

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 upcoming turn structure. In Experiment 1, we showed English and non-English conversations to English-speaking participants, finding minimal differences between the predictive looking behavior of preschoolers and adults. In Experiment 2, we phonetically controlled lexicosyntactic and prosodic cues in English-only speech, finding that children's predictive looking behavior improved from ages one to six, but that even one-year-olds made more anticipatory looks than would be expected by chance. In both experiments, children and adults anticipated more often after hearing questions. Like adults, prosody alone did not improve children's predictive gaze shifts. But, unlike adults, lexical cues alone were also not sufficient to improve prediction—children's performance was best overall with access to lexicosyntax and prosody together. Our findings support an account in which turn prediction emerges in infancy, but takes several years before becoming fully integrated with linguistic processing.

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

*Corresponding author

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
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). How does this ability
¹¹ develop?

¹² Children participate in sequential coordination with their caregivers start-
¹³ ing at three months of age—before they can rely on any linguistic cues in
¹⁴ taking turns (see, among others, Bateson, 1975; Hilbrink et al., 2015; Jaffe
¹⁵ et al., 2001; Snow, 1977). Of course, infant turn taking is different from
¹⁶ adult turn taking in several ways. Infant turn taking is heavily scaffolded
¹⁷ by caregivers, has distinct timing in comparison to adult turn taking, and
¹⁸ lacks semantic content (Hilbrink et al., 2015; Jaffe et al., 2001). However,
¹⁹ children’s early, turn-structured social interactions are presumably a critical
²⁰ precursor to their conversational turn taking: early non-verbal interactions
²¹ likely establish the protocol by which children come to use language with
²² others. How do children integrate linguistic knowledge with these preverbal
²³ turn-taking abilities?

²⁴ In this study, we investigate when children begin to make predictions
²⁵ about upcoming turn structure in conversation, and how they integrate lan-
²⁶ guage into their predictions as they grow older. In what follows, we first give
²⁷ a basic review of turn-taking research and the state of current knowledge
²⁸ about adult turn prediction. We then discuss recent work on the develop-
²⁹ ment of turn-taking skills before turning to the 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 environments (Lemasson et al., 2011; Takahashi et al., 2013). In all these
38 cases, turn taking serves as a protocol for interaction, allowing the participants
39 to coordinate with one another through sequences of contingent action.

40 Conversation distinguishes itself from non-conversational turn-taking behaviors by the complexity of the turn sequencing involved. In the examples
41 above (traffic, games, and monkeys) the set of sequence and action types is
42 far more limited and predictable than what we find in everyday talk. For
43 example, conversational turns come grouped into semantically-contingent se-
44 quences of action. The groups can span turn-by-turn exchanges (e.g., simple
45 question-response, “How are you?”–“Fine.”) or sequence-by-sequence ex-
46 changes (e.g., reciprocals, “How are you?”–“Fine, and you?”–“Great!”).

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

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

⁷⁴ Prior work suggests that adult listeners primarily use lexicosyntactic information to accurately predict upcoming turn structure (De Ruiter et al.,
⁷⁵ 2006). De Ruiter and colleagues (2006) asked participants to listen to snippets of spontaneous conversation and to press a button whenever they anticipated 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 participants'
⁸⁰ button responses was comparable to their timing with the full linguistic signal. But when no lexical information was available, participants'
⁸¹ responses were significantly earlier. The authors concluded that lexicosyntactic information¹ was necessary and possibly sufficient for turn-end
⁸² projection, while intonation was neither necessary nor sufficient. Congruent
⁸³ evidence comes from studies varying the predictability of lexicosyntactic
⁸⁴ and pragmatic content: Adults anticipate turn ends better when they can
⁸⁵ more accurately predict the exact words that will come next (Magyari and
⁸⁶ De Ruiter, 2012; see also Magyari et al., 2014). They can also identify speech
⁸⁷ acts within the first word of an utterance (Gísladóttir et al., 2015), allowing
⁸⁸ them to start planning their response at the first moment possible (Bögels
⁸⁹ et al., 2015).

⁹⁰ Despite this evidence, the role of prosody for adult turn prediction is
⁹¹ still a matter of debate. De Ruiter and colleagues' (2006) experiment focused
⁹² on the role of intonation, which is only a partial index of prosody.
⁹³ Prosodic structure is also tied closely to the syntax of an utterance, and
⁹⁴ so the two linguistic signals are difficult to control independently (Ford and
⁹⁵ Thompson, 1996). Torreira, Bögels and Levinson (2015) used a combination
⁹⁶ of button-press and verbal responses to investigate the relationship between
⁹⁷ lexicosyntactic and prosodic cues in turn-end prediction. Critically, their
⁹⁸ stimuli were cross-spliced so that each item had full prosodic cues to accom-
⁹⁹ pany the lexicosyntax. Because of the splicing, they were able to create items
¹⁰⁰ that had syntactically-complete units with no intonational phrase boundary
¹⁰¹ at the end. Participants never verbally responded or pressed the “turn-end”
¹⁰² button when hearing a syntactically-complete phrase without an intonational

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.

107 phrase boundary. And when intonational phrase boundaries were embedded
108 in multi-utterance turns, participants were tricked into pressing the “turn-
109 end” button 29% of the time. Their results suggest that listeners actually
110 do rely on prosodic cues to execute a response (see also de De Ruiter et al.
111 (2006):525). These experimental findings corroborate other corpus and ex-
112 perimental work promoting a combination of cues (lexicosyntactic, prosodic,
113 and pragmatic) as key for accurate turn-end prediction (Duncan, 1972; Ford
114 and Thompson, 1996; Hirvenkari et al., 2013). We next turn to evidence on
115 children’s developing turn-taking ability.

116 *1.2. Children’s turn prediction*

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

128 The timing of children’s responses to their caregivers’ speech shows a non-
129 linear pattern. A recent study by Hilbrink et al. (2015) finds that infants’
130 turn timing at three months is often too early or too late: They start vocal-
131 izing in overlap on 40% of their caregivers’ turns, and their non-overlapped
132 vocalizations come after an average inter-turn silent gap of 350–900 msec
133 (adult average: 200 msec). Between four and nine months, children begin
134 to reduce the number of turns happening in overlap while also improving
135 on their average response latency. But then, later on, children’s response
136 latencies slow down again, peaking at average gaps of more than 1000 msec
137 at nine months, with only very gradual improvement after that (Hilbrink
138 et al., 2015). While children’s avoidance of overlap is nearly adult-like by
139 nine months, the timing of their non-overlapped responses stays much longer
140 than the 200 msec standard for the next few years (Casillas et al., In press;
141 Garvey, 1984; Ervin-Tripp, 1979).

142 This puzzling pattern is likely due to their linguistic development: Tak-
143 ing turns on time is easier when the response is a simple vocalization rather

than a linguistic utterance. Integrating language into the turn-taking system may be one major factor in children's delayed responses (Casillas et al., In press). Consistent with this hypothesis, during the first year, children's prosodic abilities are relatively sophisticated while their lexical knowledge is limited. Infants can distinguish their native language's rhythm type from others soon after birth (Mehler et al., 1988; Nazzi and Ramus, 2003); they show preference for the typical stress patterns of their native language over others by 6–9 months (e.g., iambic vs. trochaic), and can use prosodic information to segment the speech stream into smaller chunks from 8 months onward (Johnson and Jusczyk, 2001; Morgan and Saffran, 1995). In comparison, children show at best a very limited lexical inventory before the first birthday (Bergelson and Swingley, 2013; Shi and Melancon, 2010).

If response planning (i.e., language production) is the primary hurdle in young children's spontaneous turn taking, we should find evidence that children understand turn-taking behaviors before they are able to produce the behaviors themselves. This hypothesis has been recently explored in experimental settings, but results are mixed. One study found that 12-month-olds make more predictive gaze shifts to a responder while watching human verbal conversation compared to conversation-like interactions with objects (Bakker et al., 2011), but another only found a similar effect at 36 months (von Hofsten et al., 2009). However, neither of these two studies had baselines to which the turn-relevant looking behavior could be compared. A baseline measurement is critical because there may be developmental differences in gaze shifting between conversational participants, even if the shifting is not related to turn structure. Such developmental differences could produce artefactual changes in measures of turn-contingent shifting.

Keitel and colleagues (2013) addressed the random baseline issue in a study of 6-, 12-, 24-, and 36-month-olds. They asked participants to watch short videos of conversation and tracked their eye movements at points of speaker change. They found that children's anticipatory gaze frequency was only greater than chance for 36-month-olds and adults. Their study was also the first to focus on the role of linguistic processing in children's turn predictions. They showed their participants two types of conversation videos: One normal and one with flattened pitch (i.e., with flattened intonation contours), finding that only 36-month-olds were affected by a lack of intonation contours. The adult control group made equal numbers of anticipatory looks in the videos, with and without intonation contours, consistent with prior adult findings (De Ruiter et al., 2006). Keitel and colleagues concluded that

182 children’s ability to predict upcoming turn structure relies on their ability
183 to comprehend the stimuli (emerging around 36 months), especially with re-
184 spect to semantic access. They also suggest that intonation takes a secondary
185 role in turn prediction, but only *after* children acquire more sophisticated,
186 adult-like language comprehension systems (sometime after 36 months).

187 Although the Keitel et al. (2013) study constitutes a substantial advance
188 over previous work, it has its own limitations. Because these limitations di-
189 rectly inform our own study design, we review them in some detail. First,
190 their estimates of baseline gaze frequency (“random” in their terminology)
191 were not random. Instead, they used gaze switches during ongoing speech as
192 a baseline. Ongoing speech is perhaps the period in which switching is least
193 likely to occur (Hirvenkari et al., 2013), thus maximizing chances of finding a
194 difference between gaze frequency at turn transitions and their baseline rate.
195 A more conservative baseline would be to compare participants’ looking be-
196 havior at turn transitions to their looking behavior during randomly selected
197 windows of time throughout the stimulus, including turn transitions. We
198 follow this conservative approach in our work.

199 Second, the conversation stimuli Keitel et al. (2013) used were somewhat
200 unusual. The average gap between turns was 900 msec, which is much longer
201 than typical adult timing, where gaps average around 200 msec (Stivers et al.,
202 2009). The speakers in the videos were also asked to minimize their move-
203 ments while performing a scripted and adult-directed conversation, which
204 would have created a somewhat unnatural stimulus. Additionally, in order
205 to produce more naturalistic conversation, it would have been ideal to local-
206 ize the sound sources for the two voices in the video (i.e., to have the voices
207 come out of separate left and right speakers). But both voices were recorded
208 and played back on the same audio channel, which may have made it more
209 difficult to distinguish the two talkers (again, we attempt to address these
210 issues in our current study).

211 Despite these minor methodological issues, the Keitel et al. (2013) study
212 still demonstrates intriguing age-based differences in children’s ability to pre-
213 dict upcoming turn structure, and the results suggest that both semantic and
214 intonational development *do* play a role in children’s looking patterns. Our
215 current work thus takes this paradigm as a starting point.² We report here
216 on the role of linguistic processing in children’s predictions about upcoming

²See also Casillas and Frank (2012, 2013).

217 turn structure, in particular on how children use prosodic and lexicosyntactic
218 information to make their predictions.

219 **2. Experiment 1**

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

232 *2.1. Methods*

233 *2.1.1. Participants*

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

243 *2.1.2. Materials*

244 *Video recordings.* We recorded pairs of talkers while they conversed in
245 a sound-attenuated booth (see sample frame in Figure 1). Each talker was
246 a native speaker of the language being recorded, and each talker pair was
247 male-female. Using a Marantz PMD 660 solid state field recorder, we cap-
248 tured audio from two lapel microphones, one attached to each participant,
249 while simultaneously recording video from the built-in camera of a MacBook



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

250 laptop computer. The talkers were volunteers and were acquainted with their
251 recording partner ahead of time.

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

266 After recording, we combined the audio and video files by hand, and
267 cropped each recording to the 30-second interval with the most turn activity.
268 Because we recorded the conversations in stereo, the male and female voices
269 came out of separate speakers during video playback. This gave each voice in

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.

270 the videos a localized source (from the left or right loudspeaker). We coded
271 each turn transition in the videos for language condition (English vs. non-
272 English), inter-turn gap duration (in milliseconds), and speech act (question
273 vs. non-question). The non-English stimuli were coded for speech act from
274 a monolingual English-speaker’s perspective, i.e., which turns “sound like”
275 questions, and which don’t: we asked five native American English speakers
276 to listen to the audio signal for each turn and judge whether it sounded
277 like a question. We then coded turns with at least 80% “yes” responses as
278 questions.

