

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 by
¹⁷ caregivers, has different timing from adult turn taking, and lacks semantic
¹⁸ content (Hilbrink et al., 2015; Jaffe et al., 2001). However, children’s early,
¹⁹ turn-structured social interactions are presumably a critical precursor to their
²⁰ later conversational turn taking: Early non-verbal interactions likely estab-
²¹ lish the protocol by which children come to use language with others. How
²² do children integrate linguistic knowledge with these preverbal turn-taking
²³ 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 to coordinate with one another through sequences of contingent action.
39

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. Conversational turns come grouped into semantically-contingent sequences of
43 action. The groups can span turn-by-turn exchanges (e.g., simple question–
44 response, “How are you?”–“Fine.”) or sequence-by-sequence exchanges (e.g.,
45 reciprocals, “How are you?”–“Fine, and you?”–“Great!”).

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

60 Adults have a lot of information at their disposal to help make accurate
61 predictions about upcoming turn content. Lexical, syntactic, and prosodic
62 information (e.g., *wh*- words, subject-auxiliary inversion, and list intonation)
63 can all inform addressees about upcoming linguistic structure (De Ruiter
64 et al., 2006; Duncan, 1972; Ford and Thompson, 1996; Torreira et al., 2015).
65 Non-verbal cues (e.g., gaze, posture, and pointing) often appear at turn-
66 boundaries and can sometimes act as late indicators of an upcoming speaker
67 switch (Rossano et al., 2009; Stivers and Rossano, 2010). Additionally, the
68 sequential context of a turn can make it clear what will come next: An-
69 swers after questions, thanks or denial after compliments, et cetera (Schegloff,
70 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
129 non-linear pattern. Infants’ contingent vocalizations in the first few months
130 of life show very fast timing (though with a lot of vocal overlap) that, by nine
131 months, slows down considerably, only gradually speeding up again after 12
132 months (Hilbrink et al., 2015). Taking turns with brief transitions between
133 speakers is difficult for children; while their avoidance of overlap is nearly
134 adult-like by nine months, the timing of their non-overlapped responses stays
135 much longer than the 200 msec standard for the next few years (Casillas
136 et al., In press; Garvey, 1984; Ervin-Tripp, 1979). This puzzling pattern is
137 likely due to their linguistic development: Taking turns on time is easier
138 when the response is a simple vocalization rather than a linguistic utterance.
139 Integrating language into the turn-taking system may be one major factor in
140 children’s delayed responses (Casillas et al., In press).

141 While children, like adults, could use linguistic cues in the ongoing turn to
142 make predictions about upcoming turn structure, studies of early linguistic

development point to a possible early advantage for prosody over lexicosyntax in children's turn-taking predictions. 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). Four- to five-month-olds also prefer pauses in speech to be inserted at prosodic boundaries, and by 6 months they can start using prosodic markers to pick out sub-clausal syntactic units, both of which are useful for extracting turn structure from ongoing speech (Jusczyk et al., 1995; Soderstrom et al., 2003). In comparison, children show at best a very limited lexical inventory before their first birthday (Bergelson and Swingley, 2013; Shi and Melancon, 2010).

Keitel and colleagues (2013) were one of the first to explore how children use linguistic cues to predict upcoming turn structure. They asked 6-, 12-, 24-, 36-month-old, and adult participants to watch short videos of conversation and tracked their eye movements at points of speaker change. They showed their participants two types of conversation videos—one normal and one with flattened pitch (i.e., with flattened intonation contours)—to test the role of intonation in participants' anticipatory predictions about upcoming speech. Comparing children's anticipatory gaze frequency to a random baseline, they found that only 36-month-olds and adults made anticipatory gaze switches more often than expected by chance. Among those, only 36-month-olds were affected by a lack of intonation contours, leading Keitel and colleagues to conclude that children's ability to predict upcoming turn structure relies on their ability to comprehend the stimuli lexicosemantically. They also suggest that intonation takes a secondary role in turn prediction, following their and other studies findings with adults (e.g., De Ruiter et al., 2006), but only after children acquire more sophisticated, adult-like language comprehension systems (sometime after 36 months in their data).

Although the Keitel et al. (2013) study constitutes a substantial advance over previous work in this domain, it has its own limitations. Because these limitations directly inform our own study design, we review them in some detail. First, their estimates of baseline gaze frequency ("random" in their terminology) were not random. Instead, they used gaze switches during ongoing speech as a baseline. Ongoing speech is perhaps the period in which switching is least likely to occur (Hirvenkari et al., 2013), thus maximizing

181 chances of finding a difference between gaze frequency at turn transitions
182 and their baseline rate. A more conservative baseline would be to compare
183 participants' looking behavior at turn transitions to their looking behavior
184 during randomly selected windows of time throughout the stimulus, including
185 turn transitions. We follow this conservative approach in our work.

186 Second, the conversation stimuli Keitel et al. (2013) used were somewhat
187 unusual. The average gap between turns was 900 msec, which is much longer
188 than typical adult timing, where gaps average around 200 msec (Stivers et al.,
189 2009). The speakers in the videos were also asked to minimize their move-
190 ments while performing a scripted and adult-directed conversation, which
191 would have created a somewhat unnatural stimulus. Additionally, in order
192 to produce more naturalistic conversation, it would have been ideal to local-
193 ize the sound sources for the two voices in the video (i.e., to have the voices
194 come out of separate left and right speakers). But both voices were recorded
195 and played back on the same audio channel, which may have made it more
196 difficult to distinguish the two talkers (again, we attempt to address these
197 issues in our current study). Despite these minor methodological issues, the
198 Keitel et al. (2013) study still demonstrates intriguing age-based differences
199 in children's ability to predict upcoming turn structure. Our current work
200 thus takes this paradigm as a starting point.²

201 Our goal in the current study is to find out when children begin to make
202 predictions about upcoming turn structure and to understand how their pre-
203 dictions are affected by linguistic cues across development. We present two
204 experiments in which we measure children's anticipatory gaze to respon-
205 ders while watching conversation videos with natural (people using English
206 vs. non-English; Experiment 1) and non-natural (puppets with phonetically
207 manipulated speech; Experiment 2) control over the presence of lexical and
208 prosodic cues. We tested children across a wide range of ages (Experiment
209 1: 3–5 years; Experiment 2: 1–6 years), with adult control participants in
210 each experiment.

211 Our findings suggest that children and adults use linguistic cues to make
212 predictions about upcoming turn structure, but that they primarily do so
213 to predict speaker transitions after a question. We also find that children
214 make more predictions than expected by chance, even at age XX, but that it
215 takes several years before they recruit linguistic cues to anticipate responses

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

216 more often after questions. We found no direct evidence of an early prosody
217 advantage in children’s anticipations and, further, no evidence that prosodic
218 or lexical cues alone can substitute for their combination in the full linguistic
219 signal, as is proposed for adults (De Ruiter et al., 2006).

220 2. Experiment 1

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

233 2.1. Methods

234 2.1.1. Participants

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

244 2.1.2. Materials

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



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

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

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

267 After recording, we combined the audio and video files by hand, and
268 cropped each recording to the 30-second interval with the most turn activity.
269 Because we recorded the conversations in stereo, the male and female voices

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

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

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

301 2.1.3. *Procedure*

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

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

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

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

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

342 *2.1.4. Data preparation and coding*

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

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

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

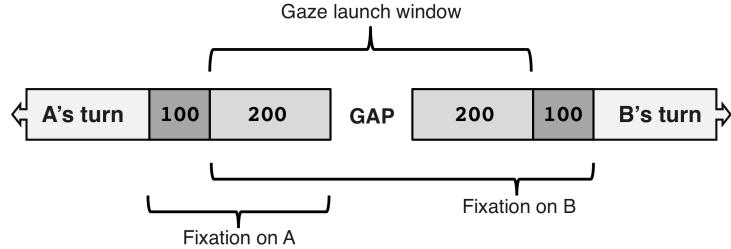


Figure 2: Schematic summary of criteria for anticipatory gaze shifts from speaker A to speaker B during a turn transition.

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

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

394 2.2. Results

395 Participants looked at the screen most of the time during video playback
 396 (81% and 91% on average for children and adults, respectively). They pri-
 397 marily kept their eyes on the person who was currently speaking in both
 398 English and non-English videos: They gazed at the current speaker between
 399 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.

400 20% of the time (Table 1). Even three-year-olds looked more at the current
 401 speaker than anything else, whether the videos were in a language they could
 402 understand or not. Children looked at the current speaker less than adults
 403 did during the non-English videos. Despite this, their looks to the addressee
 404 did not increase substantially in the non-English videos, indicating that their
 405 looks away were probably related to boredom rather than confusion about
 406 ongoing turn structure. Overall, participants' pattern of gaze to current
 407 speakers demonstrated that they performed basic turn tracking during the
 408 videos, regardless of language. Figure 3 shows participants' anticipatory gaze
 409 rates across age, language condition, and transition type, in the real data and
 410 the random baseline (described below).

411 *2.2.1. Statistical models*

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

⁵Because each non-English language was represented by a single stimulus, we cannot treat individual languages as factors. Gaze behavior might be best for non-native languages

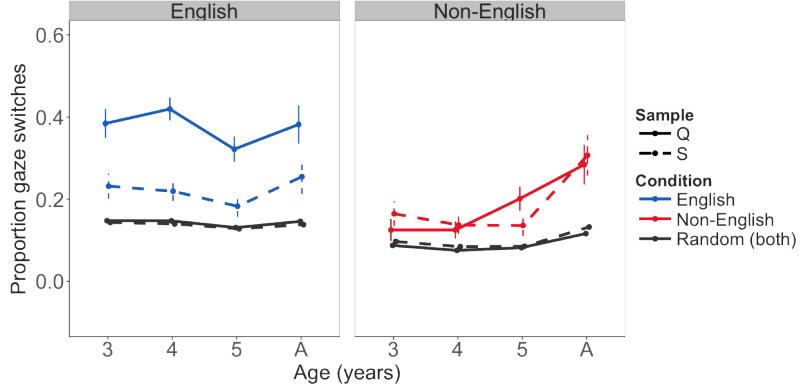


Figure 3: Anticipatory gaze rates across language condition and transition type for the real (red and blue) and randomly permuted baseline (gray). Vertical bars represent the standard error.

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

Children’s anticipatory gaze switches showed effects of language condition ($\beta=-3.29$, $SE=0.961$, $t=-3.43$, $p<.001$) and gap duration ($\beta=3.4$, $SE=1.229$,

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

⁶The models we report are all qualitatively unchanged by the exclusion of their random slopes. We have left the random slopes in because of minor participant-level variation in the predictors modeled.

<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.

431 $t=2.77$, $p<.01$) with additional effects of an age-by-language condition in-
 432 teraction ($\beta=0.52$, $SE=0.212$, $t=2.46$, $p<.05$) and a language condition-by-
 433 transition type interaction ($\beta=2.68$, $SE=1.35$, $t=1.99$, $p<.05$). There were
 434 no significant effects of age or transition type alone ($\beta=-0.18$, $SE=0.175$,
 435 $t=-1.04$, $p=.3$ and $\beta=-1.10$, $SE=0.865$, $t=-1.27$, $p=.2$, respectively).

