source (from the left or right loudspeaker).  
300 We coded each turn transition in the videos for language condition (English  
301 vs. non-English), inter-turn gap duration (in milliseconds), and transition  
302 type (question vs. non-question). Each non-English turn was coded as a

---

<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.

303 question or non-question from a monolingual English-speaker’s perspective,  
304 i.e., turns that “sound like” questions and turns that do not. We asked five  
305 native American English speakers to listen to the audio recording for each  
306 non-English turn and judge whether it sounded like a question. We marked  
307 non-English turns as questions when at least 4 of the 5 listeners (80%) said  
308 that the turn “sounded like a question”. Thus, “question” cues in the non-  
309 English condition only resembled native English question cues, and therefore  
310 were likely harder to identify than cues to questionhood in the English con-  
311 dition. However, since participants did not speak the non-English languages  
312 and would only ever treat “question-sounding” turns as questions, we pro-  
313 ceeded with these analyses to see how pervasive question effects were—could  
314 they show up even without lexical access in a foreign language? If partic-  
315 ipants primarily rely on prosodic cues to question turns, it’s possible that  
316 even non-English prosody can elicit anticipatory gaze switches for question-  
317 like turns.

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

328 Questions made up exactly half of the turn transitions in the English  
329 (N=10) and non-English (N=8) videos. In the English videos, inter-turn  
330 gaps were slightly shorter for questions (mean=310, median=293, stdev=112  
331 msec) than non-questions (mean=325, median=315, stdev=118 msec). Non-  
332 English videos did not show a large difference in transition time for questions  
333 (mean=270, median=257, stdev=116 msec) and non-questions (mean=302,

---

<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.

334 median=252, stdev=134 msec).

335 *Procedure*

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

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

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

368 All participants (child and adult) completed a five-point calibration rou-  
369 tine before the first video started. We used a dancing Elmo for the children's  
370 calibration image. During the experiment, participants watched all six 30-

371 second conversation videos. The first and last conversations were in American  
372 English and the intervening conversations were Hebrew, Japanese, German,  
373 and Korean. The presentation order of the non-English videos was shuffled  
374 into four lists, which participants were assigned to randomly. The entire  
375 experiment, including instructions, took 10–15 minutes.

376 *Data preparation and coding*

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

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

403 As mentioned, most anticipatory switches happen in the inter-turn gap,  
404 but we also allowed anticipatory gaze switches that occurred in the final syl-  
405 lables of the prior turn. Early switches are consistent with the distribution  
406 of responses in explicit turn-boundary prediction tasks. For example, in a  
407 button press task, adult participants anticipated turn ends approximately

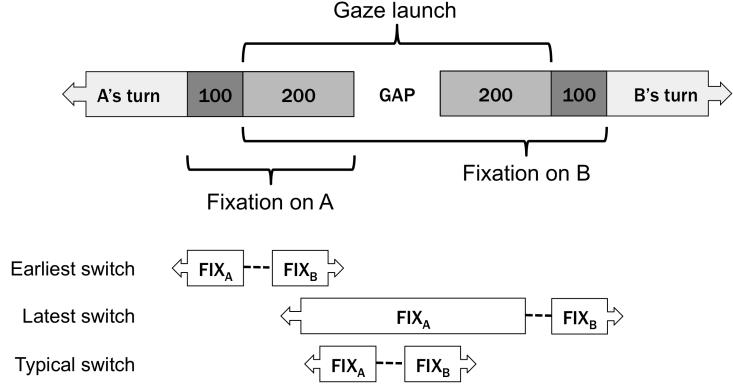


Figure 2: Schematic summary of the criteria for anticipatory gaze shifts from speaker A to speaker B during a turn transition.

408 200 msec in advance of the turn’s end, and anticipatory responses to pitch-  
 409 flattened stimuli came even earlier (De Ruiter et al., 2006). We therefore  
 410 allowed switches to occur as early as 200 msec before the end of the prior  
 411 turn. Again, because it takes 200 msec to plan an eye movement, we counted  
 412 anticipatory switches, at the latest, 200 msec after the onset of speech. There-  
 413 fore, for very early and very late switches, our requirement of 100 msec of  
 414 fixation on each speaker would sometimes extend outside of the gaze launch  
 415 window boundaries (200 msec before and after the inter-turn gap; Figure 2).  
 416 The maximally available fixation window was therefore 100 msec before and  
 417 after the earliest and latest possible switch point (300 msec before and after  
 418 the inter-turn gap). We did not count switches made during the fixation win-  
 419 dows as anticipatory. We *did* count switches made during the inter-turn gap.  
 420 The period of time from the beginning of the possible fixation window on  
 421 the prior speaker to the end of the possible fixation window on the responder  
 422 was our total analysis window (300 msec + the inter-turn gap + 300 msec).  
 423 *Predictions.* We expected participants to show greater anticipation in the  
 424 English videos than in the non-English videos because of their increased  
 425 access to linguistic information in English. We also predicted that anticipa-  
 426 tion would be greater following questions compared to non-questions; ques-  
 427 tions have early cues to upcoming turn transition (e.g., *wh*- words, subject-  
 428 auxiliary inversion), and also make a next response immediately relevant.  
 429 Our third prediction was that anticipatory looks would increase with devel-

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 across ages in Experiment 1.

opment, along with children’s increased linguistic competence. Finally, we predicted that transitions with longer inter-turn gaps would show greater anticipation because longer gaps provide (a) more time to make a gaze switch and (b) are themselves a cue to possible upcoming speaker switch.

#### Results

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

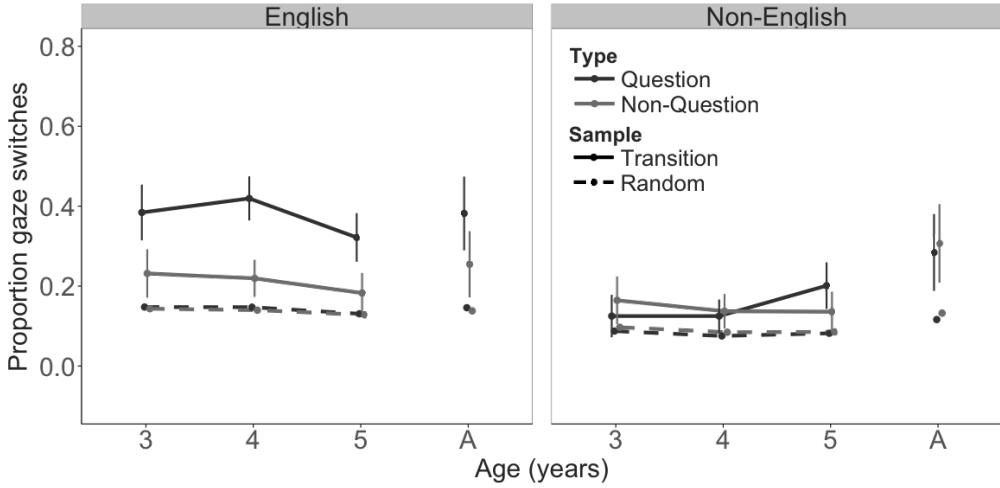


Figure 3: Anticipatory gaze rates across language condition and transition type for the real and randomly permuted datasets. Vertical bars represent 95% confidence intervals.

450 *Statistical models*

451 We identified anticipatory gaze switches for all 36 usable turn transitions,  
 452 based on the criteria outlined above, and analyzed them for effects of lan-  
 453 guage, transition type, and age with two mixed-effects logistic regressions  
 454 (Bates et al., 2014; R Core Team, 2014). We built one model each for chil-  
 455 dren and adults. We modeled children and adults separately because effects  
 456 of age are only pertinent to the children’s data.

457 The child model included condition (English vs. non-English)<sup>6</sup>, transition  
 458 type (question vs. non-question), age (3, 4, 5; numeric; intercept as age=0),  
 459 and duration of the inter-turn gap (seconds, e.g., 0.441) as predictors, with  
 460 two-way interactions between gap duration and the other simple fixed effects  
 461 (language condition, transition type, and age) and a three-way interaction

---

<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 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.



466 tion) and participant, with maximal random slopes of condition, transition  
467 type, and their interaction for participants (Barr et al., 2013).<sup>8</sup>

468 The adult model included fixed effects of condition, transition type, and  
469 their interaction, plus two-way interactions between gap duration and the  
470 other simple fixed effects (language condition and transition type, as in the  
471 child model). The adult model also included random effects of item and  
472 participant with maximal random slopes of condition, transition type, and  
473 their interaction for participant.

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

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

490 *Random baseline comparison*

491 Our primary analysis (above) makes the assumption that participants'  
492 eye movements generally follow the turn structure of the stimulus, i.e., that  
493 participants track the current speaker and switch their gaze to the upcoming  
494 speaker near turn transitions. As just described, based on this assumption,  
495 we used linear mixed effects regressions to see how anticipatory looking is  
496 affected by aspects of participant group (age) and stimulus (e.g., transition  
497 type, language condition). But what if the assumption that participants gen-

---

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.

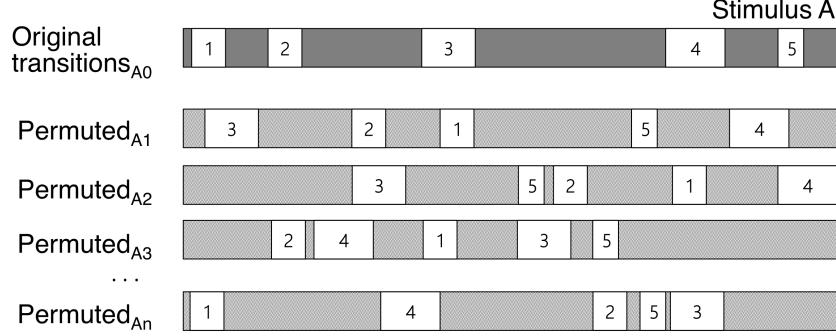


Figure 4: Example of analysis window permutations for a stimulus with five turn transitions. The windows included  $\pm 300$  msec around the inter-turn gap.

498 really track turn structure were wrong? Could these results have emerged  
 499 if participants' eye movements were *not* linked to turn structure? For ex-  
 500 ample, if participants were randomly looking back and forth between the  
 501 two speakers, we might still find some anticipatory switching by chance. To  
 502 test whether our primary results (the regression output above) could have  
 503 arisen from random switching, not linked to turn structure, we conducted  
 504 a secondary analysis comparing participants' anticipatory gaze at real and  
 505 randomly shuffled points of turn transition.

506 We conducted this analysis by running the same regression models on  
 507 participants' eye-tracking data, only this time calculating their anticipatory  
 508 gaze switches with respect to randomly permuted turn transition windows.  
 509 This process involved: (1) randomizing the order and temporal placement of  
 510 the analysis windows within each stimulus (see Figure 4; "analysis window"  
 511 as shown in Figure 2) to randomly redistribute the analysis windows across  
 512 the eye-tracking signal, (2) re-running each participant's eye tracking data  
 513 through switch identification (described above) on each of the randomly per-  
 514 muted analysis windows, and (3) modeling the anticipatory switches from the  
 515 randomly permuted data (our random baseline dataset) with the same statis-  
 516 tical models we used for the original dataset (Table 2). Importantly, although  
 517 the onset time of each transition was shuffled within the eye-tracking signal,  
 518 the other intrinsic properties of each turn transition (e.g., prior speaker iden-  
 519 tity, transition type, gap duration, language condition, etc.) stayed constant  
 520 across each permutation.

521 The random shuffling procedure de-links participants' gaze data from the

522 turn structure in the original stimulus, thereby allowing us to compare turn-  
523 related (original) and non-turn-related (randomly permuted) looking behav-  
524 ior using the same eye movement data. We created 5,000 permutations of the  
525 original turn transitions (step 1 above), thereby creating 5,000 anticipatory  
526 gaze datasets with randomly de-linked gaze data (step 2 above). Because  
527 the randomly shuffled turn transitions could occur anywhere in the stimu-  
528 lus (so long as they didn't overlap each other within a single iteration), the  
529 resulting turn-transition windows collectively covered the entire stimulus—  
530 during speech and silence, during speaker change and speaker continuation,  
531 and during all turn transitions in the stimulus, even those excluded in the  
532 original analyses (e.g., because they were overlapped).<sup>9</sup> Pooled together, the  
533 anticipatory gaze datasets yielded an average anticipatory switch rate for  
534 each participant over all possible starting points in the stimuli: a random  
535 baseline.

536 Using this technique we compared participants' anticipatory switches at  
537 turn transition windows to their anticipatory switches over the stimulus as a  
538 whole. If participants looked randomly back and forth between the speakers,  
539 we would have seen similar patterns in both cases. Rather than simply com-  
540 paring participants' overall anticipatory switch rates with real and random  
541 transition windows, we estimated the likelihood that each of the predictor  
542 effects in the original data (e.g., the effect of language condition; Table 2)  
543 could have arisen with random gaze switching: we ran identical statistical  
544 models on the real and randomly permuted data sets. This tells us not only  
545 whether participants' switches were above chance, but whether the specific  
546 underlying effects of their anticipatory gaze patterns (e.g., the effect of lan-  
547 guage condition) were above that expected by chance. Because these analyses  
548 are complex and secondary to the main results, we report their full details  
549 in Appendix A.

550 Our baseline analyses revealed that none of the significant predictors from  
551 models of the original, turn-related data can be explained by random looking.  
552 For the children's data, the original  $z$ -values for language condition, transi-  
553 tion type, the age-language condition interaction, the transition type-gap  
554 duration interaction, and the language condition-transition type interaction

---

<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, which maximized the possibility of finding a difference from the baseline measure.

