

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

Marisa Casillas^{a,*}, Michael C. Frank^b

^a*Max Planck Institute for Psycholinguistics, Nijmegen*

^b*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. Like adults, prosody alone did not improve children's predictive gaze shifts. But, unlike adults, lexical cues alone were not sufficient to improve prediction—children's performance was best overall with lexicosyntax and prosody together. Our findings support an account in which turn prediction emerges in infancy, but becomes fully integrated with linguistic processing only gradually.

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

¹ 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

*Corresponding author

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

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

30 *1.1. Adults’ turn taking*

31 Turn taking itself is not unique to conversation. Many other human activ-
32 ities are organized around sequential turns at action. Traffic intersections and
33 computer network communication both use turn-taking systems. Children’s
34 early games (e.g., give-and-take, peek-a-boo) have built-in, predictable turn
35 structure (Ratner and Bruner, 1978; Ross and Lollis, 1987). Even monkeys
36 take turns: non-human primates such as marmosets and Campbell’s monkeys
37 vocalize contingently with each other in both natural and lab-controlled en-
38 vironments (Lemasson et al., 2011; Takahashi et al., 2013). In all these cases,
39 turn taking serves as a protocol for interaction, allowing the participants to
40 coordinate with each other through sequences of contingent action.

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

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

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

72 Prior work suggests that adult listeners primarily use lexicosyntactic in-
73 formation to accurately predict upcoming turn structure (De Ruiter et al.,
74 2006). De Ruiter and colleagues (2006) asked participants to listen to snip-
75 pets of spontaneous conversation and to press a button whenever they antici-
76 pated that the current speaker was about to finish his or her turn. The speech
77 snippets were controlled for the amount of linguistic information present;
78 some were normal, but others had flattened pitch, low-pass filtered speech,

79 or further manipulations. With pitch-flattened speech, the timing of par-
80 ticipants' button responses was comparable to their timing with the full
81 linguistic signal. But when no lexical information was available, partici-
82 pants' responses were significantly earlier. The authors concluded that lex-
83 icosyntactic information¹ was necessary and possibly sufficient for turn-end
84 projection, while intonation was neither necessary nor sufficient. Congru-
85 ent evidence comes from studies varying the predictability of lexicosyntactic
86 and pragmatic content: adults anticipate turn ends better when they can
87 more accurately predict the exact words that will come next (Magyari and
88 De Ruiter, 2012; see also Magyari et al., 2014). They can also identify speech
89 acts within the first word of an utterance (Gísladóttir et al., 2015), allowing
90 them to start planning their response at the first moment possible (Bögels
91 et al., 2015).

92 Despite this body of evidence, the role of prosody for adult turn predic-
93 tion is still a matter of debate. De Ruiter and colleagues' (2006) experiment
94 focused on the role of intonation, which is only a partial index of prosody.
95 Prosody is tied closely to the syntax of an utterance, so the two linguistic
96 signals are difficult to control independently (Ford and Thompson, 1996).
97 Torreira, Bögels and Levinson (2015) used a combination of button-press
98 and verbal responses to investigate the relationship between lexicosyntac-
99 tic and prosodic cues in turn-end prediction. Critically, their stimuli were
100 cross-spliced so that each item had full prosodic cues to accompany the lex-
101 icosyntax. Because of the splicing, they were able to create items that had
102 syntactically-complete units with no intonational phrase boundary at the
103 end. Participants never verbally responded or pressed the “turn-end” but-
104 ton when hearing a syntactically-complete phrase without an intonational
105 phrase boundary. And when intonational phrase boundaries were embedded
106 in multi-utterance turns, participants were tricked into pressing the “turn-
107 end” button 29% of the time. Their results suggest that listeners actually
108 do rely on prosodic cues to execute a response (see also de De Ruiter et al.
109 (2006):525). These experimental findings corroborate other corpus and ex-
110 perimental work promoting a combination of cues (lexicosyntactic, prosodic,
111 and pragmatic) as key for accurate turn-end prediction (Duncan, 1972; Ford

¹The “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.

¹¹² and Thompson, 1996; Hirvenkari et al., 2013).

¹¹³ *1.2. Children’s turn prediction*

¹¹⁴ The majority of work on children’s early turn taking has focused on ob-
¹¹⁵ servations of spontaneous interaction. Children’s first turn-like structures
¹¹⁶ appear as early as two to three months after birth, in proto-conversation with
¹¹⁷ their caregivers (Bruner, 1975, 1985). During proto-conversations, caregivers
¹¹⁸ treat their infants as capable of making meaningful contributions: they take
¹¹⁹ every look, vocalization, arm flail, and burp as “utterances” in the joint dis-
¹²⁰ course (Bateson, 1975; Jaffe et al., 2001; Snow, 1977). Infants catch onto the
¹²¹ structure of proto-conversations quickly. By three to four months they notice
¹²² disturbances to the contingency of their caregivers’ response and, in reaction,
¹²³ change the rate and quality of their vocalizations (Bloom, 1988; Masataka,
¹²⁴ 1993).

¹²⁵ The timing of children’s responses to their caregivers’ speech shows a non-
¹²⁶ linear pattern. Infants’ contingent vocalizations in the first few months of
¹²⁷ life show very fast timing (though with a lot of vocal overlap). But by nine
¹²⁸ months, their timing slows down considerably, only to gradually speed up
¹²⁹ again after 12 months (Hilbrink et al., 2015). Taking turns with brief transi-
¹³⁰ tions between speakers is difficult for children; while their avoidance of over-
¹³¹ lap is nearly adult-like by nine months, the timing of their non-overlapped
¹³² responses stays much longer than the adult 200 msec standard for the next
¹³³ few years (Casillas et al., In press; Garvey, 1984; Garvey and Berninger,
¹³⁴ 1981; Ervin-Tripp, 1979). This puzzling pattern is likely due to their linguis-
¹³⁵ tic development: taking turns on time is easier when the response is a simple
¹³⁶ vocalization rather than a linguistic utterance. Integrating language into the
¹³⁷ turn-taking system may be one major factor in children’s delayed responses
¹³⁸ (Casillas et al., In press).

¹³⁹ While children, like adults, might use linguistic cues in the ongoing turn
¹⁴⁰ to make predictions about upcoming turn structure, studies of early linguistic
¹⁴¹ development point to a possible early advantage for prosody over lexicosyn-
¹⁴² tax in children’s turn-taking predictions. Infants can distinguish their native
¹⁴³ language’s rhythm type from others soon after birth (Mehler et al., 1988;
¹⁴⁴ Nazzi and Ramus, 2003); they show preference for the typical stress patterns
¹⁴⁵ of their native language over others by 6–9 months (e.g., iambic vs. trochaic),
¹⁴⁶ and can use prosodic information to segment the speech stream into smaller
¹⁴⁷ chunks from 8 months onward (Johnson and Jusczyk, 2001; Morgan and
¹⁴⁸ Saffran, 1995). Four- to five-month-olds also prefer pauses in speech to be

149 inserted at prosodic boundaries, and by 6 months infants can use prosodic
150 markers to pick out sub-clausal syntactic units, both of which are useful for
151 extracting turn structure from ongoing speech (Jusczyk et al., 1995; Soder-
152 strom et al., 2003). In comparison, children show at best a very limited
153 lexical inventory before their first birthday (Bergelson and Swingley, 2013;
154 Shi and Melancon, 2010).

155 Keitel and colleagues (2013) were one of the first to explore how children
156 use linguistic cues to predict upcoming turn structure. They asked 6-, 12-,
157 24-, and 36-month-old infants, and adult participants to watch short videos
158 of conversation and tracked their eye movements at points of speaker change.
159 They showed their participants two types of conversation videos—one normal
160 and one with flattened pitch (i.e., with flattened intonation contours)—to test
161 the role of intonation in participants’ anticipatory predictions about upcom-
162 ing speech. Comparing children’s anticipatory gaze frequency to a random
163 baseline, they found that only 36-month-olds and adults made anticipatory
164 gaze switches more often than expected by chance. Among those, only 36-
165 month-olds were affected by a lack of intonation contours, leading Keitel
166 and colleagues to conclude that children’s ability to predict upcoming turn
167 structure relies on their ability to comprehend the stimuli lexicosemantically.
168 They also suggest that intonation might play a secondary role in turn predic-
169 tion, but only after children acquire more sophisticated, adult-like language
170 comprehension abilities (also see Keitel and Daum, 2015).

171 Although the Keitel et al. (2013) study constitutes a substantial ad-
172 vance over previous work in this domain, it has some limitations. Because
173 these limitations directly inform our own study design, we review them in
174 some detail. First, their estimates of baseline gaze frequency (“random” in
175 their terminology) were not random. Instead, they used gaze switches dur-
176 ing ongoing speech as a baseline. But ongoing speech is the period in which
177 switching is least likely to occur (Hirvenkari et al., 2013)—their baseline thus
178 maximizes the chance of finding a difference in gaze frequency at turn transi-
179 tions compared to the baseline. A more conservative baseline would compare
180 participants’ looking behavior at turn transitions to their looking behavior
181 during randomly selected windows of time throughout the stimulus, includ-
182 ing turn transitions. We follow this conservative approach in the current
183 study.

184 Second, the conversation stimuli Keitel et al. (2013) used were somewhat
185 unusual. The average gap between turns was 900 msec, which is much longer
186 than typical adult timing, which averages around 200 msec (Stivers et al.,

2009). The speakers in the videos were also asked to minimize their movements while performing scripted, adult-directed conversation, which would have created a somewhat unnatural interaction. Additionally, to produce more naturalistic conversation, it would have been ideal to localize the sound sources for the two voices in the video (i.e., to have the voices come out of separate left and right speakers). But both voices were recorded and played back on the same audio channel, which may have made it difficult to distinguish the two talkers. Again, we attempt to address these issues in our current study. Despite these minor methodological issues, the Keitel et al. (2013) study still demonstrates intriguing age-based differences in children's ability to predict upcoming turn structure. Our current work thus takes this paradigm as a starting point.²

199 *1.3. The current study*

200 Our goal in the current study is to find out when children begin to make predictions about upcoming turn structure and to understand how their predictions are affected by linguistic cues across development. We present two experiments in which we measured children's anticipatory gaze to responders 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.

210 We highlight three primary findings: first, although children and adults use linguistic cues to make predictions about upcoming turn structure, they do so primarily in 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. Third, we find no evidence of an early prosody advantage in children's anticipations and, further, no evidence that prosodic or lexical cues alone can substitute

²But also see Casillas and Frank (2012, 2013).

221 for their combination in the full linguistic signal (as is proposed for adults;
222 De Ruiter et al., 2006); instead, anticipation is strongest for stimuli with the
223 full range of cues. In sum, our findings support an account in which turn
224 prediction emerges in infancy, but becomes fully integrated with linguistic
225 processing only gradually across development.

226 **2. Experiment 1**

227 We recorded participants' eye movements as they watched six short videos
228 of two-person (dyadic) conversation interspersed with attention-getting filler
229 videos. Each conversation video featured an improvised discourse in one
230 of five languages (English, German, Hebrew, Japanese, and Korean); par-
231 ticipants saw two videos in English and one in every other language. The
232 participants, all native English speakers, were only expected to understand
233 the two videos in English. We showed participants non-English videos to
234 limit their access to lexical information while maintaining their access to
235 other cues to turn boundaries (e.g., (non-native) prosody, gaze, inbreaths,
236 phrase final lengthening). Using this method, we compared children and
237 adult's anticipatory looks from the current speaker to the upcoming speaker
238 at points of turn transition in English and non-English videos.

239 *2.1. Methods*

240 *2.1.1. Participants*

241 We recruited 74 children between ages 3;0–5;11 and 11 undergraduate
242 adults to participate in the experiment. Our child sample included 19 three-
243 year-olds, 32 four-year-olds, and 23 five-year-olds, all enrolled in a local nurs-
244 ery school. All participants were native English speakers. Approximately
245 one-third ($N=25$) of the children's parents and teachers reported that their
246 child regularly heard a second (and sometimes third or further) language, but
247 only one child frequently heard a language that was used in our non-English
248 video stimuli, and we excluded his data from analyses. None of the adult
249 participants reported fluency in a second language.

250 *2.1.2. Materials*

251 *Video recordings.* We recorded pairs of talkers while they conversed in
252 a sound-attenuated booth (see a sample frame in Figure 1). Each talker
253 was a native speaker of the language being recorded, and each talker pair
254 was male-female. Using a Marantz PMD 660 solid state field recorder, we



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

255 captured audio from two lapel microphones, one attached to each participant,
256 while simultaneously recording video from the built-in camera of a MacBook
257 laptop computer. The talkers were volunteers and were acquainted with their
258 recording partner ahead of time.

259 Each recording session began with a 20-minute warm-up period of sponta-
260 neous conversation during which the pair talked for five minutes on four
261 topics (favorite foods, entertainment, hometown layout, and pets). Then we
262 asked talkers to choose a new topic—one relevant to young children (e.g.,
263 riding a bike, eating breakfast)—and to improvise a dialogue on that topic.
264 We asked them to speak as if they were on a children’s television show in
265 order to elicit child-directed speech toward each other. We recorded until the
266 talkers achieved at least 30 seconds of uninterrupted discourse with enthu-
267 siastic, child-directed speech. Most talker pairs took less than five minutes
268 to complete the task, usually by agreeing on a rough script at the start. We
269 encouraged talkers to ask at least a few questions to each other during the
270 improvisation. The resulting conversations were therefore not entirely spon-
271 taneous, but were as close as possible while still remaining child-oriented in
272 topic, prosodic pattern, and lexicosyntactic construction.³

273 After recording, we combined the audio and video recordings by hand, and
274 cropped each recording to the 30-second interval with the most turn activity.

³All 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.

275 Because we recorded the conversations in stereo, the male and female voices
276 came out of separate speakers during video playback. This gave each voice in
277 the videos a localized source (from the left or right loudspeaker). We coded
278 each turn transition in the videos for language condition (English vs. non-
279 English), inter-turn gap duration (in milliseconds), and speech act (question
280 vs. non-question). The non-English stimuli were coded for speech act from
281 a monolingual English-speaker’s perspective, i.e., which turns “sound like”
282 questions, and which do not: we asked five native American English speakers
283 to listen to the audio recording for each turn and judge whether it sounded
284 like a question. We then coded turns with at least 80% “yes” responses as
285 questions.

286 Because the conversational stimuli were recorded semi-spontaneously, the
287 duration of turn transitions and the number of speaker transitions in each
288 video was variable. We measured the duration of each turn transition from
289 the audio recording associated with each video. We excluded turn transi-
290 tions longer than 550 msec and shorter than 90 msec, also excluding over-
291 lapped transitions, from analysis.⁴ This left approximately equal numbers
292 of turn transitions available for analysis in the English (N=20) and non-
293 English (N=16) videos. On average, the inter-turn gaps for English videos
294 (mean=318, median=302, stdev=112 msec) were slightly longer than for non-
295 English videos (mean=286, median=251, stdev=122 msec). The longer gaps
296 in the English videos could give them a slight advantage: our definition of
297 an “anticipatory gaze shift” includes shifts that are initiated during the gap
298 between turns (Figure 2), so participants had slightly more time to make
299 anticipatory shifts in the English videos.

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

⁴Overlap 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.