436 Adults' anticipatory gaze switches shows an effect of transition type ($\beta=$
 437 4.5 , $SE=1.314$, $t=-3.42$, $p<.001$) and significant interactions between lan-
 438 guage condition and transition type ($\beta=3.3$, $SE=1.61$, $t=2.05$, $p<.05$) and
 439 transition type and gap duration ($\beta=10.51$, $SE=3.346$, $t=3.141$, $p<.01$).

440 2.2.2. Random baseline comparison

441 We estimated the probability that these patterns were the result of ran-
 442 dom looking by running the same regression models on participants' real
 443 eye-tracking data, only this time calculating their anticipatory gaze switches
 444 with respect to randomly permuted turn transition windows. This process

445 involved: (1) Randomizing the order and temporal placement of the analysis
 446 windows within each stimulus (Figure 4; “analysis window” is defined
 447 in Figure 2), thereby randomly redistributing the analysis windows across
 448 the eye-tracking signal, (2) re-running each participant’s eye tracking data
 449 through switch identification (described in Section 2.1.4), this time using
 450 the randomly permuted analysis windows, and (3) modeling the anticipatory
 451 gazes from the randomly permuted data with the same statistical models we
 452 used for the original data (Section 2.2.1; Table 2). Importantly, although
 453 the onset time of each transition was shuffled within the eye-tracking signal,
 454 the other intrinsic properties of each turn transition (e.g., prior speaker iden-
 455 tity, transition type, gap duration, language condition, etc.) stayed constant
 456 across each random permutation.

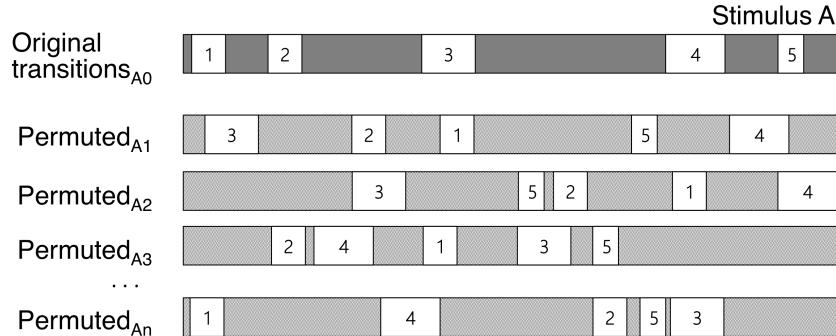


Figure 4: Example of analysis window permutations for a stimulus with five turn transitions. The windows were ± 300 msec around the inter-turn gap.

457 This procedure effectively de-links participants’ gaze data from the turn
 458 structure in the original stimulus, thereby allowing us to compare turn-
 459 related (original) and non-turn-related (randomly permuted) looking behav-
 460 ior using the same eye movements. The resulting anticipatory gazes from the
 461 randomly permuted analysis windows represent an average anticipatory gaze
 462 rate over all possible starting points: a random baseline. By running the real
 463 and randomly permuted data sets through identical statistical models, we
 464 can also estimate how likely it is that predictor effects in the original data
 465 (e.g., the effect of language condition; Table 2) arose from random looking.

466 We completed this baseline procedure on 5,000 random permutations of
 467 the original turn transition analysis windows and compared the t -values from
 468 each predictor in the original models (Table 2) to the distribution of t -values

469 for each predictor in the 5,000 models of the randomly permuted datasets.⁷
470 We could then test whether significant effects from the original statistical
471 models differed from the random baseline by calculating the proportion of
472 random data t -values exceeded by the original t -value for each predictor,
473 using the absolute value of all t -values for a two-tailed test. For example,
474 children's original "language condition" t -value was $|3.429|$, which is greater
475 than 99.9% of all $|t\text{-value}|$ estimates from the randomly-permuted data mod-
476 els (i.e., $p=.001$). This leads us to conclude that the effect of language
477 condition in the original model was highly unlikely to be the result of random
478 gaze shifting.

479 We excluded the output of random-permutation models that gave con-
480 vergence warnings in order to remove unreliable estimates from our analyses
481 (non-converging models were 22.4% and 24.4% of all models for children and
482 adults, respectively; see the Supplementary Materials for more information
483 on model exclusion).

484 Our baseline analyses revealed that none of the significant predictors from
485 models of the original, turn-related data can be explained by random look-
486 ing (see Figures A.1a and A.1b for the full distributional data). For the
487 children's data, the original t -values for language condition, gap duration,
488 the age-language condition interaction, and the language condition-transition
489 type interaction were all greater than 95% of t -values for the randomly per-
490 muted data (99.9%, 95.5%, 99.4%, and 96%, respectively). Similarly, the
491 adults' data showed significant differentiation from the randomly permuted
492 data for two of the three originally significant predictors—transition type and
493 the transition type-gap duration interaction (greater than 99.9% and 99.7%
494 of random t -values, respectively)—with marginal differentiation for the in-
495 teraction of language condition and transition type (greater than 94.7% of
496 random t -values).

497 *Developmental effects.* The original statistical model of the children's real
498 data revealed a significant interaction of age and language condition (Table
499 2) that was highly unlikely to have come from random looking (Figure 3).
500 To further explore this effect, we compared the average effect of language

⁷We report t -values rather than beta estimates because the standard errors in the ran-
domly 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 Supple-
mentary Materials.

501 condition for each age group: We extracted the average difference score for
502 the two language conditions (English minus non-English) for each subject,
503 computing an overall average for each random permutation of the data. For
504 each random permutation, we then made pairwise comparisons of the average
505 difference scores across participant age groups, finding that, while 3- and 4-
506 year olds showed similarly large effects of language condition, 5-year-olds
507 showed a significantly smaller effect of language condition, compared to both
508 younger age groups. In other words, the difference in the effect of language
509 condition for 5-year-olds compared to younger children was larger than would
510 be expected by chance in 99.52% of the randomly permuted data sets for 3-
511 year-olds and 99.96% of the data sets for 4-year-olds—differences of $p < .01$
512 and $p < .001$, respectively (see Figure B.1 for difference score distributions).

513 When does spontaneous turn prediction emerge developmentally during
514 natural speech? To test whether the youngest age group (3-year-olds) already
515 exceeded chance in their anticipatory gaze switches, we used two-tailed t -tests
516 to compare their real gaze rates to the random baseline in the most familiar
517 speech condition for our participants: English. We found that three-year-
518 olds made anticipatory gaze switches significantly above chance, both when
519 all transitions were considered ($t(22.824) = -4.147, p < .001$) and for question
520 transitions alone ($t(21.677) = -5.268, p < .001$).

521 *2.3. Discussion*

522 Children and adults spontaneously tracked the turn structure of the con-
523 versations, making anticipatory gaze switches at an above-chance rate across
524 all ages and conditions (Table 1; Figure 3). Children’s anticipatory gaze rates
525 were affected by language condition, transition type, age, and gap duration
526 (Table 2), none of which could be explained by a baseline of random gaze
527 switching (Figure A.1a).

528 Language condition (English vs. non-English) affected children’s antici-
529 pations in two ways (Table 2; Figure 3). First, children made more antici-
530 patory switches overall in English videos, compared to non-English videos.
531 This effect suggests that lexical access is important for children’s ability
532 to anticipate upcoming turn structure; children had no lexical access to the
533 speech in the non-English videos, though they did have access to (non-native)
534 prosodic cues and non-verbal behavior. This finding is consistent with prior
535 work on turn-end prediction in adults (De Ruiter et al., 2006; Magyari and
536 De Ruiter, 2012) and children (Keitel et al., 2013). Second, children system-
537 atically made more anticipatory switches after hearing a question compared

538 to a non-question, but only in the English condition, suggesting that, when
539 children have access to lexical cues, they are more likely to make an anticipa-
540 tory gaze switch if they can expect an immediate response from the addressee.
541 If so, then children's (and adults') attention to lexical cues for turn taking
542 may primarily be in monitoring the signal for cues to questionhood (e.g.,
543 subject-auxiliary inversion, *wh*-words, etc.).

544 Children's anticipatory gaze switches were also affected by their age, but
545 only in the non-English videos: 3- and 4-year-olds made many more anticipa-
546 tory switches when watching videos in English compared to non-English, but
547 this effect of language condition had attenuated significantly by age 5 (Table
548 2; Figure 3; Figure B.1). This interaction suggests that the 5-year-olds were
549 able to leverage anticipatory cues in the non-English videos in a way that
550 3- and 4-year-olds could not, possibly by shifting more attention to the non-
551 native prosodic or non-verbal cues. Prior work on children's turn-structure
552 anticipation proposed that children's turn-end predictions rely primarily on
553 lexicosyntactic structure (and not, e.g., prosody) as they get older (Keitel
554 et al., 2013). The current results suggest more flexibility in children's pre-
555 dictions; when they do not have access to lexical information, older children
556 are likely to find alternative cues to turn taking behavior.

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

560 Adults' anticipatory gaze rates were also affected by transition type, lan-
561 guage condition, and gap duration (Table 2), none of which could be easily
562 explained by a baseline of random gaze switching (Figure A.1b). Like chil-
563 dren, adults made more anticipatory switches after hearing questions com-
564 pared to non-questions, suggesting that anticipation mattered more to them
565 when an immediate response was expected. Also like children, the advantage
566 for questions was driven by lexical access such that adults must have relied
567 on lexicosyntactic cues to questionhood in picking out turns that potentially
568 require an immediate response, though this effect was only marginally di-
569 vergent from the distribution of randomly permuted data ($p=.053$; Figure
570 A.1b). Finally, adults' anticipation rates were also affected by gap dura-
571 tion, but more so for questions than non-questions (Table 2), suggesting that
572 adults were less likely overall to make switches at non-questions, and so did
573 not benefit from extra time to do so.

574 2.3.1. *Summary*

575 Children and adults' predictions alike were benefited by access to lexical
576 information (English) and speech act status (questionhood), suggesting that
577 linguistic cues, particularly lexical ones, facilitate their spontaneous predic-
578 tions about upcoming turn structure through the identification of turns with
579 immediate responses. Children's anticipatory gaze rates for questions and
580 non-questions in English was stable across ages and comparable to adult be-
581 havior (Figure 3), suggesting that they can identify questions in native stimuli
582 with adult-like competence by age three. Although participants' ability to
583 recognize questions was facilitated by lexical access (i.e., English vs. non-
584 English), the prosody in the non-English videos was non-native, and so the
585 experimental design can not conclusively show which linguistic cues children
586 relied on in the English videos to identify question turns. Relatedly, though
587 lexical access clearly facilitated participants' anticipatory gaze rate, it was
588 not necessary for participants—especially adults—in order to exceed chance
589 switching rates (Figure 3), suggesting that participants use non-lexical cues
590 (e.g., prosody, non-verbal behavior) to make anticipatory eye movements at
591 least some of the time.