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

564 *Developmental effects*

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

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

585 When does spontaneous turn prediction emerge developmentally? We

---

<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 switches vs. “boredom” gaze switches in Appendix C

586 tested whether the youngest age group (3-year-olds) already exceeded chance  
587 in their anticipatory gaze switches by comparing children's real gaze rates  
588 to the random baseline in the English condition with two-tailed *t*-tests.  
589 We used the English condition because we are most interested in finding  
590 out when children begin to make spontaneous turn predictions for natural  
591 speech. We found that three-year-olds made anticipatory gaze switches signif-  
592 icantly above chance, when all transitions were considered ( $t(22.824)=-4.147$ ,  
593  $p<.001$ ) as well as for question transitions alone ( $t(21.677)=-5.268$ ,  $p<.001$ ).

594 *Discussion*

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

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

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

619 Third, we saw that older children made anticipatory switches more re-  
620 liably than younger children, but only in the non-English videos. In the  
621 English videos, children anticipated well at all ages, especially after hear-  
622 ing questions. This interaction between age and language condition suggests

that the 5-year-olds were able to leverage anticipatory cues in the non-English videos in a way that 3- and 4-year-olds could not, possibly by shifting more attention to the non-English prosodic or non-verbal cues. Prior work on children’s turn-structure anticipation has proposed that children’s turn-end predictions rely primarily on lexicosemantic structure (and not, e.g., prosody) as they get older (Keitel et al., 2013). The current results suggest more flexibility in children’s predictions; when they do not have access to lexical information, older children and adults find alternative cues to turn taking behavior.

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

In Experiment 2, we followed up on these findings, improving on two aspects of the design: first, our language manipulation in this first experiment was too coarse to provide data regarding specific linguistic information channels (e.g., the effect of prosodic information alone). In Experiment 2, we compared lexicosyntactic and prosodic cues with phonetically altered speech and used puppets to eliminate non-verbal cues to turn taking. Second, we were not able to pinpoint the emergence of anticipatory switching because the youngest age group in our sample was already able to make anticipatory switches at above-chance rates. In Experiment 2, we explored a wider developmental range.

659 **Experiment 2**

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

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

682 *Methods*

683 *Participants*

684 We recruited 27 undergraduate adults and 129 children ages 1;0–6;11 to  
685 participate in our experiment. Adult participants were recruited again via  
686 the Stanford University Psychology participant database and were either paid  
687 or received course credit for their time. We recruited our child participants  
688 from the Children's Discovery Museum in San Jose, California<sup>11</sup>, targeting  
689 approximately 20 children for each of the six one-year age groups (range:  
690 20–23). All participants were native English speakers, though some parents

---

<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.

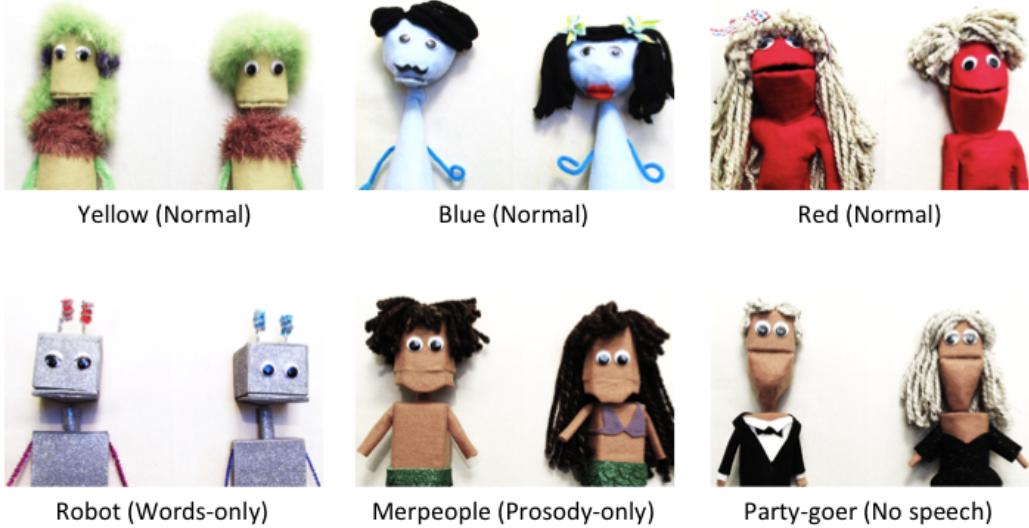


Figure 5: The six puppet pairs (and associated audio conditions). Each pair was linked to three distinct conversations from the same condition across the three experiment versions.

691 (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  
 692 language.  
 693

694 *Materials*

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

701 *Audio recordings.* The recording session was set up in the same way as  
 702 the first experiment, but with a shorter warm up period (5–10 minutes) and  
 703 a pre-determined topic for the child-friendly improvisation ('riding bikes',  
 704 'pets', 'breakfast', 'birthday cake', 'rainy days', or 'the library'). All of the  
 705 talkers were native English speakers, and were recorded in male-female pairs.  
 706 As before, we asked talkers to speak "as if they were on a children's television  
 707 show" and to ask at least a few questions during the improvisation. We cut  
 708 each audio recording down to the (approximate) 20-second interval with the

709 most turn activity. The 20-second clips were then phonetically manipulated  
710 and used in the final video stimuli.

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

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

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

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

738 Finally, because low-pass filtering removes significant acoustic energy, the  
739 *prosody only* conversations were much quieter than the other three conditions.  
740 Our last step was to downscale the intensity of the audio tracks in the three  
741 other conditions to match the volume of the *prosody only* clips. We referred  
742 to the conditions as “normal”, “robot”, “mermaid”, and “birthday party”  
743 speech when interacting with participants.

---

<sup>12</sup>We excluded hyper-lengthened words like [wau!] ‘woooow!’.

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

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

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

772        There were 85 turn transitions for analysis after excluding transitions  
773        longer than 550 msec and shorter than 90 msec. The remaining turn tran-  
774        sitions had more questions than non-questions ( $N=47$  and  $N=38$ , respec-  
775        tively), with transitions distributed somewhat evenly across conditions (keep-  
776        ing in mind that there were three *normal* videos and only one video for  
777        each partial cue condition in each experiment version): *normal* ( $N=36$ ),  
778        *words only* ( $N=13$ ), *prosody only* ( $N=17$ ), and *no speech* ( $N=19$ ). Inter-turn  
779        gaps for questions (mean=366, median=438, stdev=138 msec) were longer  
780        than those for non-questions (mean=305, median=325, stdev=94 msec) on  
781        average, but gap duration was overall comparable across conditions: *nor-*

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 in Experiment 2.

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 in Experiment 2.

782     *mal* (mean=334, median=321, stdev=130 msec), *words only* (mean=347,  
 783     median=369, stdev= 115 msec), *prosody only* (mean=365, median=369,  
 784     stdev=104 msec), and *no words* (mean=319, median=329, stdev=136 msec).

785     *Procedure*

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

796 *Data preparation and coding*

797 We coded each turn transition for its linguistic condition (*normal, words*  
 798 *only, prosody only*, and *no speech*) and transition type (question/non-question)<sup>13</sup>,  
 799 and identified anticipatory gaze switches to the upcoming speaker using the  
 800 methods from Experiment 1.

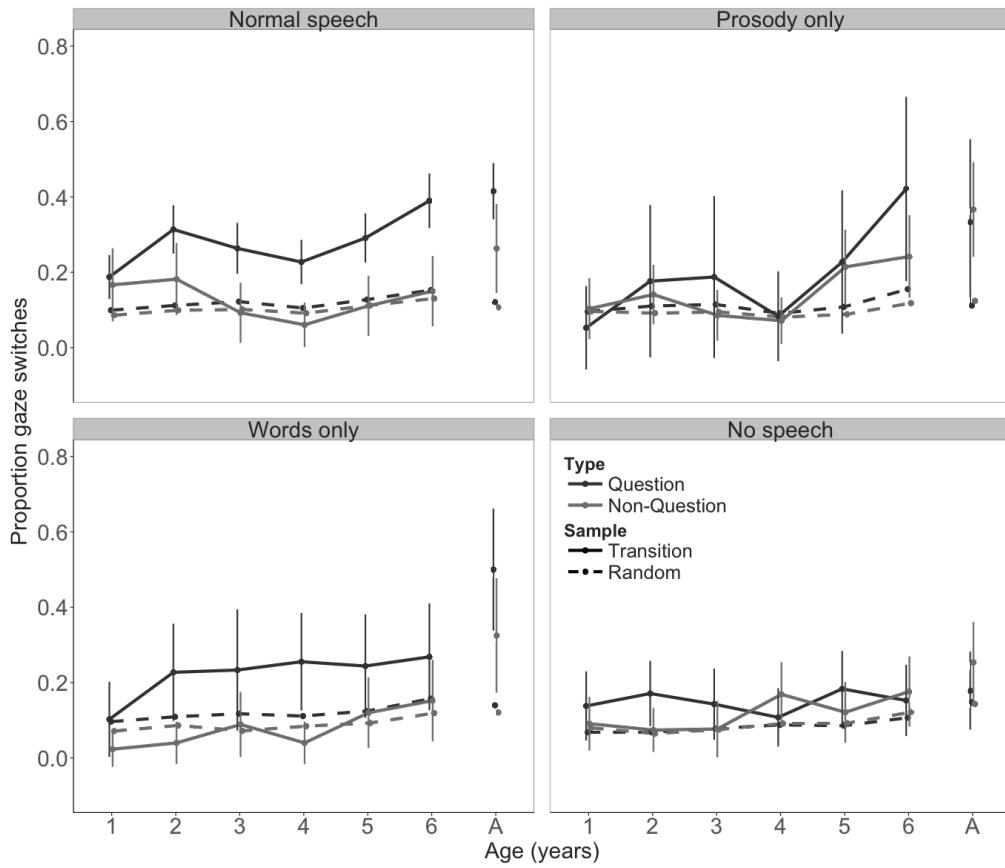


Figure 6: Anticipatory gaze rates across language condition and transition type for the real and randomly permuted datasets. Vertical bars represent 95% confidence intervals.

---

<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 do not (Hedberg et al., 2010).

801    *Results*

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

814    *Statistical models*

815    We identified anticipatory gaze switches for all 85 usable turn transi-  
816    tions, and analyzed them for effects of language condition, transition type,  
817    and age with two mixed-effects logistic regressions. We again built separate  
818    models for children and adults because effects of age were only pertinent to  
819    the children's data. The child model included condition (*normal/prosody*  
820    *only/words only/no speech*; with *no speech* as the reference level), transition  
821    type (question vs. non-question), age (1, 2, 3, 4, 5, 6; numeric, intercept as  
822    age=0), and duration of the inter-turn gap (in seconds) as predictors, with  
823    full interactions between language condition, transition type, and age and  
824    two-way interactions between gap duration and the other basic fixed effects  
825    (age, linguistic condition, and transition type). We also included random ef-  
826    fects of participant and item (turn transition), with maximal random slopes  
827    of transition type for participant. The adult model included condition, transi-  
828    tion type, their interactions, gap duration, and two-way interactions between  
829    gap duration and condition and transition type, with participant and item as  
830    random effects and maximal random slopes of condition and transition type  
831    for participant.

832    Children's anticipatory gaze switches showed an effect of gap duration  
833    ( $\beta=3.85$ ,  $SE=1.73$ ,  $z=2.22$ ,  $p<.05$ ), a two-way interaction of age and lan-  
834    guage condition (for *prosody only* speech compared to the *no speech* reference  
835    level;  $\beta=0.38$ ,  $SE=0.19$ ,  $z=1.97$ ,  $p<.05$ ), a marginal two-way interaction of  
836    language condition and gap duration (for *prosody only* speech compared to  
837    the *no speech* reference level;  $\beta=-4.77$ ,  $SE=2.63$ ,  $z=-1.82$ ,  $p=.07$ ), and a

838 three-way interaction of age, transition type, and language condition (for  
839 *normal* speech compared to the *no speech* reference level;  $\beta=-0.35$ ,  $SE=0.17$ ,  
840  $z=-2.05$ ,  $p<.05$ ). There were no significant effects of age or transition type  
841 alone (Table 5;  $\beta=-0.05$ ,  $SE=0.14$ ,  $z=-0.38$ ,  $p=.7$  and  $\beta=-1.22$ ,  $SE=0.96$ ,  
842  $z=-1.27$ ,  $p=.2$ , respectively)

843 Adults' anticipatory gaze switches showed a significant effect of language  
844 condition (for *words only* speech compared to the *no speech* reference level;  
845  $\beta=3.79$ ,  $SE=1.62$ ,  $z=2.34$ ,  $p<.05$ ) and a marginal two-way interaction be-  
846 tween language condition and transition type (for *words only* speech com-  
847 pared to the *no speech* reference level;  $\beta=-1.68$ ,  $SE=0.89$ ,  $z=-1.89$ ,  $p=.06$ ).  
848 There was no significant effect of transition type alone (Table 6;  $\beta=-0.02$ ,  
849  $SE=1.44$ ,  $z=-0.02$ ,  $p=.99$ ).