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

293 Questions made up exactly half of the turn transitions in the English
294 (N=10) and non-English (N=8) videos. In the English videos, inter-turn
295 gaps were slightly shorter for questions (mean=310, median=293, stdev=112
296 msec) than non-questions (mean=325, median=315, stdev=118 msec). Non-
297 English videos did not show a large difference in transition time for questions
298 (mean=270, median=257, stdev=116 msec) and non-questions (mean=302,
299 median=252, stdev=134 msec).

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). As such, they would still be fixated on the prior speaker when the overlap started, and then would have to switch their gaze *reactively* to the responder.

300 2.1.3. *Procedure*

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

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

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

333 All participants (child and adult) completed a five-point calibration rou-
334 tine before the first video started. We used a dancing Elmo for the children's
335 calibration image. During the experiment, participants watched all six 30-
336 second conversation videos. The first and last conversations were in American

337 English and the intervening conversations were Hebrew, Japanese, German,
338 and Korean. The presentation order of the non-English videos was shuffled
339 into four lists, which participants were assigned to randomly. The entire
340 experiment, including instructions, took 10–15 minutes.

341 *2.1.4. Data preparation and coding*

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

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

367 As mentioned, most anticipatory switches happen in the inter-turn gap,
368 but we also allowed anticipatory gaze switches that occurred in the final
369 syllables of the prior turn. Early switches are consistent with the distribution
370 of responses in explicit turn-boundary prediction tasks. For example, in
371 a button press task, adult participants anticipate turn ends approximately
372 200 msec in advance of the turn’s end, and anticipatory responses to pitch-
373 flattened stimuli come even earlier (De Ruiter et al., 2006). We therefore

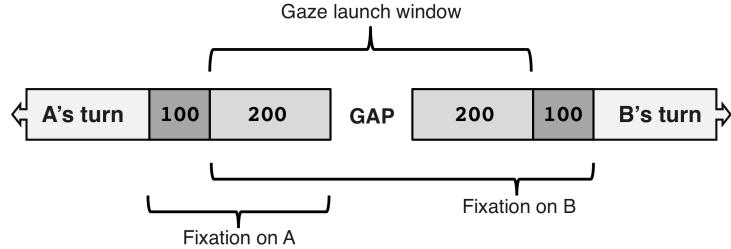


Figure 2: Schematic summary of criteria for anticipatory gaze shifts from speaker A to speaker B during a turn transition.

374 allowed switches to occur as early as 200 msec before the end of the prior turn.
 375 For very early and very late switches, our requirement for 100 msec of fixation
 376 on each speaker would sometimes extend outside of the transition window
 377 boundaries (200 msec before and after the inter-turn gap). The maximally
 378 available fixation window was 100 msec before and after the earliest and
 379 latest possible switch point (300 msec before and after the inter-turn gap).
 380 We did not count switches made during the fixation window as anticipatory.
 381 We *did* count switches made during the inter-turn gap. The period of time
 382 from the beginning of the possible fixation window on the prior speaker to the
 383 end of the possible fixation window on the responder was our total analysis
 384 window (300 msec + the inter-turn gap + 300 msec).

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

393 *2.2. Results*

394 Participants looked at the screen most of the time during video playback
 395 (81% and 91% on average for children and adults, respectively). They pri-
 396 marily kept their eyes on the person who was currently speaking in both
 397 English and non-English videos: They gazed at the current speaker between
 398 38% and 63% of the time, looking back at the addressee between 15% and

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.

399 20% of the time (Table 1). Even three-year-olds looked more at the current
 400 speaker than anything else, whether the videos were in a language they could
 401 understand or not. Children looked at the current speaker less than adults
 402 did during the non-English videos. Despite this, their looks to the addressee
 403 did not increase substantially in the non-English videos, indicating that their
 404 looks away were probably related to boredom rather than confusion about
 405 ongoing turn structure. Overall, participants’ pattern of gaze to current
 406 speakers demonstrated that they performed basic turn tracking during the
 407 videos, regardless of language.

408 *2.2.1. Statistical models*

409 We identified anticipatory gaze switches for all 36 usable turn transitions,
 410 based on the criteria outlined in Section 2.1.4, and analyzed them for effects
 411 of language, transition type, and age with two mixed-effects logistic regres-
 412 sions (Bates et al., 2014; R Core Team, 2014). We built one model each
 413 for children and adults. We modeled children and adults separately because
 414 effects of age are only pertinent to the children’s data. The child model
 415 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 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

<i>Children</i>	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.96146	0.84901	-1.132	0.257446
Age	-0.18268	0.17507	-1.043	0.296725
LgCond= <i>non-English</i>	-3.29347	0.96045	-3.429	0.000606 ***
Type= <i>non-Question</i>	-1.10129	0.86494	-1.273	0.202925
Duration	3.40169	1.22826	2.770	0.005614 **
Age*LgCond= <i>non-English</i>	0.52065	0.21190	2.457	0.014008 **
Age*TypeS= <i>non-Question</i>	-0.01628	0.19437	-0.084	0.933232
LgCond= <i>non-English</i> *	2.68166	1.35016	1.986	0.047013 *
Type= <i>non-Question</i>				
Age*LgCond= <i>non-English</i> *	-0.45632	0.30163	-1.513	0.130315
Type= <i>non-Question</i>				

<i>Adults</i>	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.1966	0.6942	-0.283	0.776988
LgCond= <i>non-English</i>	-0.8812	0.9602	-0.918	0.358754
Type= <i>non-Question</i>	-4.4953	1.3139	-3.421	0.000623 ***
Duration	-1.1227	1.9880	-0.565	0.572238
LgCond= <i>non-English</i> *	3.2972	1.6101	2.048	0.040581 *
Type= <i>non-Question</i>				
LgCond= <i>non-English</i> *	1.3626	3.0077	0.453	0.650527
Duration				
Type= <i>non-Question</i> *	10.5107	3.3459	3.141	0.001682 **
Duration				
LgCond= <i>non-English</i> *	-6.3156	4.4926	-1.406	0.159790
Type= <i>non-Question</i> *				
Duration				

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

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

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.

425 slopes of condition, transition type, and their interaction for participant.

426 Children's anticipatory gaze switches showed effects of language condition
427 ($\beta=-3.29$, $SE=0.961$, $t=-3.43$, $p<.001$) and gap duration ($\beta=3.4$, $SE=1.229$,
428 $t=2.77$, $p<.01$) with additional effects of an age-by-language condition in-
429 teraction ($\beta=0.52$, $SE=0.212$, $t=2.46$, $p<.05$) and a language condition-by-
430 transition type interaction ($\beta=2.68$, $SE=1.35$, $t=1.99$, $p<.05$). There were
431 no significant effects of age or transition type alone ($\beta=-0.18$, $SE=0.175$,
432 $t=-1.04$, $p=.3$ and $\beta=-1.10$, $SE=0.865$, $t=-1.27$, $p=.2$, respectively).

433 Adults' anticipatory gaze switches shows an effect of transition type ($\beta=$
434 4.5 , $SE=1.314$, $t=-3.42$, $p<.001$) and significant interactions between lan-
435 guage condition and transition type ($\beta=3.3$, $SE=1.61$, $t=2.05$, $p<.05$) and
436 transition type and gap duration ($\beta=10.51$, $SE=3.346$, $t=3.141$, $p<.01$).

437 *2.2.2. Random baseline comparison*

438 We estimated the probability that these patterns were the result of ran-
439 dom looking by running the same regression models on participants' real
440 eye-tracking data, only this time calculating their anticipatory gaze switches
441 with respect to randomly permuted turn transition windows. This process
442 involved: (1) randomizing the order and temporal placement of the anal-
443 ysis windows within each stimulus (Figure 3; "analysis window" is defined
444 in Figure 2), thereby randomly redistributing the analysis windows across
445 the eye-tracking signal, (2) re-running each participant's eye tracking data
446 through switch identification (described in 2.1.4), this time using the ran-
447 domly permuted analysis windows, and (3) modeling the anticipatory gazes
448 from the randomly permuted data with the same statistical models we used
449 for the original data (Section 2.2.1; Table 2). Importantly, although the onset
450 time of each transition was shuffled within the eye-tracking signal, the other
451 intrinsic properties of each turn transition (e.g., prior speaker identity, transi-
452 tion type, gap duration, language condition, etc.) stayed constant across
453 each random permutation.

454 This procedure effectively de-links participants' gaze data from the turn
455 structure in the original stimulus, thereby allowing us to compare turn-
456 related (original) and non-turn-related (randomly permuted) looking behav-
457 ior using the same eye movements. The resulting anticipatory gazes from the
458 randomly permuted analysis windows represent an average anticipatory gaze
459 rate over all possible starting points: a random baseline. By running the real
460 and randomly permuted data sets through identical statistical models, we
461 can also estimate how likely it is that predictor effects in the original data

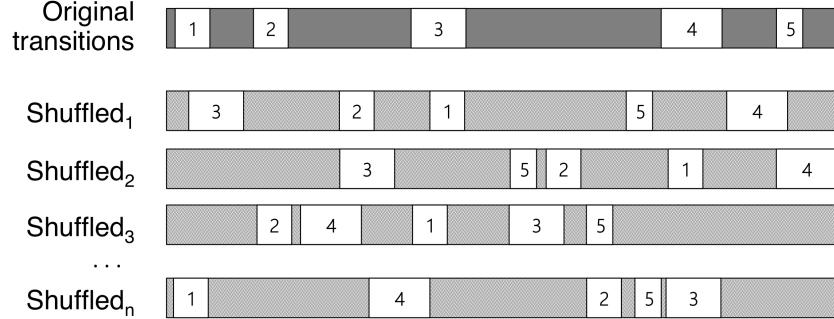


Figure 3: Example of shuffling for five turn transition analysis windows. The windows were ± 300 msec around the inter-turn gap.

(e.g., the effect of language condition; Table 2) arose from random looking.

We completed this baseline procedure on 5,000 random permutations of the original turn transition analysis windows and compared the t -values from each predictor in the original models (Table 2) to the distribution of t -values for each predictor in the 5,000 models of the randomly permuted datasets.⁷ We could then test whether significant effects from the original statistical models differed from the random baseline by calculating the proportion of random data t -values exceeded by the original t -value for each predictor, using the absolute value of all t -values for a two-tailed test. For example, children's original "language condition" t -value was $|3.429|$, which is greater than 99.9% of all $|t\text{-value}|$ estimates from the randomly-permuted data models (i.e., $p = .001$). This leads us to conclude that the effect of language condition in the original model was highly unlikely to be the result of random gaze shifting.

Our baseline analyses revealed that none of the significant predictors from models of the original, turn-related data can be explained by random looking. The children's data showed strong evidence of differentiation from the randomly permuted data for all four significant effects in the original model (Table 2: Children): the original t -values for language condition, gap duration, the age-language condition interaction, and the language condition-transition

⁷We report t -values rather than beta estimates because the standard errors in the randomly permuted data models were much higher than for the original data. For those interested, plots of the beta and standard error distributions are available in the Appendix.

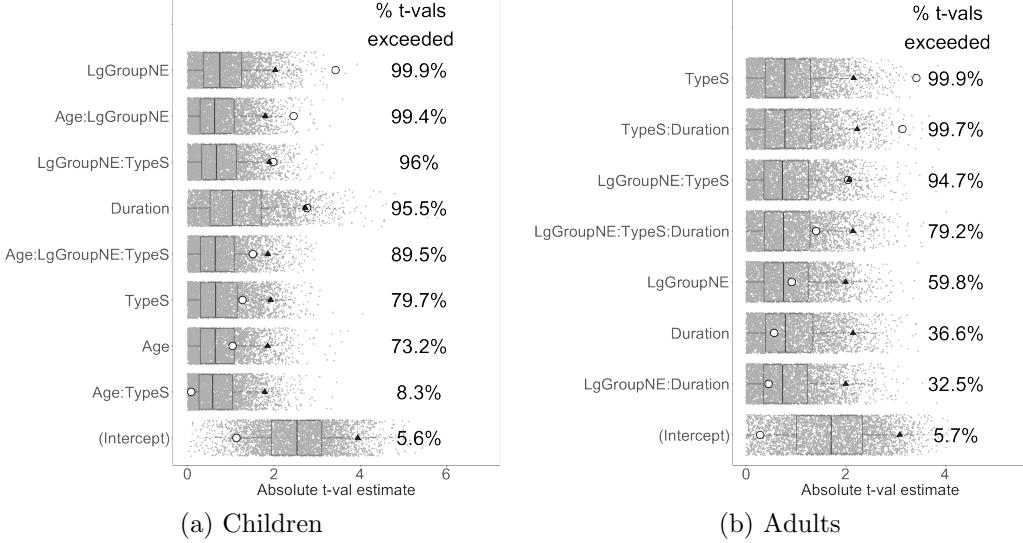


Figure 4: Random-permutation and original $|t\text{-values}|$ for predictors of children and adults' anticipatory gaze rates. Gray dots = random model estimates, White dots = original model estimates, Triangles = 95th percentile for each t -value distribution.