592 Interestingly, adults and children both were strongly affected by transition
593 type, in that they made more anticipatory switches after hearing questions,
594 compared to non-questions. Even in the English videos, when participants
595 had full access to linguistic cues, their rates of anticipation were relatively
596 low—in fact, comparable to the non-English videos—unless the turn was a
597 question (Figure 3). Prior work using online, metalinguistic tasks has shown
598 that participants can use linguistic cues to accurately predict upcoming turn
599 ends (Torreira et al., 2015; ?; De Ruiter et al., 2006). The current results
600 add a new dimension to our understanding of how listeners make predic-
601 tions about turn ends: Both children and adults spontaneously monitor the
602 linguistic structure of unfolding turns for cues to upcoming responses.

603 Children and adults behaved relatively similarly in this first experiment,
604 and our language manipulation (English vs. non-English) was too coarse to
605 comment on when children begin to use specific linguistic cues (e.g., prosody
606 vs. lexicosyntax). We would instead need to directly compare lexicosyn-
607 tactic and prosodic cues in the participants' native language, controlling for
608 the presence of non-verbal cues. To see the emergence of anticipatory gaze
609 switching we would also need to include younger children, since participants
610 already reliably made anticipatory gaze switches at age three. We follow up

611 on both of these ideas in Experiment 2.

612 **3. Experiment 2**

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

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

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

638 *3.1. Methods*

639 *3.1.1. Participants*

640 We recruited 27 undergraduate adults and 129 children between ages 1;0–
641 6;11 to participate in our experiment. We recruited our child participants
642 from the Children's Discovery Museum in San Jose, California, targeting ap-
643 proximately 20 children for each of the six 1-year age groups (range=20–23).
644 All participants were native English speakers, though some parents (N=27)
645 reported that their child heard a second (and sometimes third) language at

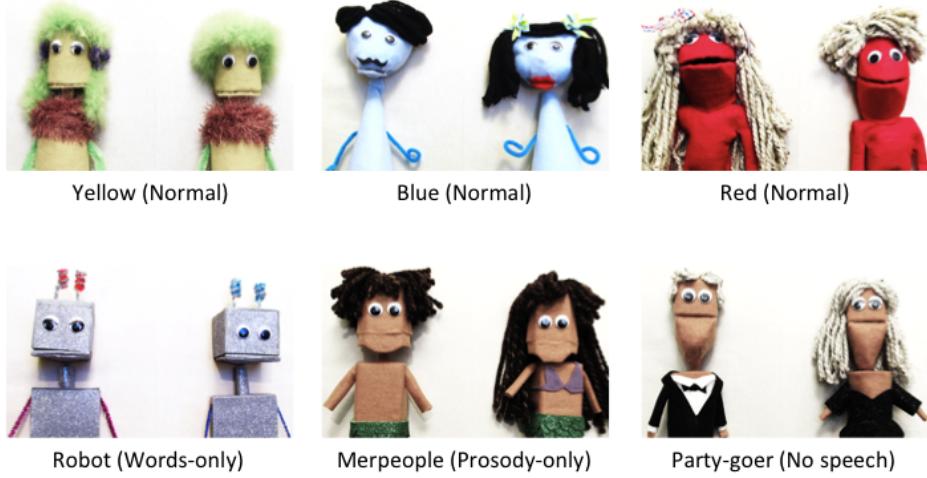


Figure 5: The six puppet pairs (and associated audio conditions). Each pair was linked to three distinct conversations from the same condition across the three experiment versions.

646 home. None of the adult participants reported fluency in a second language.
 647 We ran Experiment 2 at a local children’s museum because it gave us access
 648 to children with a more diverse range of ages.

649 *3.1.2. Materials*

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

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

664 activity. The 20-second clips were then phonetically manipulated and used
665 in the final video stimuli.

666 *Audio Manipulation.* We created four versions of each audio clip: *nor-*
667 *mal*, *words only*, *prosody only*, and *no speech*. That is, one version with a full
668 linguistic signal (*normal*), and three with incomplete linguistic information
669 (hereafter “limited cue” conditions). The *normal* clips were the unmanipu-
670 lated, original audio clips.

671 The *words only* clips were manipulated to have robot-like speech: We
672 flattened the intonation contours to each talker’s average pitch (F0) and
673 we reset the duration of every nucleus and coda to each talker’s average
674 nucleus and coda duration.⁸ We made duration and pitch manipulations
675 using PSOLA resynthesis in Praat (Boersma and Weenink, 2012). Thus,
676 the *words only* versions of the audio clips had no pitch or durational cues
677 to upcoming turn boundaries, but did have intact lexicosyntactic cues (and
678 residual phonetic correlates of prosody, e.g., intensity).

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

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

693 Finally, because low-pass filtering removes significant acoustic energy, the
694 *prosody only* clips were much quieter than the other three conditions. Our
695 last step was to downscale the intensity of the audio tracks in the three other
696 conditions to match the volume of the *prosody only* clips. We referred to the
697 conditions as “normal”, “robot”, “mermaid”, and “birthday party” speech

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

698 when interacting with participants.

699 *Video recordings.* We created puppet video recordings to match the ma-
700 nicipated 20-second audio clips. The puppets were minimally expressive;
701 the experimenter could only control the opening and closing of their mouths;
702 their head, eyes, arms, and body stayed still. Puppets were positioned look-
703 ing forward to eliminate shared gaze as a cue to turn structure (Thorgrímsson
704 et al., 2015). We took care to match the puppets' mouth movements to the
705 syllable onsets as closely as possible, specifically avoiding any mouth move-
706 ment before the onset of a turn. We then added the manipulated audio clips
707 to the puppet video recordings by hand.

708 We used three pairs of puppets used for the *normal* condition—‘red’,
709 ‘blue’ and ‘yellow’—and one pair of puppets for each limited cue condition:
710 “robots”, “merpeople”, and “party-goers” (Figure 8). We randomly assigned
711 half of the conversation topics (‘birthday cake’, ‘pets’, and ‘breakfast’) to the
712 *normal* condition, and half to the limited cue conditions (‘riding bikes’, ‘rainy
713 days’, and ‘the library’). We then created three versions of the experiment,
714 so that each of the six puppet pairs was associated with three different con-
715 versation topics across the different versions of the experiment (18 videos
716 in total). We ensured that the position of the talkers (left and right) was
717 counterbalanced in each version by flipping the video and audio channels as
718 needed.

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

726 There were 85 turn transitions for analysis after excluding transitions
727 longer than 550 msec and shorter than 90 msec. The remaining turn tran-
728 sitions had slightly more questions than non-question ($N=50$ and $N=35$, re-
729 spectively), with transitions distributed somewhat evenly across conditions
730 (keeping in mind that there were three *normal* videos and only one lim-
731 ited cue video for each experiment version): *normal* ($N=36$), *words only*
732 ($N=13$), *prosody only* ($N=17$), and *no speech* ($N=19$). Inter-turn gaps for
733 questions (mean=365, median=427) were longer than those for non-questions
734 (mean=302, median=323) on average, but gap duration was overall com-
735 parable across conditions: *normal* (mean=334, median=321), *words only*

736 (mean=347, median=369), *prosody only* (mean=365, median=369), and *no
737 words* (mean=319, median=329). The longer gaps for question transitions
738 could give them an advantage because our anticipatory measure includes
739 shifts initiated during the gap between turns (Figure 2).

740 *3.2. Procedure*

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

749 *3.2.1. Data preparation and coding*

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

754 *3.3. Results*

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

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).

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.

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.

768 3.3.1. Statistical models

769 We identified anticipatory gaze switches for all 85 usable turn transitions,
770 and analyzed them for effects of language condition, transition type,
771 and age with two mixed-effects logistic regressions. We again built separate
772 models for children and adults because effects of age were only pertinent to
773 the children’s data. The child model included condition (normal/prosody
774 only/words only/no speech; with no speech as the reference level), transition
775 type (question vs. non-question), age (1, 2, 3, 4, 5, 6; numeric), and duration
776 of the inter-turn gap (in seconds) as predictors, with full interactions between
777 language condition, transition type, and age. We again included the dura-
778 tion of the inter-turn gap as a control predictor and added random effects of
779 item (turn transition) and participant, with random slopes of transition type
780 for participants. The adult model included condition, transition type, their
781 interactions, and duration as a control predictor, with participant and item
782 included as random effects and random slopes of condition and transition

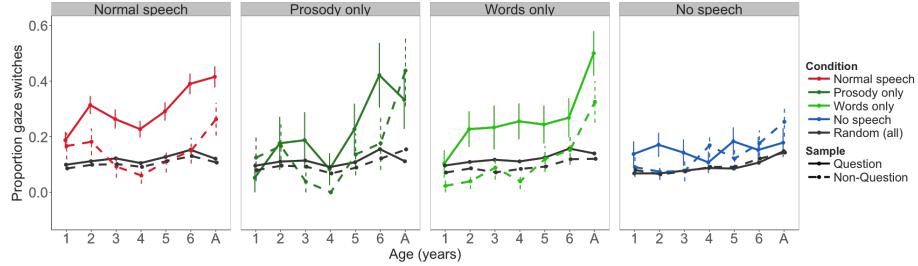


Figure 6: Anticipatory gaze rates across language condition and transition type for the real (blue, dark green, light green, and red) and randomly permuted baseline (gray). Vertical bars represent the standard error.

783 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=non-Question	-0.81873	0.59985	-1.365	0.1723
Duration	4.17672	0.62446	6.689	2.25e-11 ***
Age*Type=non-Question	0.15116	0.13643	1.108	0.2679
Condition=normal	0.36710	0.43296	0.848	0.3965
Age*Condition=normal	0.12919	0.10227	1.263	0.2065
Condition=normal*	0.91059	0.72095	1.263	0.2066
Type=non-Question				
Age*Condition=normal/*	-0.37542	0.16963	-2.213	0.0269 *
Type=non-Question				
Condition=muffled	-1.63429	0.86390	-1.892	0.0585 .
Age*Condition=muffled	0.39317	0.18907	2.080	0.0376 *
Condition=muffled*	1.77190	1.24864	1.419	0.1559
Type=non-Question				
Age*Condition=muffled*	-0.47057	0.28703	-1.639	0.1011
Type=non-Question				
Condition=robot	-0.26741	0.59071	-0.453	0.6508
Age*Condition=robot	0.13740	0.13568	1.013	0.3112
Condition=robot*	-1.02193	1.01227	-1.010	0.3127
Type=non-Question				
Age*Condition=robot*	0.08946	0.22349	0.400	0.6890
Type=non-Question				
<i>Adults</i>				
	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.4557	0.7199	-4.800	1.58e-06 ***
Type=non-Question	0.4292	0.6089	0.705	0.480916
Duration	4.7500	1.2480	3.806	0.000141 ***
Condition=normal	1.2556	0.5633	2.229	0.025805 *
Condition=normal*	-0.9452	0.7631	-1.239	0.215475
Type=non-Question				

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.