306 median=252, stdev=134 msec).

307 *2.1.3. Procedure*

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

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

324 Children were less inclined to simply sit and watch videos of conversation
325 in languages they didn't speak, so we used a different procedure to keep them
326 engaged: the experimenter started each session by asking the child about
327 what languages he or she could speak, and about what other languages he
328 or she had heard of. Then the experimenter expressed her own enthusiasm
329 for learning about new languages, and invited the child to watch a video
330 about "new and different languages" together. If the child agreed to watch,
331 the experimenter and the child sat together in front of the display, with
332 the child centered in front of the tracker and the experimenter off to the
333 side. Each conversation video was preceded and followed by a 15–30 second
334 attention-getting filler video (e.g., running puppies, singing muppets, flying
335 bugs). If the child began to look bored, the experimenter would talk during
336 the fillers, either commenting on the previous conversation ("That was a neat
337 language!") or giving the language name for the next conversation ("This
338 next one is called Hebrew. Let's see what it's like.") The experimenter's
339 comments reinforced the video-watching as a joint task.

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

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

348 *2.1.4. Data preparation and coding*

349 To determine whether participants predicted upcoming turn transitions,
350 we needed to define a set of criteria for what counted as an anticipatory gaze
351 shift. Prior work using similar experimental procedures has found that adults
352 and children make anticipatory gaze shifts to upcoming talkers within a wide
353 time frame; the earliest shifts occur before the end of the prior turn, and the
354 latest occur after the onset of the response turn, with most shifts occurring
355 in the inter-turn gap (Keitel et al., 2013; Hirvenkari, 2013; Tice and Henetz,
356 2011). Following prior work, we measured how often our participants shifted
357 their gaze from the prior to the upcoming speaker *before* the shift in gaze
358 could have been initiated in reaction to the onset of the speaker’s response.
359 In doing so, we assumed that it takes participants 200 msec to plan an eye
360 movement, following standards from adult anticipatory processing studies
361 (e.g., Kamide et al., 2003).

362 We checked each participant’s gaze at each turn transition for three char-
363 acteristics (Figure 2): (1) that the participant fixated on the prior speaker for
364 at least 100 msec at the end of the prior turn, (2) that sometime thereafter
365 the participant switched to fixate on the upcoming speaker for at least 100
366 ms, and (3) that the switch in gaze was initiated within the first 200 msec of
367 the response turn, or earlier. These criteria guarantee that we only counted
368 gaze shifts when: (1) participants were tracking the previous speaker, (2)
369 switched their gaze to track the upcoming speaker, and (3) did so before
370 they could have simply reacted to the onset of speech in the response. Under
371 this assumption, a gaze shift that was initiated within the first 200 msec of
372 the response (or earlier) was planned *before* the child could react to the onset
373 of speech itself.

374 As mentioned, most anticipatory switches happen in the inter-turn gap,
375 but we also allowed anticipatory gaze switches that occurred in the final
376 syllables of the prior turn. Early switches are consistent with the distribution
377 of responses in explicit turn-boundary prediction tasks. For example, in
378 a button press task, adult participants anticipate turn ends approximately
379 200 msec in advance of the turn’s end, and anticipatory responses to pitch-

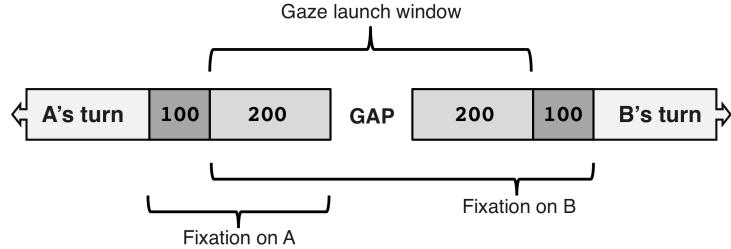


Figure 2: Schematic summary of criteria for anticipatory gaze shifts from speaker A to speaker B during a turn transition.

380 flattened stimuli come even earlier (De Ruiter et al., 2006). We therefore
 381 allowed switches to occur as early as 200 msec before the end of the prior turn.
 382 For very early and very late switches, our requirement for 100 msec of fixation
 383 on each speaker would sometimes extend outside of the transition window
 384 boundaries (200 msec before and after the inter-turn gap). The maximally
 385 available fixation window was 100 msec before and after the earliest and
 386 latest possible switch point (300 msec before and after the inter-turn gap).
 387 We did not count switches made during the fixation window as anticipatory.
 388 We *did* count switches made during the inter-turn gap. The period of time
 389 from the beginning of the possible fixation window on the prior speaker to the
 390 end of the possible fixation window on the responder was our total analysis
 391 window (300 msec + the inter-turn gap + 300 msec).

392 *Predictions.* We expected participants to show greater anticipation in the
 393 English videos than in the non-English videos because of their increased
 394 access to linguistic information in English. We also predicted that anticipa-
 395 tion would be greater following questions compared to non-questions; ques-
 396 tions have early cues to upcoming turn transition (e.g., *wh*- words, subject-
 397 auxiliary inversion), and also make a next response immediately relevant.
 398 Our third prediction was that anticipatory looks would increase with devel-
 399 opment, along with children’s increased linguistic competence.

400 *2.2. Results*

401 Participants looked at the screen most of the time during video playback
 402 (81% and 91% on average for children and adults, respectively). They pri-
 403 marily kept their eyes on the person who was currently speaking in both
 404 English and non-English videos: they gazed at the current speaker between

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.

405 38% and 63% of the time, looking back at the addressee between 15% and
 406 20% of the time (Table 1). Even three-year-olds looked more at the current
 407 speaker than anything else, whether the videos were in a language they could
 408 understand or not. Children looked at the current speaker less than adults
 409 did during the non-English videos. Despite this, their looks to the addressee
 410 did not increase substantially in the non-English videos, indicating that their
 411 looks away were probably related to boredom rather than confusion about
 412 ongoing turn structure. Overall, participants' pattern of gaze to current
 413 speakers demonstrated that they performed basic turn tracking during the
 414 videos, regardless of language. Figure 3 shows participants' anticipatory gaze
 415 rates across age, language condition, and transition type.

416 *2.2.1. Statistical models*

417 We identified anticipatory gaze switches for all 36 usable turn transitions,
 418 based on the criteria outlined in Section 2.1.4, and analyzed them for effects
 419 of language, transition type, and age with two mixed-effects logistic regres-
 420 sions (Bates et al., 2014; R Core Team, 2014). We built one model each
 421 for children and adults. We modeled children and adults separately because
 422 effects of age are only pertinent to the children's data. The child model
 423 included condition (English vs. non-English)⁵, transition type (question vs.

⁵Because 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

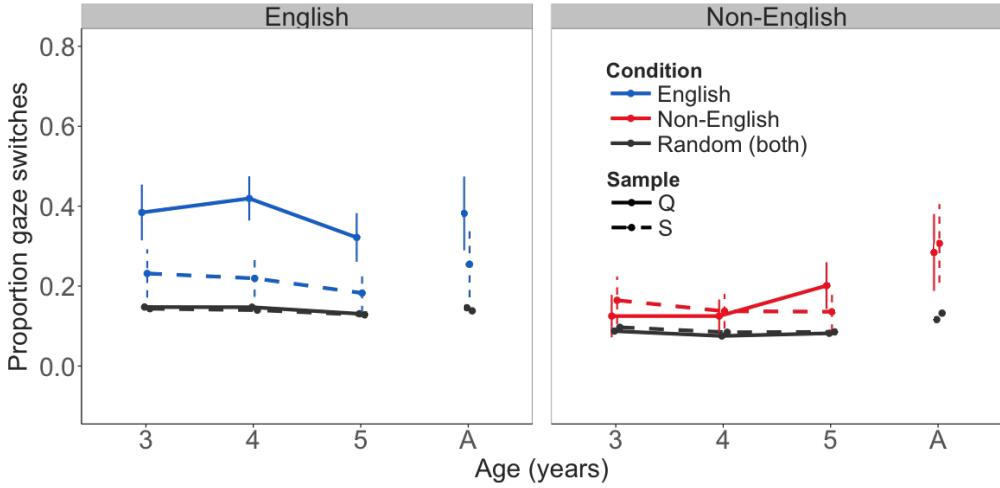


Figure 3: Anticipatory gaze rates across language condition and transition type for the real (red and blue) and randomly permuted baseline (gray). Vertical bars represent 95% confidence intervals.

non-question), age (3, 4, 5; numeric), and duration of the inter-turn gap (seconds, e.g., 0.441) as predictors, with full interactions between condition, transition type, and age. We included the duration of the inter-turn gap as a control predictor since longer gaps provide more opportunities to make anticipatory switches (Figure 2). We additionally included random effects of item (turn transition) and participant, with random slopes of condition, transition type, and their interaction for participants (Barr et al., 2013).⁶ The adult model included condition, transition type, duration, and their interactions as predictors with participant and item included as random effects and random slopes of condition, transition type, and their interaction for participant.

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.

⁶The models we report 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.

Children

	Estimate	Std. Error	<i>z</i> value	Pr(> <i>z</i>)
(Intercept)	-0.96145	0.84915	-1.132	0.257531
Age	-0.18268	0.17509	-1.043	0.296764
LgCond= <i>non-English</i>	-3.29349	0.96055	-3.429	0.000606 ***
Type= <i>non-Question</i>	-1.10131	0.86520	-1.273	0.203055
Duration	3.40171	1.22878	2.768	0.005634 **
Age*LgCond= <i>non-English</i>	0.52066	0.21192	2.457	0.014015 *
Age*TypeS= <i>non-Question</i>	-0.01628	0.19442	-0.084	0.933262
LgCond= <i>non-English</i> *	2.68171	1.35045	1.986	0.047057 *
Type= <i>non-Question</i>				
Age*LgCond= <i>non-English</i> *	-0.45633	0.30168	-1.513	0.130378
Type= <i>non-Question</i>				

Adults

	Estimate	Std. Error	<i>z</i> value	Pr(> <i>z</i>)
(Intercept)	-0.1966	0.6945	-0.283	0.777062
LgCond= <i>non-English</i>	-0.8812	0.9607	-0.917	0.359028
Type= <i>non-Question</i>	-4.4953	1.3147	-3.419	0.000628 ***
Duration	-1.1227	1.9889	-0.565	0.572414
LgCond= <i>non-English</i> *	3.2972	1.6115	2.046	0.040747 *
Type= <i>non-Question</i>				
LgCond= <i>non-English</i> *	1.3625	3.0097	0.453	0.650749
Duration				
Type= <i>non-Question</i> *	10.5107	3.3482	3.139	0.001694 **
Duration				
LgCond= <i>non-English</i> *	-6.3156	4.4969	-1.404	0.160191
Type= <i>non-Question</i> *				
Duration				

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

434 Children's anticipatory gaze switches showed effects of language condition
 435 ($\beta=-3.29$, $SE=0.961$, $z=-3.43$, $p<.001$) and gap duration ($\beta=3.4$, $SE=1.229$,
 436 $z=2.77$, $p<.01$) with additional effects of an age-by-language condition in-
 437 teraction ($\beta=0.52$, $SE=0.212$, $z=2.46$, $p<.05$) and a language condition-by-
 438 transition type interaction ($\beta=2.68$, $SE=1.35$, $z=1.99$, $p<.05$). There were
 439 no significant effects of age or transition type alone ($\beta=-0.18$, $SE=0.175$,

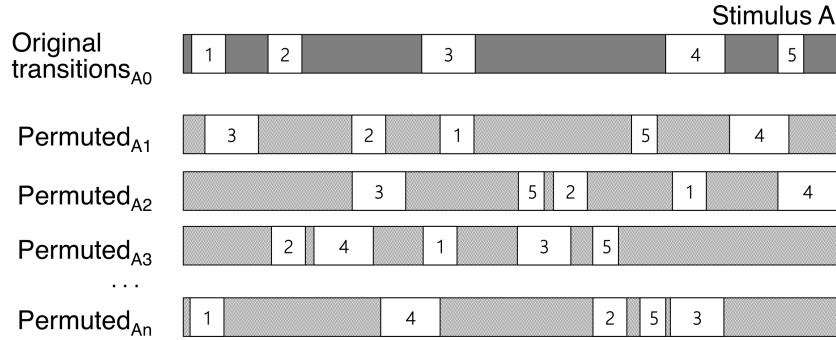


Figure 4: Example of analysis window permutations for a stimulus with five turn transitions. The windows were ± 300 msec around the inter-turn gap.

⁴⁴⁰ $z=-1.04$, $p=.3$ and $\beta=-1.10$, $SE=0.865$, $z=-1.27$, $p=.2$, respectively).

⁴⁴¹ Adults' anticipatory gaze switches showed an effect of transition type
⁴⁴² ($\beta=-4.5$, $SE=1.315$, $z=-3.42$, $p<.001$) and significant interactions between
⁴⁴³ language condition and transition type ($\beta=3.3$, $SE=1.61$, $z=2.05$, $p<.05$)
⁴⁴⁴ and transition type and gap duration ($\beta=10.51$, $SE=3.348$, $z=3.139$, $p<.01$).

⁴⁴⁵ *2.2.2. Random baseline comparison*

⁴⁴⁶ We estimated the probability that these patterns were the result of ran-
⁴⁴⁷ dom looking by running the same regression models on participants' real
⁴⁴⁸ eye-tracking data, only this time calculating their anticipatory gaze switches
⁴⁴⁹ with respect to randomly permuted turn transition windows. This process
⁴⁵⁰ involved: (1) randomizing the order and temporal placement of the anal-
⁴⁵¹ ysis windows within each stimulus (Figure 4; "analysis window" is defined
⁴⁵² in Figure 2) to randomly redistribute the analysis windows across the eye-
⁴⁵³ tracking signal, (2) re-running each participant's eye tracking data through
⁴⁵⁴ switch identification (described in Section 2.1.4) on each of the randomly per-
⁴⁵⁵ mated analysis windows, and (3) modeling the anticipatory switches from the
⁴⁵⁶ randomly permuted data with the same statistical models we used for the
⁴⁵⁷ original data (Section 2.2.1; Table 2). Importantly, although the onset time
⁴⁵⁸ of each transition was shuffled within the eye-tracking signal, the other in-
⁴⁵⁹ trin-
⁴⁶⁰ sic properties of each turn transition (e.g., prior speaker identity, transition
⁴⁶¹ type, gap duration, language condition, etc.) stayed constant across each
⁴⁶² random permutation.

⁴⁶² This procedure effectively de-links participants' gaze data from the turn

463 structure in the original stimulus, thereby allowing us to compare turn-
464 related (original) and non-turn-related (randomly permuted) looking behav-
465 ior using the same eye movement data. The resulting anticipatory gazes from
466 the randomly permuted analysis windows represent an average anticipatory
467 gaze rate over all possible starting points: a random baseline.

468 By running the real and randomly permuted data sets through identical
469 statistical models, we can estimate how likely it is that predictor effects in
470 the original data (e.g., the effect of language condition; Table 2) arose from
471 random looking. Because these analyses are complex, we report their full
472 details in Appendix A.