482 type interaction were all greater than 95% of t -values for the randomly
 483 permuted data (99.9%, 95.5%, 99.4%, and 96%, respectively; Figure 4a).
 484 Similarly, the adults' data showed significant differentiation from the ran-
 485 domly permuted data for two of the three originally significant predictors—
 486 transition type and the transition type-gap duration interaction (greater than
 487 99.9% and 99.7% of random t -values, respectively)—with marginal differen-
 488 tiation for the interaction of language condition and transition type (greater
 489 than 94.7% of random t -values; Figure 4b). The effects of language condi-
 490 tion and transition type for the real and randomly permuted data can also
 491 be observed in Figure 5. We excluded the output of random-permutation
 492 models that did not ultimately converge to remove unreliable model results
 493 from our percentile calculations below (78% and 76% of models for children
 494 and adults, respectively).

495 *Developmental effects.* The model of the children's data revealed a significant
 496 interaction of age and language condition (Table 2) that was highly unlikely
 497 to have derived from random looking (Figure 5). To further explore this
 498 effect, we compared the average effect of language condition across all ages:

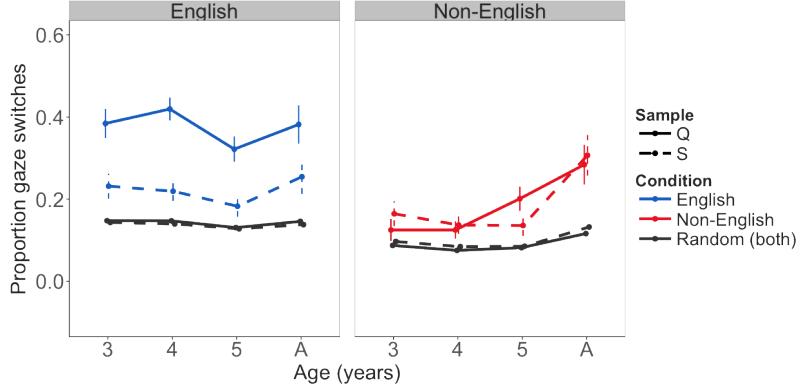


Figure 5: Anticipatory gaze rates across language condition and transition type for the real (red and blue) and randomly permuted (gray) data. Vertical bars represent the standard error.

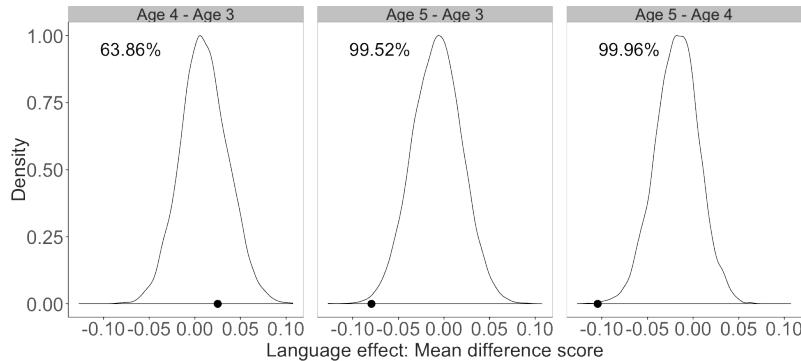


Figure 6: Pairwise comparisons of the language condition effect across ages for the original data (black dots) and the 5,000 randomly permuted datasets (distribution).

499 we extracted the average difference score for the two language conditions
 500 (English minus non-English) for each subject, computing an overall average
 501 for each random permutation of the data. For each random permutation,
 502 we then made pairwise comparisons of the average difference scores across
 503 ages. Figure 6 plots the real-data difference scores against the random-data
 504 difference score distribution for each pairwise age comparison, showing that
 505 3- and 4-year olds were affected equally by language condition, but that 5-

506 year-olds affected less than both 3- and 4-year-olds (with 99.52% and 99.96%
507 of difference scores greater than the randomly permuted data, respectively,
508 i.e., differences of $p < .01$ and $p < .001$).

509 *2.3. Discussion*

510 Children and adults spontaneously tracked the turn structure of the con-
511 versations, making anticipatory gaze switches at an above-chance rate across
512 all ages and conditions (Table 1; Figure 5). Children’s anticipatory gaze rates
513 were affected by language condition, transition type, age, and gap duration
514 (Table 2), none of which could be explained by a baseline of random gaze
515 switching (Figure 4a).

516 Language condition (English vs. non-English) affected children’s antici-
517 pations in two ways (Table 2; Figure 5). First, children made more antici-
518 patory switches overall in English videos, compared to non-English videos. This
519 effect suggests that lexical access is important for children’s ability to antici-
520 pate upcoming turn structure; children had no lexical access to the speech in
521 the non-English videos, though they did have access to (non-native) prosodic
522 cues and non-verbal behavior, consistent with prior work on turn-end pre-
523 diction in adults (De Ruiter et al., 2006; Magyari and De Ruiter, 2012) and
524 children (Keitel et al., 2013). Second, children systematically made more an-
525 ticipatory switches after hearing a question compared to a non-question, but
526 only in the English condition, suggesting that, when children have access to
527 lexical cues, they are more likely to make an anticipatory gaze switch if they
528 can expect an immediate response from the addressee. If so, then children’s
529 attention to lexical cues for turn taking may primarily be in monitoring the
530 signal for cues to questionhood (e.g., subject-auxiliary inversion, *wh*-words,
531 etc.).

532 Children’s anticipatory gaze switches were also affected by their age, but
533 only in the non-English videos: 3- and 4-year-olds made many more antici-
534 patory switches when watching videos in English compared to non-English, but
535 this effect of language condition had attenuated significantly by age 5 (Table
536 2; Figure 5; Figure 6). This interaction suggests that the 5-year-olds were
537 able to leverage anticipatory cues in the non-English videos in a way that
538 3- and 4-year-olds could not, possibly by shifting more attention to the non-
539 native prosodic or non-verbal cues. Prior work on children’s turn-structure
540 anticipation proposed that children’s turn-end predictions rely primarily on
541 lexicosyntactic structure (and not, e.g., prosody) as they get older (Keitel

542 et al., 2013). The current results suggest more flexibility in children’s pre-
543 dictions; when they do not have access to lexical information, older children
544 are likely to find alternative cues to turn taking behavior.

545 Finally, children showed an effect of gap duration (Table 2). This effect
546 is straightforward: longer gaps resulted in longer analysis windows, yielding
547 more time for children to make an anticipatory gaze.

548 Adults’ anticipatory gaze rates were also affected by transition type, lan-
549 guage condition, and gap duration (Table 2), none of which could be easily
550 explained by a baseline of random gaze switching (Figure 4b). Like children,
551 adults made more anticipatory switches after hearing questions compared
552 to non-questions, suggesting that anticipation mattered more to them when
553 an immediate response was expected. Also like children, the advantage for
554 questions was driven by lexical access such that adults must have relied on
555 lexicosyntactic cues to questionhood in picking out turns that potentially
556 require an immediate response, though this effect was only marginally di-
557 vergent from the distribution of randomly permuted data ($p = .053$; Figure
558 4b). Finally, adults’ anticipation rates were also affected by gap duration,
559 but more so for questions than non-questions (Table 2). This interaction
560 suggests that adults were less likely overall to make anticipatory switches at
561 non-questions (as is evident for adults and children in Figure 5), and so did
562 not benefit from extra time to do so compared to long gaps for questions.

563 2.3.1. Summary

564 Children and adults’ predictions alike were benefited by access to lexical
565 information (English) and speech act status (questionhood), suggesting that
566 linguistic cues, particularly lexical ones, facilitate their spontaneous predic-
567 tions about upcoming turn structure through the identification of turns with
568 immediate responses. Children’s anticipatory gaze rates for questions and
569 non-questions in English was stable across ages and comparable to adult be-
570 havior (Figure 5), suggesting that they can identify questions in native stimuli
571 with adult-like competence by age three. Although participants’ ability to
572 recognize questions was facilitated by lexical access (i.e., English vs. non-
573 English), the prosody in the non-English videos was non-native, and so the
574 experimental design can not conclusively show which linguistic cues children
575 relied on in the English videos to identify question turns. Relatedly, though
576 lexical access clearly facilitated participants’ anticipatory gaze rate, it was
577 not necessary for participants—especially adults—in order to exceed chance
578 switching rates (Figure 5), suggesting that participants use non-lexical cues

579 (e.g., prosody, non-verbal behavior) to make anticipatory eye movements at
580 least some of the time.

581 Interestingly, adults and children both were strongly affected by transition
582 type, in that they made more anticipatory switches after hearing questions,
583 compared to non-questions. Even in the English videos, when participants
584 had full access to linguistic cues, their rates of anticipation were relatively
585 low—in fact, comparable to the non-English videos—unless the turn was a
586 question (Figure 5). Prior work using online, metalinguistic tasks has shows
587 that participants can use linguistic cues to accurately predict upcoming turn
588 ends. The current results suggest that, in their spontaneous predictions
589 about third-party conversation, both children and adults monitor the lin-
590 guistic structure of unfolding turns for cues to upcoming responses.

591 Children and adults generally behaved relatively similarly in this first ex-
592 periment and our language manipulation (English vs. non-English) was too
593 coarse to comment on when children begin to use different types of native
594 linguistic cues (e.g., prosody vs. lexicosyntax); we would instead need to di-
595 rectly compare lexicosyntactic and prosodic cues in the participants' native
596 language, controlling for the presence of non-verbal cues. To see the emer-
597 gence of anticipatory gaze switching we would also need to include younger
598 children, since participants already reliably made anticipatory gaze switches
599 at age three. Experiment 1 thus lays the analytic groundwork for a method
600 that allows for greater experimental control, which we introduce in Experi-
601 ment 2.

602 3. Experiment 2

603 We improved our design by using native-language stimuli, controlling for
604 lexical and prosodic information, eliminating non-verbal cues, and testing
605 children from a wider age range. All of the videos in Experiment 2 were in
606 the participants' native language (American English). To tease apart the
607 role of lexical and prosodic information, we phonetically manipulated the
608 speech signal for pitch, syllable duration, and lexical access. By testing one-
609 to six-year-olds we hoped to find the developmental onset of turn-predictive
610 gaze. We also hoped to measure changes in the relative roles of prosody and
611 lexicosyntax across development.

612 Non-verbal cues in Experiment 1 (e.g., gaze and gesture) could have
613 helped participants make predictions about upcoming turn structure (Rossano
614 et al., 2009; Stivers and Rossano, 2010). Since our focus is on linguistic cues,

615 we eliminated all gaze and gestural signals in Experiment 2 by replacing
616 the videos of human actors with videos of puppets. Puppets are less real-
617 istic and expressive than human actors, but they create a natural context
618 for having somewhat motionless talkers in the videos (thereby allowing us
619 to eliminate gestural and gaze cues). Additionally, the prosody-controlled
620 condition included small but global changes to syllable duration that would
621 have required complex video manipulation or precise re-enactment with hu-
622 man talkers, neither of which was feasible. For these reasons, we decided to
623 substitute puppet videos for human videos in the final stimuli.

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

628 *3.1. Methods*

629 *3.1.1. Participants*

630 We recruited 27 undergraduate adults and 129 children between ages 1;0–
631 6;11 to participate in our experiment. We recruited our child participants
632 from the Children's Discovery Museum in San Jose, California, targeting ap-
633 proximately 20 children for each of the six 1-year age groups (range=20–23).
634 All participants were native English speakers, though some parents (N=27)
635 reported that their child heard a second (and sometimes third) language at
636 home. None of the adult participants reported fluency in a second language.
637 We ran Experiment 2 at a local children's museum because it gave us access
638 to children with a more diverse range of ages.