850 *Random baseline comparison*

851 Using the same technique described in Experiment 1, we created and  
852 modeled random permutations of participants' anticipatory gaze switches.  
853 These analyses revealed that the significant predictors from models of the  
854 original, turn-related data were unlikely to be explained by random looking.  
855 In the children's data, the original model's  $z$ -values for gap duration, the  
856 two-way interaction of age and language condition (*prosody only*) and the  
857 three-way interaction of age, transition type, and language condition (*normal*  
858 speech) were all greater than 93% of the randomly permuted  $z$ -values (95.6%,  
859 94%, and 93.3%, respectively,  $p=.04$ , .06, and .07). Similarly, the adults'  
860 data showed significant differentiation from the randomly permuted data for  
861 the effect of language condition (*words only* speech; greater than 98.3% of  
862 random  $z$ -values,  $p<.02$ ). See Section Appendix A for more information on  
863 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
TType= <i>non-Question</i>	-1.217	0.958	-1.27	0.204
GapDuration	3.852	1.735	2.221	0.026 *
Age*TType= <i>non-Question</i>	0.152	0.141	1.081	0.28
Age*GapDuration	0.214	0.266	0.805	0.421
TType= <i>non-Question*</i> GapDuration	0.995	2.134	0.466	0.641
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> TType= <i>non-Question</i>	0.908	0.748	1.215	0.224
Age*Condition= <i>normal*</i> TType= <i>non-Question</i>	-0.355	0.173	-2.051	0.04 *
Condition= <i>normal*</i> GapDuration	-0.431	1.67	-0.258	0.797
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> TType= <i>non-Question</i>	1.076	1.105	0.974	0.33
Age*Condition= <i>prosody*</i> TType= <i>non-Question</i>	-0.296	0.235	-1.257	0.209
Condition= <i>prosody*</i> GapDuration	-4.767	2.625	-1.816	0.069 (.)
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> TType= <i>non-Question</i>	-1.244	1.031	-1.207	0.228
Age*Condition= <i>words*</i> TType= <i>non-Question</i>	0.111	0.225	0.495	0.621
Condition= <i>words*</i> GapDuration	-2.285	2.232	-1.024	0.306

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 **
TType= <i>non-Question</i>	-0.022	1.44	-0.015	0.988
GapDuration	4.073	2.947	1.382	0.167
TType= <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
TType= <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
TType= <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 (.)
TType= <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.

866 *Developmental effects*

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

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

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

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

919 (99.5%, 97.36%, and 95.04%), all significantly different from chance ( $p < .05$ ).

920 See Figure B.3 for these *normal* speech difference score distributions.

921 *Discussion*

922 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  
923 to find the age of emergence for spontaneous turn anticipation. Many of our  
924 results replicate the findings from Experiment 1: participants often made  
925 more anticipatory switches when they had access to linguistic information  
926 and, when they did, tended to make more anticipatory switches for questions  
927 compared to non-questions.

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

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

951 **General Discussion**

952 Children begin to develop conversational turn-taking skills long before  
953 their first words emerge (Bateson, 1975; Hilbrink et al., 2015; Jaffe et al.,

954 2001; Snow, 1977). As they acquire language, they also acquire the information needed to make accurate predictions about upcoming turn structure.  
955 Until recently, we have had very little data on how children weave language  
956 into their already-existing turn-taking behaviors. In two experiments investigating  
957 children's anticipatory gaze to upcoming speakers, we found evidence  
958 that turn prediction develops early in childhood and that, when spontaneous  
959 predictions begin, they are primarily driven by participants' expectation of  
960 an immediate response in the next turn (e.g., after questions). In making  
961 predictions about upcoming turn structure, children used a combination of  
962 lexical and prosodic cues; neither signal alone was sufficient to support increased  
963 anticipatory gaze. We also found no early advantage for prosody over lexisyntax; children's anticipatory switch rates in the *prosody only*  
964 condition were initially low, but showed significant gains by age five. We dis-  
965 cuss these findings with respect to the role of linguistic cues and inter-turn  
966 silence for predicting upcoming turn structure, the importance of questions  
967 in predictions about conversation, and children's developing competence as  
968 conversationalists.

971 *Predicting upcoming turn structure*

972 Prior work with adults has found a consistent role for lexisyntax in pre-  
973 dicting upcoming turn structure (De Ruiter et al., 2006; Magyari & De Ruiter,  
974 2012), whereas the role of prosody is still under debate (Duncan, 1972; Ford  
975 & Thompson, 1996; Torreira et al., 2015). Knowing that children com-  
976 prehend more about prosody than lexisyntax early on (see Speer & Ito, 2009  
977 for a review), we thought it possible that young children would instead show  
978 an advantage for prosody in their predictions about turn structure in con-  
979 versation. Our results suggest that, on the contrary, exclusively presenting  
980 prosodic information to children limits their spontaneous predictions about  
981 upcoming turn structure until age five.

982 Thus, using prosody alone to accurately predict turn boundaries in con-  
983 versation appears to be difficult for adults and children. Prosodic informa-  
984 tion is continuous, multidimensional, and indexically complex—it encodes  
985 syntactic structure, speech act, and extralinguistic information, sometimes  
986 simultaneously, without clear one-to-one mappings between form and mean-  
987 ing (Cutler et al., 1997; Shriberg et al., 1998; Lammertink et al., 2015). For  
988 these reasons, prosodic information alone may not be enough to both (a)  
989 make precise temporal predictions about turn structure, and (b) identify  
990 question turns, which otherwise appear to drive anticipatory gaze switching.

991 Therefore, although children show early facility with prosodic discrimination  
992 (Nazzi & Ramus, 2003; Soderstrom et al., 2003; Johnson & Jusczyk,  
993 2001; Jusczyk et al., 1995; Morgan & Saffran, 1995; Mehler et al., 1988),  
994 using prosodic knowledge for turn prediction may generally be too difficult  
995 without additional information from lexical or syntactic cues.

996 Our findings suggest that there is one prosodic cue that is an exception  
997 to this rule: inter-turn silence. Generally speaking, participants showed a  
998 greater anticipatory switches for longer inter-turn gaps, but the effect of  
999 inter-turn gap duration is strongest in our data when upcoming responses  
1000 are less predictable, whether due to the asymmetrical response expectations  
1001 for questions vs. non-questions (Experiment 1) or the lack of non-verbal  
1002 cues and any linguistic information (Experiment 2). Notably, there were  
1003 no significant interactions of gap duration with participant age. This pat-  
1004 tern of results suggests that, when predictive information about upcoming  
1005 responses is absent, long silences may increase participants' expectation for  
1006 a speaker change and promote more anticipatory gaze switches. Pauses are  
1007 detected and related to phrasal structure from early on; 5-month-old infants  
1008 use pauses to parse intonational phrases (Männel & Friederici, 2009). The  
1009 lack of interactions between age and gap duration suggests that the use of  
1010 inter-turn silence remains important for older speakers and the interactions  
1011 between transition type and gap duration (Experiment 1) and condition and  
1012 gap duration (Experiment 2; marginal), suggest that this effect is simply  
1013 the result of having more time to make a gaze switch. These findings thus  
1014 suggest that silence is an early and lasting cue for identifying turn structure  
1015 online when other predictive information is not adequate.

1016 Notably, many other non-linguistic cues encode information about trans-  
1017 sition type, including gaze and gesture. We did not systematically test those  
1018 cues here but, like inter-turn silence, they may play a critical role in parsing  
1019 and making predictions about turn structure when other linguistic informa-  
1020 tion is not sufficient to make accurate predictions.

1021 Perhaps surprisingly, we found no evidence that lexical information alone  
1022 is equivalent to the full linguistic signal in driving children's predictions, as  
1023 has been shown previously for adults (Magyari & De Ruiter, 2012; De Ruiter  
1024 et al., 2006) and as is replicated with adult participants in the current study.  
1025 Unlike prosodic cues, lexicosyntactic cues are discreet and have much clearer  
1026 form-to-meaning mappings, with clear lexicosyntactic cues to questionhood  
1027 that occur early within turns (e.g., *wh*-words, *do*-insertion, and subject-  
1028 auxiliary inversion). That said, children's lexical and syntactic knowledge is

1029 limited for quite some time (Tomasello & Brooks, 1999, but see also Bergel-  
1030 son & Swingley, 2013; Shi & Melancon, 2010). Our stimuli were made in  
1031 a child-friendly style, they are still other-directed and fairly complex, with  
1032 20–30 seconds of continuous conversational speech.

1033 It is perhaps for this reason that children’s performance was always best  
1034 with the full signal, where lexicosyntactic information was supported by  
1035 prosodic information and vice versa. Even in adults, Torreira and colleagues  
1036 (2015) showed that the trade-off in informativity between lexical and prosodic  
1037 cues is more subtle in semi-natural speech. The present findings are the first  
1038 to show evidence of a similar effect developmentally.

1039 *The question effect*

1040 In both experiments, anticipatory looking was primarily driven by ques-  
1041 tion transitions, a pattern that has not been previously reported in other an-  
1042 ticipatory gaze studies, on children or adults (Keitel et al., 2013; Hirvenkari,  
1043 2013; Tice and Henetz, 2011). Questions make an upcoming speaker switch  
1044 immediately relevant, helping the listener to predict with high certainty what  
1045 will happen next (i.e., an answer from the addressee), and are often easily  
1046 identifiable by overt prosodic and lexicosyntactic cues.

1047 Prior work on children’s acquisition of questions indicates that they may  
1048 already have some knowledge of question-answer sequences by the time they  
1049 begin to speak: questions make up approximately one third of the utter-  
1050 ances children hear, before and after the onset of speech, and even into  
1051 their preschool years, though the type and complexity of questions changes  
1052 throughout development (Casillas et al., In press; Fitneva, 2012; Henning  
1053 et al., 2005; Shatz, 1979).<sup>14</sup> For the first few years, many of the questions  
1054 directed to children are “test” questions—questions that the caregiver al-  
1055 ready has the answer to (e.g., “What does a cat say?”), but this changes as  
1056 children get older. Questions help caregivers to get their young children’s  
1057 attention and to ensure that information is in common ground, even if the  
1058 responses are non-verbal or infelicitous (Bruner, 1985; Fitneva, 2012; Snow,  
1059 1977). Moreover, because of their high frequency and relatively limited num-  
1060 ber of formats, questions, especially *wh*-questions, may be more identifiable  
1061 and predictable compared to other types of speech acts. So, in addition to

---

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

1062 having a special interactive status, questions are a frequent, predictable and  
1063 core characteristic of many caregiver-child interactions, motivating a general  
1064 benefit for questions in turn structure anticipation.

1065 Two important questions for future work are then: (1) how does children's  
1066 ability to monitor for questions in conversation relate to their prior experi-  
1067 ence with questions? and (2) what is it about questions that makes children  
1068 and adults more likely to anticipatorily switch their gaze to addressees? If  
1069 this "question" effect exists for all turns that require an immediate response  
1070 ("adjacency pairs"; Schegloff, 2007), other turn types, such as imperatives,  
1071 compliments, and complaints should show similar patterns. If the effect is  
1072 instead about overall predictability of the syntactic frame, children would  
1073 instead show similar patterns for other frequent frames from child-directed  
1074 speech (e.g., "Look at the X"; Mintz, 2003). The recognizability and pre-  
1075 dictability of syntactic frames is likely to play a role in turn prediction as  
1076 children become more sophisticated language users, even if the effect is truly  
1077 about adjacency pairs; for example, rhetorical and tag questions take a very  
1078 similar form to prototypical polar questions, but usually do not require an  
1079 answer. So, though it is clear that adults and children anticipate responses  
1080 more often for questions than non-questions, we do not yet know whether  
1081 their predictive action is limited to turns formatted as questions, turns with  
1082 high recognizability and predictability, or turns that project an immediate  
1083 response from the addressee.

1084 The question effect suggests that participants' spontaneous predictions  
1085 may be driven by what lies *beyond* the end of the current turn—not just  
1086 by the upcoming end of the turn itself, as has been focused on in prior  
1087 work (Torreira et al., 2015; Keitel et al., 2013; Magyari & De Ruiter, 2012;  
1088 De Ruiter et al., 2006). In future work, it will be crucial to measure prediction  
1089 from a first-person perspective to find out what kinds of predictions are most  
1090 relevant to addressees in conversation (see also Holler & Kendrick, 2015).

1091 One possible scenario is that listeners in spontaneous, first-person con-  
1092 versation use multiple strategies to make predictions about upcoming turn  
1093 structure: they could semi-passively attend to incoming speech for cues to  
1094 upcoming speaker transition (e.g., questions and other adjacency pairs) and,  
1095 when possible upcoming transition is detected, switch into a more precise  
1096 turn-end prediction mode (à la De Ruiter et al., 2006). A flexible prediction  
1097 system like this one allows listeners to continuously monitor ongoing conver-  
1098 sation for turn-related cues at a low cost while still managing to plan their  
1099 responses and come in quickly when needed.

1100 To test this hypothesis, we would need to look at prediction from a first-  
1101 person perspective, which very little work so far has accomplished (present  
1102 work included). Although third-party measures enable us to measure partic-  
1103 ipants' predictions without any interference from language production, they  
1104 also limit our knowledge about how the need to give a response might it-  
1105 self play an important role in addressees' prediction strategies. Recent work  
1106 has shown that shifts in addressee gaze similar to those measured here in-  
1107 deed occur in spontaneous conversation (Holler & Kendrick, 2015), but much  
1108 more work is needed to determine how participants make predictions about  
1109 turn structure in first-person contexts and whether those mechanisms shift  
1110 at points of imminent speaker change.

1111 *Early competence for turn taking?*

1112 One of the core aims of our study was to test whether children show an  
1113 early competence for turn taking, as is proposed by studies of spontaneous  
1114 mother-infant proto-conversation and theories about the mechanisms under-  
1115 lying human interaction in general (Hilbrink et al., 2015; Levinson, 2006).  
1116 We found evidence that young children make spontaneous predictions about  
1117 upcoming turn structure: definitely at age two and marginally at age one.

1118 These results contrast with Keitel and colleagues' (2013) finding that chil-  
1119 dren cannot anticipate upcoming turn structure at above-chance rates until  
1120 age three. The current study used an appreciably more conservative random  
1121 baseline than the one used in Keitel and colleagues' study. Therefore, this  
1122 difference in age of emergence more likely stems from our use of a more en-  
1123 gaging speech style, stereo speech playback, and more typical turn transition  
1124 durations. The child-friendly style of speech in particular may have helped in  
1125 two ways: keeping children more engaged with the stimuli and using less syn-  
1126 tactically complex and more prosodically exaggerated speech (Fernald et al.,  
1127 1989; Werker & McLeod, 1989; Snow, 1977) compared to what they would  
1128 get with adult-adult conversation.