784 Children's anticipatory gaze switches showed an effect of gap duration
 785 ($\beta=4.18$, $SE=0.624$, $t=6.689$, $p<.001$), a two-way interaction of age and
 786 language condition (for *prosody only* speech compared to the *no speech* refer-
 787 ence level; $\beta=0.393$, $SE=0.189$, $t=2.08$, $p<.05$), and a three-way interaction
 788 of age, transition type, and language condition (for *normal* speech compared
 789 to the *no speech* reference level; $\beta=-0.375$, $SE=0.17$, $t=-2.213$, $p<.05$). There
 790 were no significant effects of age or transition type alone (Table 3.3.1), with
 791 only a marginal effect of language condition (for *prosody only* compared to
 792 the *no speech* reference level; $\beta=-1.634$, $SE=0.864$, $t=-1.89$, $p=.06$)

793 Adults' anticipatory gaze switches showed effects of gap duration ($\beta=4.75$,
 794 $SE=1.248$, $t=3.806$, $p<.001$) and language condition (for *normal* speech
 795 $\beta=1.256$, $SE=0.563$, $t=2.229$, $p<.05$ and *words only* speech $\beta=1.594$, $SE=0.721$,
 796 $t=2.211$, $p<.05$ compared to the *no speech* reference level). There were no
 797 effects of transition type ($\beta=0.429$, $SE=0.609$, $t=0.705$, $p=.48$).

798 3.3.2. Random baseline comparison

799 Using the same technique described in Experiment 1 (Section 2.2.2), we
 800 created and modeled 5,000 random permutations of participants' antici-
 801 patory gaze (see Figures A.4a and A.4b for the full distributional data). Our
 802 baseline analyses revealed that none of the significant predictors from mod-
 803 els of the original, turn-related data (Table 5: Children) can be explained by
 804 random looking. In the children's data, the original model's *t*-values for lan-
 805 guage condition (*prosody only*), gap duration, the two-way interaction of age
 806 and language condition (*prosody only*) and the three-way interaction of age,
 807 transition type, and language condition (*normal* speech) were all greater than
 808 95% of the randomly permuted *t*-values (96.2%, 100%, 95.1%, and 95.1%,
 809 respectively). Similarly, the adults' data showed significant differentiation
 810 from the randomly permuted data for all originally significant predictors: gap
 811 duration and language condition for *normal* speech and words-only speech
 812 (greater than 100%, 96.8%, and 98.7% of random *t*-values, respectively). The

813 effects of language condition and transition type for the real and randomly
814 permuted data can also be observed in Figure 3.

815 As before, we excluded the output of random-permutation models that
816 resulted in convergence warnings in order to remove unreliable model esti-
817 mates from our analyses (non-converging models made up 69% and 70% of
818 models for children and adults, respectively; see the Supplementary Materials
819 for more information on model exclusion).

820 *Developmental effects.* The model of the children’s data revealed two signif-
821 icant interactions with age, neither of which derived from random looking
822 (Table 5; Figure 6). The first was a significant interaction of age and lan-
823 guage condition (for *prosody only* compared to the *no speech* reference level),
824 suggesting a different age effect between the two linguistic conditions. As
825 in Experiment 1, we further explored each age interaction by extracting an
826 average difference score over subjects for the effect of language condition (*no*
827 *speech* vs. *prosody only*) within each random permutation of the data, mak-
828 ing pairwise comparisons between the six age groups. These tests revealed
829 that children’s anticipation in the *prosody only* condition significantly im-
830 proved at ages five and six (difference scores greater than 95% of the random
831 data scores; $p < .05$; see Figure B.2 for the full distributions).

832 The second age-based interaction was a three-way interaction of age, transi-
833 tion type, and language condition (for *normal* speech compared to the *no*
834 *speech* baseline). We again created pairwise comparisons of the average dif-
835 ference scores for the transition type-language condition interaction across
836 age groups in each random permutation of the data, finding that the effect
837 of transition type in the *normal* speech condition became larger with age,
838 with significant improvements by age 4 over ages 1 and 2 (99.9% and 98.86%,
839 respectively), by age 5 over age 4 (97.54%), and by age 6 over ages 1, 2, and 5
840 (99.5%, 97.36%, and 95.04%), all significantly different from chance ($p < .05$;
841 see Figure B.3 for the full distributions).

842 Our main goal in extending the age range to 1- and 2-year-olds in Ex-
843 periment 2 was to find the age of emergence for spontaneous turn structure
844 predictions. As in Experiment 1, we used two-tailed *t*-tests to compare chil-
845 dren’s real gaze rates to the random baseline in the *normal* speech condition,
846 in which the speech stimulus is most like what children hear every day. We
847 tested real gaze rates against baseline for three age groups: ages one, two,
848 and three. Two- and three-year-old children made anticipatory gaze switches
849 significantly above chance both when all transitions were considered (2-year-

olds: $t(26.193)=-4.137$, $p<.001$; 3-year-olds: $t(22.757)=-2.662$, $p<.05$) and for question transitions alone (2-year-olds: $t(25.345)=-4.269$, $p<.001$; 3-year-olds: $t(21.555)=-3.03$, $p<.01$). One-year-olds, however, made anticipatory gaze shifts that were only marginally above chance for turn transitions overall and for question turns alone (overall: $t(24.784)=-2.049$, $p=.051$; questions: $t(25.009)=-2.03$, $p=.053$).

3.4. Discussion

As in Experiment 1, children and adults spontaneously tracked the turn structure of the conversations, making anticipatory gaze switches at an above-chance rate across all ages, but not in all conditions (Table 3; Figure 6). However, when children had access to full linguistic information (in the *normal* speech condition), they made more anticipatory gaze switches than expected by chance at age two (and even marginally at age one). Across conditions, children's anticipatory gaze rates were affected by gap duration, plus interactions of language condition, transition type and age (Table 5), none of which could be explained by a baseline of random gaze switching (Figure A.4a).

Two of the three linguistic conditions affected children's anticipatory gaze switches, compared to the no-speech reference level: prosody-only speech and *normal* speech. *Prosody only* speech resulted in fewer anticipatory gazes overall, compared to the no-speech baseline, though the effect was only marginal ($p=.06$). However, prosody-only speech *did* significantly differ from the *no speech* condition in its interaction with age: Whereas age gains were negligible in the *no speech* condition, 5- and 6-year-olds in the prosody-only condition showed significantly more anticipatory gaze switches than younger children (Figure B.2), going from at-chance anticipatory gaze rates to rates well above chance (Figure 6).

Anticipatory gaze for questions increased significantly more with age in the *normal* speech condition compared to the no-speech baseline. As in Experiment 1, children consistently made more anticipatory gazes after hearing questions when they had access to lexical material. But now, with a larger age span in Experiment 2, we can begin to see a developmental path for the question effect. At least for *normal* speech, greater anticipatory looking for questions is not present from the start: 4-year-olds are the first to make significant gains on 1- and 2-year-olds, but then there are further significant gains at age 5, and again age 6 (Figure B.3). While children at ages five and six showed adult-like differentiation of questions and non-questions in the *normal* speech condition, 1-year-olds had nearly identical switch rates

887 for the two transition types. This suggests that the participants' tendency
888 to make more anticipatory switches for questions emerges after age one and
889 continues developing, with rapid gains in ages three through age six. Fi-
890 nally, children showed a straightforward effect of gap duration (Table 5), as
891 in Experiment 1.

892 Adults' anticipatory gaze rates were affected by gap duration and two
893 of the language conditions (Table 5), none of which could be explained by
894 a baseline of random gaze switching (Figure A.4b). Adults made more an-
895 ticipatory switches overall for the *normal* speech and words-only conditions
896 compared to the no-speech condition, falling in-line with past work showing
897 that adults primarily use lexical information in making predictions about up-
898 coming turn structure (De Ruiter et al., 2006). Though adults did make more
899 anticipatory switches for questions than non-questions on average in the two
900 lexical conditions (Figure 6), the effect of transition type was not significant
901 in either one (Table 5). Like children, adults also showed a straightforward
902 effect of gap duration.

903 *3.4.1. Summary*

904 Children and adults both showed more anticipatory gaze switches with
905 increased linguistic information, but only for a subset of the linguistic con-
906 ditions and transition types. We had expected to see the most anticipatory
907 switches in the *normal* condition and the least anticipatory switches in the
908 *no speech* condition because they contained the most and least linguistic in-
909 formation, respectively. We had also expected to replicate our finding from
910 Experiment 1 that questions result in more anticipatory switches than non-
911 questions, with the added hypothesis that the question effect is driven by
912 lexicosyntactic cues. We additionally anticipated an overall increase in an-
913 ticipatory switches with age. Finally, since the development of prosodic skills
914 partially precedes the development of lexicosyntax, we expected to see more
915 switches in the *prosody only* condition compared to the *words only* condition
916 in the youngest age groups.

917 In fact, children and adults did show more anticipatory gaze switching
918 in the *normal* speech condition, compared to the no-speech condition, but
919 for children this effect only emerged in the form of anticipations following
920 question turns, which increased with age significantly faster than it did when
921 there were no linguistic cues present at all. Adults also showed significantly
922 higher anticipatory switch rates in the *words only* condition but no effects
923 of transition type, alone or within linguistic conditions. Taken together,

924 these results partly replicate the findings from Experiment 1: Participants
925 make more anticipatory switches when they have access to lexical information
926 and, when they do, tend to make more anticipatory switches for questions
927 compared to non-questions.

928 We had anticipated significant gains in anticipatory switching with age,
929 but children only showed significant developmental increases in the *prosody*
930 *only* condition and the *normal* condition (for question transitions). Rather
931 than showing an early advantage for prosody over lexical information, chil-
932 dren's anticipation were at chance with prosody-only speech and did not show
933 significant improvement until age five. In contrast, their anticipatory gaze
934 switches were already marginally above chance at age one and significantly
935 so at age two for *normal* speech. These findings therefore do not support an
936 early role for prosody in children's spontaneous turn structure predictions.
937 On the contrary, children's predictions were best when lexical information
938 was present, especially when following a question.

939 However, these results do not support the idea that lexical information is
940 sufficient for children, as has been proposed for adults previously De Ruiter
941 et al., 2006 and replicated in the present study. On the contrary, children
942 showed significant gains in the *normal* speech condition, but *not* in the *words*
943 *only* condition. Notably, the *normal* speech condition, in addition to having
944 both lexical and prosodic information, is also the one most likely to occur
945 in our participants' daily lives (compared to muffled or robotic speech) and
946 therefore may have an extra advantage of familiarity over the other condi-
947 tions.

948 Finally, though participants were often able to make anticipatory gaze
949 switches more often than would be expected by chance with no linguistic in-
950 formation (in the *no speech* condition), the advantage for lexical cues emerg-
951 ing in the children's data with age and in the adult's data overall suggests
952 that participants were using linguistic information in making their predic-
953 tions when it was available.