473 Our baseline analyses revealed that none of the significant predictors
474 from models of the original, turn-related data can be explained by random
475 looking. For the children’s data, the original z -values for language condi-
476 tion, gap duration, the age-language condition interaction, and the language
477 condition-transition type interaction were all greater than 95% of z -values
478 for the randomly permuted data (99.9%, 95.5%, 99.4%, and 96%, respec-
479 tively, all $p < .05$). Similarly, the adults’ data showed significant differen-
480 tiation from the randomly permuted data for two of the three originally
481 significant predictors—transition type and the transition type-gap duration
482 interaction (greater than 99.9% and 99.7% of random z -values, respectively,
483 all $p < .01$)—with marginal differentiation for the interaction of language con-
484 dition and transition type (greater than 94.6% of random z -values; $p = .054$).

485 2.2.3. Developmental effects

486 The models reported above revealed a significant interaction of age and
487 language condition (Table 2) that was unlikely be due to random looking
488 (Figure 3). To further explore this effect, we compared the effect of language
489 condition across age groups: Using the permutation analyses above, we ex-
490 tracted the average difference score for the two language conditions (English
491 minus non-English) for each participant, computing an overall average for
492 each random permutation of the data. Then, within each permutation, we
493 made pairwise comparisons of the average difference scores across participant
494 age groups. This process yielded a distribution of random permutation-based
495 difference scores that we could then compare to the difference score in the
496 actual data. Details are given in Appendix B.

497 These analyses revealed that, while 3- and 4-year olds showed similarly-
498 sized effects of language condition, 5-year-olds had a significantly smaller
499 effect of language condition, compared to both younger age groups. The

500 difference in the language condition effect between 5-year-olds and 3-year-
501 olds was greater than would be expected by chance (99.52% of the randomly
502 permuted data sets; $p < .01$). Similarly, the difference in the language con-
503 dition effect between 5-year-olds and 4-year-olds was greater than would be
504 expected by chance (99.96% of the data sets; $p < .001$). See Figure B.1 for
505 each difference score distribution.

506 When does spontaneous turn prediction emerge developmentally? To
507 test whether the youngest age group (3-year-olds) already exceeded chance
508 in their anticipatory gaze switches, we compared children's real gaze rates
509 to the random baseline in the English condition with two-tailed t -tests.
510 We used the English condition because we are most interested in finding
511 out when children begin to make spontaneous turn predictions for natu-
512 ral speech. We found that three-year-olds made anticipatory gaze switches
513 significantly above chance, when all transitions were considered ($t(22.824) =$
514 4.147 , $p < .001$) as well as for question transitions alone ($t(21.677) = -5.268$,
515 $p < .001$).

516 2.3. Discussion

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

524 First, both adults' and children's anticipations were strongly affected by
525 transition type. Both groups made more anticipatory switches after hear-
526 ing questions, compared to non-questions. Even in the English videos, when
527 participants had full access to linguistic cues, their rates of anticipation were
528 relatively low—comparable to the non-English videos—unless the turn was a
529 question. Prior work using online, metalinguistic tasks has shown that partic-
530 ipants can use linguistic cues to accurately predict upcoming turn ends (Tor-
531 reira et al., 2015; Magyari and De Ruiter, 2012; De Ruiter et al., 2006). The
532 current results add a new dimension to our understanding of how listeners
533 make predictions about turn ends: both children and adults spontaneously
534 monitor the linguistic structure of unfolding turns for cues to imminent re-
535 sponses.

536 Second, children made more anticipatory switches overall in English videos,
537 compared to non-English videos. This effect suggests that lexical access is
538 important for children’s ability to anticipate upcoming turn structure, con-
539 sistent with prior work on turn-end prediction in adults (De Ruiter et al.,
540 2006; Magyari and De Ruiter, 2012) and children (Keitel et al., 2013).

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

554 In Experiment 2 we follow up on these findings, improving on two as-
555 pects of the design: first, our language manipulation in this first experiment
556 was too coarse to provide data regarding specific linguistic channels (e.g.,
557 prosody vs. lexicosyntax). In Experiment 2 we compared lexicosyntactic
558 and prosodic cues with phonetically altered speech and used puppets to elim-
559 inate non-verbal cues to turn taking. Second, we were not able to pinpoint
560 the emergence of anticipatory switching because the youngest age group in
561 our sample was already able to make anticipatory switches at above chance
562 rates. In Experiment 2 we explore a wider developmental range.

563 3. Experiment 2

564 Experiment 2 used native-language stimuli, controlled for lexical and
565 prosodic information, eliminating non-verbal cues, and tested children from a
566 wider age range. To tease apart the role of lexical and prosodic information,
567 we phonetically manipulated the speech signal for pitch, syllable duration,
568 and lexical access. By testing 1- to 6-year-olds we hoped to find the devel-
569 opmental onset of turn-predictive gaze. We also hoped to measure changes
570 in the relative roles of prosody and lexicosyntax across development.

571 Non-verbal cues in Experiment 1 (e.g., gaze and gesture) could have
572 helped participants make predictions about upcoming turn structure (Rossano
573 et al., 2009; Stivers and Rossano, 2010). Since our focus was on linguistic
574 cues, we eliminated all gaze and gestural signals in Experiment 2 by replacing
575 the videos of human actors with videos of puppets. Puppets are less real-
576 istic and expressive than human actors, but they create a natural context
577 for having somewhat motionless talkers in the videos (thereby allowing us
578 to eliminate gestural and gaze cues). Additionally, the prosody-controlled
579 condition included small but global changes to syllable duration that would
580 have required complex video manipulation or precise re-enactment with hu-
581 man talkers, neither of which was feasible. For these reasons, we decided to
582 substitute puppet videos for human videos in the final stimuli.

583 As in the first experiment, we recorded participants' eye movements as
584 they watched six short videos of dyadic conversation, and then analyzed
585 their anticipatory glances from the current speaker to the upcoming speaker
586 at points of turn transition.

587 *3.1. Methods*

588 *3.1.1. Participants*

589 We recruited 27 undergraduate adults and 129 children between ages 1;0–
590 6;11 to participate in our experiment. We recruited our child participants
591 from the Children's Discovery Museum of San Jose, California⁷, targeting
592 approximately 20 children for each of the six one-year age groups (range:
593 20–23). All participants were native English speakers, though some parents
594 ($N=27$) reported that their child heard a second (and sometimes third) lan-
595 guage at home. None of the adult participants reported fluency in a second
596 language.

597 *3.1.2. Materials*

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

⁷We ran Experiment 2 at a local children's museum because it gave us access to children with a more diverse range of ages.

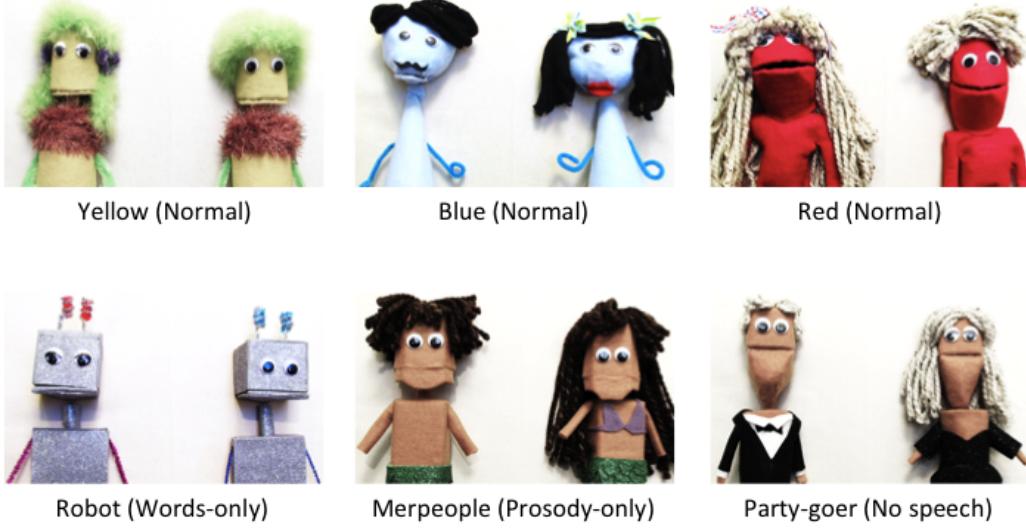


Figure 5: The six puppet pairs (and associated audio conditions). Each pair was linked to three distinct conversations from the same condition across the three experiment versions.

604 *Audio recordings.* The recording session was set up in the same way as
 605 the first experiment, but with a shorter warm up period (5–10 minutes) and
 606 a pre-determined topic for the child-friendly improvisation ('riding bikes',
 607 'pets', 'breakfast', 'birthday cake', 'rainy days', or 'the library'). All of the
 608 talkers were native English speakers, and were recorded in male-female pairs.
 609 As before, we asked talkers to speak "as if they were on a children's television
 610 show" and to ask at least a few questions during the improvisation. We cut
 611 each audio recording down to the 20-second interval with the most turn
 612 activity. The 20-second clips were then phonetically manipulated and used
 613 in the final video stimuli.

614 *Audio Manipulation.* We created four versions of each audio clip: *nor-*
 615 *mal*, *words only*, *prosody only*, and *no speech*. That is, one version with a full
 616 linguistic signal (*normal*), and three with incomplete linguistic information
 617 (hereafter "partial cue" conditions). The *normal* clips were the unmanipu-
 618 lated, original audio clips.

619 The *words only* clips were manipulated to have robot-like speech: we
 620 flattened the intonation contours to each talker's average pitch (F_0) and
 621 we reset the duration of every nucleus and coda to each talker's average

622 nucleus and coda duration.⁸ We made duration and pitch manipulations
623 using PSOLA resynthesis in Praat (Boersma and Weenink, 2012). Thus,
624 the *words only* versions of the audio clips had no pitch or durational cues
625 to upcoming turn boundaries, but did have intact lexicosyntactic cues (and
626 residual phonetic correlates of prosody, e.g., intensity).

627 We created the *prosody only* clips by low-pass filtering the original record-
628 ing at 500 Hz with a 50 Hz Hanning window (following de Ruiter et al., 2006).
629 This manipulation creates a “muffled speech” effect because low-pass filter-
630 ing removes most of the phonetic information used to distinguish between
631 phonemes. The *prosody only* versions of the audio clips lacked lexical infor-
632 mation, but retained their intonational and rhythmic cues to upcoming turn
633 boundaries.

634 The *no speech* condition served as a non-linguistic baseline. For this
635 condition, we replaced the original clip with multi-talker babble: we overlaid
636 different child-oriented conversations (not including the original one), and
637 then cropped the result to the duration of the original video. Thus, the
638 *no speech* audio clips lacked any linguistic information to upcoming turn
639 boundaries—the only cue to turn taking was the opening and closing of the
640 puppets’ mouths.

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

647 *Video recordings.* We created puppet video recordings to match the ma-
648 nipulated 20-second audio clips. The puppets were minimally expressive;
649 so that the puppeteer could only control the opening and closing of their
650 mouths; the puppets’ heads, eyes, arms, and bodies stayed still. Puppets
651 were positioned looking forward to eliminate shared gaze as a cue to turn
652 structure (Thorgrímsson et al., 2015). We took care to match the puppets’
653 mouth movements to the syllable onsets as closely as possible, specifically
654 avoiding any mouth movement before the onset of a turn. We then added
655 the manipulated audio clips to the puppet video recordings by hand.

⁸We excluded hyper-lengthened words like [wau] ‘woooow!’. These were rare in the clips.

656 We used three pairs of puppets used for the *normal* condition—‘red’,
657 ‘blue’ and ‘yellow’—and one pair of puppets for each partial cue condition:
658 “robots”, “merpeople”, and “party-goers” (Figure 8). We randomly assigned
659 half of the conversation topics (‘birthday cake’, ‘pets’, and ‘breakfast’) to the
660 *normal* condition, and half to the partial cue conditions (‘riding bikes’, ‘rainy
661 days’, and ‘the library’). We then created three versions of the experiment,
662 so that each of the six puppet pairs was associated with three different con-
663 versation topics across the different versions of the experiment (18 videos
664 in total). We ensured that the position of the talkers (left and right) was
665 counterbalanced in each version by flipping the video and audio channels as
666 needed.

667 The duration of turn transitions and the number of speaker changes
668 across videos was variable because the conversations were recorded semi-
669 spontaneously. We measured turn transitions from the audio signal of the
670 *normal*, *words only*, and *prosody only* conditions. There was no audio from
671 the original conversation in the *no speech* condition videos, so we measured
672 turn transitions from puppets’ mouth movements in the video signal, using
673 ELAN video annotation software (Wittenburg et al., 2006).

674 There were 85 turn transitions for analysis after excluding transitions
675 longer than 550 msec and shorter than 90 msec. The remaining turn transi-
676 tions had slightly more questions than non-questions ($N=50$ and $N=35$,
677 respectively), with transitions distributed somewhat evenly across condi-
678 tions (keeping in mind that there were three *normal* videos and only one
679 partial cue video for each experiment version): *normal* ($N=36$), *words only*
680 ($N=13$), *prosody only* ($N=17$), and *no speech* ($N=19$). Inter-turn gaps for
681 questions (mean=365, median=427) were longer than those for non-questions
682 (mean=302, median=323) on average, but gap duration was overall com-
683 parable across conditions: *normal* (mean=334, median=321), *words only*
684 (mean=347, median=369), *prosody only* (mean=365, median=369), and *no
685 words* (mean=319, median=329). The longer gaps for question transitions
686 could give them an advantage because our anticipatory measure includes
687 shifts initiated during the gap between turns (Figure 2).

688 3.2. Procedure

689 We used the same experimental apparatus and procedure as in the first
690 experiment. Each participant watched six puppet videos in random order,
691 with five 15–30 second filler videos placed in-between (e.g., running puppies,
692 moving balls, flying bugs). Three of the puppet videos had *normal* audio

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.

693 while the other three had *words only*, *prosody only*, and *no speech* audio.
 694 This experiment required no special instructions so, as before, the exper-
 695 imenter immediately began each session with calibration and then stimulus
 696 presentation. The entire experiment took less than five minutes.

697 *3.2.1. Data preparation and coding*

698 We coded each turn transition for its linguistic condition (*normal*, *words*
 699 *only*, *prosody only*, and *no speech*) and transition type (question/non-question)⁹,
 700 and identified anticipatory gaze switches to the upcoming speaker using the
 701 methods from Experiment 1.

702 *3.3. Results*

703 Participants' pattern of gaze indicated that they performed basic turn
 704 tracking across all ages and in all conditions. Participants looked at the
 705 screen most of the time during video playback (82% and 86% average for
 706 children and adults, respectively), primarily looking at the person who was
 707 currently speaking (Table 2). They tracked the current speaker in every
 708 condition—even one-year-olds looked more at the current speaker than at
 709 anything else in the three partial cue conditions (40% for *words only*, 43%
 710 for *prosody only*, and 39% for *no speech*). There was a steady overall increase
 711 in looks to the current speaker with age and added linguistic information

⁹We coded *wh*-questions as “non-questions” for the *prosody only* videos. Polar questions often have a final rising intonational contour, but *wh*-questions often do not (Hedberg et al., 2010).

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.

(Tables 3 and 4). Looks to the addressee also decreased with age, but the change was minimal. Figure 6 shows participants' anticipatory gaze rates across age, the four language conditions, and transition type.