1129 To be clear, young children's "above chance" performance was often still  
1130 far from adult-like predictive behavior—children at ages one and two were  
1131 still very close to chance in their anticipations and, even at age six, chil-  
1132 dren were not fully adult-like in their predictions. This may indicate that  
1133 young children rely primarily on non-verbal cues (like inter-turn silence) in  
1134 anticipating turn transitions while adults use both verbal and non-verbal  
1135 cues to make predictions. Relatedly, adults may be more expert in flexibly

1136 adapting to the turn-relevant cues present at any moment, e.g., responding  
1137 to non-English prosodic cues in Experiment 1.

1138 Taken together, our data suggest that turn-taking skills do begin to  
1139 emerge in infancy, but that children cannot consistently make effective pre-  
1140 dictions until they can identify and react to question turns. This finding  
1141 leads us to wonder how participant role (first- instead of third-person) and  
1142 differences in early interactional experience (e.g., frequent vs. infrequent  
1143 question-asking from caregivers) feed into this early predictive skill. It also  
1144 bridges prior work showing a predisposition for turn taking in infancy (e.g.,  
1145 Bateson, 1975; Hilbrink et al., 2015; Jaffe et al., 2001; Snow, 1977) with  
1146 children's apparently *late* acquisition of adult-like competence for turn tak-  
1147 ing in actual conversation (Casillas et al., In press; Garvey, 1984; Garvey  
1148 & Berninger, 1981; Ervin-Tripp, 1979) and reinforces the idea that it takes  
1149 children several years to fully integrate linguistic information into their turn-  
1150 taking systems (Casillas et al., In press; Garvey & Berninger, 1981).

1151 *Limitations and future work*

1152 There are at least two major limitations to our work: speech naturalness  
1153 and participant role. Following prior work (De Ruiter et al., 2006; Keitel  
1154 et al., 2013), we used phonetically manipulated speech in Experiment 2.  
1155 This decision resulted in speech sounds that children don't usually hear in  
1156 their natural environment. Many prior studies have used phonetically-altered  
1157 speech with infants and young children (cf. Jusczyk, 2000), but few of them  
1158 have done so in a conversational context. Future work could instead carefully  
1159 script speech or cross-splice sub-parts of turns to control for the presence of  
1160 linguistic cues for turn transition (see, e.g., Torreira et al., 2015).

1161 The prediction measure used in our studies is based on an observer's view  
1162 of conversation but, because participants' role in the interaction could affect  
1163 their online predictions about turn taking, an ideal measure would instead  
1164 capture first-person predictions. If conversational participants' predictions  
1165 are partly shaped by their need to respond, first-person measures of spon-  
1166 taneous turn prediction will be key to revealing how participants distribute  
1167 their attention over verbal and non-verbal cues while taking part in everyday  
1168 interaction, the implications of which relate to theories of online language  
1169 processing for both language learning and everyday talk.

1170 That said, the third-person paradigm used in the present study still has  
1171 much to tell us about turn prediction. The task is natural and intuitive  
1172 in that no instruction is required, which means that it captures spontaneous

1173 predictive behavior and can be used with participants of all ages. Frequencies  
1174 of anticipatory gaze switching appear to be stable across language commu-  
1175 nities where similar tasks have been tested (Keitel et al., 2013; Keitel &  
1176 Daum, 2015; Holler & Kendrick, 2015; Hirvenkari et al., 2013)—even from a  
1177 first-person perspective—so the task is one that measures robust predictive  
1178 behavior relevant to conversational processing across languages. It also lends  
1179 itself to many possibilities for controlling the presence of individual verbal  
1180 and non-verbal cues and has a clear method for assessing random switch-  
1181 ing baselines across the entire stimulus. Also, if it is the case that response  
1182 preparation interferes with our ability to see prediction at the ends of in-  
1183 coming turns, third-person paradigms are one of the only ways to measure  
1184 prediction processes in isolation.

1185 The current findings also make predictions about what we would see in  
1186 first-person paradigms. For example, a focus on possible upcoming speaker  
1187 transitions is even more important when the participants themselves may  
1188 need to respond; we would thus expect question-like effects to occur in first-  
1189 person paradigms, and perhaps even be amplified compared to third-person  
1190 paradigms. If so, participants' use of linguistic information would still sub-  
1191 serve this goal, with prediction at a premium. Regarding development, the  
1192 same facts about the complexity of prosody-based prediction and children's  
1193 initial limited lexical inventories would still hold, as would the use of silence  
1194 and non-verbal cues to assess and predict turn structure in the absence of  
1195 clear predictive linguistic information. The paradigm presented here thus has  
1196 important contributions to make in our understanding of how participants  
1197 attend to and make predictions about conversational interaction.

### 1198 *Conclusions*

1199 Conversation plays a central role in children's language learning. It is  
1200 the driving force behind what children say and what they hear. Adults use  
1201 linguistic information to accurately predict turn structure in conversation,  
1202 which facilitates their online comprehension and allows them to respond rel-  
1203 evantly and on time. The present study offers new findings regarding the  
1204 role of speech acts and linguistic processing in online turn prediction, and  
1205 has given evidence that turn prediction emerges by age two, increases with  
1206 age, and is driven by question turns. However, turn prediction is not fully  
1207 integrated with linguistic cues until much later and, in the absence of pre-  
1208 dictive linguistic cues, children and adults alike rely on retroactive cues such  
1209 as inter-turn silence to predict upcoming speaker change. Using language to

1210 make predictions about upcoming interactive content takes time to develop  
1211 and, for participants of all ages appears to be primarily driven by partic-  
1212 ipants' expectations about what will happen next, beyond the end of the  
1213 current turn.

1214 **Acknowledgements**

1215 We gratefully acknowledge the parents and children at Bing Nursery  
1216 School and the Children's Discovery Museum of San Jose. This work was  
1217 supported by an ERC Advanced Grant to Stephen C. Levinson (269484-  
1218 INTERACT), an NSF Graduate Research Fellowship and NSF Dissertation  
1219 Improvement Grant to MC, and a Merck Foundation fellowship to MCF.  
1220 Earlier versions of these data and analyses were presented to conference au-  
1221 diences (Casillas & Frank, 2012, 2013). We also thank Tania Henetz, Fran-  
1222 cisco Torreira, Stephen C. Levinson, and Eve V. Clark for their feedback on  
1223 earlier versions of this work. The analysis code for this project can be found  
1224 on GitHub at [https://github.com/langcog/turn\\_taking/](https://github.com/langcog/turn_taking/).

1225 **References**

- 1226 Allison, P. D. (2004). Convergence problems in logistic regression. In M. Alt-  
1227 man, J. Gill, & M. McDonald (Eds.), *Numerical Issues in Statistical Com-*  
1228 *puting for the Social Scientist* (pp. 247–262). Wiley-Interscience: New  
1229 York, NY.
- 1230 Allison, P. D. (2012). *Logistic Regression Using SAS: Theory and Applica-*  
1231 *tion*. SAS Institute.
- 1232 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects  
1233 structure for confirmatory hypothesis testing: Keep it maximal. *Journal*  
1234 *of Memory and Language*, 68, 255–278.
- 1235 Bates, D., Maechler, M., Bolker, B., & Walker, S. (2014).  
1236 *lme4: Linear mixed-effects models using Eigen and S4*. URL:  
1237 <https://github.com/lme4/lme4>/<http://lme4.r-forge.r-project.org/>  
1238 [Computer program] R package version 1.1-7.
- 1239 Bateson, M. C. (1975). Mother-infant exchanges: The epigenesis of conver-  
1240 sational interaction. *Annals of the New York Academy of Sciences*, 263,  
1241 101–113.

- 1242 Bergelson, E., & Swingley, D. (2013). The acquisition of abstract words by  
1243 young infants. *Cognition*, 127, 391–397.
- 1244 Bloom, K. (1988). Quality of adult vocalizations affects the quality of infant  
1245 vocalizations. *Journal of Child Language*, 15, 469–480.
- 1246 Boersma, P., & Weenink, D. (2012). *Praat: doing phonetics by computer*.  
1247 URL: <http://www.praat.org> [Computer program] Version 5.3.16.
- 1248 Bögels, S., Magyari, L., & Levinson, S. C. (2015). Neural signatures of  
1249 response planning occur midway through an incoming question in conver-  
1250 sation. *Scientific Reports*, 5.
- 1251 Bruner, J. (1985). Child's talk: Learning to use language. *Child Language  
Teaching and Therapy*, 1, 111–114.
- 1253 Bruner, J. S. (1975). The ontogenesis of speech acts. *Journal of Child Lan-  
guage*, 2, 1–19.
- 1255 Carlson, R., Hirschberg, J., & Swerts, M. (2005). Cues to upcoming swedish  
1256 prosodic boundaries: Subjective judgment studies and acoustic correlates.  
1257 *Speech Communication*, 46, 326–333.
- 1258 Casillas, M., Bobb, S. C., & Clark, E. V. (In press). Turn taking, timing,  
1259 and planning in early language acquisition. *Journal of Child Language*, .
- 1260 Casillas, M., & Frank, M. C. (2012). Cues to turn boundary prediction in  
1261 adults and preschoolers. In *Proceedings of SemDial* (pp. 61–69).
- 1262 Casillas, M., & Frank, M. C. (2013). The development of predictive processes  
1263 in children's discourse understanding. In *Proceedings of the 35th Annual  
1264 Meeting of the Cognitive Science Society* (pp. 299–304).
- 1265 Cutler, A., Dahan, D., & Van Donselaar, W. (1997). Prosody in the com-  
1266 prehension of spoken language: A literature review. *Language and speech*,  
1267 40, 141–201.
- 1268 De Ruiter, J. P., Mitterer, H., & Enfield, N. J. (2006). Projecting the end of  
1269 a speaker's turn: A cognitive cornerstone of conversation. *Language*, 82,  
1270 515–535.

- 1271 De Vos, C., Torreira, F., & Levinson, S. C. (2015). Turn-timing in signed  
1272 conversations: coordinating stroke-to-stroke turn boundaries. *Frontiers in*  
1273 *Psychology*, 6.
- 1274 Dingemanse, M., Torreira, F., & Enfield, N. (2013). Is “Huh?” a univer-  
1275 sal word? Conversational infrastructure and the convergent evolution of  
1276 linguistic items. *PloS one*, 8, e78273.
- 1277 Duncan, S. (1972). Some signals and rules for taking speaking turns in  
1278 conversations. *Journal of Personality and Social Psychology*, 23, 283.
- 1279 Ervin-Tripp, S. (1979). Children’s verbal turn-taking. In E. Ochs, & B. B.  
1280 Schieffelin (Eds.), *Developmental Pragmatics* (pp. 391–414). Academic  
1281 Press, New York.
- 1282 Fernald, A., Taeschner, T., Dunn, J., Papousek, M., de Boysson-Bardies,  
1283 B., & Fukui, I. (1989). A cross-language study of prosodic modifications  
1284 in mothers’ and fathers’ speech to preverbal infants. *Journal of Child*  
1285 *Language*, 16, 477–501.
- 1286 Fitneva, S. (2012). Beyond answers: questions and children’s learning. In  
1287 J.-P. De Ruiter (Ed.), *Questions: Formal, Functional, and Interactional*  
1288 *Perspectives* (pp. 165–178). Cambridge University Press, Cambridge, UK.
- 1289 Ford, C. E., & Thompson, S. A. (1996). Interactional units in conversation:  
1290 Syntactic, intonational, and pragmatic resources for the management of  
1291 turns. *Studies in Interactional Sociolinguistics*, 13, 134–184.
- 1292 Garvey, C. (1984). *Children’s Talk* volume 21. Harvard University Press.
- 1293 Garvey, C., & Berninger, G. (1981). Timing and turn taking in children’s  
1294 conversations 1. *Discourse Processes*, 4, 27–57.
- 1295 Gísladóttir, R., Chwilla, D., & Levinson, S. C. (2015). Conversation electri-  
1296 fied: ERP correlates of speech act recognition in underspecified utterances.  
1297 *PloS one*, 10, e0120068.
- 1298 Griffin, Z. M., & Bock, K. (2000). What the eyes say about speaking. *Psy-  
1299 chological science*, 11, 274–279.

- 1300 Hart, B., & Risley, T. R. (1992). American parenting of language-learning  
1301 children: Persisting differences in family-child interactions observed in nat-  
1302 ural home environments. *Developmental Psychology*, 28, 1096.
- 1303 Hedberg, N., Sosa, J. M., Görgülü, E., & Mameni, M. (2010). The prosody  
1304 and meaning of Wh-questions in American English. In *Speech Prosody*  
1305 2010 (pp. 100045:1–4).
- 1306 Henning, A., Striano, T., & Lieven, E. V. (2005). Maternal speech to infants  
1307 at 1 and 3 months of age. *Infant Behavior and Development*, 28, 519–536.
- 1308 Hilbrink, E., Gattis, M., & Levinson, S. C. (2015). Early developmental  
1309 changes in the timing of turn-taking: A longitudinal study of mother-  
1310 infant interaction. *Frontiers in Psychology*, 6.
- 1311 Hirvenkari, L., Ruusuvuori, J., Saarinen, V.-M., Kivioja, M., Peräkylä, A.,  
1312 & Hari, R. (2013). Influence of turn-taking in a two-person conversation  
1313 on the gaze of a viewer. *PloS one*, 8, e71569.
- 1314 Holler, J., & Kendrick, K. H. (2015). Unaddressed participants' gaze in  
1315 multi-person interaction. *Frontiers in Psychology*, 6.
- 1316 Jaffé, J., Beebe, B., Feldstein, S., Crown, C. L., Jasnow, M. D., Rochat,  
1317 P., & Stern, D. N. (2001). *Rhythms of dialogue in infancy: Coordinated  
1318 timing in development*. Monographs of the Society for Research in Child  
1319 Development. JSTOR.
- 1320 Johnson, E. K., & Jusczyk, P. W. (2001). Word segmentation by 8-  
1321 month-olds: When speech cues count more than statistics. *Journal of  
1322 Memory and Language*, 44, 548–567.
- 1323 Jusczyk, P. W. (2000). *The Discovery of Spoken Language*. MIT press.
- 1324 Jusczyk, P. W., Hohne, E., Mandel, D., & Strange, W. (1995). Picking up  
1325 regularities in the sound structure of the native language. *Speech perception  
1326 and linguistic experience: Theoretical and methodological issues in cross-  
1327 language speech research*, (pp. 91–119).
- 1328 Kamide, Y., Altmann, G., & Haywood, S. L. (2003). The time-course of  
1329 prediction in incremental sentence processing: Evidence from anticipatory  
1330 eye movements. *Journal of Memory and Language*, 49, 133–156.