954 The core aims of Experiment 2 were to gain better traction on the indi-
955 vidual roles of prosody and lexicosyntax in children's turn predictions, and
956 to find the age of emergence for spontaneous turn anticipation. Replicating
957 many of our findings from Experiment 1, we found that children and adults
958 both use linguistic cues, in children's case, primarily to anticipate responses
959 after question turns and primarily when the full speech signal was available.
960 Adults, on the other hand, showed broadly increased anticipation for lexical
961 conditions, with no overall effects of transition type.

962 Rather than finding an early advantage for prosodic cues, we found that
963 children struggled to make above-chance anticipations in the *prosody only*
964 condition until age five. They didn't show significant overall effects for lexical
965 information alone either, suggesting that children need the full linguistic
966 signal to extract predictive benefits from ongoing conversational speech. This
967 finding is surprising given children's earlier development of prosodic skills
968 but unsurprising in the light of their overall question bias—compared to
969 prosodic cues, lexicosyntactic cues to questionhood are categorical units, have
970 clearer form-to-meaning mappings, and occur early in the utterance. This
971 finding has implications for our conceptions about adults and children's use
972 of linguistic cues for spontaneous turn prediction, which we expand on in the
973 General Discussion.

974 4. General Discussion

975 Children begin to develop conversational turn-taking skills long before
976 their first words (Bateson, 1975; Hilbrink et al., 2015; Jaffe et al., 2001;
977 Snow, 1977). As they acquire language, they also acquire the information
978 needed to make accurate predictions about upcoming turn structure. Until
979 recently, we have had very little data on how children weave language into
980 their already-existing turn-taking behaviors.

981 In two experiments investigating children's anticipatory gaze to upcom-
982 ing speakers, we found evidence that turn prediction develops early in child-
983 hood and that spontaneous predictions are primarily driven by participants'
984 expectation of an immediate response in the next turn. In making predic-
985 tions about upcoming turn structure, children used a combination of lexical
986 and prosodic cues; neither lexical nor prosodic cues alone were sufficient to
987 support increased anticipatory gaze. We also found no early advantage for
988 prosody over lexicosyntax, and instead found that children were unable to
989 make above-chance anticipatory gazes in the *prosody only* condition until age
990 five. We discuss these findings with respect to the role of linguistic cues in
991 predictions about upcoming turn structure, the importance of questions in
992 spontaneous predictions about conversation, and children's developing com-
993 petence as conversationalists.

994 4.1. Predicting turn structure with linguistic cues

995 Prior work with adults has found a consistent and critical role for lex-
996 icosyntax in predicting upcoming turn structure (De Ruiter et al., 2006;

997 Magyari and De Ruiter, 2012), with the role of prosody still under debate
998 (Duncan, 1972; Ford and Thompson, 1996; Torreira et al., 2015). Knowing
999 that children comprehend more about prosody than lexicosyntax early on
1000 (see introduction; also see Speer and Ito, 2009 for a review), we thought it
1001 possible that young children would instead show an advantage for prosody in
1002 their predictions about turn structure in conversation. Our results suggest
1003 that, on the contrary, when presented with *only* prosodic information, chil-
1004 dren's spontaneous predictions about upcoming turn structure are limited
1005 until age five.

1006 Importantly, we also found no evidence that lexical information alone is
1007 equivalent to full linguistic information for children, as has been shown be-
1008 fore Magyari and De Ruiter (2012); De Ruiter et al. (2006) and replicated
1009 in the current study for adult participants: In both experiments, children's
1010 performance was best in conditions when they had access to the full linguistic
1011 signal. Adults on the other hand, showed significant gains in anticipatory
1012 gaze switching in both conditions with lexical cues. If this effect arose in
1013 our data because children do not make predictions based on phonetically
1014 manipulated speech in conversational contexts, the current findings can be
1015 overturned by follow-up work controlling linguistic information through other
1016 means. But if not, there may be something specially informative about the
1017 combined prosodic and lexical cues to questionhood that boosts children's an-
1018 ticipations before they can use these cues separately. Even in adults, Torreira
1019 and colleagues (2015) were able to show that the trade-off in informativity
1020 between lexical and prosodic cues is more subtle in semi-natural (spliced)
1021 speech. The present findings are the first to show evidence of a similar effect
1022 developmentally.

1023 4.1.1. *The question effect*

1024 In both experiments, anticipatory looking was primarily driven by ques-
1025 tion transitions, a pattern that had not been previously reported in other an-
1026 ticipatory gaze studies, on children or adults (Keitel et al., 2013; Hirvenkari,
1027 2013; Tice and Henetz, 2011). Questions make an upcoming speaker switch
1028 immediately relevant, helping the listener to predict with high certainty what
1029 will happen next (i.e., an answer from the addressee), and are often easily
1030 identifiable by overt prosodic and lexicosyntactic cues.

1031 Compared to prosodic cues (e.g., final rising intonation), lexicosyntactic
1032 cues (e.g., *wh*-words, *do*-insertion, and subject-auxiliary inversion) were fre-
1033 quent, categorical, and early-occurring in the utterance. Children may have

1034 therefore had an easier time picking out and interpreting lexical cues to ques-
1035 tionhood. The question effect showed its first significant gains between ages
1036 three and four in the *normal* speech condition of Experiment 2, by which
1037 time children frequently hear and use a variety of polar *wh*-questions (Clark,
1038 2009). Furthermore, while lexicosyntactic question cues were available on
1039 every instance of *wh*- and *yes/no* questions in our stimuli, prosodic question
1040 cues were only salient on *yes/no* questions and, even then, the mapping of
1041 prosodic contour to speech act (e.g., high final rises for polar questions) is
1042 far from one-to-one.

1043 Prior work on children’s acquisition of questions indicates that they may
1044 already have some understanding about question-answer sequences by the
1045 time they begin to speak: Questions make up approximately one third of
1046 the utterances children hear, before and after the onset of speech, and even
1047 into their preschool years, even though the types and complexity of questions
1048 change throughout development (Casillas et al., In press; Fitneva, 2012; Hen-
1049 ning et al., 2005; Shatz, 1979).¹⁰ For the first few years, many of the questions
1050 directed to children are “test” questions—questions that the caregiver already
1051 has the answer to (e.g., “What does a cat say?”), but this changes as children
1052 get older. Questions help caregivers to get their young children’s attention
1053 and to ensure that information is in common ground, even if the responses are
1054 non-verbal or infelicitous (Bruner, 1985; Fitneva, 2012; Snow, 1977). So, in
1055 addition to having a special interactive status (for adults and children alike),
1056 questions are a core characteristic of many caregiver-child interactions, mot-
1057 ivating a general benefit for questions in turn structure anticipation.

1058 Two important questions for future work are then: (a) How does chil-
1059 dren’s ability to monitor for questions in conversation relate to their prior
1060 experience with questions? and (b) what is it about questions that makes
1061 children and adults more likely to anticipatorily switch their gaze to ad-
1062 dressees? Other request formats, such as imperatives, compliments, and
1063 complaints make a response from the addressee highly likely in the next turn
1064 (Schegloff, 2007). Rhetorical and tag questions, on the other hand, take a
1065 similar form to prototypical polar questions, but often do not require an an-
1066 swer. So, though it is clear that adults and children anticipated responses
1067 more often for questions than non-questions, we do not yet know whether

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

1068 their predictive action is limited to turns formatted as questions or is generally applicable to turn structures that project an immediate response from
1069 the addressee.
1070

1071 The question effect itself has implications for our current theories about
1072 turn prediction. While participants may always use linguistic information
1073 to predict upcoming speaker changes (as they have in the two studies pre-
1074 sented here), they might not always use linguistic information to predict
1075 upcoming turn ends, as is assumed in metalinguistic measures of turn-end
1076 prediction (e.g., pressing a button while listening to speech: Torreira et al.,
1077 2015; Magyari and De Ruiter, 2012; De Ruiter et al., 2006) and as has been
1078 the focus of other turn-end prediction studies Ford and Thompson, 1996;
1079 Duncan, 1972. Our results rather suggest that participants' *spontaneous*
1080 predictions, at least while viewing third-party conversation, are the results
1081 of question-monitoring, presumably achieved by recruiting linguistic cues to
1082 questionhood from the unfolding signal. In other words, our results suggest
1083 that predictions are driven by what is *beyond* the end of the current turn—
1084 that questions, not lexical cues, are sufficient for prediction at the first level
1085 of conversation monitoring.

1086 There is at least one clear hypothesis that can bridge these apparently
1087 conflicting results: Listeners in spontaneous, first-person conversation may
1088 use multiple strategies in making predictions about upcoming turn structure,
1089 using more passive prediction to detect upcoming speaker transition (e.g.,
1090 questions), and then switching into precise turn-end prediction mode (à la
1091 De Ruiter et al., 2006) when necessary. A flexible prediction system like this
1092 one allows listeners to continuously monitor ongoing conversation at a low
1093 cost while still managing to plan their responses and come in quickly when
1094 needed.

1095 To test this hypothesis, which integrates findings from turn-structure pre-
1096 diction with multiple measures, age groups, and styles of linguistic control,
1097 it will be crucial to look at prediction from a first-person perspective. The
1098 results we present here are based on predictions about third-party conver-
1099 sation, which enables participants to follow interactions with no chance of
1100 actually participating. Although recent work has shown that similar antici-
1101 patory eye gazes do occur in spontaneous conversation (Holler and Kendrick,
1102 2015), more work is needed to determine if the same question advantage oc-
1103 curs, which linguistic cues seem to drive it, and whether participants switch
1104 into a “precision” mode when they detect imminent speaker change.

1105 4.1.2. *Early competence for turn taking?*

1106 One of the core aims of our study was to test whether children show an
1107 early competence for turn taking, as is proposed by studies of spontaneous
1108 mother-infant proto-conversation and theories about the mechanisms under-
1109 lying human interaction in general (Hilbrink et al., 2015; Levinson, 2006). We
1110 did find evidence that young children already make spontaneous predictions
1111 about upcoming turn structure, definitely at age two and even marginally
1112 at age one. However, “above chance” performance was far from adult-like
1113 predictive behavior, and children in our studies did not show adult-like com-
1114 petency in their predictions, even at age six. This may indicate that children
1115 rely more on non-verbal cues in anticipating turn transitions or, alternatively,
1116 that adults are better at flexibly adapting to the turn-relevant cues present
1117 at any moment.

1118 Taken together, the data suggest that turn-taking skills do begin to
1119 emerge in infancy, but that their predictions don’t help them anticipate much
1120 until they have acquired the ability to pick out question turns. This finding
1121 leads us to wonder how participant role (first- instead of third-person) and
1122 cultural differences (e.g., high vs. low parent-infant interaction styles) might
1123 feed into this early predictive skill. It also bridges the prior work showing a
1124 predisposition for turn taking in infancy (e.g., Hilbrink et al., 2015) but late
1125 acquisition of adult-like competence when it comes to integrating linguistic
1126 information into turn-taking behaviors (Casillas et al., In press; Garvey, 1984;
1127 Ervin-Tripp, 1979).