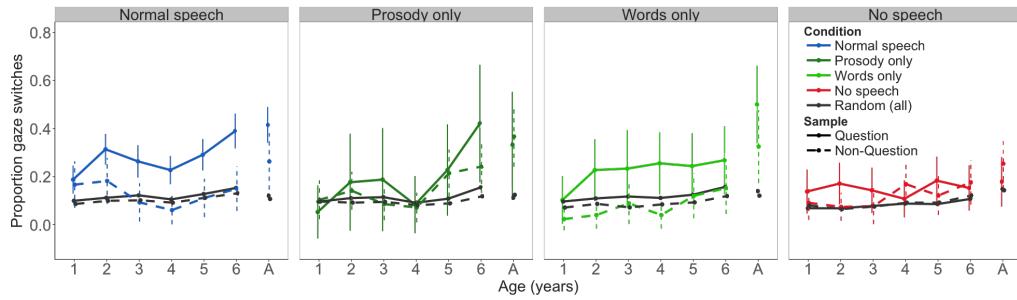


Figure 6: Anticipatory gaze rates across language condition and transition type for the real (blue, dark green, light green, and red) and randomly permuted baseline (gray). Vertical bars represent 95% confidence intervals.

3.3.1. Statistical models

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

language condition, transition type, and age. We again included the duration of the inter-turn gap as a control predictor and added random effects of item (turn transition) and participant, with random slopes of transition type for participants. The adult model included condition, transition type, their interactions, and duration as a control predictor, with participant and item included as random effects and random slopes of condition and transition type.

Children's anticipatory gaze switches showed an effect of gap duration ($\beta=3.95$, $SE=0.617$, $z=6.401$, $p<.001$), a two-way interaction of age and language condition (for *prosody only* speech compared to the *no speech* reference level; $\beta=0.393$, $SE=0.189$, $z=2.08$, $p<.05$), and a three-way interaction of age, transition type, and language condition (for *normal* speech compared to the *no speech* reference level; $\beta=-0.38$, $SE=0.17$, $z=-2.229$, $p<.05$). There were no significant effects of age or transition type alone (Table 3.3.1), with only a marginal effect of language condition (for *prosody only* compared to the *no speech* reference level; $\beta=-1.607$, $SE=0.867$, $z=-1.85$, $p=.064$)

Adults' anticipatory gaze switches showed effects of gap duration ($\beta=4.7$, $SE=1.18$, $z=3.978$, $p<.001$) and language condition (for *normal* speech $\beta=1.203$, $SE=0.536$, $z=2.245$, $p<.05$ and *words only* speech $\beta=1.561$, $SE=0.709$, $z=2.203$, $p<.05$ compared to the *no speech* reference level). There were no effects of transition type ($\beta=0.437$, $SE=0.585$, $z=0.747$, $p=.45$).

3.3.2. Random baseline comparison

Using the same technique described in Experiment 1 (Section 2.2.2), we created and modeled random permutations of participants' anticipatory gaze. These analyses revealed that none of the significant predictors from models of the original, turn-related data could be explained by random looking. In the children's data, the original model's z -values for language condition (*prosody only*), gap duration, the two-way interaction of age and language condition (*prosody only*) and the three-way interaction of age, transition type, and language condition (*normal* speech) were all greater than 95% of the randomly permuted z -values (96.5%, 100%, 95.6%, and 95.5%, respectively, all $p<.05$). Similarly, the adults' data showed significant differentiation from the randomly permuted data for all originally significant predictors: gap duration and language condition for *normal* speech and *words only* speech (greater than 100%, 97.8%, and 99.1% of random z -values, respectively, all $p<.05$).

<i>Children</i>	Estimate	Std. Error	<i>z</i> value	Pr(> <i>z</i>)
(Intercept)	-3.49403	0.48454	-7.211	5.55e-13 ***
Age	0.02436	0.10249	0.238	0.8121
Type= <i>non-Question</i>	-0.88900	0.61192	-1.453	0.1463
Duration	3.94743	0.61668	6.401	1.54e-10 ***
Age*Type= <i>non-Question</i>	0.15359	0.13996	1.097	0.2725
Condition= <i>normal</i>	0.37337	0.43421	0.860	0.3899
Age*Condition= <i>normal</i>	0.12950	0.10217	1.267	0.2050
Condition= <i>normal</i> *	0.91074	0.72581	1.255	0.2096
Type= <i>non-Question</i>				
Age*Condition= <i>normal</i> *	-0.37965	0.17031	-2.229	0.0258 *
Type= <i>non-Question</i>				
Condition= <i>prosody</i>	-1.60734	0.86680	-1.854	0.0637 .
Age*Condition= <i>prosody</i>	0.39271	0.18905	2.077	0.0378 *
Condition= <i>prosody</i> *	1.68552	1.05414	1.599	0.1098
Type= <i>non-Question</i>				
Age*Condition= <i>prosody</i> *	-0.32360	0.23229	-1.393	0.1636
Type= <i>non-Question</i>				
Condition= <i>words</i>	-0.26996	0.59313	-0.455	0.6490
Age*Condition= <i>words</i>	0.14044	0.13565	1.035	0.3005
Condition= <i>words</i> *	-1.03066	1.01610	-1.014	0.3104
Type= <i>non-Question</i>				
Age*Condition= <i>words</i> *	0.08829	0.22387	0.394	0.6933
Type= <i>non-Question</i>				
<i>Adults</i>	Estimate	Std. Error	<i>z</i> value	Pr(> <i>z</i>)
(Intercept)	-3.3811	0.6884	-4.912	9.03e-07 ***
Type= <i>non-Question</i>	0.4375	0.5854	0.747	0.4549
Duration	4.6961	1.1804	3.978	6.94e-05 ***
Condition= <i>normal</i>	1.2033	0.5359	2.245	0.0247 *
Condition= <i>normal</i> *	-0.9627	0.7358	-1.308	0.1907
Type= <i>non-Question</i>				
Condition= <i>prosody</i>	0.2407	0.8011	0.301	0.7638
Condition= <i>prosody</i> *	0.5525	0.9374	0.589	0.5556
Type= <i>non-Question</i>				
Condition= <i>words</i>	1.5613	0.7087	2.203	0.0276 *
Condition= <i>words</i> *	-1.1557	0.8854	-1.305	0.1918
Type= <i>non-Question</i>				

Table 5: Model output for children and adults' anticipatory gaze switches.

761 3.3.3. *Developmental effects*

762 Our main goal in extending the age range to 1- and 2-year-olds in Experiment 2 was to find the age of emergence for spontaneous predictions about
763 upcoming turn structure. As in Experiment 1, we used two-tailed *t*-tests to
764 compare children's real gaze rates to the random baseline rates in the *normal*
765 speech condition, in which the speech stimulus is most like what children hear
766 every day. We tested real gaze rates against baseline for three age groups:
767 ages one, two, and three. Two- and three-year-old children made anticipatory
768 gaze switches significantly above chance both when all transitions were
769 considered (2-year-olds: $t(26.193)=-4.137$, $p<.001$; 3-year-olds: $t(22.757)=-$
770 2.662, $p<.05$) and for question transitions alone (2-year-olds: $t(25.345)=-$
771 4.269, $p<.001$; 3-year-olds: $t(21.555)=-3.03$, $p<.01$). One-year-olds, how-
772 ever, only made anticipatory gaze shifts marginally above chance for turn
773 transitions overall and for question turns alone (overall: $t(24.784)=-2.049$,
774 $p=.051$; questions: $t(25.009)=-2.03$, $p=.053$).

775 The regression models for the children's data also revealed two significant
776 interactions with age. The first was a significant interaction of age and
777 language condition (for *prosody only* compared to the *no speech* reference
778 level), suggesting a different age effect between the two linguistic conditions.
779 As in Experiment 1, we explored each age interaction by extracting an average
780 difference score over participants for the effect of language condition
781 (*no speech* vs. *prosody only*) within each random permutation of the data,
782 making pairwise comparisons between the six age groups. These tests re-
783 vealed that children's anticipation in the *prosody only* condition significantly
784 improved at ages five and six (with difference scores greater than 95% of the
785 random data scores; $p<.05$). See Figure B.2 for these *prosody only* difference
786 score distributions.

787 The second age-based interaction was a three-way interaction of age, transi-
788 tion type, and language condition (for *normal* speech compared to the *no*
789 *speech* baseline). We again created pairwise comparisons of the average dif-
790 ference scores for the transition type-language condition interaction across
791 age groups in each random permutation of the data, finding that the effect
792 of transition type in the *normal* speech condition became larger with age,
793 with significant improvements by age 4 over ages 1 and 2 (99.9% and 98.86%,
794 respectively), by age 5 over age 4 (97.54%), and by age 6 over ages 1, 2, and 5
795 (99.5%, 97.36%, and 95.04%), all significantly different from chance ($p<.05$).
796 See Figure B.3 for these *normal* speech difference score distributions.

798 3.4. Discussion

799 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
800 to find the age of emergence for spontaneous turn anticipation. Many of our
801 results replicate the findings from Experiment 1: participants made more
802 anticipatory switches when they had access to lexical information and, when
803 they did, tended to make more anticipatory switches for questions compared
804 to non-questions.

805 As in Experiment 1, children and adults spontaneously tracked the turn
806 structure of the conversations, making anticipatory gaze switches at above-
807 chance rates across all ages when listening to natural speech. They also made
808 far more anticipations for questions than for non-question turns—at least for
809 children two years olds and older. But these effects were different for the two
810 conditions with partial linguistic information: *prosody only* and *words only*.
811 In the *prosody only* condition, performance was low for younger children and
812 increased significantly with age (especially for questions). In the *words only*
813 condition, children age two and older showed robust anticipatory switching
814 for questions (much like in *normal* speech), but never rose above chance for
815 non-question turns. These findings do not therefore support an early role
816 for prosody in children’s spontaneous turn structure predictions. There is
817 also no evidence that lexical information is sufficient on its own to support
818 children’s anticipatory switching.

820 4. General Discussion

821 Children begin to develop conversational turn-taking skills long before
822 their first words emerge (Bateson, 1975; Hilbrink et al., 2015; Jaffe et al.,
823 2001; Snow, 1977). As they acquire language, they also acquire the infor-
824 mation needed to make accurate predictions about upcoming turn structure.
825 Until recently, we have had very little data on how children weave language
826 into their already-existing turn-taking behaviors. In two experiments inves-
827 tigating children’s anticipatory gaze to upcoming speakers, we found evi-
828 dence that turn prediction develops early in childhood and that spontaneous
829 predictions are primarily driven by participants’ expectation of an imme-
830 diate response in the next turn (e.g., after questions). In making predic-
831 tions about upcoming turn structure, children used a combination of lexical
832 and prosodic cues; neither lexical nor prosodic cues alone were sufficient to
833 support increased anticipatory gaze. We also found no early advantage for

834 prosody over lexicosyntax, and instead found that children were unable to
835 make above-chance anticipatory gazes in the *prosody only* condition until age
836 five. We discuss these findings with respect to the role of linguistic cues in
837 predictions about upcoming turn structure, the importance of questions in
838 spontaneous predictions about conversation, and children's developing com-
839 petence as conversationalists.

840 *4.1. Predicting turn structure with linguistic cues*

841 Prior work with adults has found a consistent role for lexicosyntax in
842 predicting upcoming turn structure (De Ruiter et al., 2006; Magyari and
843 De Ruiter, 2012), whereas the role of prosody still under debate (Duncan,
844 1972; Ford and Thompson, 1996; Torreira et al., 2015). Knowing that chil-
845 dren comprehend more about prosody than lexicosyntax early on (Section
846 1; also see Speer and Ito, 2009 for a review), we thought it possible that
847 young children would instead show an advantage for prosody in their predic-
848 tions about turn structure in conversation. Our results suggest that, on the
849 contrary, exclusively presenting prosodic information to children limits their
850 spontaneous predictions about upcoming turn structure until age five.

851 Perhaps surprisingly, we also found no evidence that lexical information
852 alone is equivalent to the full linguistic signal in driving children's predic-
853 tions, as has been shown previously for adults (Magyari and De Ruiter, 2012;
854 De Ruiter et al., 2006) and as is replicated with adult participants in the cur-
855 rent study. That said, our findings point more toward early lexical effects
856 in children's turn anticipations than prosodic ones: children's performance
857 in both experiments was consistently better when they had access to lexical
858 information, especially after question turns. And although the *words only*
859 condition in Experiment 2 was not significantly different from the baseline *no*
860 *speech* condition, children's anticipations trended toward above-chance rates.

861 Above all, children and adults anticipated best with they had access to
862 the full linguistic signal. There may be something informative about com-
863 bined prosodic and lexical cues to questionhood that helps to boost children's
864 anticipations before they can use these cues separately. Even in adults, Tor-
865 reira and colleagues (2015) showed that the trade-off in informativity between
866 lexical and prosodic cues is more subtle in semi-natural speech. The present
867 findings are the first to show evidence of a similar effect developmentally.

868 4.2. *The question effect*

869 In both experiments, anticipatory looking was primarily driven by ques-
870 tion transitions, a pattern that has not been previously reported in other an-
871 ticipatory gaze studies, on children or adults (Keitel et al., 2013; Hirvenkari,
872 2013; Tice and Henetz, 2011). Questions make an upcoming speaker switch
873 immediately relevant, helping the listener to predict with high certainty what
874 will happen next (i.e., an answer from the addressee), and are often easily
875 identifiable by overt prosodic and lexicosyntactic cues.

876 Compared to prosodic cues (e.g., final rising intonation), lexicosyntactic
877 cues to questionhood (e.g., *wh*-words, *do*-insertion, and subject-auxiliary in-
878 version) are categorical, and early-occurring in the utterance. Children may
879 have therefore had an easier time picking out and interpreting lexical cues to
880 questionhood. The question effect showed its first significant gains between
881 ages three and four in the *normal* speech condition of Experiment 2 (Figure
882 B.3), by which time children frequently hear and use a variety of polar and
883 *wh*-questions (Clark, 2009). Furthermore, while lexicosyntactic question cues
884 were available on every instance of *wh*- and *yes/no* questions in our stimuli,
885 prosodic question cues were only salient on *yes/no* questions. Finally, the
886 mapping of prosodic contour to speech act (e.g., high final rises for polar
887 questions) is far from one-to-one, leaving substantial room for uncertainty in
888 prosodic contour interpretation in general.

889 Prior work on children’s acquisition of questions indicates that they may
890 already have some knowledge of question-answer sequences by the time they
891 begin to speak: questions make up approximately one third of the utter-
892 ances children hear, before and after the onset of speech, and even into
893 their preschool years, though the types and complexity of questions change
894 throughout development (Casillas et al., In press; Fitneva, 2012; Henning
895 et al., 2005; Shatz, 1979).¹⁰ For the first few years, many of the questions
896 directed to children are “test” questions—questions that the caregiver al-
897 ready has the answer to (e.g., “What does a cat say?”), but this changes as
898 children get older. Questions help caregivers to get their young children’s
899 attention and to ensure that information is in common ground, even if the
900 responses are non-verbal or infelicitous (Bruner, 1985; Fitneva, 2012; Snow,
901 1977). So, in addition to having a special interactive status (for adults and

¹⁰There is substantial variation question frequency by individual and socioeconomic class (Hart and Risley, 1992).

902 children alike), questions are a core characteristic of many caregiver-child
903 interactions, motivating a general benefit for questions in turn structure an-
904 ticipation.