- 1331 Keitel, A., & Daum, M. M. (2015). The use of intonation for turn anticipation  
1332 in observed conversations without visual signals as source of information.  
1333 *Frontiers in Psychology*, 6.
- 1334 Keitel, A., Prinz, W., Friederici, A. D., Hofsten, C. v., & Daum, M. M.  
1335 (2013). Perception of conversations: The importance of semantics and  
1336 intonation in childrens development. *Journal of Experimental Child Psychology*, 116, 264–277.
- 1338 Lammertink, I., Casillas, M., Benders, T., Post, B., & Fikkert, P. (2015).  
1339 Dutch and english toddlers' use of linguistic cues in predicting upcoming  
1340 turn transitions. *Frontiers in Psychology*, 6.
- 1341 Lemasson, A., Glas, L., Barbu, S., Lacroix, A., Guilloux, M., Remeuf, K., &  
1342 Koda, H. (2011). Youngsters do not pay attention to conversational rules:  
1343 is this so for nonhuman primates? *Nature Scientific Reports*, 1.
- 1344 Levelt, W. J. (1989). *Speaking: From intention to articulation*. MIT press.
- 1345 Levinson, S. C. (2006). On the human “interaction engine”. In N. Enfield,  
1346 & S. Levinson (Eds.), *Roots of Human Sociality: Culture, Cognition and*  
1347 *Interaction* (pp. 39–69). Oxford: Berg.
- 1348 Levinson, S. C. (2013). Action formation and ascriptions. In T. Stivers, &  
1349 J. Sidnell (Eds.), *The Handbook of Conversation Analysis* (pp. 103–130).  
1350 Wiley-Blackwell, Malden, MA.
- 1351 Levinson, S. C. (2016). Turn-taking in Human Communication – Origins  
1352 and Implications for Language Processing. *Trends in Cognitive Sciences*,  
1353 20, 6–14.
- 1354 Magyari, L., Bastiaansen, M. C. M., De Ruiter, J. P., & Levinson, S. C.  
1355 (2014). Early anticipation lies behind the speed of response in conversation.  
1356 *Journal of Cognitive Neuroscience*, 26, 2530–2539.
- 1357 Magyari, L., & De Ruiter, J. P. (2012). Prediction of turn-ends based on  
1358 anticipation of upcoming words. *Frontiers in Psychology*, 3:376, 1–9.
- 1359 Männel, C., & Friederici, A. D. (2009). Pauses and intonational phrasing:  
1360 ERP studies in 5-month-old German infants and adults. *Journal of Cognitive Neuroscience*, 21, 1988–2006.

- 1362 Masataka, N. (1993). Effects of contingent and noncontingent maternal stimulation  
1363 on the vocal behaviour of three-to four-month-old Japanese infants.  
1364 *Journal of Child Language*, 20, 303–312.
- 1365 Mehler, J., Jusczyk, P., Lambertz, G., Halsted, N., Bertoni, J., & Amiel-  
1366 Tison, C. (1988). A precursor of language acquisition in young infants.  
1367 *Cognition*, 29, 143–178.
- 1368 Mintz, T. H. (2003). Frequent frames as a cue for grammatical categories in  
1369 child directed speech. *Cognition*, 90, 91–117.
- 1370 Morgan, J. L., & Saffran, J. R. (1995). Emerging integration of sequential  
1371 and suprasegmental information in preverbal speech segmentation. *Child  
1372 Development*, 66, 911–936.
- 1373 Nazzi, T., & Ramus, F. (2003). Perception and acquisition of linguistic  
1374 rhythm by infants. *Speech Communication*, 41, 233–243.
- 1375 Nomikou, I., & Rohlfing, K. J. (2011). Language does something: Body  
1376 action and language in maternal input to three-month-olds. *IEEE Transactions  
1377 on Autonomous Mental Development*, 3, 113–128.
- 1378 R Core Team (2014). *R: A Language and Environment for Statistical Com-  
1379 puting*. R Foundation for Statistical Computing Vienna, Austria. URL:  
1380 <http://www.R-project.org> [Computer program] Version 3.1.1.
- 1381 Ratner, N., & Bruner, J. (1978). Games, social exchange and the acquisition  
1382 of language. *Journal of Child Language*, 5, 391–401.
- 1383 Reddy, V., Markova, G., & Wallot, S. (2013). Anticipatory adjustments to  
1384 being picked up in infancy. *PloS one*, 8, e65289.
- 1385 Ross, H. S., & Lollis, S. P. (1987). Communication within infant social games.  
1386 *Developmental Psychology*, 23, 241.
- 1387 Rossano, F., Brown, P., & Levinson, S. C. (2009). Gaze, questioning and cul-  
1388 ture. In J. Sidnell (Ed.), *Conversation Analysis: Comparative Perspectives*  
1389 (pp. 187–249). Cambridge University Press, Cambridge.
- 1390 Sacks, H., Schegloff, E. A., & Jefferson, G. (1974). A simplest systematics for  
1391 the organization of turn-taking for conversation. *Language*, 50, 696–735.

- 1392 Schegloff, E. A. (2007). *Sequence organization in interaction: Volume 1: A*  
1393 *primer in conversation analysis*. Cambridge University Press.
- 1394 Shatz, M. (1978). On the development of communicative understandings:  
1395 An early strategy for interpreting and responding to messages. *Cognitive*  
1396 *Psychology*, 10, 271–301.
- 1397 Shatz, M. (1979). How to do things by asking: Form-function pairings in  
1398 mothers' questions and their relation to children's responses. *Child Develop-*  
1399 *ment*, 50, 1093–1099.
- 1400 Shi, R., & Melancon, A. (2010). Syntactic categorization in French-learning  
1401 infants. *Infancy*, 15, 517–533.
- 1402 Shriberg, E., Stolcke, A., Jurafsky, D., Coccato, N., Meteer, M., Bates, R.,  
1403 Taylor, P., Ries, K., Martin, R., & Van Ess-Dykema, C. (1998). Can  
1404 prosody aid the automatic classification of dialog acts in conversational  
1405 speech? *Language and Speech*, 41, 443–492.
- 1406 Snow, C. E. (1977). The development of conversation between mothers and  
1407 babies. *Journal of Child Language*, 4, 1–22.
- 1408 Soderstrom, M., Seidl, A., Kemler Nelson, D. G., & Jusczyk, P. W. (2003).  
1409 The prosodic bootstrapping of phrases: Evidence from prelinguistic in-  
1410 fants. *Journal of Memory and Language*, 49, 249–267.
- 1411 Speer, S. R., & Ito, K. (2009). Prosody in first language acquisition—  
1412 Acquiring intonation as a tool to organize information in conversation.  
1413 *Language and Linguistics Compass*, 3, 90–110.
- 1414 Stivers, T., Enfield, N. J., Brown, P., Englert, C., Hayashi, M., Heinemann,  
1415 T., Hoymann, G., Rossano, F., De Ruiter, J. P., Yoon, K.-E. et al. (2009).  
1416 Universals and cultural variation in turn-taking in conversation. *Proceed-  
1417 ings of the National Academy of Sciences*, 106, 10587–10592.
- 1418 Stivers, T., & Rossano, F. (2010). Mobilizing response. *Research on Language*  
1419 *and Social Interaction*, 43, 3–31.
- 1420 Takahashi, D. Y., Narayanan, D. Z., & Ghazanfar, A. A. (2013). Coupled  
1421 oscillator dynamics of vocal turn-taking in monkeys. *Current Biology*, 23,  
1422 2162–2168.

- 1423 Thorgrímsson, G., Fawcett, C., & Liszkowski, U. (2015). 1- and 2-year-olds'  
1424 expectations about third-party communicative actions. *Infant Behavior*  
1425 and *Development*, 39, 53–66.
- 1426 Tice (Casillas), M., & Henetz, T. (2011). Turn-boundary projection: Looking  
1427 ahead. In *Proceedings of the 33rd Annual Meeting of the Cognitive Science*  
1428 *Society* (pp. 838–843).
- 1429 Toda, S., & Fogel, A. (1993). Infant response to the still-face situation at 3  
1430 and 6 months. *Developmental Psychology*, 29, 532.
- 1431 Tomasello, M., & Brooks, P. J. (1999). Early syntactic development: A  
1432 construction grammar approach. In M. Barrett (Ed.), *The development of*  
1433 *language* (pp. 161–190). Psychology Press.
- 1434 Torreira, F., Bögels, S., & Levinson, S. C. (2015). Intonational phrasing is  
1435 necessary for turn-taking in spoken interaction. *Journal of Phonetics*, 52,  
1436 46–57.
- 1437 Weisleder, A. (2012). *Richer language experience leads to faster understanding: Links between language input, processing efficiency, and vocabulary growth*. Ph.D. thesis Stanford University.
- 1440 Werker, J. F., & McLeod, P. J. (1989). Infant preference for both male and  
1441 female infant-directed talk: A developmental study of attentional and af-  
1442 ffective responsiveness. *Canadian Journal of Psychology/Revue Canadienne*  
1443 *de Psychologie*, 43, 230.
- 1444 Wittenburg, P., Brugman, H., Russel, A., Klassmann, A., & Sloetjes, H.  
1445 (2006). Elan: a professional framework for multimodality research. In  
1446 *Proceedings of LREC*.

1447 **Appendix A. Permutation Analyses**

1448 How can we be sure that our primary dependent measure (anticipatory  
1449 gaze switching) actually relates to turn transitions? Even if children were  
1450 gazing back and forth randomly during the experiment, we would have still  
1451 captured some false hits—switches that ended up in the turn-transition win-  
1452 dows by chance.

1453 We estimated the baseline probability of making an anticipatory switch  
1454 by randomly permuting the placement of the transition windows within each  
1455 stimulus (Figure 4). We then used the switch identification procedure from  
1456 Experiments 1 and 2 to find out how often participants made “anticipatory”  
1457 switches within these randomly permuted windows. This procedure de-links  
1458 participants’ gaze data from turn structure by randomly re-assigning the on-  
1459 set time of each turn-transition in each permutation. We created 5,000 of  
1460 these permutations for each experiment to get an anticipatory switch base-  
1461 lines over all possible starting points.

1462 Importantly, the randomized windows were not allowed to overlap with  
1463 each other, keeping true to the original stimuli. We also made sure that the  
1464 properties of each turn transition stayed constant across permutations. So,  
1465 while “transition window A” might start 2 seconds into Random Permu-  
1466 tation 1 and 17 seconds into Random Permutation 2, it maintained the same  
1467 prior speaker identity, transition type, gap duration, language condition, etc.,  
1468 across both permutations.

1469 We then re-ran the statistical models from the original data on each of the  
1470 random permutations, e.g., using Experiment 1’s original model to analyze  
1471 the anticipatory switches from each random permutation of the Experiment  
1472 1 looking data. We could then calculate the proportion of random data  
1473  $z$ -values exceeded by the original  $z$ -value for each predictor. We used the  
1474 absolute value of all  $z$ -values to conduct a two-tailed test. If the original  
1475 effect of a predictor exceeded 95% of the random model effects for that same  
1476 predictor, we deemed that predictor’s effect to be significantly different from  
1477 the random baseline (i.e.,  $p < .05$ ).

1478 For example, children’s “language condition” effect from Experiment 1  
1479 had a  $z$ -value of  $|3.65|$ , which is greater than 99.3% of all  $|z\text{-value}|$  estimates  
1480 from Experiment 1’s random permutation models (i.e.,  $p = .007$ ). It is there-  
1481 fore highly unlikely that the effect of language condition in the original model  
1482 came from random gaze shifting.

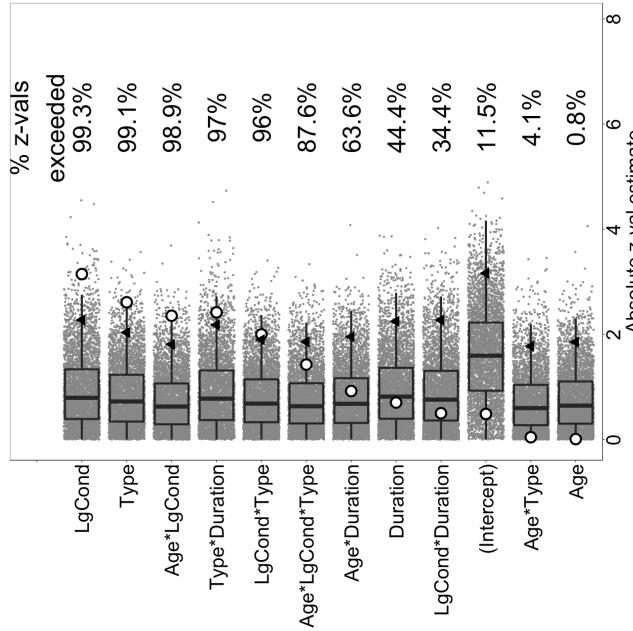
1483 We used this procedure to derive the random-baseline comparison values

1484 in the main text (above). However, we ran into two issues along the way:  
1485 first, we had to report  $z$ -values rather than beta estimates of each effect.  
1486 Second, we had to exclude a substantial portion of the models, especially in  
1487 Experiment 2 because of model non-convergence. We address each of these  
1488 issues below.