1128 4.2. *Limitations and future work*

1129 There are at least two major limitations to our work: Speech naturalness
1130 and participant role. Following prior work (De Ruiter et al., 2006; Keitel et al.,
1131 2013), we used phonetically manipulated speech in Experiment 2, resulting in
1132 speech sounds that children don’t usually hear in their natural environment.
1133 Many prior studies have used phonetically-altered speech with infants and
1134 young children (cf. Jusczyk, 2000), but almost none of them have done so
1135 in a conversational context. Future work could instead carefully script or
1136 cross-splice parts of turns to control for the presence or absence of linguistic
1137 cues for turn transition.

1138 The prediction measure used in these studies is based on an observer’s
1139 view of third-party conversation but, because participants’ role in the inter-
1140 action could affect their online predictions about turn taking, an ideal exper-
1141 imental measure would capture first-person behavior. First-person measures

1142 of spontaneous turn prediction will be key to revealing how participants dis-
1143 tribute their attention over linguistic and non-verbal cues while taking part
1144 in everyday interaction, the implications of which relate to theories of online
1145 language processing for both language learning and everyday talk.

1146 *4.3. Conclusions*

1147 Conversation plays a central role in children’s language learning. It is
1148 the driving force behind what children say and what they hear. Adults use
1149 linguistic information to accurately predict turn structure in conversation,
1150 which facilitates their online comprehension and allows them to respond rel-
1151 evantly and on time. The present study offers new findings regarding the
1152 role of speech acts and linguistic processing in online turn prediction, and
1153 has given evidence that turn prediction emerges by age two is not integrated
1154 with linguistic cues until much later. Using language to make predictions
1155 about upcoming interactive content takes time and, for both children and
1156 adults, is primarily driven by participants’ orientation to what will happen
1157 beyond the end of the current turn.

1158 **Acknowledgements**

1159 We gratefully acknowledge the parents and children at Bing Nursery
1160 School and the Children’s Discovery Museum of San Jose. This work was
1161 supported by an ERC Advanced Grant to Stephen C. Levinson (269484-
1162 INTERACT), NSF graduate research and dissertation improvement fellow-
1163 ships to the first author, and a Merck Foundation fellowship to the second
1164 author. Earlier versions of these data and analyses were presented to confer-
1165 ence audiences (Casillas and Frank, 2012, 2013). We also thank Tania Henetz,
1166 Francisco Torreira, Stephen C. Levinson, Eve V. Clark, and the First Lan-
1167 guage Acquisition group at Radboud University for their feedback on earlier
1168 versions of this work. The analysis code and raw data for this project can
1169 be found on GitHub at https://github.com/langcog/turn_taking/.

1170 **References**

- 1171 Barr, D.J., Levy, R., Scheepers, C., Tily, H.J., 2013. Random effects structure
1172 for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory*
1173 and *Language* 68, 255–278.

- 1174 Bates, D., Maechler, M., Bolker, B., Walker, S., 2014. lme4:
1175 Linear mixed-effects models using Eigen and S4. URL:
1176 <https://github.com/lme4/lme4><http://lme4.r-forge.r-project.org/>.
1177 [Computer program] R package version 1.1-7.
- 1178 Bateson, M.C., 1975. Mother-infant exchanges: The epigenesis of conver-
1179 sational interaction. Annals of the New York Academy of Sciences 263,
1180 101–113.
- 1181 Bergelson, E., Swingley, D., 2013. The acquisition of abstract words by young
1182 infants. Cognition 127, 391–397.
- 1183 Bloom, K., 1988. Quality of adult vocalizations affects the quality of infant
1184 vocalizations. Journal of Child Language 15, 469–480.
- 1185 Boersma, P., Weenink, D., 2012. Praat: doing phonetics by computer. URL:
1186 <http://www.praat.org>. [Computer program] Version 5.3.16.
- 1187 Bögels, S., Magyari, L., Levinson, S.C., 2015. Neural signatures of response
1188 planning occur midway through an incoming question in conversation. Sci-
1189 entific Reports 5.
- 1190 Bruner, J., 1985. Child's talk: Learning to use language. Child Language
1191 Teaching and Therapy 1, 111–114.
- 1192 Bruner, J.S., 1975. The ontogenesis of speech acts. Journal of Child Language
1193 2, 1–19.
- 1194 Carlson, R., Hirschberg, J., Swerts, M., 2005. Cues to upcoming swedish
1195 prosodic boundaries: Subjective judgment studies and acoustic correlates.
1196 Speech Communication 46, 326–333.
- 1197 Casillas, M., Bobb, S.C., Clark, E.V., In press. Turn taking, timing, and
1198 planning in early language acquisition. Journal of Child Language .
- 1199 Casillas, M., Frank, M.C., 2012. Cues to turn boundary prediction in adults
1200 and preschoolers. Proceedings of SemDial .
- 1201 Casillas, M., Frank, M.C., 2013. The development of predictive processes
1202 in children's discourse understanding, in: Proceedings of the 35th Annual
1203 Meeting of the Cognitive Science Society.

- 1204 Clark, E.V., 2009. First language acquisition. Cambridge University Press.
- 1205 De Ruiter, J.P., Mitterer, H., Enfield, N.J., 2006. Projecting the end of
1206 a speaker's turn: A cognitive cornerstone of conversation. *Language* 82,
1207 515–535.
- 1208 De Vos, C., Torreira, F., Levinson, S.C., 2015. Turn-timing in signed con-
1209 versations: coordinating stroke-to-stroke turn boundaries. *Frontiers in*
1210 *Psychology* 6.
- 1211 Dingemanse, M., Torreira, F., Enfield, N., 2013. Is “Huh?” a universal word?
1212 Conversational infrastructure and the convergent evolution of linguistic
1213 items. *PloS one* 8, e78273.
- 1214 Duncan, S., 1972. Some signals and rules for taking speaking turns in con-
1215 versations. *Journal of Personality and Social Psychology* 23, 283.
- 1216 Ervin-Tripp, S., 1979. Children's verbal turn-taking, in: Ochs, E., Schieffelin,
1217 B.B. (Eds.), *Developmental Pragmatics*. Academic Press, New York, pp.
1218 391–414.
- 1219 Fitneva, S., 2012. Beyond answers: questions and children's learning, in:
1220 de Ruiter, J.P. (Ed.), *Questions: Formal, Functional, and Interactional*
1221 *Perspectives*. Cambridge University Press, Cambridge, UK, pp. 165–178.
- 1222 Ford, C.E., Thompson, S.A., 1996. Interactional units in conversation: Syn-
1223 tactic, intonational, and pragmatic resources for the management of turns.
1224 *Studies in Interactional Sociolinguistics* 13, 134–184.
- 1225 Garvey, C., 1984. Children's Talk. volume 21. Harvard University Press.
- 1226 Gísladóttir, R., Chwilla, D., Levinson, S.C., 2015. Conversation electrified:
1227 ERP correlates of speech act recognition in underspecified utterances. *PloS*
1228 *one* 10, e0120068.
- 1229 Griffin, Z.M., Bock, K., 2000. What the eyes say about speaking. *Psycho-*
1230 *logical science* 11, 274–279.
- 1231 Hart, B., Risley, T.R., 1992. American parenting of language-learning chil-
1232 dren: Persisting differences in family-child interactions observed in natural
1233 home environments. *Developmental Psychology* 28, 1096.

- 1234 Hedberg, N., Sosa, J.M., Görgülü, E., Mameni, M., 2010. The prosody and
1235 meaning of Wh-questions in American English, in: Speech Prosody 2010–
1236 Fifth International Conference.
- 1237 Henning, A., Striano, T., Lieven, E.V., 2005. Maternal speech to infants at
1238 1 and 3 months of age. *Infant Behavior and Development* 28, 519–536.
- 1239 Hilbrink, E., Gattis, M., Levinson, S.C., 2015. Early developmental changes
1240 in the timing of turn-taking: A longitudinal study of mother-infant inter-
1241 action. *Frontiers in Psychology* 6.
- 1242 Hirvenkari, L., Ruusuvuori, J., Saarinen, V.M., Kivioja, M., Peräkylä, A.,
1243 Hari, R., 2013. Influence of turn-taking in a two-person conversation on
1244 the gaze of a viewer. *PloS one* 8, e71569.
- 1245 Holler, J., Kendrick, K.H., 2015. Unaddressed participants' gaze in multi-
1246 person interaction. *Frontiers in Psychology* 6.
- 1247 Jaffe, J., Beebe, B., Feldstein, S., Crown, C.L., Jasnow, M.D., Rochat, P.,
1248 Stern, D.N., 2001. Rhythms of dialogue in infancy: Coordinated timing in
1249 development. *Monographs of the Society for Research in Child Develop-
1250 ment*. JSTOR.
- 1251 Johnson, E.K., Jusczyk, P.W., 2001. Word segmentation by 8-month-olds:
1252 When speech cues count more than statistics. *Journal of Memory and
1253 Language* 44, 548–567.
- 1254 Jusczyk, P.W., 2000. *The Discovery of Spoken Language*. MIT press.
- 1255 Jusczyk, P.W., Hohne, E., Mandel, D., Strange, W., 1995. Picking up reg-
1256 ularities in the sound structure of the native language. *Speech perception*
1257 and linguistic experience: Theoretical and methodological issues in cross-
1258 language speech research , 91–119.
- 1259 Kamide, Y., Altmann, G., Haywood, S.L., 2003. The time-course of predic-
1260 tion in incremental sentence processing: Evidence from anticipatory eye
1261 movements. *Journal of Memory and Language* 49, 133–156.
- 1262 Keitel, A., Prinz, W., Friederici, A.D., Hofsten, C.v., Daum, M.M., 2013.
1263 Perception of conversations: The importance of semantics and intonation
1264 in childrens development. *Journal of Experimental Child Psychology* 116,
1265 264–277.