905 Two important questions for future work are then: (1) how does children's
906 ability to monitor for questions in conversation relate to their prior experience
907 with questions? and (2) what is it about questions that makes children and
908 adults more likely to anticipatorily switch their gaze to addressees? Other
909 request formats, such as imperatives, compliments, and complaints make a
910 response from the addressee highly likely in the next turn (Schegloff, 2007).
911 Rhetorical and tag questions, on the other hand, take a similar form to pro-
912 totypical polar questions, but often do not require an answer. So, though it
913 is clear that adults and children anticipated responses more often for ques-
914 tions than non-questions, we do not yet know whether their predictive action
915 is limited to turns formatted as questions or is generally applicable to turn
916 structures that project an immediate response from the addressee.

917 More broadly, our results suggest that participants' spontaneous predic-
918 tions, at least while viewing third-party conversation, are driven by what
919 lies *beyond* the end of the current turn—not by the upcoming end of the
920 turn itself, as has been focused on in prior work (Torreira et al., 2015; Keitel
921 et al., 2013; Magyari and De Ruiter, 2012; De Ruiter et al., 2006). In future
922 work, it will be crucial to measure prediction from a first-person perspective
923 to resolve this apparent contradiction (see also Holler and Kendrick, 2015).

924 4.3. Early competence for turn taking?

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

931 These results contrast with Keitel and colleagues' (2013) finding that chil-
932 dren cannot anticipate upcoming turn structure at above-chance rates until
933 age three. The current study used an appreciably more conservative random
934 baseline than the one used in Keitel and colleagues' study. Therefore, this
935 difference in age of emergence more likely stems from our use of a more en-
936 gaging speech style, stereo speech playback, and more typical turn transition
937 durations.

938 To be clear, young children’s “above chance” performance was often still
939 far from adult-like predictive behavior—children at ages one and two were
940 still very close to chance in their anticipations and, even at age six, children
941 were not fully adult-like in their predictions. This may indicate that young
942 children rely more on non-verbal cues in anticipating turn transitions or,
943 alternatively, that adults are better at flexibly adapting to the turn-relevant
944 cues present at any moment.

945 Taken together, our data suggest that turn-taking skills do begin to
946 emerge in infancy, but that children cannot make effective predictions until
947 they can pick out question turns. This finding leads us to wonder how partic-
948 ipant role (first- instead of third-person) and cultural differences (e.g., high
949 vs. low parent-infant interaction styles) feed into this early predictive skill.
950 It also bridges prior work showing a predisposition for turn taking in infancy
951 (e.g., Bateson, 1975; Hilbrink et al., 2015; Jaffe et al., 2001; Snow, 1977) but
952 late acquisition of adult-like competence for turn-taking in children’s conver-
953 sations (Casillas et al., In press; Garvey, 1984; Garvey and Berninger, 1981;
954 Ervin-Tripp, 1979).

955 *4.4. Limitations and future work*

956 There are at least two major limitations to our work: speech naturalness
957 and participant role. Following prior work (De Ruiter et al., 2006; Keitel
958 et al., 2013), we used phonetically manipulated speech in Experiment 2.
959 This decision resulted in speech sounds that children don’t usually hear in
960 their natural environment. Many prior studies have used phonetically-altered
961 speech with infants and young children (cf. Jusczyk, 2000), but almost none
962 of them have done so in a conversational context. Future work could instead
963 carefully script or cross-splice sub-parts of turns to control for the presence
964 of linguistic cues for turn transition (see, e.g., Torreira et al., 2015).

965 The prediction measure used in our studies is based on an observer’s view
966 of third-party conversation but, because participants’ role in the interaction
967 could affect their online predictions about turn taking, a better measure
968 would instead capture first-person behavior. First-person measures of spon-
969 taneous turn prediction will be key to revealing how participants distribute
970 their attention over linguistic and non-linguistic cues while taking part in
971 everyday interaction, the implications of which relate to theories of online
972 language processing for both language learning and everyday talk.

973 *4.5. Conclusions*

974 Conversation plays a central role in children’s language learning. It is
975 the driving force behind what children say and what they hear. Adults use
976 linguistic information to accurately predict turn structure in conversation,
977 which facilitates their online comprehension and allows them to respond rel-
978 evantly and on time. The present study offers new findings regarding the
979 role of speech acts and linguistic processing in online turn prediction, and
980 has given evidence that turn prediction emerges by age two, but is not inte-
981 grated with linguistic cues until much later. Using language to make predic-
982 tions about upcoming interactive content takes time and, for both children
983 and adults, is primarily driven by participants’ orientation to what will hap-
984 pen beyond the end of the current turn.

985 **Acknowledgements**

986 We gratefully acknowledge the parents and children at Bing Nursery
987 School and the Children’s Discovery Museum of San Jose. This work was
988 supported by an ERC Advanced Grant to Stephen C. Levinson (269484-
989 INTERACT), an NSF graduate research and dissertation improvement fel-
990 lowship to the first author, and a Merck Foundation fellowship to the second
991 author. Earlier versions of these data and analyses were presented to con-
992 ference audiences (Casillas and Frank, 2012, 2013). We also thank Tania
993 Henetz, Francisco Torreira, Stephen C. Levinson, Eve V. Clark, and the
994 First Language Acquisition group at Radboud University for their feedback
995 on earlier versions of this work. The analysis code for this project can be
996 found on GitHub at https://github.com/langcog/turn_taking/.

997 **References**

- 998 Allison, P.D., 2004. Convergence problems in logistic regression, in: Alt-
999 man, M., Gill, J., McDonald, M. (Eds.), Numerical Issues in Statistical
1000 Computing for the Social Scientist. Wiley-Interscience: New York, NY,
1001 pp. 247–262.
- 1002 Allison, P.D., 2012. Logistic Regression Using SAS: Theory and Application.
1003 SAS Institute.

- 1004 Barr, D.J., Levy, R., Scheepers, C., Tily, H.J., 2013. Random effects structure
1005 for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory*
1006 and Language
- 1007 68, 255–278.
- 1008 Bates, D., Maechler, M., Bolker, B., Walker, S., 2014. lme4:
1009 Linear mixed-effects models using Eigen and S4. URL:
1010 <https://github.com/lme4/lme4><http://lme4.r-forge.r-project.org/>.
1011 [Computer program] R package version 1.1-7.
- 1012 Bateson, M.C., 1975. Mother-infant exchanges: The epigenesis of conver-
1013 sational interaction. *Annals of the New York Academy of Sciences* 263,
1014 101–113.
- 1015 Bergelson, E., Swingley, D., 2013. The acquisition of abstract words by young
1016 infants. *Cognition* 127, 391–397.
- 1017 Bloom, K., 1988. Quality of adult vocalizations affects the quality of infant
1018 vocalizations. *Journal of Child Language* 15, 469–480.
- 1019 Boersma, P., Weenink, D., 2012. Praat: doing phonetics by computer. URL:
<http://www.praat.org>. [Computer program] Version 5.3.16.
- 1020 Bögels, S., Magyari, L., Levinson, S.C., 2015. Neural signatures of response
1021 planning occur midway through an incoming question in conversation. *Sci-
1022 entific Reports* 5.
- 1023 Bruner, J., 1985. Child's talk: Learning to use language. *Child Language*
1024 Teaching and Therapy 1, 111–114.
- 1025 Bruner, J.S., 1975. The ontogenesis of speech acts. *Journal of Child Language*
1026 2, 1–19.
- 1027 Carlson, R., Hirschberg, J., Swerts, M., 2005. Cues to upcoming swedish
1028 prosodic boundaries: Subjective judgment studies and acoustic correlates.
1029 *Speech Communication* 46, 326–333.
- 1030 Casillas, M., Bobb, S.C., Clark, E.V., In press. Turn taking, timing, and
1031 planning in early language acquisition. *Journal of Child Language* .
- 1032 Casillas, M., Frank, M.C., 2012. Cues to turn boundary prediction in adults
1033 and preschoolers, in: *Proceedings of SemDial*, pp. 61–69.

- 1034 Casillas, M., Frank, M.C., 2013. The development of predictive processes
1035 in children's discourse understanding, in: Proceedings of the 35th Annual
1036 Meeting of the Cognitive Science Society, pp. 299–304.
- 1037 Clark, E.V., 2009. First language acquisition. Cambridge University Press.
- 1038 De Ruiter, J.P., Mitterer, H., Enfield, N.J., 2006. Projecting the end of
1039 a speaker's turn: A cognitive cornerstone of conversation. Language 82,
1040 515–535.
- 1041 De Vos, C., Torreira, F., Levinson, S.C., 2015. Turn-timing in signed con-
1042 versations: coordinating stroke-to-stroke turn boundaries. Frontiers in
1043 Psychology 6.
- 1044 Dingemanse, M., Torreira, F., Enfield, N., 2013. Is "Huh?" a universal word?
1045 Conversational infrastructure and the convergent evolution of linguistic
1046 items. PloS one 8, e78273.
- 1047 Duncan, S., 1972. Some signals and rules for taking speaking turns in con-
1048 versations. Journal of Personality and Social Psychology 23, 283.
- 1049 Ervin-Tripp, S., 1979. Children's verbal turn-taking, in: Ochs, E., Schieffelin,
1050 B.B. (Eds.), Developmental Pragmatics. Academic Press, New York, pp.
1051 391–414.
- 1052 Fitneva, S., 2012. Beyond answers: questions and children's learning, in:
1053 De Ruiter, J.P. (Ed.), Questions: Formal, Functional, and Interactional
1054 Perspectives. Cambridge University Press, Cambridge, UK, pp. 165–178.
- 1055 Ford, C.E., Thompson, S.A., 1996. Interactional units in conversation: Syn-
1056 tactic, intonational, and pragmatic resources for the management of turns.
1057 Studies in Interactional Sociolinguistics 13, 134–184.
- 1058 Garvey, C., 1984. Children's Talk. volume 21. Harvard University Press.
- 1059 Garvey, C., Berninger, G., 1981. Timing and turn taking in children's con-
1060 versations 1. Discourse Processes 4, 27–57.
- 1061 Gísladóttir, R., Chwilla, D., Levinson, S.C., 2015. Conversation electrified:
1062 ERP correlates of speech act recognition in underspecified utterances. PloS
1063 one 10, e0120068.

- 1064 Griffin, Z.M., Bock, K., 2000. What the eyes say about speaking. Psychological science 11, 274–279.
- 1065
- 1066 Hart, B., Risley, T.R., 1992. American parenting of language-learning children: Persisting differences in family-child interactions observed in natural home environments. Developmental Psychology 28, 1096.
- 1067
- 1068
- 1069 Hedberg, N., Sosa, J.M., Görgülü, E., Mameni, M., 2010. The prosody and meaning of Wh-questions in American English, in: Speech Prosody 2010, pp. 100045:1–4.
- 1070
- 1071
- 1072 Henning, A., Striano, T., Lieven, E.V., 2005. Maternal speech to infants at 1 and 3 months of age. Infant Behavior and Development 28, 519–536.
- 1073
- 1074 Hilbrink, E., Gattis, M., Levinson, S.C., 2015. Early developmental changes in the timing of turn-taking: A longitudinal study of mother-infant interaction. Frontiers in Psychology 6.
- 1075
- 1076
- 1077 Hirvenkari, L., Ruusuvuori, J., Saarinen, V.M., Kivioja, M., Peräkylä, A., Hari, R., 2013. Influence of turn-taking in a two-person conversation on the gaze of a viewer. PloS one 8, e71569.
- 1078
- 1079
- 1080 Holler, J., Kendrick, K.H., 2015. Unaddressed participants' gaze in multi-person interaction. Frontiers in Psychology 6.
- 1081
- 1082 Jaffe, J., Beebe, B., Feldstein, S., Crown, C.L., Jasnow, M.D., Rochat, P., Stern, D.N., 2001. Rhythms of dialogue in infancy: Coordinated timing in development. Monographs of the Society for Research in Child Development. JSTOR.
- 1083
- 1084
- 1085
- 1086 Johnson, E.K., Jusczyk, P.W., 2001. Word segmentation by 8-month-olds: When speech cues count more than statistics. Journal of Memory and Language 44, 548–567.
- 1087
- 1088
- 1089 Jusczyk, P.W., 2000. The Discovery of Spoken Language. MIT press.
- 1090 Jusczyk, P.W., Hohne, E., Mandel, D., Strange, W., 1995. Picking up regularities in the sound structure of the native language. Speech perception and linguistic experience: Theoretical and methodological issues in cross-language speech research , 91–119.
- 1091
- 1092
- 1093

- 1094 Kamide, Y., Altmann, G., Haywood, S.L., 2003. The time-course of prediction
1095 in incremental sentence processing: Evidence from anticipatory eye
1096 movements. *Journal of Memory and Language* 49, 133–156.
- 1097 Keitel, A., Daum, M.M., 2015. The use of intonation for turn anticipation
1098 in observed conversations without visual signals as source of information.
1099 *Frontiers in Psychology* 6.
- 1100 Keitel, A., Prinz, W., Friederici, A.D., Hofsten, C.v., Daum, M.M., 2013.
1101 Perception of conversations: The importance of semantics and intonation
1102 in childrens development. *Journal of Experimental Child Psychology* 116,
1103 264–277.
- 1104 Lemasson, A., Glas, L., Barbu, S., Lacroix, A., Guilloux, M., Remeuf, K.,
1105 Koda, H., 2011. Youngsters do not pay attention to conversational rules:
1106 is this so for nonhuman primates? *Nature Scientific Reports* 1.
- 1107 Levelt, W.J., 1989. Speaking: From intention to articulation. MIT press.
- 1108 Levinson, S.C., 2006. On the human “interaction engine”, in: Enfield, N.,
1109 Levinson, S. (Eds.), *Roots of Human Sociality: Culture, Cognition and*
1110 *Interaction*. Oxford: Berg, pp. 39–69.
- 1111 Levinson, S.C., 2013. Action formation and ascriptions, in: Stivers, T., Sid-
1112 nell, J. (Eds.), *The Handbook of Conversation Analysis*. Wiley-Blackwell,
1113 Malden, MA, pp. 103–130.
- 1114 Magyari, L., Bastiaansen, M.C.M., De Ruiter, J.P., Levinson, S.C., 2014.
1115 Early anticipation lies behind the speed of response in conversation. *Journal*
1116 *of Cognitive Neuroscience* 26, 2530–2539.
- 1117 Magyari, L., De Ruiter, J.P., 2012. Prediction of turn-ends based on antici-
1118 pation of upcoming words. *Frontiers in Psychology* 3:376, 1–9.
- 1119 Masataka, N., 1993. Effects of contingent and noncontingent maternal stimu-
1120 lation on the vocal behaviour of three-to four-month-old Japanese infants.
1121 *Journal of Child Language* 20, 303–312.
- 1122 Mehler, J., Jusczyk, P., Lambertz, G., Halsted, N., Bertoni, J., Amiel-
1123 Tison, C., 1988. A precursor of language acquisition in young infants.
1124 *Cognition* 29, 143–178.