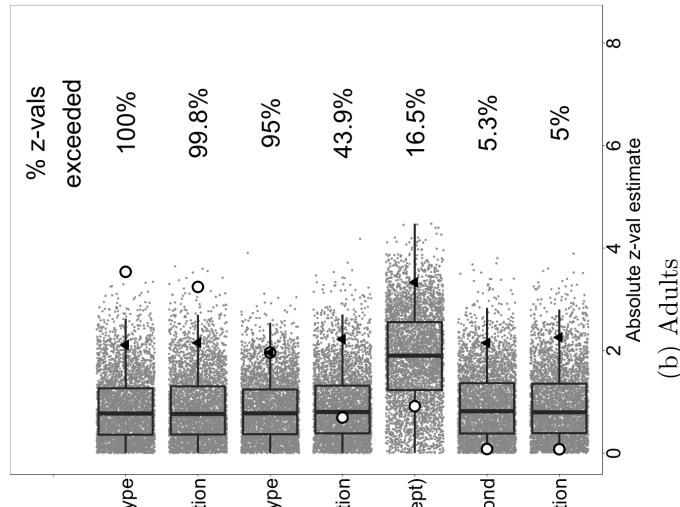
1489 *Appendix A.1. Beta, standard error, and  $z$  estimates*

1490 We reported  $z$ -values in the main text rather than beta estimates because  
1491 the standard errors in the randomly permuted data models were much higher  
1492 than for the original data. The distributions for each predictor's beta esti-  
1493 mate, standard error, and  $z$ -value for adults and children in each experiment  
1494 are shown in the graphs below (Figures A.1a–A.6b). In each plot, the gray  
1495 dots represent the absolute value of the 5,000 randomly permuted model es-  
1496 timates for the estimate type plotted (beta, standard error, or  $z$ ), the white  
1497 circles represent the model estimates from the original data, and the black  
1498 triangles represent the 95th percentile for each random distribution.

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



53

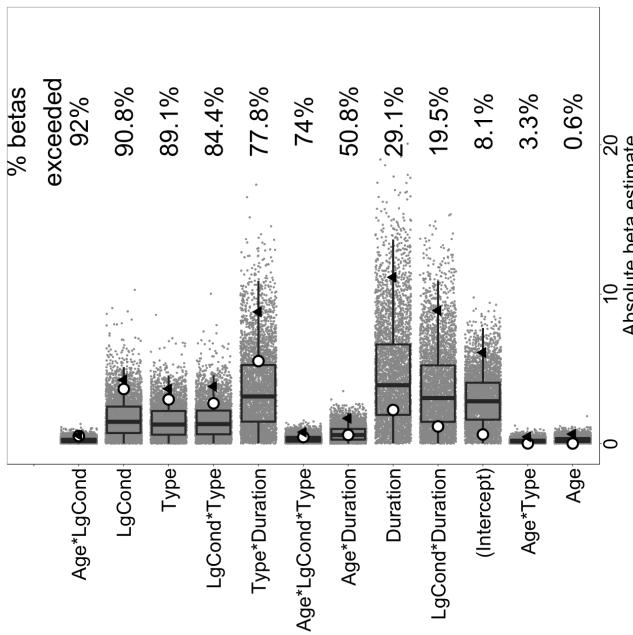


(a) Children

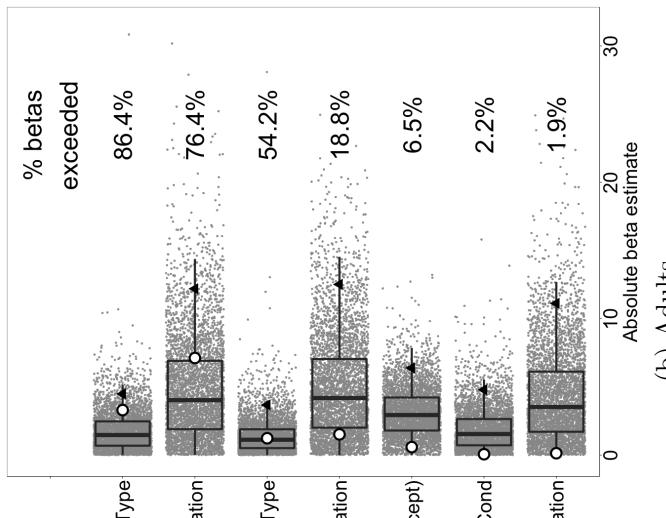
(b) Adults

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

### Experiment 1: $\beta$ estimates



(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

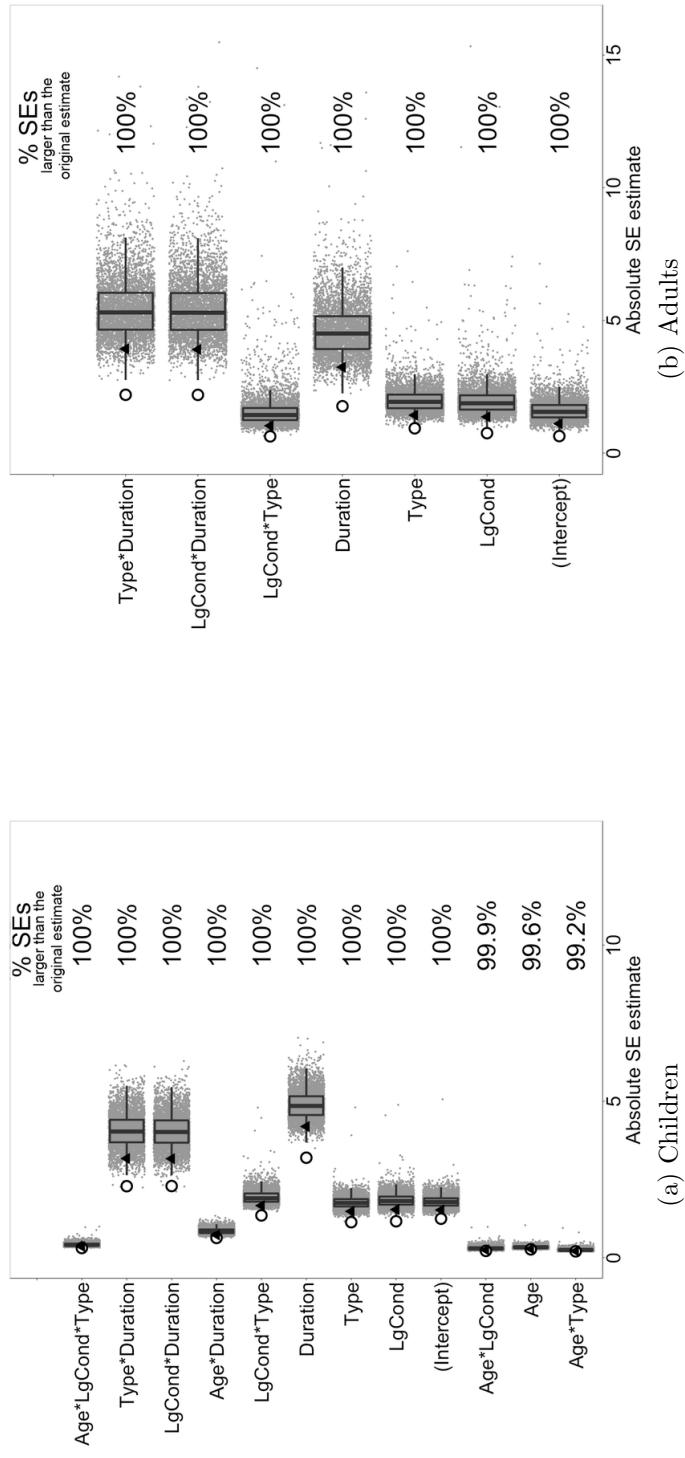


Figure A.3: Random-permutation and original *SE*-values for predictors of anticipatory gaze rates in Experiment 1.

## Experiment 2: $z$ estimates

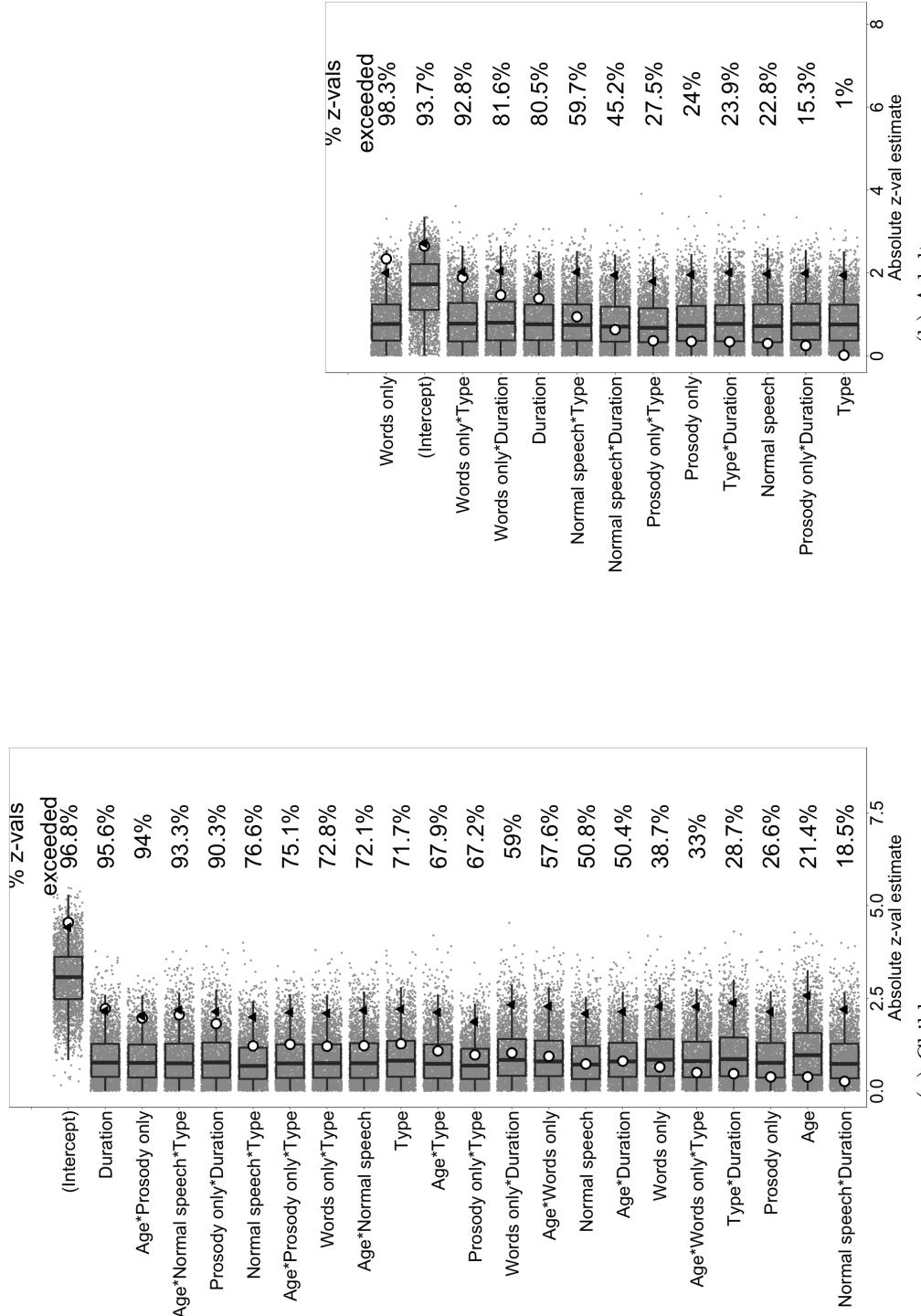


Figure A.4: Random-permutation and original  $|z\text{-values}|$  for predictors of anticipatory gaze rates in Experiment 2.