- 1266 Lemasson, A., Glas, L., Barbu, S., Lacroix, A., Guilloux, M., Remeuf, K.,
1267 Koda, H., 2011. Youngsters do not pay attention to conversational rules:
1268 is this so for nonhuman primates? *Nature Scientific Reports* 1.
- 1269 Levelt, W.J., 1989. Speaking: From intention to articulation. MIT press.
- 1270 Levinson, S.C., 2006. On the human “interaction engine”, in: Enfield, N.,
1271 Levinson, S. (Eds.), *Roots of human sociality: Culture, cognition and*
1272 *interaction*. Oxford: Berg, pp. 39–69.
- 1273 Levinson, S.C., 2013. Action formation and ascriptions, in: Stivers, T., Sid-
1274 nell, J. (Eds.), *The Handbook of Conversation Analysis*. Wiley-Blackwell,
1275 Malden, MA, pp. 103–130.
- 1276 Magyari, L., Bastiaansen, M.C.M., De Ruiter, J.P., Levinson, S.C., 2014.
1277 Early anticipation lies behind the speed of response in conversation. *Jour-*
1278 *nal of Cognitive Neuroscience* 26, 2530–2539.
- 1279 Magyari, L., De Ruiter, J.P., 2012. Prediction of turn-ends based on antici-
1280 pation of upcoming words. *Frontiers in Psychology* 3:376, 1–9.
- 1281 Masataka, N., 1993. Effects of contingent and noncontingent maternal stimu-
1282 lation on the vocal behaviour of three-to four-month-old Japanese infants.
1283 *Journal of Child Language* 20, 303–312.
- 1284 Mehler, J., Jusczyk, P., Lambertz, G., Halsted, N., Bertoni, J., Amiel-
1285 Tison, C., 1988. A precursor of language acquisition in young infants.
1286 *Cognition* 29, 143–178.
- 1287 Morgan, J.L., Saffran, J.R., 1995. Emerging integration of sequential and
1288 suprasegmental information in preverbal speech segmentation. *Child De-*
1289 *velopment* 66, 911–936.
- 1290 Nazzi, T., Ramus, F., 2003. Perception and acquisition of linguistic rhythm
1291 by infants. *Speech Communication* 41, 233–243.
- 1292 R Core Team, 2014. R: A Language and Environment for Statistical Com-
1293 puting. R Foundation for Statistical Computing. Vienna, Austria. URL:
1294 <http://www.R-project.org>. [Computer program] Version 3.1.1.
- 1295 Ratner, N., Bruner, J., 1978. Games, social exchange and the acquisition of
1296 language. *Journal of Child Language* 5, 391–401.

- 1297 Ross, H.S., Lollis, S.P., 1987. Communication within infant social games.
1298 *Developmental Psychology* 23, 241.
- 1299 Rossano, F., Brown, P., Levinson, S.C., 2009. Gaze, questioning and culture,
1300 in: Sidnell, J. (Ed.), *Conversation Analysis: Comparative Perspectives*.
1301 Cambridge University Press, Cambridge, pp. 187–249.
- 1302 Sacks, H., Schegloff, E.A., Jefferson, G., 1974. A simplest systematics for the
1303 organization of turn-taking for conversation. *Language* 50, 696–735.
- 1304 Schegloff, E.A., 2007. Sequence organization in interaction: Volume 1: A
1305 primer in conversation analysis. Cambridge University Press.
- 1306 Shatz, M., 1979. How to do things by asking: Form-function pairings in
1307 mothers' questions and their relation to children's responses. *Child Develop-*
1308 *ment* 50, 1093–1099.
- 1309 Shi, R., Melancon, A., 2010. Syntactic categorization in French-learning
1310 infants. *Infancy* 15, 517–533.
- 1311 Snow, C.E., 1977. The development of conversation between mothers and
1312 babies. *Journal of Child Language* 4, 1–22.
- 1313 Soderstrom, M., Seidl, A., Kemler Nelson, D.G., Jusczyk, P.W., 2003. The
1314 prosodic bootstrapping of phrases: Evidence from prelinguistic infants.
1315 *Journal of Memory and Language* 49, 249–267.
- 1316 Speer, S.R., Ito, K., 2009. Prosody in first language acquisition—Acquiring
1317 intonation as a tool to organize information in conversation. *Language and*
1318 *Linguistics Compass* 3, 90–110.
- 1319 Stivers, T., Enfield, N.J., Brown, P., Englert, C., Hayashi, M., Heinemann,
1320 T., Hoymann, G., Rossano, F., De Ruiter, J.P., Yoon, K.E., et al., 2009.
1321 Universals and cultural variation in turn-taking in conversation. *Proceed-*
1322 *ings of the National Academy of Sciences* 106, 10587–10592.
- 1323 Stivers, T., Rossano, F., 2010. Mobilizing response. *Research on Language*
1324 and Social Interaction 43, 3–31.
- 1325 Takahashi, D.Y., Narayanan, D.Z., Ghazanfar, A.A., 2013. Coupled oscillator
1326 dynamics of vocal turn-taking in monkeys. *Current Biology* 23, 2162–2168.

1327 Thorgrímsson, G., Fawcett, C., Liszkowski, U., 2015. 1- and 2-year-olds'
1328 expectations about third-party communicative actions. Infant Behavior
1329 and Development 39, 53–66.

1330 Tice (Casillas), M., Henetz, T., 2011. Turn-boundary projection: Looking
1331 ahead, in: Proceedings of the 33rd Annual Meeting of the Cognitive Science
1332 Society.

1333 Torreira, F., Bögels, S., Levinson, S.C., 2015. Intonational phrasing is neces-
1334 sary for turn-taking in spoken interaction. Journal of Phonetics 52, 46–57.

1335 Wittenburg, P., Brugman, H., Russel, A., Klassmann, A., Sloetjes, H., 2006.
1336 Elan: a professional framework for multimodality research, in: Proceedings
1337 of LREC.

1338 **Appendix A. Real vs. randomly permuted model outcomes: t , β ,**

1339 and SE

1340 In all of the following plots, the gray dots represent the randomly per-
1341 mitted data's model estimates for the value listed (beta or standard error),
1342 the white dots represent the model estimates from the original data, and the
1343 triangles represent the 95th percentile for each distribution being shown.

1344 *Appendix A.1. Experiment 1*

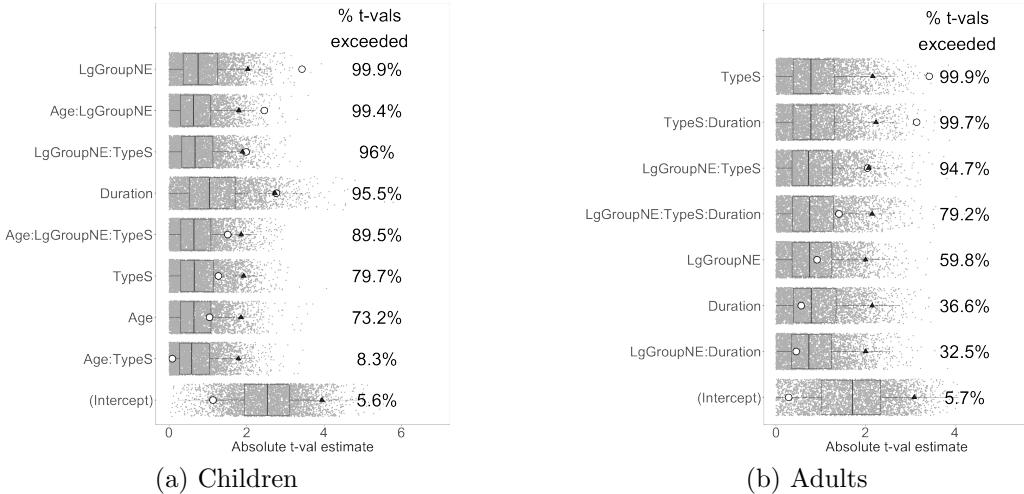


Figure A.1: Random-permutation and original $|t\text{-values}|$ for predictors of anticipatory gaze rates in Experiment 1.

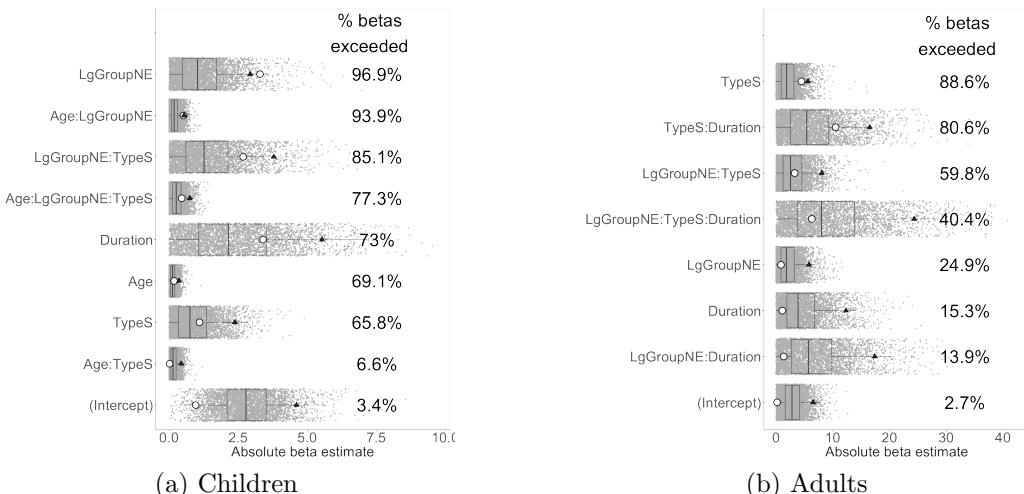


Figure A.2: Random-permutation and original $|\beta\text{-values}|$ for predictors of gaze rates in Experiment 1.

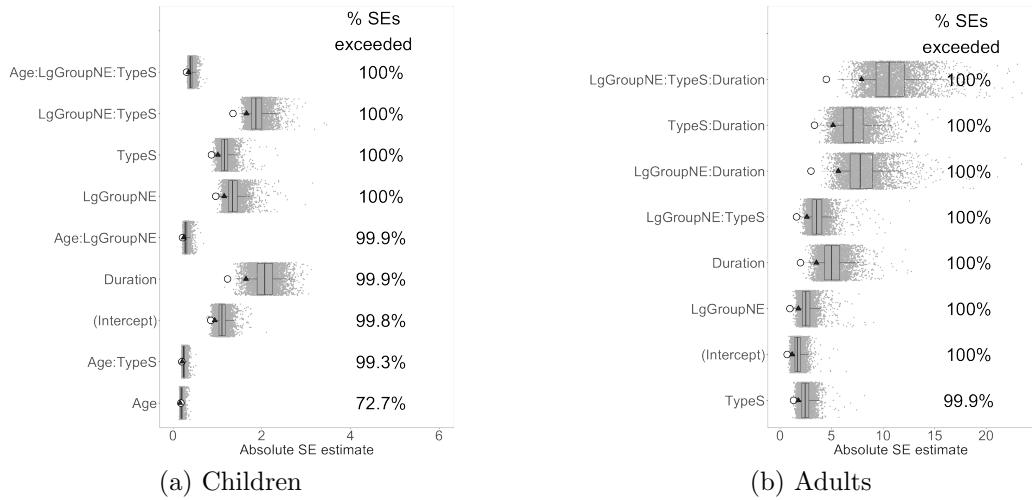


Figure A.3: Random-permutation and original $|SE\text{-values}|$ for predictors of anticipatory gaze rates in Experiment 1.