639 *3.1.2. Materials*

640 We created 18 short videos of improvised, child-friendly conversation (Fig-
641 ure 7). To eliminate non-verbal cues to turn transition and to control the
642 types of linguistic information available in the stimuli we first audio-recorded
643 improvised conversations, then phonetically manipulated those recordings to
644 limit the availability of prosodic and lexical information, and finally recorded
645 video to accompany the manipulated audio, featuring puppets as talkers.

646 *Audio recordings.* The recording session was set up in the same way as
647 the first experiment, but with a shorter warm up period (5–10 minutes) and
648 a pre-determined topic for the child-friendly improvisation ('riding bikes',
649 'pets', 'breakfast', 'birthday cake', 'rainy days', or 'the library'). All of the
650 talkers were native English speakers, and were recorded in male-female pairs.

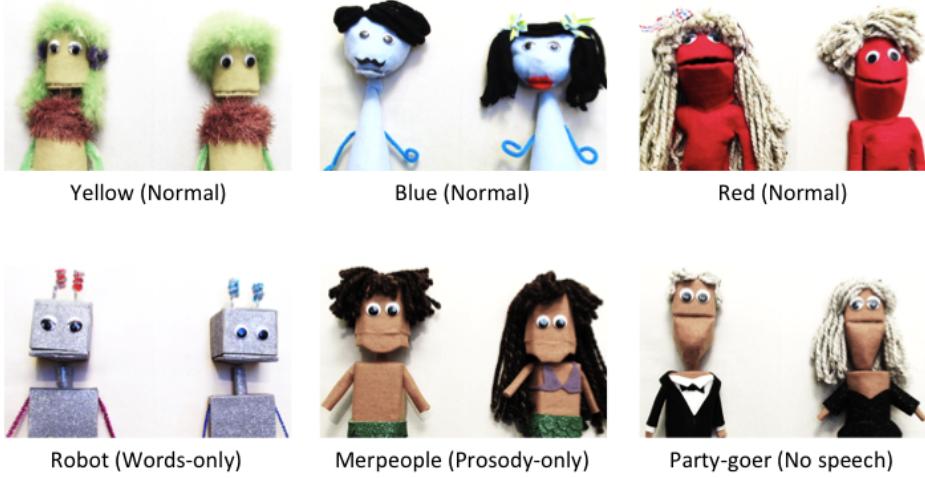


Figure 7: The six puppet pairs (and associated audio conditions). Each pair was linked to three distinct conversations from the same condition across the three experiment versions.

651 As before, we asked talkers to speak “as if they were on a children’s television
 652 show” and to ask at least a few questions during the improvisation. We cut
 653 each audio recording down to the 20-second interval with the most turn
 654 activity. The 20-second clips were then phonetically manipulated and used
 655 in the final video stimuli.

656 *Audio Manipulation.* We created four versions of each audio clip: *Normal*,
 657 *words only*, *prosody only*, and *no speech*. That is, one version with a full
 658 linguistic signal (*normal*), and three with incomplete linguistic information
 659 (hereafter “limited cue” conditions). The *normal* clips were the unmanipu-
 660 lated, original audio clips.

661 The *words only* clips were manipulated to have robot-like speech: We
 662 flattened the intonation contours to each talker’s average pitch (F0) and
 663 we reset the duration of every nucleus and coda to each talker’s average
 664 nucleus and coda duration.⁸ We made duration and pitch manipulations
 665 using PSOLA resynthesis in Praat (Boersma and Weenink, 2012). Thus,
 666 the *words only* versions of the audio clips had no pitch or durational cues
 667 to upcoming turn boundaries, but did have intact lexicosyntactic cues (and

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

668 residual phonetic correlates of prosody, e.g., intensity).

669 We created the *prosody only* clips by low-pass filtering the original recording
670 at 500 Hz with a 50 Hz Hanning window (following de Ruiter et al., 2006).
671 This manipulation creates a “muffled speech” effect because low-pass filtering
672 removes most of the phonetic information used to distinguish between
673 phonemes. The *prosody only* versions of the audio clips lacked lexical infor-
674 mation, but retained their intonational and rhythmic cues to upcoming turn
675 boundaries.

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

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

689 *Video recordings.* We created puppet video recordings to match the ma-
690 nipulated 20-second audio clips. The puppets were minimally expressive;
691 the experimenter could only control the opening and closing of their mouths;
692 their head, eyes, arms, and body stayed still. Puppets were positioned look-
693 ing forward to eliminate shared gaze as a cue to turn structure (Thorgrímsson
694 et al., 2015). We took care to match the puppets’ mouth movements to the
695 syllable onsets as closely as possible, specifically avoiding any mouth move-
696 ment before the onset of a turn. We then added the manipulated audio clips
697 to the puppet video recordings by hand.

698 We used three pairs of puppets used for the *normal* condition—‘red’,
699 ‘blue’ and ‘yellow’—and one pair of puppets for each limited cue condition:
700 “robots”, “merpeople”, and “party-goers” (Figure 8). We randomly assigned
701 half of the conversation topics (‘birthday cake’, ‘pets’, and ‘breakfast’) to the
702 *normal* condition, and half to the limited cue conditions (‘riding bikes’, ‘rainy
703 days’, and ‘the library’). We then created three versions of the experiment,
704 so that each of the six puppet pairs was associated with three different con-
705 versation topics across the different versions of the experiment (18 videos

706 in total). We ensured that the position of the talkers (left and right) was
707 counterbalanced in each version by flipping the video and audio channels as
708 needed.

709 The duration of turn transitions and the number of speaker changes
710 across videos was variable because the conversations were recorded semi-
711 spontaneously. We measured turn transitions from the audio recording of
712 the *normal*, *words only*, and *prosody only* conditions. There was no audio
713 from the original conversation in the *no speech* condition videos, so we mea-
714 sured turn transitions from the video recording, using ELAN video editing
715 software (Wittenburg et al., 2006).

716 There were 85 turn transitions for analysis after excluding transitions
717 longer than 550 msec and shorter than 90 msec. The remaining turn tran-
718 sitions had slightly more questions than non-question ($N=50$ and $N=35$, re-
719 spectively), with transitions distributed somewhat evenly across conditions
720 (keeping in mind that there were three *normal* videos and only one lim-
721 ited cue video for each experiment version): *Normal* ($N=36$), *words only*
722 ($N=13$), *prosody only* ($N=17$), and *no speech* ($N=19$). Inter-turn gaps for
723 questions (mean=365, median=427) were longer than those for non-questions
724 (mean=302, median=323) on average, but gap duration was overall com-
725 parable across conditions: *Normal* (mean=334, median=321), *words only*
726 (mean=347, median=369), *prosody only* (mean=365, median=369), and *no*
727 *words* (mean=319, median=329). The longer gaps for question transitions
728 could give them an advantage because our anticipatory measure includes
729 shifts initiated during the gap between turns (Figure 2).

730 *3.2. Procedure*

731 We used the same experimental apparatus and procedure as in the first
732 experiment. Each participant watched six puppet videos in random order,
733 with five 15–30 second filler videos placed in-between (e.g., running pup-
734 pies, moving balls, flying bugs). Three of the puppet videos had *normal*
735 audio while the other three had *words only*, *prosody only*, and *no speech* au-
736 dio. This experiment required no special instructions so the experimenter
737 immediately began each session with calibration (same as before) and then
738 stimulus presentation. The entire experiment took less than five minutes.

739 *3.2.1. Data preparation and coding*

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

744 *3.3. Results*

745 Participants' pattern of gaze indicated that they performed basic turn
746 tracking across all ages and in all conditions. Participants looked at the
747 screen most of the time during video playback (82% and 86% average for
748 children and adults, respectively). Children and adults primarily kept their
749 eyes on the person who was currently speaking: They gazed at the current
750 speaker between 44% and 69% of the time, looking back at the addressee
751 between 11% and 14% of the time (Table 2). They tracked the current
752 speaker in every condition—even one-year-olds looked more at the current
753 speaker than at anything else in the three limited cue conditions (40% for
754 *words only*, 43% for *prosody only*, and 39% for *no speech*). There was a steady
755 overall increase in looks to the current speaker with age and added linguistic
756 information (Tables 3 and 4). Looks to the addressee also decreased with
757 age, but the change was minimal.

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.

⁹We coded *wh*-questions as “non-questions” for the *prosody only* videos. Polar questions had a final rising prosodic contour, but *wh*-questions did 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.

758 *3.3.1. Statistical models*

759 We identified anticipatory gaze switches for all 85 usable turn transitions,
 760 and analyzed them for effects of language condition, transition type, and age
 761 with two mixed-effects logistic regressions (Bates et al., 2014; R Core Team,
 762 2014). We again built separate models for children and adults because effects
 763 of age were only pertinent to the children’s data. The child model included
 764 condition (normal/prosody only/words only/no speech; with no speech as
 765 the reference level), transition type (question vs. non-question), age (1, 2, 3,
 766 4, 5, 6), and duration of the inter-turn gap (in seconds) as predictors, with
 767 full interactions between language condition, transition type, and age. We
 768 again included the duration of the inter-turn gap as a control predictor and
 769 added random effects of item (turn transition) and participant, with random
 770 slopes of transition type for participants (Barr et al., 2013). The adult model
 771 included condition, transition type, their interactions, and duration as a
 772 control predictor, with participant and item included as random effects and
 773 random slopes of condition and transition type.

<i>Children</i>				
	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.57414	0.48576	-7.358	1.87e-13 ***
Age	0.02543	0.10260	0.248	0.8042
Type= <i>non-Question</i>	-0.81873	0.59985	-1.365	0.1723
Duration	4.17672	0.62446	6.689	2.25e-11 ***
Age*Type= <i>non-Question</i>	0.15116	0.13643	1.108	0.2679
Condition= <i>normal</i>	0.36710	0.43296	0.848	0.3965
Age*Condition= <i>normal</i>	0.12919	0.10227	1.263	0.2065
Condition= <i>normal</i> *	0.91059	0.72095	1.263	0.2066
Type= <i>non-Question</i>				
Age*Condition= <i>normal</i> * Type= <i>non-Question</i>	-0.37542	0.16963	-2.213	0.0269 *
Condition= <i>muffled</i>	-1.63429	0.86390	-1.892	0.0585 .
Age*Condition= <i>muffled</i>	0.39317	0.18907	2.080	0.0376 *
Condition= <i>muffled</i> *	1.77190	1.24864	1.419	0.1559
Type= <i>non-Question</i>				

Age*Condition= <i>muffled</i> *	-0.47057	0.28703	-1.639	0.1011
Type= <i>non-Question</i>				
Condition= <i>robot</i>	-0.26741	0.59071	-0.453	0.6508
Age*Condition= <i>robot</i>	0.13740	0.13568	1.013	0.3112
Condition= <i>robot</i> *	-1.02193	1.01227	-1.010	0.3127
Type= <i>non-Question</i>				
Age*Condition= <i>robot</i> *	0.08946	0.22349	0.400	0.6890
Type= <i>non-Question</i>				
<hr/>				
Adults				
	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.4557	0.7199	-4.800	1.58e-06 ***
Type= <i>non-Question</i>	0.4292	0.6089	0.705	0.480916
Duration	4.7500	1.2480	3.806	0.000141 ***
Condition= <i>normal</i>	1.2556	0.5633	2.229	0.025805 *
Condition= <i>normal</i> *	-0.9452	0.7631	-1.239	0.215475
Type= <i>non-Question</i>				
Condition= <i>muffled</i>	0.3349	0.8965	0.374	0.708692
Condition= <i>muffled</i> *	0.6627	1.2138	0.546	0.585108
Type= <i>non-Question</i>				
Condition= <i>robot</i>	1.5938	0.7208	2.211	0.027023 *
Condition= <i>robot</i> *	-1.1265	0.9109	-1.237	0.216201
Type= <i>non-Question</i>				

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

774 Children's anticipatory gaze switches showed an effect of gap duration
 775 ($\beta=4.18$, $SE=0.624$, $t=6.689$, $p<.001$), a two-way interaction of age and lan-
 776 guage condition (for prosody only speech compared to the no speech reference
 777 level; $\beta=0.393$, $SE=0.189$, $t=2.08$, $p<.05$), and a three-way interaction of
 778 age, transition type, and language condition (for normal speech compared to
 779 the no speech reference level; $\beta=-0.375$, $SE=0.17$, $t=-2.213$, $p<.05$). There
 780 were no significant effects of age or transition type alone (Table 3.3.1), with
 781 only a marginal effect of language condition (for prosody only compared to
 782 the no speech reference level; $\beta=-1.634$, $SE=0.864$, $t=-1.89$, $p=.06$)

783 Adults' anticipatory gaze switches showed effects of gap duration ($\beta=4.75$,
 784 $SE=1.248$, $t=3.806$, $p<.001$) and language condition (for normal speech
 785 $\beta=1.256$, $SE=0.563$, $t=2.229$, $p<.05$. and words only speech $\beta=1.594$, $SE=0.721$,
 786 $t=2.211$, $p<.05$ compared to the no speech reference level). There were no
 787 effects of transition type ($\beta=0.429$, $SE=0.609$, $t=0.705$, $p=.48$).