- 1125 Morgan, J.L., Saffran, J.R., 1995. Emerging integration of sequential and
1126 suprasegmental information in preverbal speech segmentation. *Child De-*
1127 *velopment* 66, 911–936.
- 1128 Nazzi, T., Ramus, F., 2003. Perception and acquisition of linguistic rhythm
1129 by infants. *Speech Communication* 41, 233–243.
- 1130 R Core Team, 2014. R: A Language and Environment for Statistical Com-
1131 puting. R Foundation for Statistical Computing. Vienna, Austria. URL:
1132 <http://www.R-project.org>. [Computer program] Version 3.1.1.
- 1133 Ratner, N., Bruner, J., 1978. Games, social exchange and the acquisition of
1134 language. *Journal of Child Language* 5, 391–401.
- 1135 Ross, H.S., Lollis, S.P., 1987. Communication within infant social games.
1136 *Developmental Psychology* 23, 241.
- 1137 Rossano, F., Brown, P., Levinson, S.C., 2009. Gaze, questioning and culture,
1138 in: Sidnell, J. (Ed.), *Conversation Analysis: Comparative Perspectives*.
1139 Cambridge University Press, Cambridge, pp. 187–249.
- 1140 Sacks, H., Schegloff, E.A., Jefferson, G., 1974. A simplest systematics for the
1141 organization of turn-taking for conversation. *Language* 50, 696–735.
- 1142 Schegloff, E.A., 2007. Sequence organization in interaction: Volume 1: A
1143 primer in conversation analysis. Cambridge University Press.
- 1144 Shatz, M., 1979. How to do things by asking: Form-function pairings in
1145 mothers' questions and their relation to children's responses. *Child Devel-*
1146 *opment* 50, 1093–1099.
- 1147 Shi, R., Melancon, A., 2010. Syntactic categorization in French-learning
1148 infants. *Infancy* 15, 517–533.
- 1149 Snow, C.E., 1977. The development of conversation between mothers and
1150 babies. *Journal of Child Language* 4, 1–22.
- 1151 Soderstrom, M., Seidl, A., Kemler Nelson, D.G., Jusczyk, P.W., 2003. The
1152 prosodic bootstrapping of phrases: Evidence from prelinguistic infants.
1153 *Journal of Memory and Language* 49, 249–267.

- 1154 Speer, S.R., Ito, K., 2009. Prosody in first language acquisition—Acquiring
1155 intonation as a tool to organize information in conversation. *Language and*
1156 *Linguistics Compass* 3, 90–110.
- 1157 Stivers, T., Enfield, N.J., Brown, P., Englert, C., Hayashi, M., Heinemann,
1158 T., Hoymann, G., Rossano, F., De Ruiter, J.P., Yoon, K.E., et al., 2009.
1159 Universals and cultural variation in turn-taking in conversation. *Proceedings of the National Academy of Sciences* 106, 10587–10592.
- 1160
- 1161 Stivers, T., Rossano, F., 2010. Mobilizing response. *Research on Language*
1162 and Social Interaction 43, 3–31.
- 1163 Takahashi, D.Y., Narayanan, D.Z., Ghazanfar, A.A., 2013. Coupled oscillator
1164 dynamics of vocal turn-taking in monkeys. *Current Biology* 23, 2162–2168.
- 1165 Thorgrímsson, G., Fawcett, C., Liszkowski, U., 2015. 1- and 2-year-olds'
1166 expectations about third-party communicative actions. *Infant Behavior*
1167 and Development 39, 53–66.
- 1168 Tice (Casillas), M., Henetz, T., 2011. Turn-boundary projection: Looking
1169 ahead, in: *Proceedings of the 33rd Annual Meeting of the Cognitive Science*
1170 Society, pp. 838–843.
- 1171 Torreira, F., Bögels, S., Levinson, S.C., 2015. Intonational phrasing is neces-
1172 sary for turn-taking in spoken interaction. *Journal of Phonetics* 52, 46–57.
- 1173 Wittenburg, P., Brugman, H., Russel, A., Klassmann, A., Sloetjes, H., 2006.
1174 Elan: a professional framework for multimodality research, in: *Proceedings*
1175 of LREC.

1176 **Appendix A. Permutation Analyses**

1177 How can we be sure that our primary dependent measure (anticipatory
1178 gaze switching) actually relates to turn transitions? Even if children were
1179 gazing back and forth randomly during the experiment, we would have still
1180 captured some false hits—switches that ended up in the turn-transition win-
1181 dows by chance.

1182 We estimated the baseline probability of making an anticipatory switch
1183 by randomly permuting the placement of the transition windows within each
1184 stimulus (Figure 4). We then used the switch identification procedure from
1185 Experiments 1 and 2 (Section 2.1.4) to find out how often participants made
1186 “anticipatory” switches within these randomly permuted windows. This pro-
1187 cedure de-links participants’ gaze data from turn structure by randomly re-
1188 assigning the onset time of each turn-transition in each permutation. We
1189 created 5,000 of these permutations for each experiment to get an anticipa-
1190 tory switch baselines over all possible starting points.

1191 Importantly, the randomized windows were not allowed to overlap with
1192 each other, keeping true to the original stimuli. We also made sure that the
1193 properties of each turn transition stayed constant across permutations. So,
1194 while “transition window A” might start 2 seconds into Random Permu-
1195 tation 1 and 17 seconds into Random Permutation 2, it maintained the same
1196 prior speaker identity, transition type, gap duration, language condition, etc.,
1197 across both permutations.

1198 We then re-ran the statistical models from the original data on each of the
1199 random permutations, e.g., using Experiment 1’s original model to analyze
1200 the anticipatory switches from each random permutation of the Experiment
1201 1 looking data. We could then calculate the proportion of random data
1202 z -values exceeded by the original z -value for each predictor. We used the
1203 absolute value of all z -values to conduct a two-tailed test. If the original
1204 effect of a predictor exceeded 95% of the random model effects for that same
1205 predictor, we deemed that predictor’s effect to be significantly different from
1206 the random baseline (i.e., $p < .05$).

1207 For example, children’s “language condition” effect from Experiment 1
1208 had a z -value of $|3.429|$, which is greater than 99.9% of all $|z\text{-value}|$ esti-
1209 mates from Experiment 1’s random permutation models (i.e., $p = .001$). It is
1210 therefore highly unlikely that the effect of language condition in the original
1211 model derived from random gaze shifting.

1212 We used this procedure to derive the random-baseline comparison values

1213 in the main text (above). However, we ran into two issues along the way:
1214 first, we had to report z -values rather than beta estimates. Second, we had
1215 to exclude a substantial portion of the models, especially in Experiment 2
1216 because of model non-convergence. We address each of these issues below.

1217 *Appendix A.1. Beta, standard error, and z estimates*

1218 We reported z -values in the main text rather than beta estimates because
1219 the standard errors in the randomly permuted data models were much higher
1220 than for the original data. The distributions of each predictor's beta estimate,
1221 standard error, and z -value for adults and children in each experiment are
1222 shown in the graphs below (Figures A.1a–A.6b). In each plot, the gray dots
1223 represent randomly permuted model estimates for the value listed (beta,
1224 standard error, or z), the white dots represent the model estimates from the
1225 original data, and the triangles represent the 95th percentile for each random
1226 distribution.

Experiment 1: z -value estimates

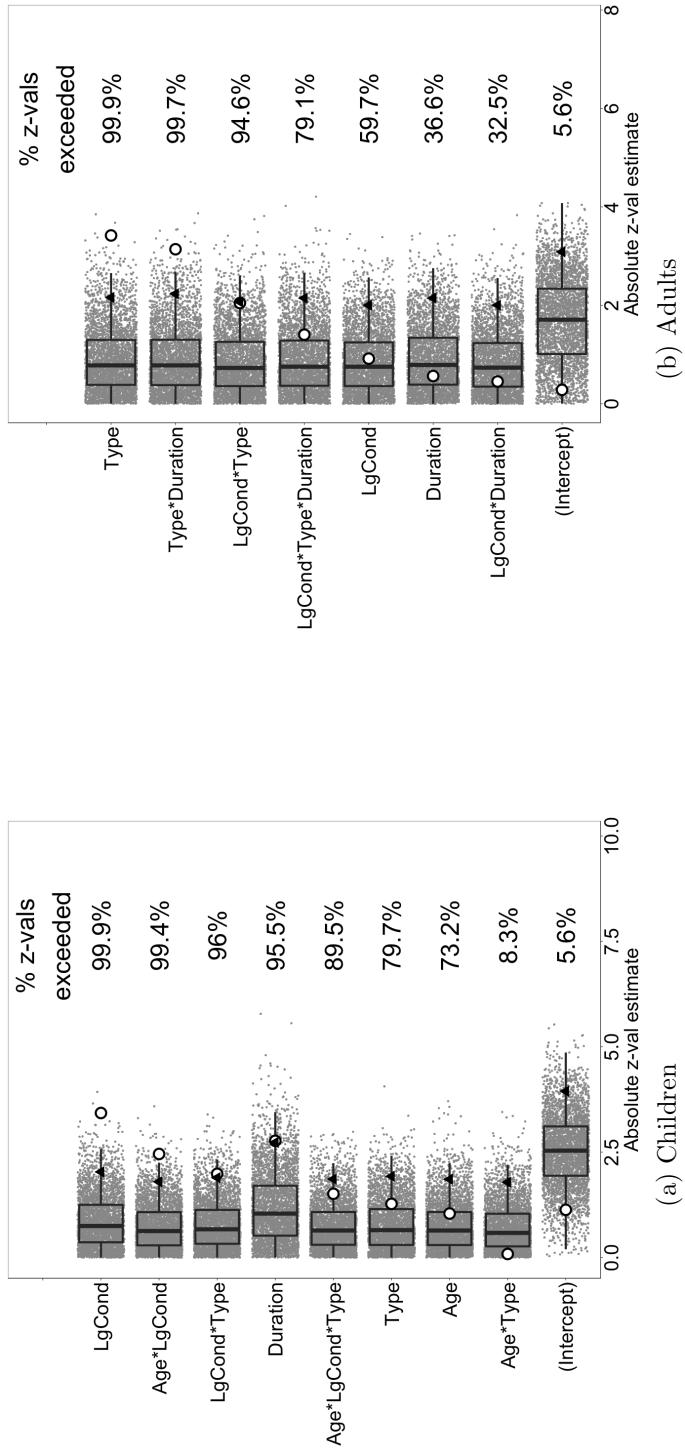
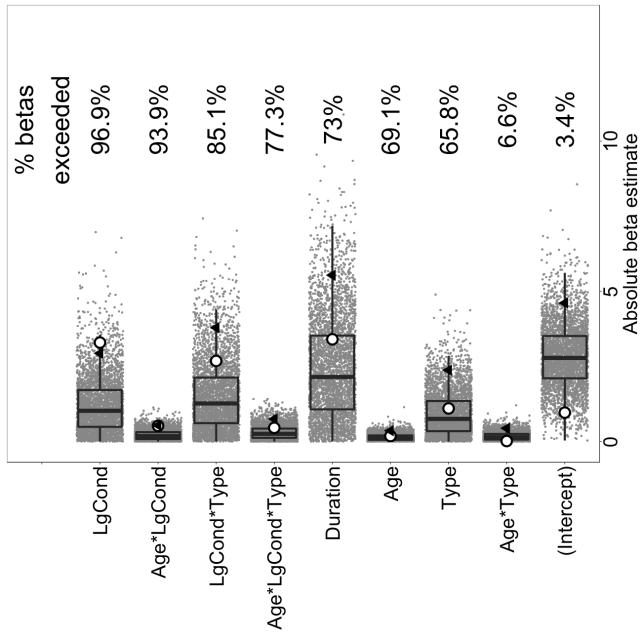


Figure A.1: Random-permutation and original $|z$ -values for predictors of anticipatory gaze rates in Experiment 1.

Experiment 1: β estimates



(a) Children

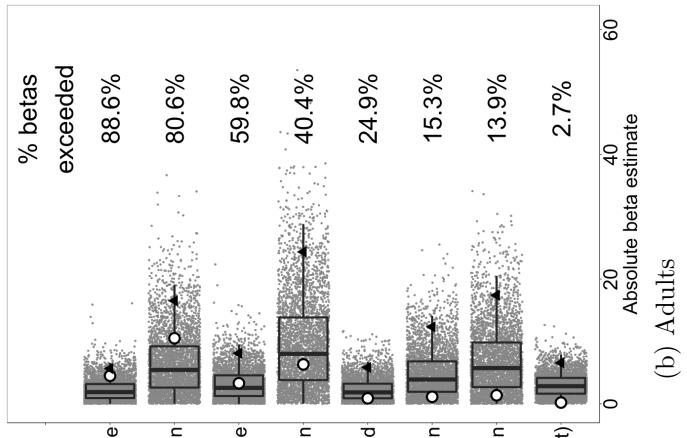


Figure A.2: Random-permutation and original $|\beta\text{-values}|$ for predictors of gaze rates in Experiment 1.

Experiment 1: *SE* estimates

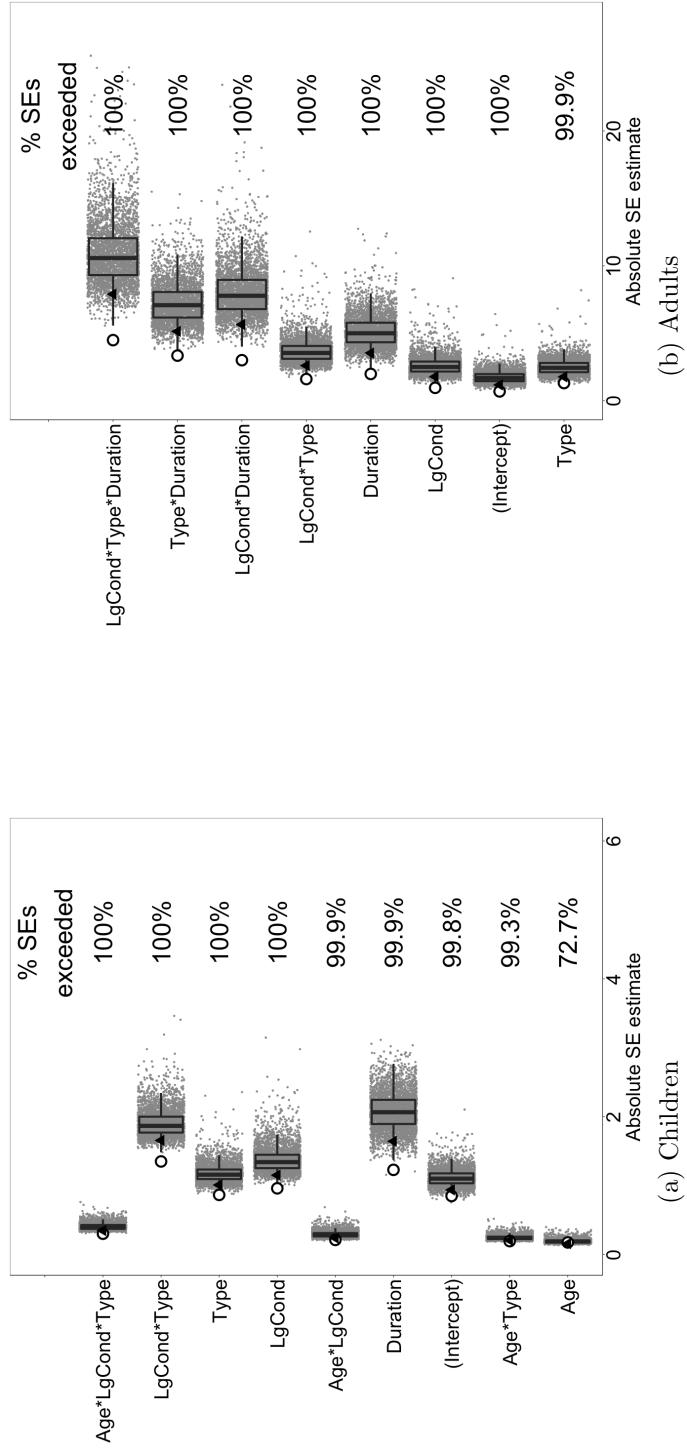


Figure A.3: Random-permutation and original $|SE\text{-values}|$ for predictors of anticipatory gaze rates in Experiment 1.