1345 *Appendix A.2. Experiment 2*

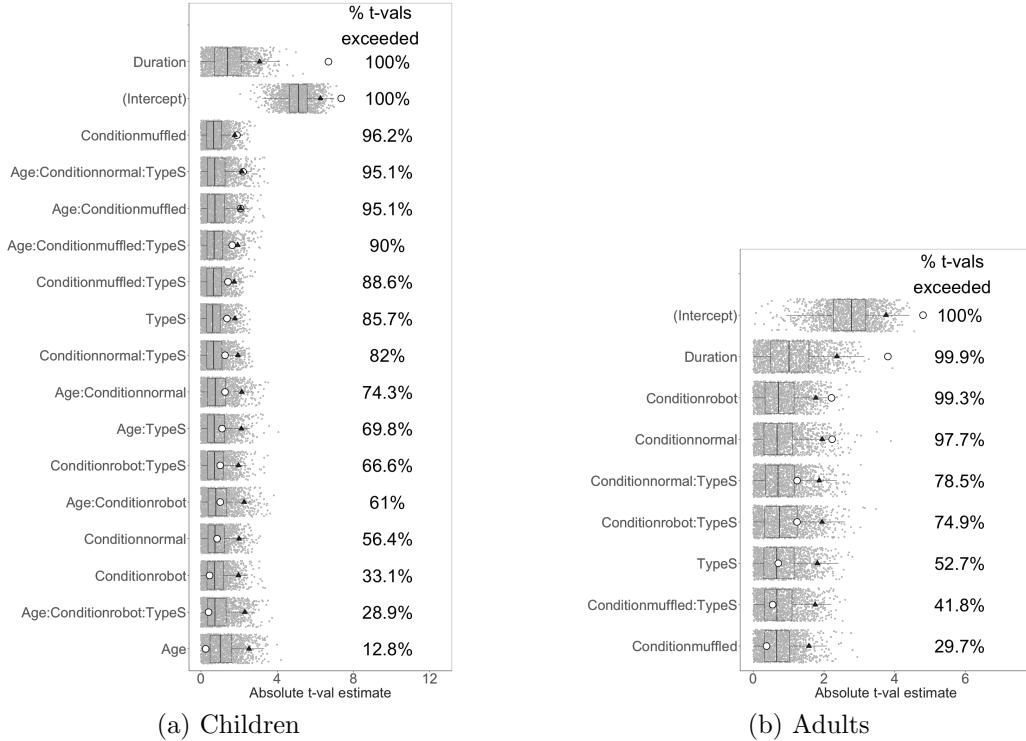


Figure A.4: Random-permutation and original $|t\text{-values}|$ for predictors of anticipatory gaze rates in Experiment 2.

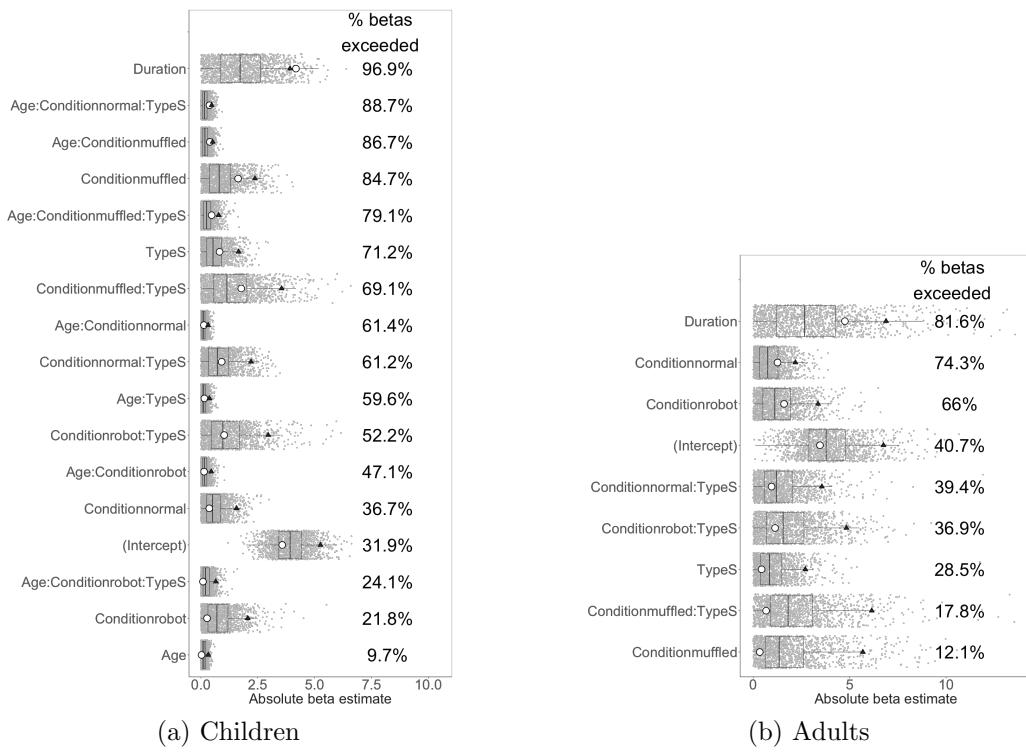


Figure A.5: Random-permutation and original $|\beta\text{-values}|$ for predictors of anticipatory gaze rates in Experiment 2.

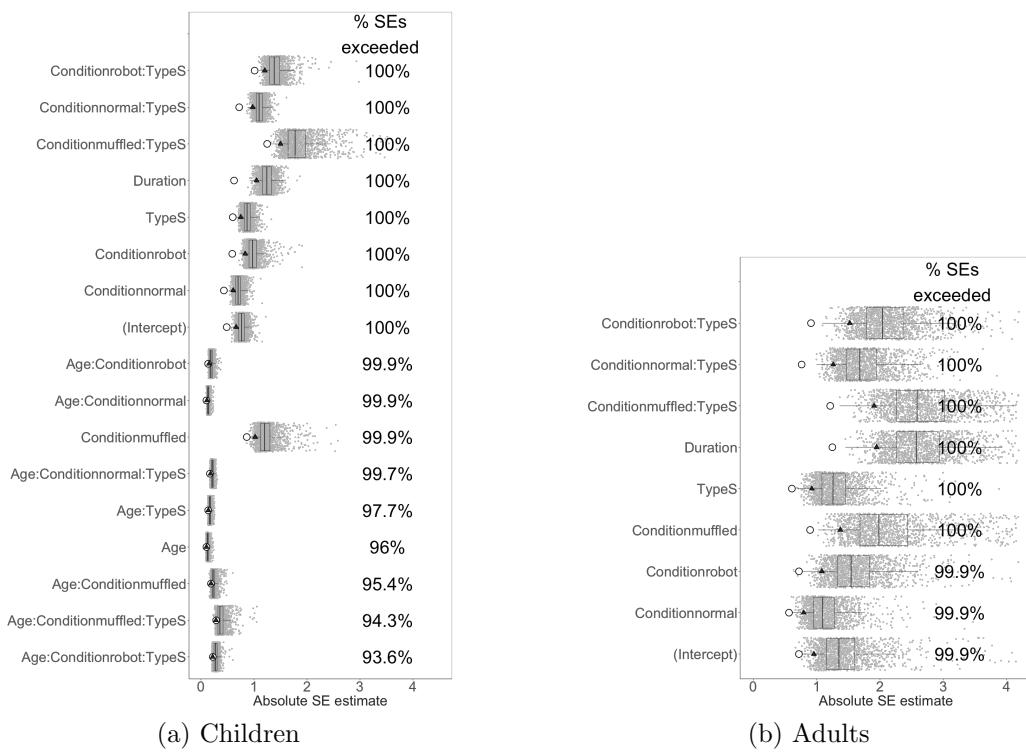


Figure A.6: Random-permutation and original $|SE\text{-values}|$ for predictors of anticipatory gaze rates in Experiment 2.

1346 **Appendix B. Pairwise developmental tests: Real vs. randomly
1347 permuted effects**

1348 In each of the plots below, the dot represents the original data value for
1349 the effect and the 5,000 randomly permuted data effect sizes are shown in the
1350 distribution. The percentage shown is the percentage of random permutation
1351 values exceeded by the original data value (taking the absolute value of all
1352 data points for a two-tailed test.)

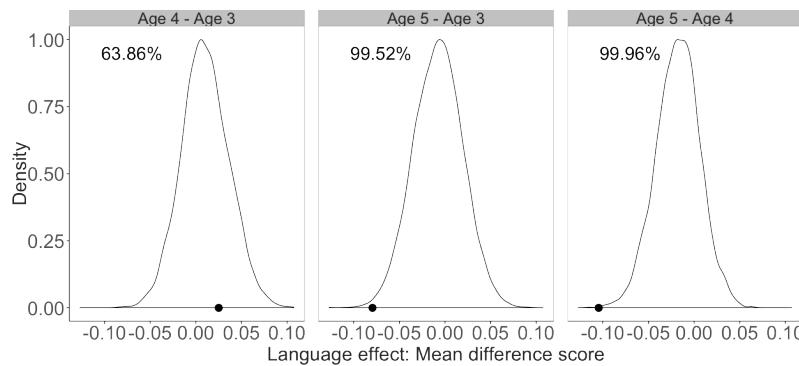


Figure B.1: Pairwise comparisons of the language condition effect across ages in Experiment 1.

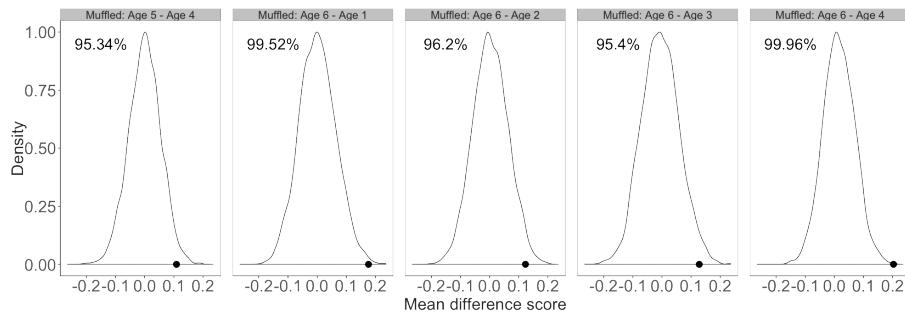


Figure B.2: Significant pairwise comparisons of the *prosody only-no speech* linguistic condition effect, across ages in Experiment 2

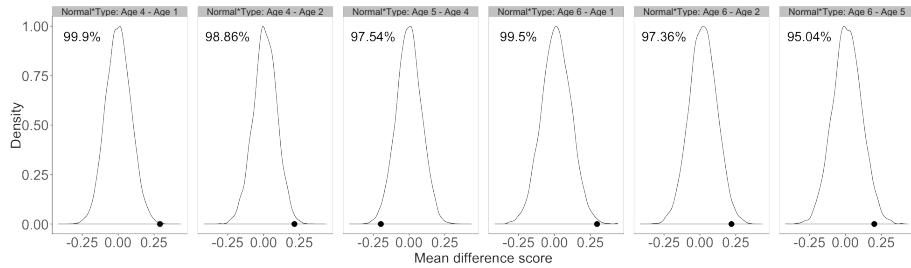


Figure B.3: Significant pairwise comparisons of the *normal speech-no speech* language condition effect for questions status, across ages, in Experiment 2.

1353 **Appendix C. Non-convergent models**