## Experiment 2: $\beta$ estimates

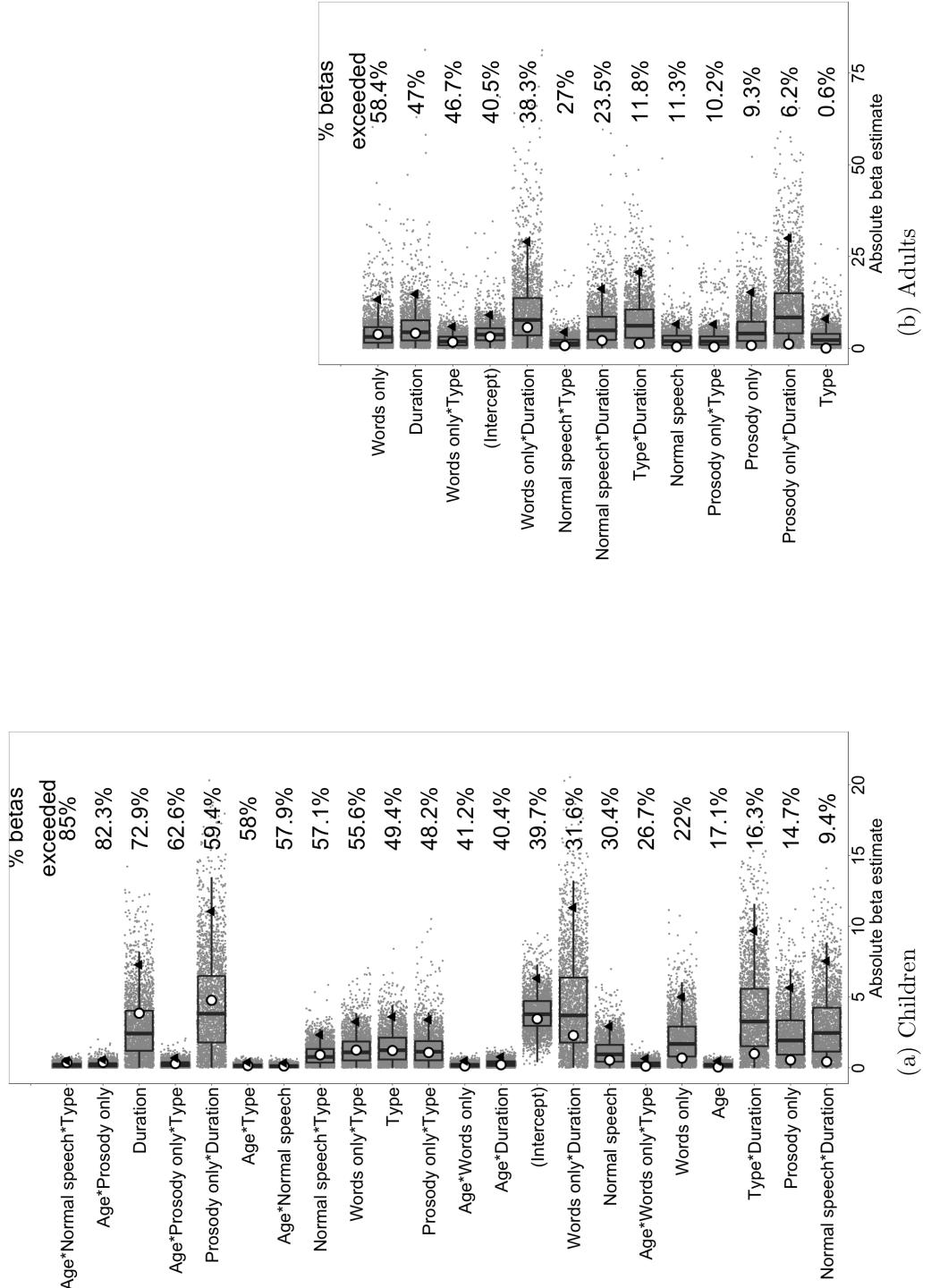


Figure A.5: Random-permutation and original  $|\beta\text{-values}|$  for predictors of anticipatory gaze rates in Experiment 2.

## Experiment 2: SE estimates

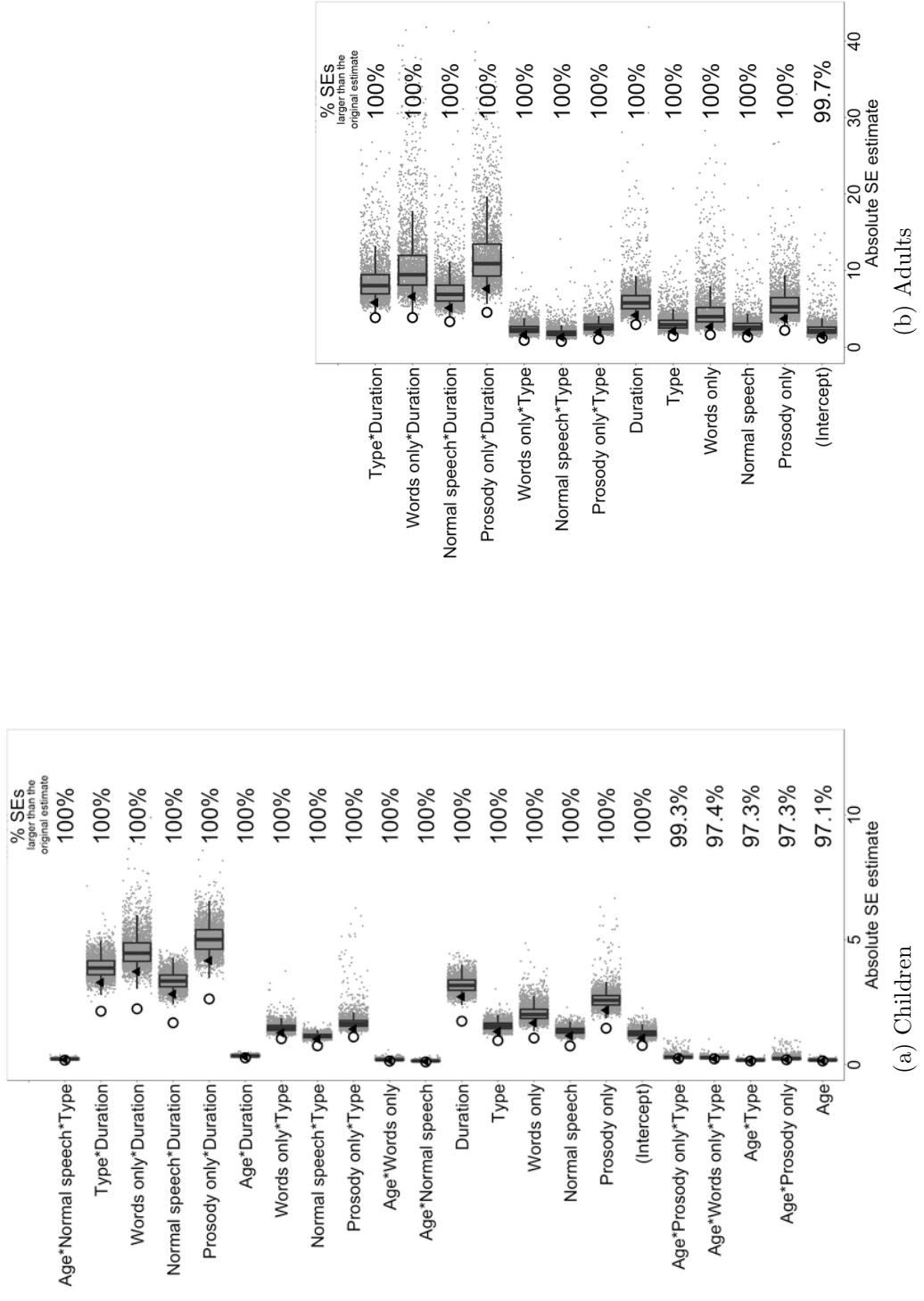


Figure A.6: Random-permutation and original SE-values for predictors of anticipatory gaze rates in Experiment 2.

1499 *Appendix A.2. Non-convergent models*

1500 In comparing the real and randomly permuted datasets, we excluded the  
1501 output of random-permutation models that gave convergence warnings to  
1502 remove erratic model estimates from our analyses. Non-convergent models  
1503 made up 22.4–24.4% of the random permutation models in Experiment 1 and  
1504 69–70% of the random permutation models in Experiment 2. The  $z$ -values for  
1505 each predictor in the converging and non-converging models from Experiment  
1506 1 are shown in Table A.1.

1507 Although many of the non-converging models show estimates within range  
1508 of the converging models (e.g., with a mean difference of only 0.09 in median  
1509  $z$ -value across predictors), they also show many radically outlying estimates  
1510 (e.g., showing a mean difference of 237.3 in mean  $z$ -value across predictors).  
1511 Similar patterns were obtained in the non-converging models for Experiment  
1512 2 and persisted across multiple attempts with different optimizers.

1513 We suspect that the issue derives from data sparsity in some of the ran-  
1514 dom permutations. This problem is known to occur when there are limited  
1515 numbers of binary observations in each of a design matrix’s bins (Allison,  
1516 2004). We could instead use zero-inflated poisson or negative binomial re-  
1517 gression models to allow for overdispersion in our data (Allison, 2012). How-  
1518 ever, these would give us baselines for the normal, convergent model, which  
1519 is not the aim of this analysis.



1520 **Appendix B. Pairwise developmental tests**

1521 Experiments 1 and 2 both showed effects of age in interaction with lin-  
1522 guistic condition and transition type. To explore these effects in more depth,  
1523 in each permutation we recorded the average difference score for each par-  
1524 ticipant, for each predictor that interacted with age (e.g., English minus  
1525 non-English anticipatory switches for each participant). We then used these  
1526 values to compute an average difference score over the participants in each  
1527 age group (e.g., age 3, 4, and 5) within each random permutation. This  
1528 averaging process produces 5,000 baseline-derived difference scores for each  
1529 age group.

1530 We then made pairwise age comparisons of the difference scores (e.g.,  
1531 the linguistic condition effect in 3-year-olds vs. 4-year-olds), computing the  
1532 percent of random-permutation difference scores exceeded by the real-data  
1533 difference score. If the real-data difference score exceeded 95% of the random-  
1534 data age difference scores, we deemed it to be an age effect significantly  
1535 different from chance, e.g., a significant difference between ages three and  
1536 four in the effect of linguistic condition. This procedure is essentially a two-  
1537 tailed  $t$ -test, adapted for use with the randomly permuted baseline data.

1538 In each of the plots below, the black dot represents the real data value  
1539 for the effect being shown. The effect sizes from the 5,000 randomly per-  
1540 muted data sets are shown as a distribution. The percentage displayed is the  
1541 percentage of random permutation difference scores exceeded by the original  
1542 data differences score (taking the absolute value of all data points for a two-  
1543 tailed test). Comparisons marked with 95% or higher are significant at the  
1544  $p < 0.05$  level.

Experiment 1: Age and linguistic condition

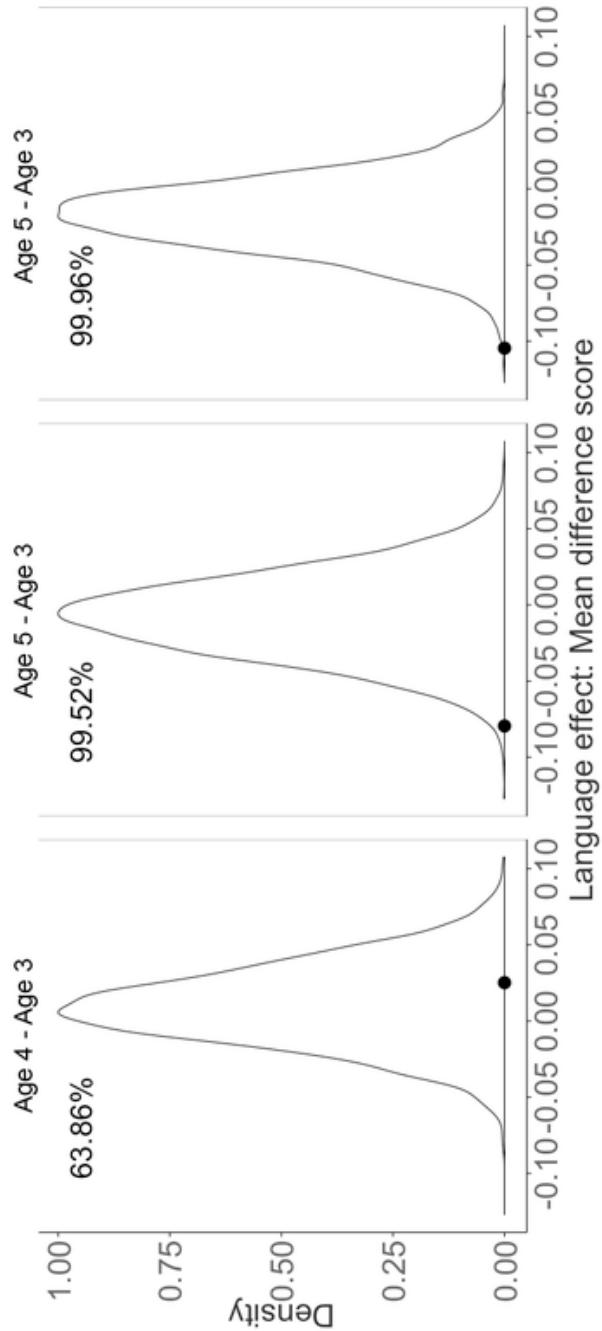


Figure B.1: Pairwise comparisons of the language condition effect across ages in Experiment 1.

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

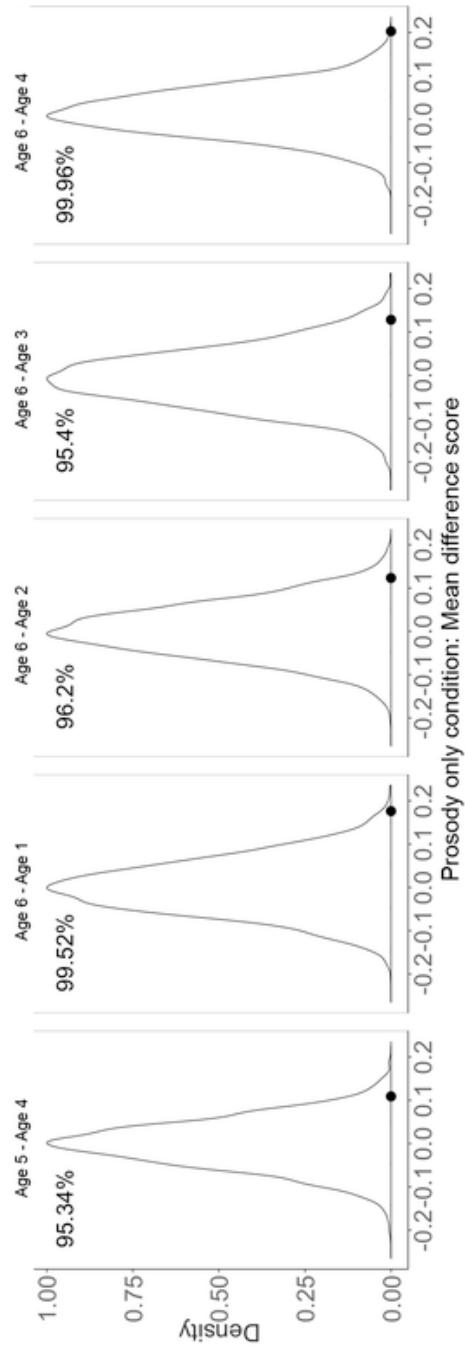


Figure B.2: Significant pairwise comparisons of the *prosody only-no speech* linguistic condition effect, across ages in Experiment 2. Non-significant comparisons are not shown here.