788 3.3.2. Random baseline comparison

789 Using the same technique described in experiment 1 (Section 2.2.2), we
 790 created and modeled 5,000 random permutations of participants' antici-
 791 patory gaze. Our baseline analyses revealed that none of the significant pre-
 792 dictors from models of the original, turn-related data (Table 5: Children)

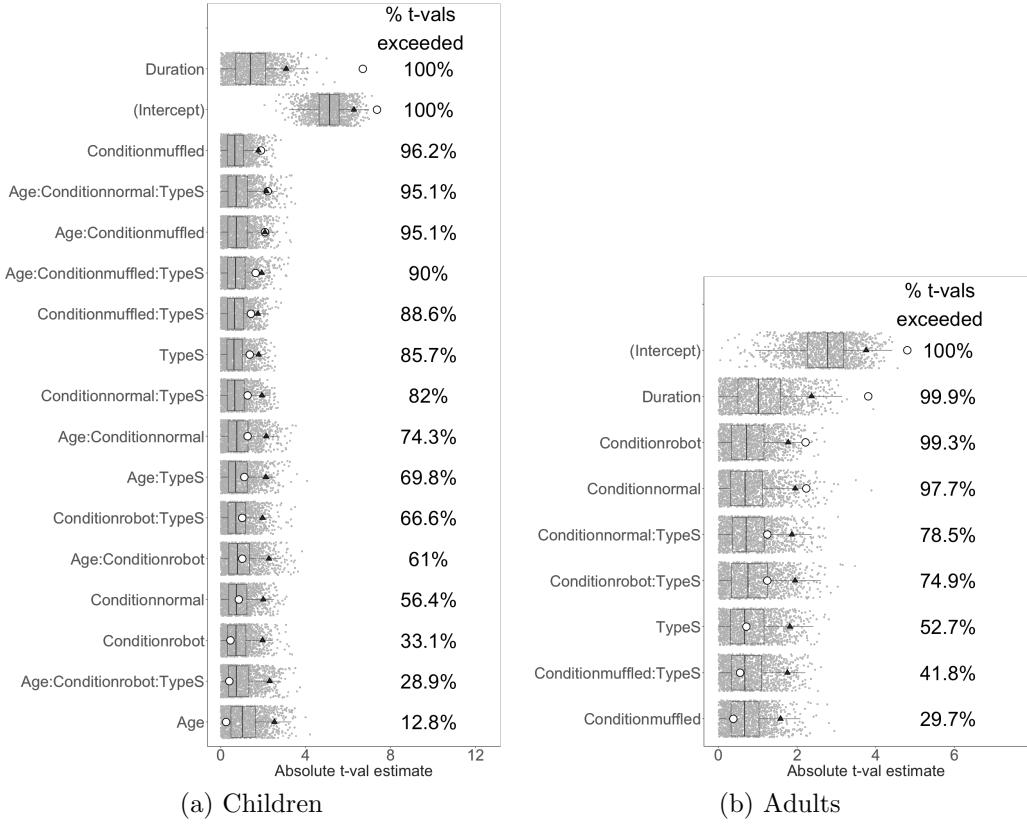


Figure 8: Random-permutation and original $|t\text{-values}|$ for predictors of children and adults' anticipatory gaze rates. Gray dots = random model estimates, White dots = original model estimates, Triangles = 95th percentile for each $t\text{-value}$ distribution.

793 can be explained by random looking. In the children's data, the original
 794 model's t -values for language condition (prosody only), gap duration, the
 795 two-way interaction of age and language condition (prosody only) and the
 796 three-way interaction of age, transition type, and language condition (normal
 797 speech) were all greater than 95% of the randomly permuted t -values (96.2%,
 798 100%, 95.1%, and 95.1%, respectively; Figure 8a). Similarly, the adults' data
 799 showed significant differentiation from the randomly permuted data for all
 800 originally significant predictors: gap duration and language condition for
 801 normal speech and words-only speech (greater than 100%, 96.8%, and 98.7%
 802 of random t -values, respectively; Figure 8b). The effects of language condi-
 803 tion and transition type for the real and randomly permuted data can also
 804 be observed in Figure 5. We excluded the output of random-permutation

models that did not ultimately converge to remove unreliable model results from our percentile calculations below (69% and 70% of models for children and adults, respectively).

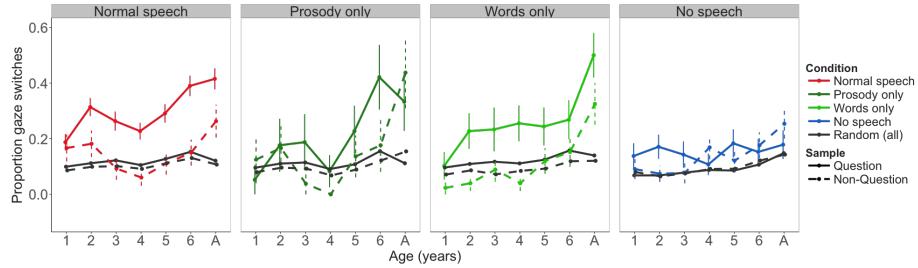


Figure 9: Anticipatory gaze rates across language condition and transition type for the real (blue, dark green, light green, and red) and randomly permuted (gray) data. Vertical bars represent the standard error.

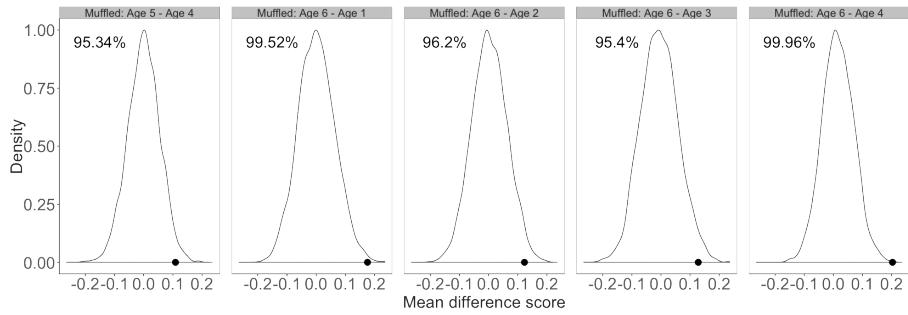


Figure 10: Significant pairwise comparisons of the prosody only-no speech linguistic condition effect, across ages, for the original data (black dots) and the 5,000 randomly permuted datasets (distribution).

Developmental effects. The model of the children’s data revealed two significant interactions with age, neither of which derived from random looking (Table 5; Figure 9). The first was a significant interaction of age and language condition (for prosody only compared to the no speech reference level), suggesting a different age effect between the two linguistic conditions. As in Experiment 1, we further explored the age effect by extracting the average

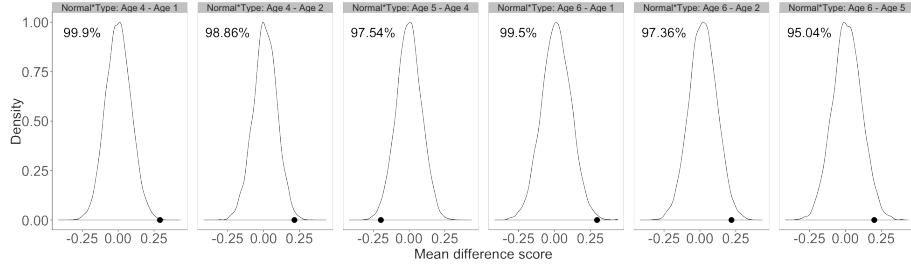


Figure 11: Significant pairwise comparisons of the normal speech-no speech language condition effect for questions status, across ages, for the original data (black dots) and the 5,000 randomly permuted datasets (distribution).

814 difference score over subjects for this age-language condition interaction in
 815 each random permutation of the data, making pairwise comparisons between
 816 the six age groups. Figure 10 shows the comparisons that demonstrate a sig-
 817 nificant difference in age for the no speech baseline vs. the prosody only
 818 condition, revealing that anticipatory gaze in the prosody only condition sig-
 819 nificantly improves with age, especially at ages 5 and 6 (all with difference
 820 scores greater than 95% of the random data scores; $p < .05$).

821 The second interaction with age was a three-way interaction of age, trans-
 822 sition type, and language condition (for normal speech compared to the no
 823 speech baseline). We again created pairwise comparisons of the average dif-
 824 ference scores for this age-transition type-language condition interaction in
 825 each random permutation of the data, with significant differences shown in
 826 Figure 11. These pairwise comparisons showed that the effect of transition
 827 type in the normal speech condition became larger with age, with significant
 828 improvements by age 4 over ages 1 and 2 (99.9% and 98.86%, respectively),
 829 by age 5 over age 4 (97.54%), and by age 6 over ages 1, 2, and 5 (99.5%,
 830 97.36%, and 95.04%), all significantly different from chance ($p < .05$).

831 3.4. Discussion

832 As in Experiment 1, children and adults spontaneously tracked the turn
 833 structure of the conversations, making anticipatory gaze switches at an above-
 834 chance rate across all ages but not in all conditions (Table ??; Figure 9).
 835 In the normal speech condition, however, when children had access to full
 836 linguistic information, they made more anticipatory gaze switches than ex-
 837 pected by chance even at age one. Children’s anticipatory gaze rates were

838 affected by gap duration, plus interactions of language condition, transition
839 type and age (Table 5), none of which could be explained by a baseline of
840 random gaze switching (Figure 8a).

841 Two of the three linguistic conditions affected children's anticipatory gaze
842 switches, compared to the no-speech reference level: prosody-only speech and
843 normal speech. Prosody only speech resulted in fewer anticipatory gazes over-
844 all, compared to the no-speech baseline, though the effect was only marginal
845 ($p=.06$). However, prosody-only speech *did* significantly differ from the no
846 speech condition in its interaction with age: whereas age gains were negligible
847 in the no speech condition, 5- and 6-year-olds in the prosody-only condition
848 showed significantly more anticipatory gaze switches than younger children
849 (Figure 10), going from below- and at-chance anticipatory gaze rates to being
850 well above chance (Figure 9).

851 The normal speech condition showed a significantly more gains in antici-
852 patory gaze for questions with age compared to the no-speech condition. As
853 in Experiment 1, children consistently made more anticipatory gazes after
854 hearing questions when they had access to lexical material. But now, with
855 a larger age span in Experiment 2, we can begin to see a developmental
856 path for the question effect. At least for normal speech, greater anticipatory
857 looking for questions is not present from the start: 4-year-olds are the first
858 to make significant gains on 1- and 2-year-olds, but then there are further
859 significant gains at age 5, and again age 6 (Figure 11). While children at ages
860 five and six showed adult-like differentiation of questions and non-questions
861 in the normal speech condition, 1-year-olds had nearly identical switch rates
862 for the two transition types. This suggests that the participants' tendency
863 to make more anticipatory switches for questions emerges after age one and
864 continues developing, with rapid gains in ages three through age six. Fi-
865 nally, children showed a straightforward effect of gap duration (Table 5), as
866 in Experiment 1.

867 Adults' anticipatory gaze rates were affected by gap duration and two
868 of the language conditions (Table 5), none of which could be explained by
869 a baseline of random gaze switching (Figure 8b). Adults made more antici-
870 patory switches overall for the normal speech and words-only conditions
871 compared to the no-speech condition, falling in-line with past work showing
872 that adults primarily use lexical information in making predictions about
873 upcoming turn structure (De Ruiter et al., 2006). Though adults did make
874 more anticipatory switches for questions than non-questions on average in
875 the two lexical conditions (Figure 9), the effect of transition type was not

876 significant in either (Table 5), unlike Experiment 1. Like children, adults
877 also showed a straightforward effect of gap duration.