Experiment 2: z estimates

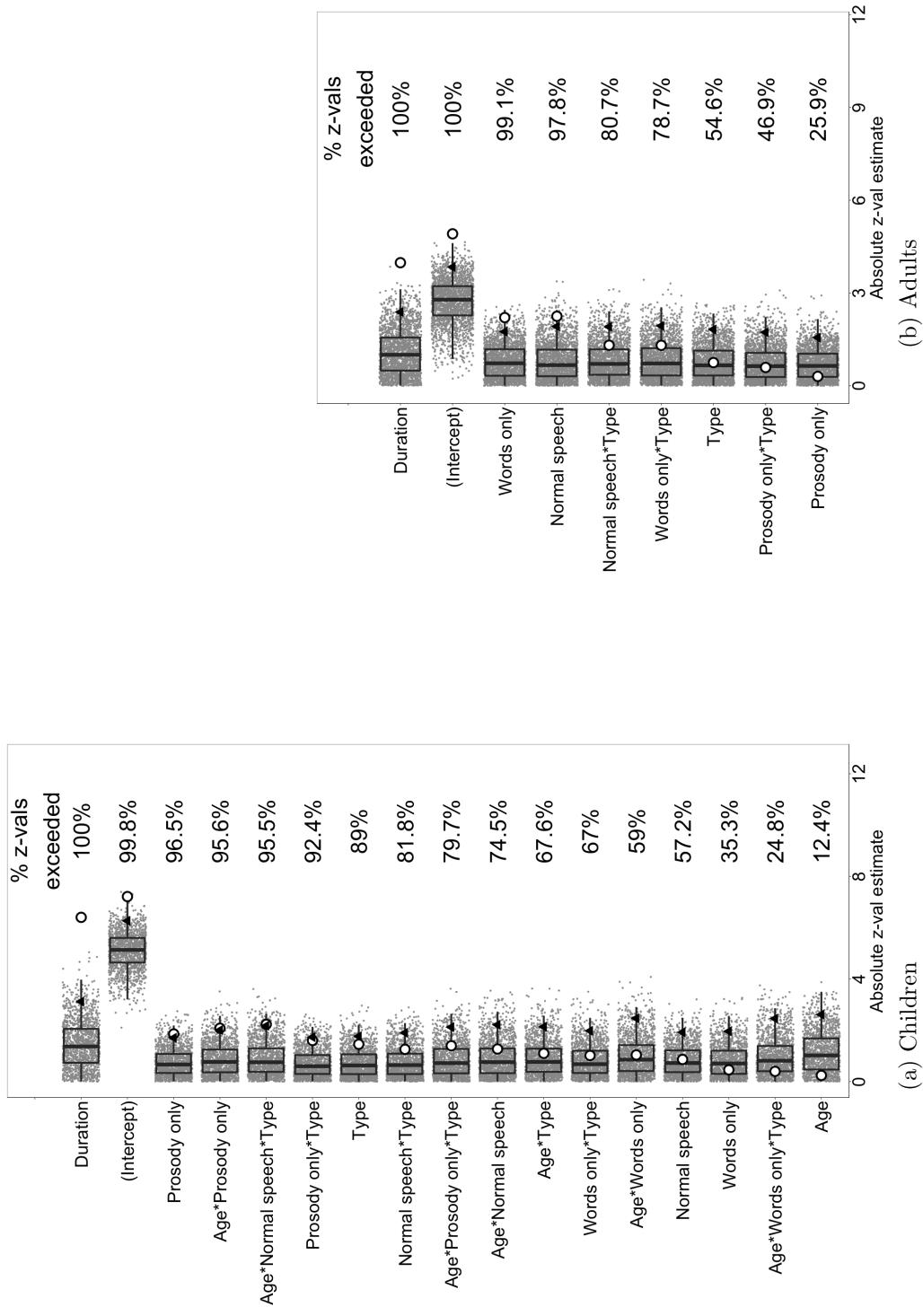


Figure A.4: Random-permutation and original $|z$ -values| for predictors of anticipatory gaze rates in Experiment 2.

Experiment 2: β estimates

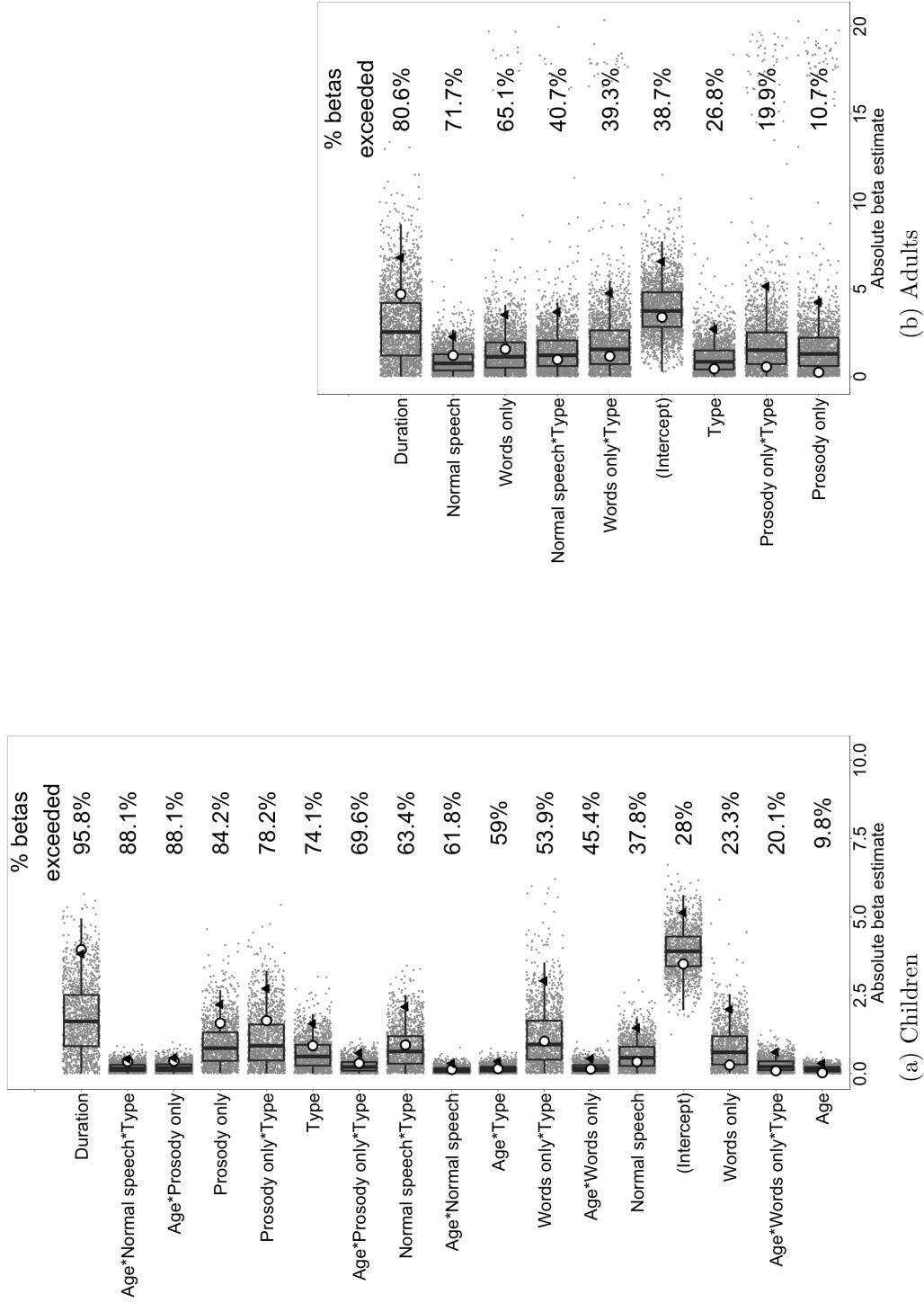


Figure A.5: Random-permutation and original $|\beta\text{-values}|$ for predictors of anticipatory gaze rates in Experiment 2.

Experiment 2: SE estimates

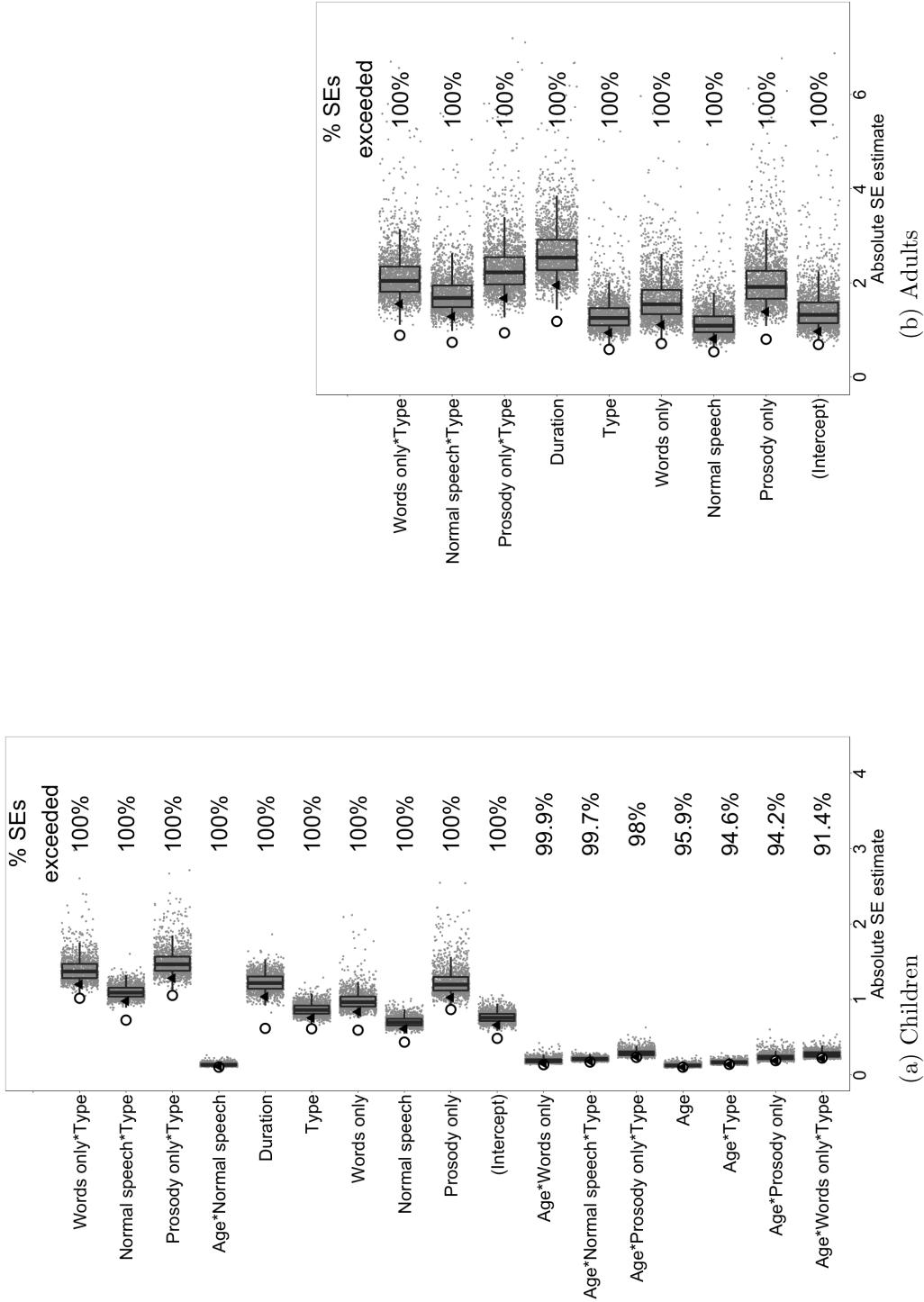


Figure A.6: Random-permutation and original $|SE\text{-values}|$ for predictors of anticipatory gaze rates in Experiment 2.

1227 *Appendix A.2. Non-convergent models*

1228 In comparing the real and randomly permuted datasets, we excluded the
1229 output of random-permutation models that gave convergence warnings to
1230 remove erratic model estimates from our analyses. Non-convergent models
1231 made up 22.4–24.4% of the random permutation models in Experiment 1
1232 and 69–70% of the random permutation models in Experiment 2. The z -
1233 values for each predictor in the converging and non-converging models from
1234 Experiment 1 are shown in table A.1.

1235 Although many of the non-converging models show estimates within range
1236 of the converging models (e.g., with a mean difference of only 0.09 in median
1237 z -value across predictors), they also show many radically outlying estimates
1238 (e.g., showing a mean difference of 237.3 in mean z -value across predictors).
1239 Similar patterns were obtained in the non-converging models for Experiment
1240 2 and persisted even when we tried other optimizers.

1241 We suspect that the issue derives from data sparsity in some of the ran-
1242 dom permutations. This problem is known to occur when there are limited
1243 numbers of binary observations in each of a design matrix’s bins (Allison,
1244 2004). We could instead use zero-inflated poisson or negative binomial re-
1245 gression models to allow for overdispersion in our data (Allison, 2012). How-
1246 ever, these would give us baselines for the normal, convergent model, which
1247 is not the aim of this analysis.

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

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

1248 **Appendix B. Pairwise developmental tests**

1249 Experiments 1 and 2 both showed effects of age in interaction with lin-
1250 guistic condition and transition type (e.g., English vs. non-English). To
1251 explore these effects in more depth, we recorded the average difference score
1252 for the predictor that interacted with age for each participant (e.g., English
1253 minus non-English anticipatory switches), using these values to compute an
1254 average difference score over participants in each age group (e.g., age 3, 4,
1255 and 5) within each random permutation. That averaging process produces
1256 5,000 baseline-derived difference scores for each age group.

1257 We then made pairwise age comparisons of these difference scores (e.g.,
1258 the linguistic condition effect in 3-year-olds vs. 4-year-olds), computing the
1259 percent of random-permutation difference scores exceeded by the real-data
1260 difference score. If the real-data difference score exceeded 95% of the random-
1261 data age difference scores, we deemed it to be an age effect significantly
1262 different from chance—e.g., a significant difference between ages three and
1263 four in the effect of linguistic condition. This procedure is essentially a two-
1264 tailed t -test, adapted for use with the randomly permuted baseline data.

1265 In each of the plots below, the dot represents the real data value for the
1266 effect being shown. The effect sizes from the 5,000 randomly permuted data
1267 sets are shown in the distribution. The percentage displayed is the percentage
1268 of random permutation values exceeded by the original data value (taking the
1269 absolute value of all data points for a two-tailed test). Comparisons marked
1270 with 95% or higher are significant at the $p < 0.05$ level.