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

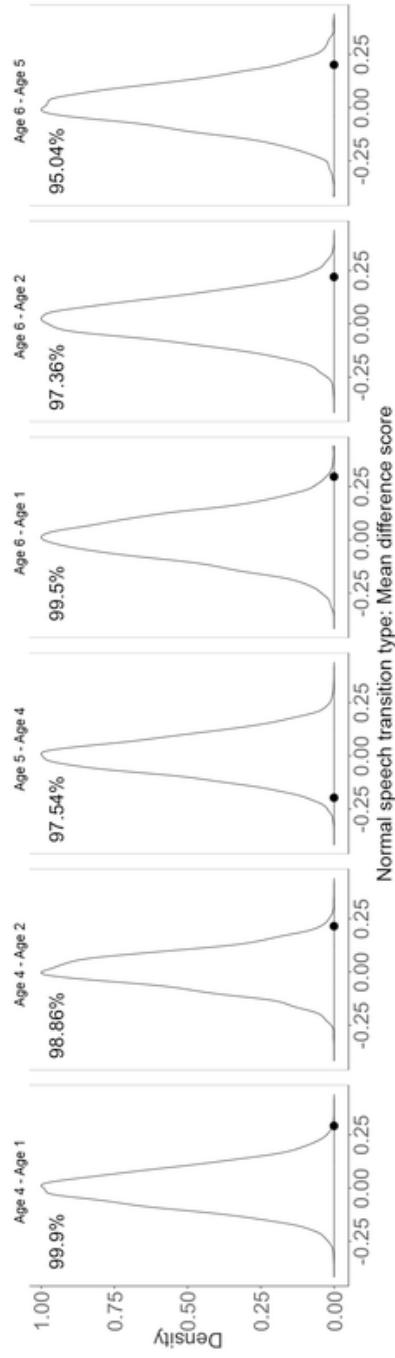


Figure B.3: Significant pairwise comparisons of the *normal speech-no speech* language condition effect for transition type, across ages, in Experiment 2. Non-significant comparisons are not shown here.

1545 **Appendix C. Boredom-driven anticipatory looking**

1546 One alternative hypothesis for children’s anticipatory gazes is that they  
1547 look at the current speaker at the start of each turn, but then grow bored  
1548 and start looking away at a constant rate. Even though this alternative  
1549 hypothesis does not predict the primary effects in our data (e.g., the difference  
1550 between questions and non-questions), we cannot rule out the possibility that  
1551 a portion of participants’ saccades come from boredom.

1552 The data plotted here show a hypothetical group of boredom-driven par-  
1553 ticipants (gray dots) and participants from the actual data in Experiment 2  
1554 (black dots). The hypothetical boredom-driven participants look away from  
1555 the current speaker at a linear rate, beginning one second after the start of  
1556 a turn.

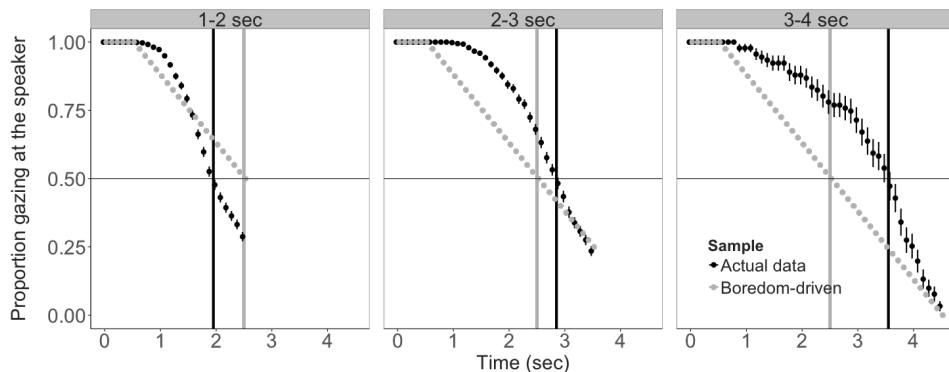


Figure C.1: Proportion of participants (hypothetical boredom-driven=gray; actual Ex-  
periment 2=black) looking at the current speaker, split by turn duration. Vertical bars  
indicate standard error in the experimental data.

1557 If children’s switches away from the current speaker were purely driven  
1558 by boredom, they would switch away equally quickly on long and short turns.  
1559 Therefore, their crossover point—the point in time at which 50% of the chil-  
1560 dren have switched away from the current speaker—would be the same for  
1561 all turns, no matter the length of the turn. This pattern is demonstrated  
1562 in the hypothetical boredom-driven crossover points, which always occur 2.5  
1563 seconds after the start of speech (gray vertical lines; Figure C.1).

1564 In children’s *actual* looking data we see that crossover points increase with  
1565 turn duration: 2.0, 2.9, and 3.6 seconds after the start of speech for turns

1566 with durations of 1–2, 2–3, and 3–4 seconds, respectively (black vertical lines;  
1567 Figure C.1). This pattern suggests that, though children do look away as  
1568 the turn is unfolding, their looks away are not simply driven by boredom.

1569 Are the looks away in Figure C.1 still too early to count as “turn-transition”  
1570 anticipation? It is true that children start looking away after one second  
1571 has passed, but then only gradually. Some of these early looks away may be  
1572 boredom-driven, but it is equally plausible that some of them are turn-driven.  
1573 Early predictive behavior is common in turn-taking studies with adults,  
1574 in both constrained turn-taking tasks (De Ruiter et al., 2006; Gísladóttir  
1575 et al., 2015; Bögels et al., 2015) and in spontaneous conversation (Holler &  
1576 Kendrick, 2015; Torreira et al., 2015). Although this same pattern has yet to  
1577 be established for children’s turn predictions, the looking behavior here is at  
1578 least consistent with adult response patterns in previous work. Additionally,  
1579 because our analysis windows in the main study only overlapped with the  
1580 pre-gap utterance by 300 msec (Figure 2), our primary results are unlikely to  
1581 capture any of these very early or early boredom-driven gaze switches, which  
1582 makes them unproblematic either way in the current analysis.

1583 We therefore conclude that the boredom-driven effects in our data are  
1584 unlikely to change our primary results, though we acknowledge that characterizing  
1585 different gaze switching strategies in this kind of data is an important  
1586 avenue for future work.

1587 **Appendix D. Puppet pair and linguistic condition**

1588     The design for Experiment 2 does not fully cross puppet pair (e.g., robots,  
1589     blue puppets) with linguistic condition (e.g., *words only* and *no speech*). Even  
1590     though each puppet pair is associated with different conversation clips across  
1591     children (e.g., robots talking about kitties, birthday parties, and pancakes),  
1592     the robot puppets themselves were exclusively associated with the *words only*  
1593     condition. Similarly, merpeople were exclusively associated with *prosody only*  
1594     speech, and the puppets wearing dressy clothes were exclusively associated  
1595     with the *no speech* condition. We designed the experiment this way to in-  
1596     crease its pragmatic felicity for older children (i.e., robots make robot sounds,  
1597     merpeople’s voices are muffled under the water, the party-going puppets are  
1598     in a ‘party’ room with many other voices). There is therefore a confound  
1599     between linguistic condition and puppet pair; for example, children could  
1600     have made fewer anticipatory switches in the *prosody only* condition because  
1601     the puppets were less interesting. To test whether puppet pair drove the  
1602     condition-based differences found in Experiment 2, we ran a short follow-up  
1603     study.

1604 **Methods**

1605 We recruited 30 children between ages 3;0 and 5;11 from the Children’s Dis-  
1606 covery Museum of San Jose, California to participate in our experiment. All  
1607 participants were native English speakers. Children were randomly assigned  
1608 to one of six videos (five children per video).

1609 *Materials.* We created 6 short videos from the stimulus recordings made for  
1610 Experiment 2. Each video featured a puppet pair (red/blue/yellow/robot/  
1611 merpeople/party-goer; Figure 5). Puppets in all six videos performed the  
1612 exact same conversation recording (‘birthday party’; Experiment 2) with  
1613 normal, unmanipulated speech. This experiment therefore holds all things  
1614 constant across stimuli except for the appearance of the puppets.

1615 *Procedure.* We used the same experimental apparatus and procedure as in  
1616 Experiments 1 and 2. Each participant was randomly assigned to watch only  
1617 one of the six puppet videos. Five children watched each video. As in Experi-  
1618 ment 2, the experimenter immediately began each session with calibration  
1619 and then stimulus presentation because no special instructions were required.  
1620 The entire experiment took less than three minutes.

1621 *Data preparation.* We identified anticipatory gaze switches to the upcoming  
1622 speaker using the same method as in Experiments 1 and 2.

1623 **Results and discussion**

1624 We modeled children’s anticipatory switches (yes or no at each transition)  
1625 with mixed effects logistic regression, including puppet pair (robots/mer-  
1626 people/party-goers/other-3) as a fixed effect and participant and turn tran-  
1627 sition as random effects. We grouped the red, blue, and yellow puppets  
1628 together because they collectively represented the puppets used in the *nor-*  
1629 *mal* speech condition—this follow-up experiment is meant to test whether  
1630 the condition-based differences from Experiment 2 arose from the puppets  
1631 used in each condition.

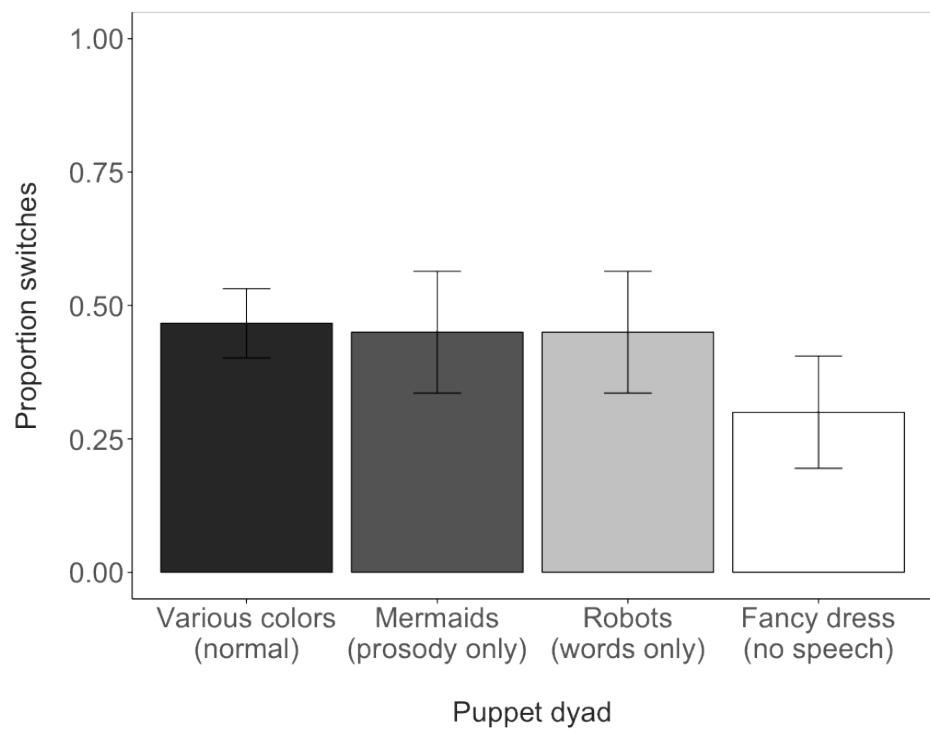


Figure D.1: Proportion gaze switches across puppet pairs when linguistic condition and conversation are held constant.

	Estimate	Std. Error	<i>z</i> value	Pr(>  <i>z</i>  )
<i>Reference level: normal-condition puppets</i>				
(Intercept)	-0.148	0.328	-0.451	0.652
Puppets= <i>mermaid</i>	-0.076	0.655	-0.116	0.908
Puppets= <i>robot</i>	-0.071	0.653	-0.109	0.913
Puppets= <i>party</i>	-0.782	0.687	-1.138	0.255
<i>Reference level: mer-puppets</i>				
(Intercept)	-0.224	0.568	-0.394	0.694
Puppets= <i>robot</i>	0.0048	0.801	0.006	0.995
Puppets= <i>party</i>	-0.706	0.827	-0.854	0.393
<i>Reference level: robot puppets</i>				
(Intercept)	-0.219	0.566	-0.387	0.699
Puppets= <i>party</i>	-0.711	0.827	-0.860	0.390
<i>Reference level: party-goer puppets</i>				
(Intercept)	-0.93	0.607	-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.

1632 In four versions of this model, we systematically varied the reference level  
 1633 of the puppet pair to check for any cross-condition differences. We found no  
 1634 significant effects of puppet pair on switching rate (all  $p > 0.25$ ; Table D.2).

1635 We take this finding as evidence that our decision to not fully cross puppet  
 1636 pairs and linguistic conditions in Experiment 2 was unlikely to have affected  
 1637 children’s anticipatory gaze rates above and beyond the intended effects of  
 1638 linguistic condition.