878 *3.4.1. Summary*

879 Children and adults both showed more anticipatory gaze switches with
880 increased linguistic information, but only for a subset of the linguistic con-
881 ditions and transition types. We had expected to see the most anticipatory
882 switches in the *normal* condition and the least anticipatory switches in the
883 *no speech* condition because they contained the most and least linguistic in-
884 formation, respectively. We had also expected to replicate our finding from
885 Experiment 1 that questions result in more anticipatory switches than non-
886 questions, with the added hypothesis that the question effect is driven by
887 lexicosyntactic cues. We additionally anticipated an overall increase in an-
888 ticipatory switches with age. Finally, since the development of prosodic skills
889 partially precedes the development of lexicosyntax, we expected to see more
890 switches in the *prosody only* condition compared to the *words only* condition
891 in the youngest age groups.

892 In fact, children and adults did show more anticipatory gaze switching
893 in the normal condition compared to the no-speech condition, but for chil-
894 dren this effect only emerged in the form of anticipations following question
895 turns, which increased with age faster than it did with no linguistic cues at
896 all (i.e., in the no-speech condition). Adults also showed significantly higher
897 anticipatory switch rates in the words only condition but no effects of transi-
898 tion type, alone or within linguistic conditions. Taken together, these results
899 partly replicate the findings from Experiment 1: participants make more an-
900 ticipatory switches when they have access to lexical information and, when
901 they do, tend to make more anticipatory switches for questions compared to
902 non-questions.

903 We had anticipated significant gains in anticipatory switching with age,
904 but children only showed significant developmental increases in the prosody
905 only condition and the normal condition (for question transitions). Rather
906 than showing an early advantage for prosody over lexical information, chil-
907 dren did not show significant improvement (or even above-chance perfor-
908 mance) until age five. In contrast, their anticipatory gaze switches were, on
909 average, already above chance at age one in the normal condition (with both
910 lexical and prosodic information) and at age three in the words only condi-
911 tion (with lexical information). These findings do not support an early role
912 for prosody in children’s spontaneous predictions about upcoming turn struc-

ture. On the contrary, their predictions were best when lexical information was present, especially when following a question.

However, these results do not support the idea that lexical information is sufficient for children (as is proposed for adults; De Ruiter et al., 2006). On the contrary, children showed significant gains in the normal speech condition, but *not* in the words only condition. Notably, the normal speech condition, in addition to having both lexical and prosodic information, is also the one most likely to occur in our participants' daily lives (compared to muffled or robotic speech) and therefore may have an extra advantage over the other conditions.

Finally, on average, participants at all ages made anticipatory gaze switches more often than would be expected by chance in the *no speech* condition. One interpretation of this finding is that participants were not using linguistic information to predict upcoming turns at all—the only cue to turn taking in the *no speech* condition was the alternating mouth movements of the conversing puppets. But this interpretation is not compatible with the effects of linguistic information that do emerge in the adult and child data: An advantage for lexical cues shaped by transition type and age.

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 to expand our age range to capture more developmental change. We found that effects of linguistic processing and age *were* present in the dataset, but were primarily related to the question effect, and primarily occur in the normal speech condition. The relation between prosody and lexicosyntax in young children's anticipations is still not clear: The results of the prosody only manipulation suggest that children do not reliably use low-pass filtered speech to anticipate upcoming turn structure until age five, but on the other hand, the words only condition did not show significant overall differences from the no-speech baseline, including its interactions with age or transition type. If anything, the results do not support the idea that there is a strong early advantage for prosody in turn predictions.

4. General Discussion

Children begin to develop conversational turn-taking skills long before their first words (Bateson, 1975; Hilbrink et al., 2015; Jaffé et al., 2001; Snow, 1977). As they acquire language, they also acquire the information needed to make accurate predictions about upcoming turn structure. Until

949 recently, we have had very little data on how children weave language into
950 their already-existing turn-taking behaviors.

951 In two experiments investigating children’s anticipatory gaze to upcom-
952 ing speakers, we found evidence that turn prediction develops early in child-
953 hood and that spontaneous predictions are primarily driven by participants’
954 expectation of an immediate response in the next turn. In making predic-
955 tions about upcoming turn structure, children used a combination lexical
956 and prosodic cues; neither prosodic nor lexical cues alone were sufficient to
957 support increased anticipatory gaze. We also found no early advantage for
958 prosody over lexisyntax, and instead found that children were unable to
959 make above-chance anticipatory gazes in the prosody only condition until age
960 five. We discuss these findings with respect to the role of linguistic cues in
961 predictions about upcoming turn structure, the importance of questions in
962 spontaneous predictions about conversation, and children’s developing com-
963 petence as conversationalists.

964 *4.1. Predicting turn structure with linguistic cues*

965 Prior work with adults has found a consistent and critical role for lex-
966 icosyntax in predicting upcoming turn structure (De Ruiter et al., 2006;
967 Magyari and De Ruiter, 2012), with the role of prosody still under debate
968 (Duncan, 1972; Ford and Thompson, 1996; Torreira et al., 2015). Knowing
969 that children comprehend more about prosody than lexisyntax early on
970 (see introduction; also see Speer and Ito, 2009 for a review), we thought it
971 possible that young children would instead show an advantage for prosody
972 over lexisyntax. Our results suggest that, on the contrary, when presented
973 with *only* prosodic information, children’s spontaneous predictions about up-
974 coming turn structure are limited until age five. Do children instead rely on
975 lexisyntax, as adults are proposed to do?

976 We found no evidence that, for children, lexisyntax alone is “suffi-
977 cient” (equal to full linguistic information) for spontaneous turn prediction
978 (De Ruiter et al., 2006, pg. 531): In both experiments, children’s perfor-
979 mance was best in conditions when they had access to the full linguistic
980 signal. Adults on the other hand, showed significant gains in anticipatory
981 gaze switching in both conditions with lexical cues, compared to a no speech
982 baseline.

983 Participants appeared to use linguistic cues to make more anticipatory
984 gaze switches in both experiments. Questions are often marked both prosod-
985 ically and lexisyntactically for their speech act status. The close link be-

986 tween prosodic and syntactic structure makes it difficult to tease apart how
987 predictive processing in one domain is distinct from the other in turn pre-
988 diction. That being said, compared to prosodic contours (e.g., final rising
989 intonation), lexicosyntactic cues like *wh*-words, *do*-insertion, and subject-
990 auxiliary inversion are frequent, categorical, and early-occurring in the utter-
991 ance. Children may therefore have an easier time picking out and interpreting
992 lexical cues to questionhood on the fly. Question turns started yielding sig-
993 nificantly more anticipatory gazes by age 3–4 in the normal speech condition
994 of Experiment 2, by which time children frequently hear and use a variety of
995 polar *wh*-questions (Clark, 2009). Furthermore, while lexicosyntactic ques-
996 tion cues were available on every instance of *wh*- and *yes/no* questions in our
997 stimuli, prosodic question cues were only salient on *yes/no* questions and,
998 even then, the mapping of prosodic contour to speech act (e.g., high final
999 rises for polar questions) is far from one-to-one.

1000 *4.1.1. The question effect*

1001 In both experiments, anticipatory looking was primarily driven by ques-
1002 tion transitions, a pattern that had not been previously reported in other
1003 observer gaze studies, on children or adults (Keitel et al., 2013; Hirvenkari,
1004 2013; Tice and Henetz, 2011). Questions make a speaker switch immediately
1005 relevant, helping the listener to predict with high certainty what will happen
1006 next (i.e., an answer from the addressee). As mentioned, linguistic cues to
1007 questionhood also often occur early in the utterance, giving observers time
1008 to plan a gaze switch. Because the form of a question can constrain the
1009 type of response that will come next (e.g., a location after a *where* question),
1010 questions can even help listeners predict specific upcoming content in the
1011 next turn.

1012 Our results suggest that question turns start significantly affecting chil-
1013 dren’s predictions between ages three and four, but further testing with finer-
1014 grained age samples with stimuli focused on specific communicative acts and
1015 linguistic cues is needed to better sketch out the developmental trajectory of
1016 this effect. For example, the effect size and age of emergence might differ by
1017 question type (e.g., *wh*- vs. *yes-no*) or the location of question-identifying
1018 cues within the unfolding utterance (early vs. late), but we would not be
1019 able to see it given the current design.

1020 Prior work on children’s acquisition of questions indicates that they may
1021 already have some understanding about question-answer sequences by the
1022 time they begin to speak: Questions make up approximately one third of

1023 the utterances children hear, before and after the onset of speech, and even
1024 into their preschool years, even though the types and complexity of questions
1025 change throughout development (Casillas et al., In press; Fitneva, 2012; Hen-
1026 ning et al., 2005; Shatz, 1979).¹⁰ For the first few years, many of the questions
1027 directed to children are “test” questions—questions that the caregiver already
1028 has the answer to (e.g., “What does a cat say?”), but this changes as children
1029 get older. Questions help caregivers to get their young children’s attention
1030 and to ensure that information is in common ground, even if the responses are
1031 non-verbal or infelicitous (Bruner, 1985; Fitneva, 2012; Snow, 1977). So, in
1032 addition to having a special interactive status (for adults and children alike),
1033 questions are a core characteristic of many caregiver-child interactions, mo-
1034 tivating a general benefit for questions in turn structure anticipation.

1035 All that being said, our current data do not tell us what it is about ques-
1036 tions that makes children and adults more likely to anticipatorily switch their
1037 gaze to addressees. Other request formats, such as imperatives, compliments,
1038 and complaints make a response from the addressee highly likely in the next
1039 turn (Schegloff, 2007). Rhetorical and tag questions, on the other hand, take
1040 a similar form to prototypical polar questions, but often do not require an
1041 answer. So, though it is clear that participants anticipated responses more
1042 often for questions than non-questions, we do not yet know whether their
1043 predictive action is limited to turns formatted as questions or is generally
1044 applicable to turn structures that project an immediate response from the
1045 addressee.

1046 Much recent work on prediction during turn taking has focused on partic-
1047 ipants’ use of linguistic cues to predict the end of the current turn (Torreira
1048 et al., 2015; Magyari and De Ruiter, 2012; De Ruiter et al., 2006; Ford and
1049 Thompson, 1996; Duncan, 1972), in some cases finding that lexical informa-
1050 tion was sufficient for prediction. Our current results suggest that, sponta-
1051 neous predictions are instead driven by predictions about what is *beyond* the
1052 end of the current turn—that questions are sufficient for prediction.

1053 To integrate these results and understand how listeners make predictions
1054 for turn taking, it is crucial to account for the participants’ role in the in-
1055 teraction. The results we present here are based on predictions about third-
1056 party conversation, which enables participants to follow interactions with no

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

chance of actually participating. Although recent work has shown that similar anticipatory eye gazes do occur in spontaneous conversation (Holler and Kendrick, 2015), we do not yet know if the same question advantage occurs, or which linguistic cues seem to drive it. It is possible that participants track conversation for cues to upcoming opportunities/obligations to speak and, when one is found, focus more attention on predicting the precise timing of the current turn's end. To answer these questions we will need to innovate first-person prediction measures that can be used in real-time interaction.

4.1.2. Early competence for turn taking?

One of the core aims of our study was to test whether children show an early competence for turn taking, as is proposed by studies of spontaneous mother-infant interaction and theories about the mechanisms underlying human interaction (Hilbrink et al., 2015; Levinson, 2006). Although children and adults' gaze patterns were quite similar in Experiment 1, children in Experiment two showed change with age in their anticipatory looking, even in the most prototypical speech condition ("normal" speech), and 6-year-olds still did not achieve adult-like levels of anticipatory looking on average. This may indicate that children rely more in non-verbal cues in anticipating turn transitions or, alternatively, that adults are better at flexibly adapting to the turn-relevant cues present at any moment. When they *did* have full linguistic information, children on average still made above-chance anticipatory gazes to responders, suggesting that at least some turn-taking competence for third-party conversation is present in infancy, though it is a weak effect.