Experiment 1: Age and linguistic condition

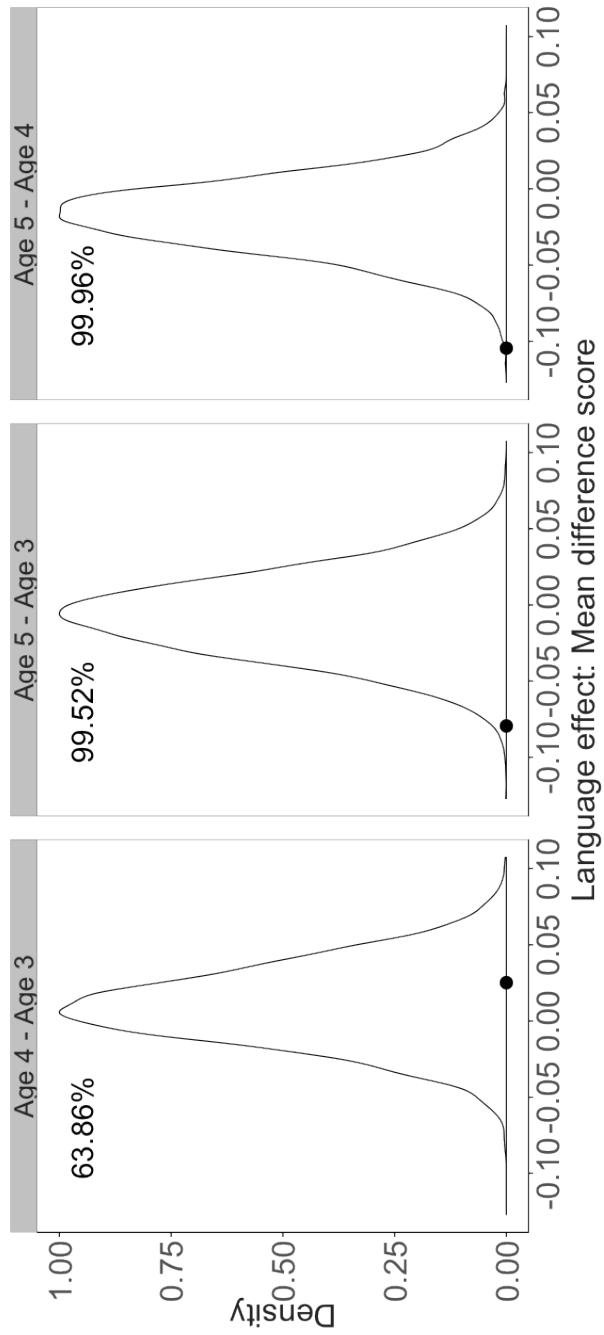


Figure B.1: Pairwise comparisons of the language condition effect across ages in Experiment 1.

Experiment 2: Age and the *prosody only* condition

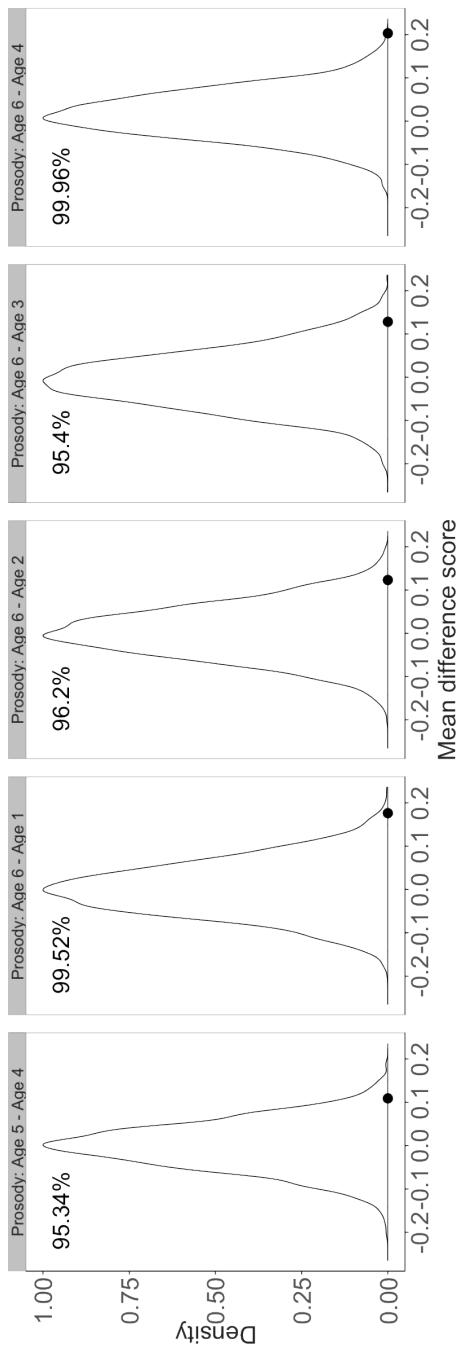


Figure B.2: Significant pairwise comparisons of the *prosody only-no speech* linguistic condition effect, across ages in Experiment 2

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

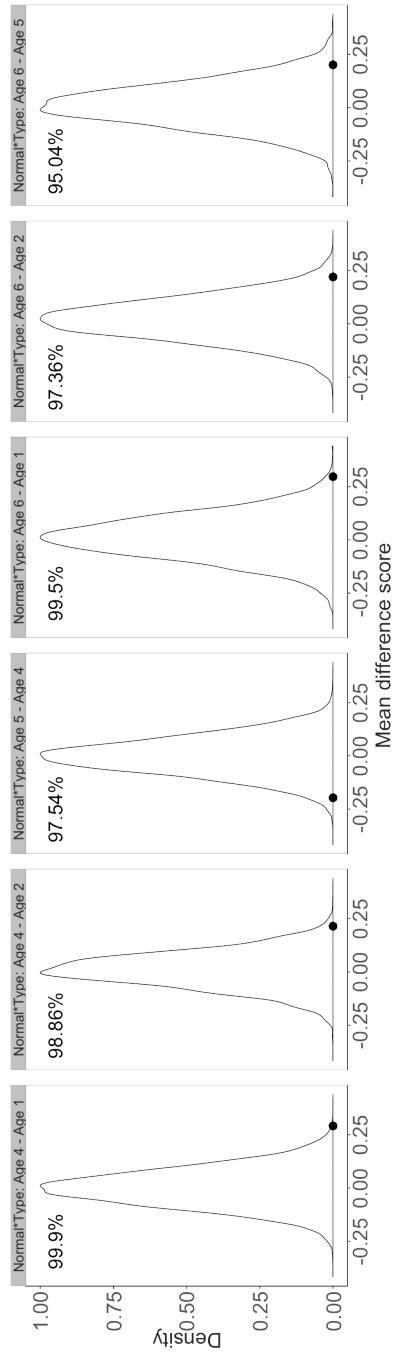


Figure B.3: Significant pairwise comparisons of the *normal speech-no speech* language condition effect for transition type, across ages, in Experiment 2.

1271 **Appendix C. Boredom-driven anticipatory looking**

1272 One alternative hypothesis for children’s anticipatory gazes is that they
1273 simply grow bored and start looking away at a constant rate after a turn
1274 begins. This data plotted here show a hypothetical group of boredom-driven
1275 participants (gray dots) compared to participants from the actual data in
1276 Experiment 2 (black dots). The hypothetical boredom-driven participants
1277 look away from the current speaker at a linear rate, beginning one second
1278 after the start of a turn.

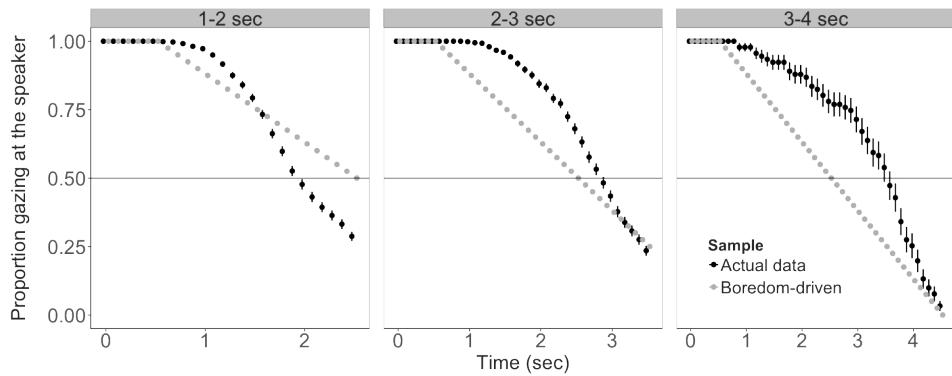


Figure C.1: Proportion of participants (hypothetical boredom-driven=gray; actual Experiment 2=black) looking at the current speaker, split by turn duration. Vertical bars indicate standard error in the experimental data.

1279 If children’s switches away from the current speaker were driven by bore-
1280 dom, they would switch away equally quickly on long and short turns. How-
1281 ever, as turn duration increased, children in the experiment looked at the
1282 current speaker for a longer period before looking away, suggesting that they
1283 were not switching away at a constant rate. We can see this pattern most
1284 clearly in the time at which 50% of the children had switched away from
1285 the current (indicated by the horizontal line in Figure C.1): children’s gaze
1286 crossed this 50% threshold at 2.0, 2.9, and 3.6 seconds after the start of speech
1287 for turns with durations of 1–2, 2–3, and 3–4 seconds, respectively. In con-
1288 trast, the hypothetical boredom-driven children always hit the 50% threshold
1289 at 2.5 seconds after the start of speech. This pattern suggests that, though
1290 children do look away with time, their looks away are not simply driven by
1291 boredom.

1292 **Appendix D. Puppet pair and linguistic condition**

1293 The design for Experiment 2 does not fully cross puppet pair (e.g., robots,
1294 blue puppets) with linguistic condition (e.g., *words only* and *no speech*). Even
1295 though each puppet pair is associated with different conversation clips across
1296 children (e.g., robots talking about kitties, birthday parties, and pancakes),
1297 the robot puppets themselves were exclusively associated with the *words only*
1298 condition. Similarly, merpeople were exclusively associated with *prosody only*
1299 speech, and the puppets wearing dress clothes were exclusively associated
1300 with the *no speech* condition. We designed the experiment this way to in-
1301 crease its pragmatic felicity for older children (i.e., robots make robot sounds,
1302 merpeople’s voices are muffled under the water, the party-going puppets are
1303 in a ‘party’ room with many other voices). There is therefore a confound
1304 between linguistic condition and puppet pair; for example, children could
1305 have made fewer anticipatory switches in the *prosody only* condition because
1306 the puppets were less interesting. To test whether puppet pair drove the
1307 condition-based differences found in Experiment 2, we ran a short follow-up
1308 study.

1309 **Methods**

1310 We recruited 30 children between ages 3;0 and 5;11 from the Children’s Dis-
1311 covery Museum of San Jose, California to participate in our experiment. All
1312 participants were native English speakers. Children were randomly assigned
1313 to one of six videos (five children per video).

1314 *Materials.* We created 6 short videos from the stimulus recordings made for
1315 Experiment 2. Each video featured a puppet pair (red/blue/yellow/robot/
1316 merpeople/party-goer; Figure 5). Puppets in all six videos performed the
1317 exact same conversation recording ('birthday party'; Experiment 2) with
1318 normal, unmanipulated speech; this experiment therefore holds all things
1319 constant across stimuli except for the appearance of the puppets.

1320 *Procedure.* We used the same experimental apparatus and procedure as in
1321 Experiments 1 and 2. Each participant was randomly assigned to watch only
1322 one of the six puppet videos. Five children watched each video. As in Exper-
1323 iment 2, the experimenter immediately began each session with calibration
1324 and then stimulus presentation because no special instructions were required.
1325 The entire experiment took less than three minutes.

1326 *Data preparation.* We identified anticipatory gaze switches to the upcoming
1327 speaker using the same method as in Experiments 1 and 2.

1328 **Results and discussion**

1329 We modeled children’s anticipatory switches (yes or no at each transition)
1330 with mixed effects logistic regression, including puppet pair (robots/mer-
1331 people/party-goers/other-3) as a fixed effect and participant and turn tran-
1332 sition as random effects. We grouped the red, blue, and yellow puppets
1333 together because they collectively represented the puppets used in the *nor-*
1334 *mal* speech condition—this follow-up experiment is meant to test whether
1335 the condition-based differences from Experiment 2 arose from the puppets
1336 used in each condition.

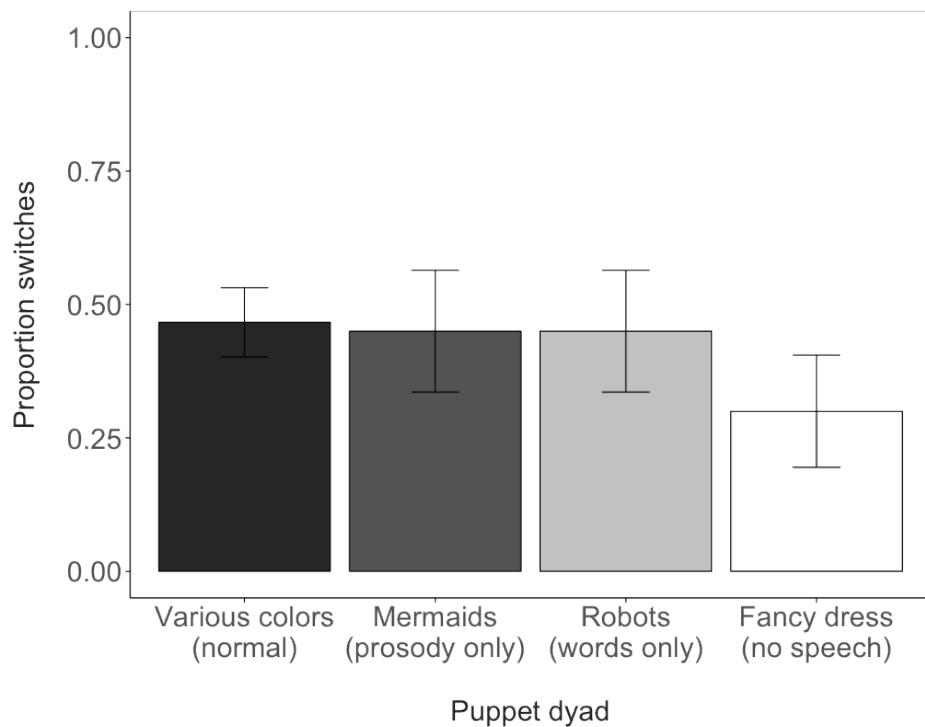


Figure D.1: Proportion gaze switches across puppet pairs when linguistic condition and conversation are held constant.

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

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

1337 In four versions of this model, we systematically varied the reference level
 1338 of the puppet pair to check for any cross-condition differences. We found no
 1339 significant effects of puppet pair on switching rate (all $p > 0.25$; Table D.2).

1340 We take this finding as evidence that our decision to not fully cross puppet
 1341 pairs and linguistic conditions in Experiment 2 was unlikely to have strongly
 1342 affected children’s anticipatory gaze rates above and beyond the intended
 1343 effects of linguistic condition.