Taken together, the data suggest that turn-taking skills do begin to emerge in infancy, but that adult-like competence is not achieved until much later, mirroring results from children's spontaneous interactions with their caregivers (infants: Hilbrink et al., 2015; Jaffe et al., 2001; Snow, 1977; older children: Casillas et al., In press; Garvey, 1984; Ervin-Tripp, 1979). It is possible, however, that first-person measures of children's predictions would show more frequent turn structure anticipations at younger ages.

4.2. Limitations and future work

Although Experiments 1 and 2 have offered new findings regarding the relation to speech act in online turn predictions and the emergence of those predictions in the first six years, there remain at least two major limitations to our work: speech naturalness and participant role.

1092 Following prior work (De Ruiter et al., 2006; Keitel et al., 2013), we used
1093 phonetically manipulated speech in Experiment 2, resulting in speech sounds
1094 that children don't usually hear in their natural environment. Many prior
1095 studies have used phonetically-altered speech with infants and young children
1096 (cf. Jusczyk, 2000), but almost none of them have done so in a conversational
1097 context. Children could have had trouble processing the *words only* and
1098 *prosody only* conditions because they were unfamiliar, and not just because
1099 they had less linguistic information available. Future work could instead
1100 carefully script or cross-splice parts of turns to control for the presence or
1101 absence of linguistic cues for turn transition.

1102 The prediction measure we present is based on an observer's view of
1103 third-party conversation but, because participants' role in the interaction
1104 could affect their online predictions about turn taking, an ideal experimen-
1105 tal measure would capture first-person behavior. First-person measures of
1106 spontaneous turn prediction will be key to revealing how participants dis-
1107 tribute their attention over linguistic and non-verbal cues while taking part
1108 in everyday interaction, the implications of which relate to theories of online
1109 language processing for both language learning and everyday talk.

1110 4.3. Conclusions

1111 Conversation plays a central role in children's language learning. It is
1112 the driving force behind what children say and what they hear. Adults use
1113 linguistic information to accurately predict turn structure in conversation,
1114 which facilitates their online comprehension and allows them to respond rel-
1115 evantly and on time. In the current study we have investigated how children's
1116 predictions about turn structure changes as their linguistic skills develop in
1117 the first six years. We found that, although some basic knowledge about turn
1118 taking exists at age one, the integration of linguistic cues into children's pre-
1119 dictions about turn structure takes time and is, like adults, primarily driven
1120 by sequences of action—in our case, by questions and their answers.

1121 Acknowledgements

1122 We gratefully acknowledge the parents and children at Bing Nursery
1123 School and the Children's Discovery Museum of San Jose. This work was
1124 supported by an ERC Advanced Grant to Stephen C. Levinson (269484-
1125 INTERACT), NSF graduate research and dissertation improvement fellow-
1126 ships to the first author, and a Merck Foundation fellowship to the second

1127 author. Earlier versions of these data and analyses were presented to con-
1128 ference audiences (Casillas and Frank, 2012, 2013). We also thank Tania
1129 Henetz, Francisco Torreira, Stephen C. Levinson, Eve V. Clark, and the
1130 First Language Acquisition group at Radboud University for their feedback
1131 on earlier versions of this work.

1132 **References**

- 1133 Bakker, M., Kochukhova, O., von Hofsten, C., 2011. Development of social
1134 perception: A conversation study of 6-, 12-and 36-month-old children.
1135 *Infant Behavior and Development* 34, 363–370.
- 1136 Barr, D.J., Levy, R., Scheepers, C., Tily, H.J., 2013. Random effects structure
1137 for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory*
1138 and *Language* 68, 255–278.
- 1139 Bates, D., Maechler, M., Bolker, B., Walker, S., 2014. lme4:
1140 Linear mixed-effects models using Eigen and S4. URL:
1141 <https://github.com/lme4/lme4>/<http://lme4.r-forge.r-project.org/>.
1142 [Computer program] R package version 1.1-7.
- 1143 Bateson, M.C., 1975. Mother-infant exchanges: The epigenesis of conver-
1144 sational interaction. *Annals of the New York Academy of Sciences* 263,
1145 101–113.
- 1146 Bergelson, E., Swingley, D., 2013. The acquisition of abstract words by young
1147 infants. *Cognition* 127, 391–397.
- 1148 Bloom, K., 1988. Quality of adult vocalizations affects the quality of infant
1149 vocalizations. *Journal of Child Language* 15, 469–480.
- 1150 Boersma, P., Weenink, D., 2012. Praat: doing phonetics by computer. URL:
1151 <http://www.praat.org>. [Computer program] Version 5.3.16.
- 1152 Bögels, S., Magyari, L., Levinson, S.C., 2015. Neural signatures of response
1153 planning occur midway through an incoming question in conversation. *Sci-
1154 entific Reports* 5.
- 1155 Bruner, J., 1985. Child's talk: Learning to use language. *Child Language
1156 Teaching and Therapy* 1, 111–114.

- 1157 Bruner, J.S., 1975. The ontogenesis of speech acts. *Journal of Child Language*
1158 2, 1–19.
- 1159 Carlson, R., Hirschberg, J., Swerts, M., 2005. Cues to upcoming swedish
1160 prosodic boundaries: Subjective judgment studies and acoustic correlates.
1161 *Speech Communication* 46, 326–333.
- 1162 Casillas, M., Bobb, S.C., Clark, E.V., In press. Turn taking, timing, and
1163 planning in early language acquisition. *Journal of Child Language* .
- 1164 Casillas, M., Frank, M.C., 2012. Cues to turn boundary prediction in adults
1165 and preschoolers. *Proceedings of SemDial* .
- 1166 Casillas, M., Frank, M.C., 2013. The development of predictive processes
1167 in children's discourse understanding, in: *Proceedings of the 35th Annual*
1168 *Meeting of the Cognitive Science Society*.
- 1169 Clark, E.V., 2009. First language acquisition. Cambridge University Press.
- 1170 De Ruiter, J.P., Mitterer, H., Enfield, N.J., 2006. Projecting the end of
1171 a speaker's turn: A cognitive cornerstone of conversation. *Language* 82,
1172 515–535.
- 1173 De Vos, C., Torreira, F., Levinson, S.C., 2015. Turn-timing in signed con-
1174 versations: coordinating stroke-to-stroke turn boundaries. *Frontiers in*
1175 *Psychology* 6.
- 1176 Dingemanse, M., Torreira, F., Enfield, N., 2013. Is “Huh?” a universal word?
1177 Conversational infrastructure and the convergent evolution of linguistic
1178 items. *PloS one* 8, e78273.
- 1179 Duncan, S., 1972. Some signals and rules for taking speaking turns in con-
1180 versations. *Journal of Personality and Social Psychology* 23, 283.
- 1181 Ervin-Tripp, S., 1979. Children's verbal turn-taking, in: Ochs, E., Schieffelin,
1182 B.B. (Eds.), *Developmental Pragmatics*. Academic Press, New York, pp.
1183 391–414.
- 1184 Fitneva, S., 2012. Beyond answers: questions and children's learning, in:
1185 de Ruiter, J.P. (Ed.), *Questions: Formal, Functional, and Interactional*
1186 *Perspectives*. Cambridge University Press, Cambridge, UK, pp. 165–178.

- 1187 Ford, C.E., Thompson, S.A., 1996. Interactional units in conversation: Syntactic, intonational, and pragmatic resources for the management of turns.
1188 Studies in Interactional Sociolinguistics 13, 134–184.
- 1190 Garvey, C., 1984. Children's Talk. volume 21. Harvard University Press.
- 1191 Gísladóttir, R., Chwilla, D., Levinson, S.C., 2015. Conversation electrified: ERP correlates of speech act recognition in underspecified utterances. PloS one 10, e0120068.
- 1194 Griffin, Z.M., Bock, K., 2000. What the eyes say about speaking. Psychological science 11, 274–279.
- 1196 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.
- 1199 Hedberg, N., Sosa, J.M., Görgülü, E., Mameni, M., 2010. The prosody and meaning of Wh-questions in American English, in: Speech Prosody 2010—Fifth International Conference.
- 1202 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.
- 1204 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.
- 1207 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.
- 1210 von Hofsten, C., Uhlig, H., Adell, M., Kochukhova, O., 2009. How children with autism look at events. Research in Autism Spectrum Disorders 3, 556–569.
- 1213 Holler, J., Kendrick, K.H., 2015. Unaddressed participants' gaze in multi-person interaction. Frontiers in Psychology 6.
- 1215 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

- 1217 development. Monographs of the Society for Research in Child Develop-
1218 ment. JSTOR.
- 1219 Johnson, E.K., Jusczyk, P.W., 2001. Word segmentation by 8-month-olds:
1220 When speech cues count more than statistics. *Journal of Memory and*
1221 *Language* 44, 548–567.
- 1222 Jusczyk, P.W., 2000. *The Discovery of Spoken Language*. MIT press.
- 1223 Kamide, Y., Altmann, G., Haywood, S.L., 2003. The time-course of predic-
1224 tion in incremental sentence processing: Evidence from anticipatory eye
1225 movements. *Journal of Memory and Language* 49, 133–156.
- 1226 Keitel, A., Prinz, W., Friederici, A.D., Hofsten, C.v., Daum, M.M., 2013.
1227 Perception of conversations: The importance of semantics and intonation
1228 in childrens development. *Journal of Experimental Child Psychology* 116,
1229 264–277.
- 1230 Lemasson, A., Glas, L., Barbu, S., Lacroix, A., Guilloux, M., Remeuf, K.,
1231 Koda, H., 2011. Youngsters do not pay attention to conversational rules:
1232 is this so for nonhuman primates? *Nature Scientific Reports* 1.
- 1233 Levelt, W.J., 1989. *Speaking: From intention to articulation*. MIT press.
- 1234 Levinson, S.C., 2006. On the human “interaction engine”, in: Enfield, N.,
1235 Levinson, S. (Eds.), *Roots of human sociality: Culture, cognition and*
1236 *interaction*. Oxford: Berg, pp. 39–69.
- 1237 Levinson, S.C., 2013. Action formation and ascriptions, in: Stivers, T., Sid-
1238 nell, J. (Eds.), *The Handbook of Conversation Analysis*. Wiley-Blackwell,
1239 Malden, MA, pp. 103–130.
- 1240 Magyari, L., Bastiaansen, M.C.M., De Ruiter, J.P., Levinson, S.C., 2014.
1241 Early anticipation lies behind the speed of response in conversation. *Jour-*
1242 *nal of Cognitive Neuroscience* 26, 2530–2539.
- 1243 Magyari, L., De Ruiter, J.P., 2012. Prediction of turn-ends based on antici-
1244 pation of upcoming words. *Frontiers in Psychology* 3:376, 1–9.
- 1245 Masataka, N., 1993. Effects of contingent and noncontingent maternal stimu-
1246 lation on the vocal behaviour of three-to four-month-old Japanese infants.
1247 *Journal of Child Language* 20, 303–312.

- 1248 Mehler, J., Jusczyk, P., Lambertz, G., Halsted, N., Bertoni, J., Amiel-
1249 Tison, C., 1988. A precursor of language acquisition in young infants.
1250 *Cognition* 29, 143–178.
- 1251 Morgan, J.L., Saffran, J.R., 1995. Emerging integration of sequential and
1252 suprasegmental information in preverbal speech segmentation. *Child De-
1253 velopment* 66, 911–936.
- 1254 Nazzi, T., Ramus, F., 2003. Perception and acquisition of linguistic rhythm
1255 by infants. *Speech Communication* 41, 233–243.
- 1256 R Core Team, 2014. R: A Language and Environment for Statistical Com-
1257 puting. R Foundation for Statistical Computing. Vienna, Austria. URL:
1258 <http://www.R-project.org>. [Computer program] Version 3.1.1.
- 1259 Ratner, N., Bruner, J., 1978. Games, social exchange and the acquisition of
1260 language. *Journal of Child Language* 5, 391–401.
- 1261 Ross, H.S., Lollis, S.P., 1987. Communication within infant social games.
1262 *Developmental Psychology* 23, 241.
- 1263 Rossano, F., Brown, P., Levinson, S.C., 2009. Gaze, questioning and culture,
1264 in: Sidnell, J. (Ed.), *Conversation Analysis: Comparative Perspectives*.
1265 Cambridge University Press, Cambridge, pp. 187–249.
- 1266 Sacks, H., Schegloff, E.A., Jefferson, G., 1974. A simplest systematics for the
1267 organization of turn-taking for conversation. *Language* 50, 696–735.
- 1268 Schegloff, E.A., 2007. Sequence organization in interaction: Volume 1: A
1269 primer in conversation analysis. Cambridge University Press.
- 1270 Shatz, M., 1979. How to do things by asking: Form-function pairings in
1271 mothers' questions and their relation to children's responses. *Child Devel-
1272 opment* 50, 1093–1099.
- 1273 Shi, R., Melancon, A., 2010. Syntactic categorization in French-learning
1274 infants. *Infancy* 15, 517–533.
- 1275 Snow, C.E., 1977. The development of conversation between mothers and
1276 babies. *Journal of Child Language* 4, 1–22.

- 1277 Speer, S.R., Ito, K., 2009. Prosody in first language acquisition—Acquiring
1278 intonation as a tool to organize information in conversation. *Language and*
1279 *Linguistics Compass* 3, 90–110.
- 1280 Stivers, T., Enfield, N.J., Brown, P., Englert, C., Hayashi, M., Heinemann,
1281 T., Hoymann, G., Rossano, F., De Ruiter, J.P., Yoon, K.E., et al., 2009.
1282 Universals and cultural variation in turn-taking in conversation. *Proceedings of the National Academy of Sciences* 106, 10587–10592.
- 1284 Stivers, T., Rossano, F., 2010. Mobilizing response. *Research on Language*
1285 and Social Interaction 43, 3–31.
- 1286 Takahashi, D.Y., Narayanan, D.Z., Ghazanfar, A.A., 2013. Coupled oscillator
1287 dynamics of vocal turn-taking in monkeys. *Current Biology* 23, 2162–2168.
- 1288 Thorgrímsson, G., Fawcett, C., Liszkowski, U., 2015. 1- and 2-year-olds'
1289 expectations about third-party communicative actions. *Infant Behavior*
1290 and Development 39, 53–66.
- 1291 Tice (Casillas), M., Henetz, T., 2011. Turn-boundary projection: Looking
1292 ahead, in: *Proceedings of the 33rd Annual Meeting of the Cognitive Science*
1293 *Society*.
- 1294 Torreira, F., Bögels, S., Levinson, S.C., 2015. Intonational phrasing is neces-
1295 sary for turn-taking in spoken interaction. *Journal of Phonetics* 52, 46–57.
- 1296 Wittenburg, P., Brugman, H., Russel, A., Klassmann, A., Sloetjes, H., 2006.
1297 Elan: a professional framework for multimodality research, in: *Proceedings*
1298 of LREC.
- 1299 In all of the following plots, the gray dots represent the randomly per-
1300 muted data's model estimates for the value listed (beta or standard error),
1301 the white dots represent the model estimates from the original data, and the
1302 triangles represent the 95th percentile for each distribution being shown.

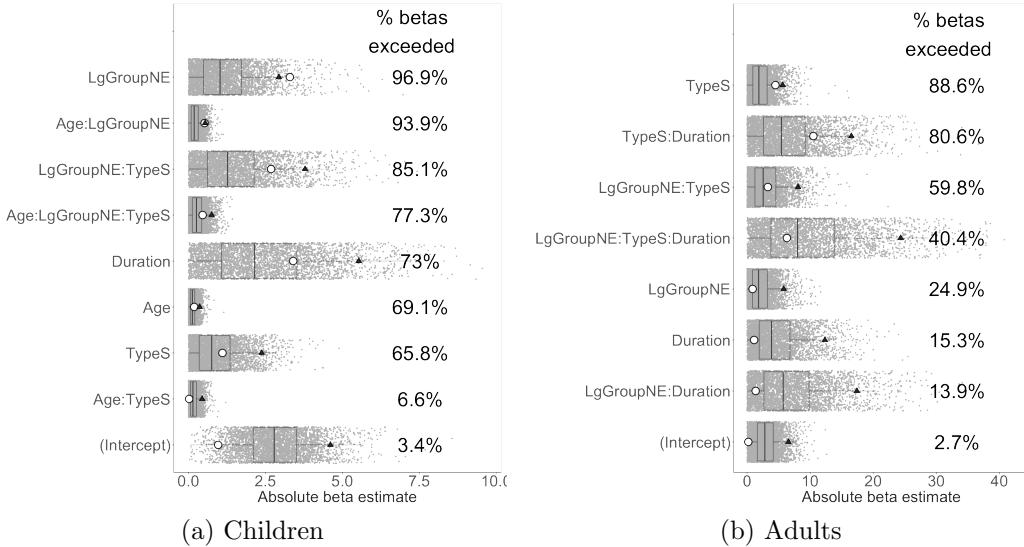


Figure .1: Random-permutation and original $|\beta\text{-values}|$ for predictors of children and adults' anticipatory gaze rates in Experiment 1.

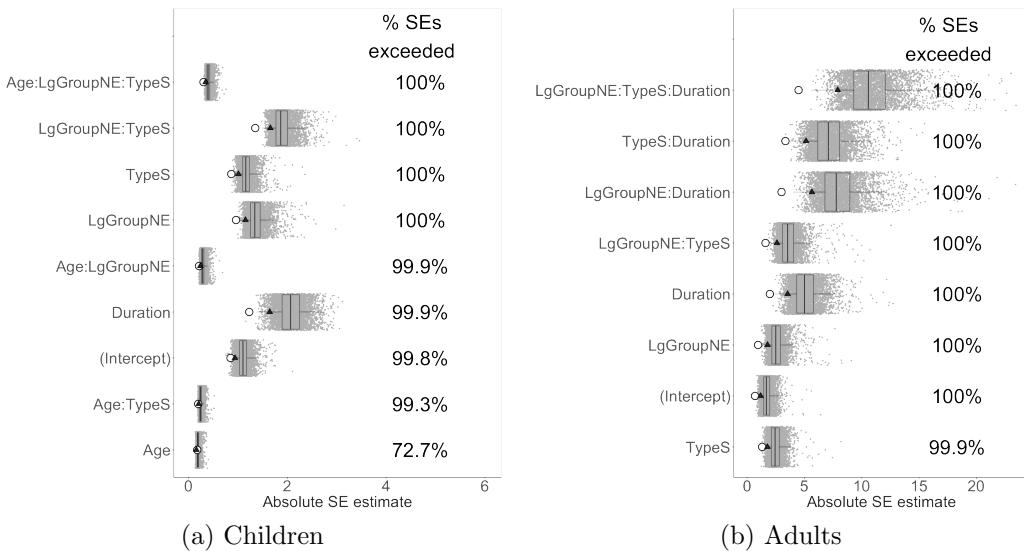


Figure .2: Random-permutation and original $|SE\text{-values}|$ for predictors of children and adults' anticipatory gaze rates in Experiment 1.

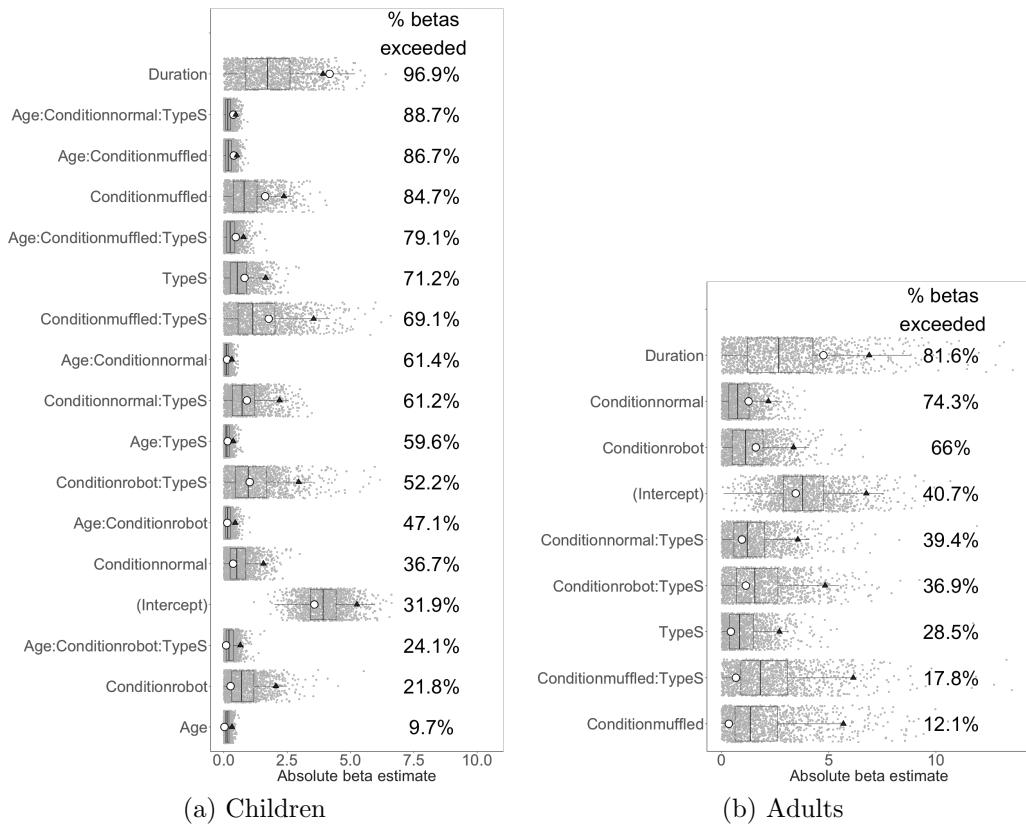


Figure .3: Random-permutation and original $|\beta\text{-values}|$ for predictors of children and adults' anticipatory gaze rates in Experiment 2.

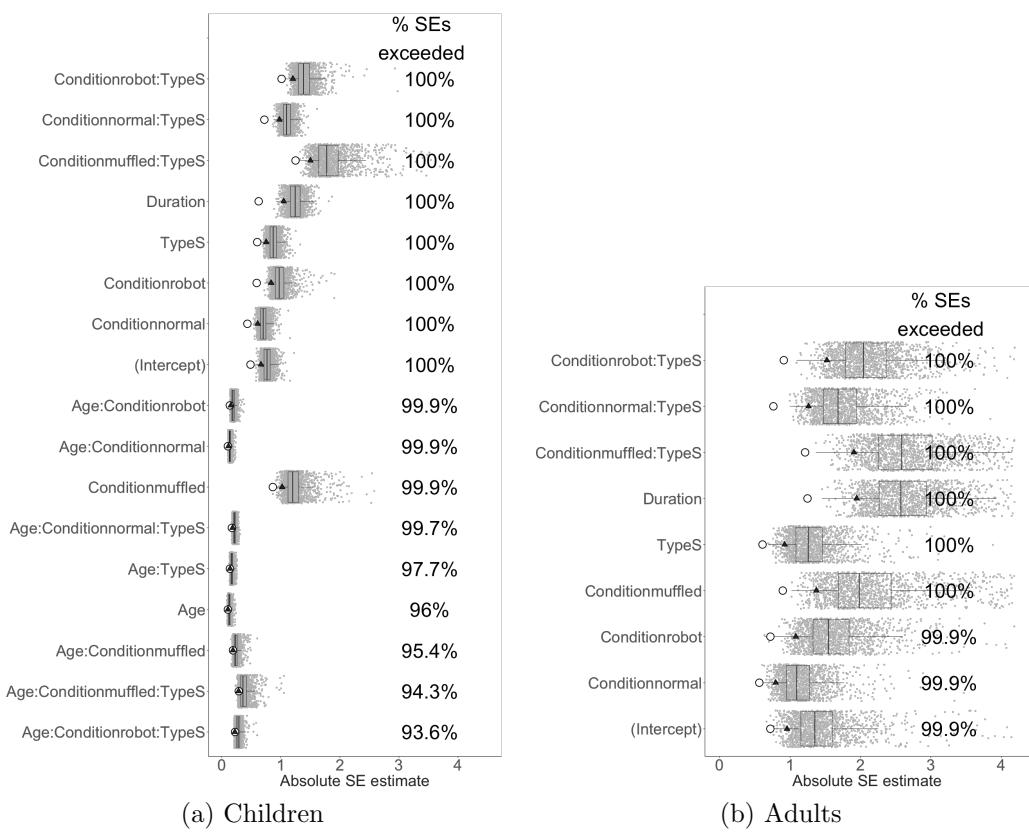


Figure .4: Random-permutation and original $|SE\text{-values}|$ for predictors of children and adults' anticipatory gaze rates in Experiment 2..