

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

We investigate language acquisition through the lens of a fundamental conversational skill: Turn taking. 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, with children across a wide developmental range, we measured spontaneous predictions about upcoming speaker change while controlling the amount of linguistic information available. We found that children already predicted upcoming turns at age one, but that they integrated linguistic cues differently at different ages. Children under three showed a general advantage for prosody over lexicosyntax, contrary to prior findings for adults. Children two and older showed more predictive switches for questions than non-questions, but only when lexical information was available. We found no evidence that lexicosyntax alone guides turn prediction—instead, participants' performance was best overall with access to lexicosyntax and prosody together. Our results point to the importance of studying children's linguistic processing in the context of conversation.

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

1. Introduction

Spontaneous conversation is a universal context for using and learning language. Like other types of human interaction, it is organized at its core

*Corresponding author

by the roles and goals of its participants. But what sets conversation apart is its structure: Sequences of interconnected, communicative actions that take place across alternating turns at talk. Sequential, turn-based structures in conversation are strikingly uniform across language communities and linguistic modalities. Turn-taking behaviors are also cross-culturally consistent in their basic features and the details of their implementation (De Vos et al., 2015; Dingemanse et al., 2013; Stivers et al., 2009). How does this ability develop?

Children participate in sequential coordination with their caregivers starting at three months of age—before they can rely on any linguistic cues in taking turns (see, among others, Bateson, 1975; Hilbrink et al., 2015; Jaffe et al., 2001; Snow, 1977). Of course, infant turn taking is different from adult turn taking in several ways. Infant turn taking is heavily scaffolded by caregivers, has distinct timing in comparison to adult turn taking, and lacks semantic content (Hilbrink et al., 2015; Jaffe et al., 2001). However, children’s early, turn-structured social interactions are presumably a critical precursor to their conversational turn taking. Along these lines, early non-verbal interactions establish the protocol by which children come to use language with others. And as they acquire language, they also come to integrate it into the preverbal turn-taking systems.

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

1.1. Turn taking

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

Conversation distinguishes itself from non-conversational turn-taking behaviors by the complexity of the turn sequencing involved. In the examples above (traffic, games, and monkeys) the set of sequence and action types is far more limited and predictable than what we find in everyday talk. For example, conversational turns come grouped into semantically-contingent sequences of action. The groups can span turn-by-turn exchanges (e.g., simple question–response, “How are you?”–“Fine.”) or sequence-by-sequence exchanges (e.g., reciprocals, “How are you?”–“Fine, and you?”–“Great!”). Sequences of action drive the conversation forward into the next, relevant sequences of talk (e.g., “And you?”–“Great!”–“Why’s that?”; Schegloff, 2007). To take a turn, participants need to make predictions about what conversational content will be relevant next. In some cases, relevant next turns are somewhat obvious (e.g., question–response) while, in other cases, there are multiple relevant next actions to choose from, or no obvious next action at all (e.g., after a closing).

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

1.2. Adults’ turn prediction

Adults have a lot of information at their disposal to help make accurate predictions about upcoming turn content. Lexical, syntactic, and prosodic information (e.g., *wh*- words, subject-auxiliary inversion, and list intonation) can all inform addressees about upcoming linguistic structure (De Ruiter et al., 2006; Duncan, 1972; Ford and Thompson, 1996; Torreira et al., 2015). Non-verbal cues (e.g., gaze, posture, and pointing) often appear at turn-boundaries and can sometimes act as late indicators of an upcoming speaker switch (Rossano et al., 2009; Stivers and Rossano, 2010). Additionally, the

sequential context of a turn can make it clear what will come next: Answers after questions; thanks or denial after compliments, et cetera (Schegloff, 2007).

Prior work suggests that adult listeners primarily use lexicosyntactic information to accurately predict upcoming turn structure (De Ruiter et al., 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 snippets were controlled for the amount of linguistic information present; some were normal, but others had flattened pitch, low-pass filtered speech, or further manipulations. De Ruiter and colleagues found that, with pitch-flattened speech, the timing of participants’ button responses was comparable to their timing with the full linguistic signal. But, when no lexical information was available, participants’ responses were significantly earlier. The authors concluded that lexicosyntactic information¹ was necessary and possibly sufficient for turn-end projection, while intonation was neither necessary nor sufficient. Congruent evidence comes from studies varying the predictability of lexicosyntactic and pragmatic content: Adults anticipate turn ends better when they can more accurately predict the exact words that will come next (Magyari and De Ruiter, 2012; see also Magyari et al., 2014). They can also identify speech acts within the first word of an utterance (Gísladóttir et al., 2015), allowing them to start planning their response at the first moment possible (Bögels et al., 2015).

The role of prosody for adult turn-prediction is still a matter of debate. De Ruiter and colleagues’ (2006) experiment focused on the role of intonation, which is only a partial index of prosody. Prosodic structure is also tied closely to the syntax of an utterance, and so the two linguistic signals are difficult to control independently (Ford and Thompson, 1996). Torreira, Bögels and Levinson (2015) used a combination of button-press and verbal responses to investigate the relationship between lexicosyntactic and prosodic cues in turn-end prediction. Critically, their stimuli were cross-spliced so that each item had full prosodic cues to accompany the lexicosyntax. Because of the splicing, they were able to create items that had syntactically-complete units

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.

with no intonational phrase boundary at the end. Participants never verbally responded or pressed the “turn-end” button when hearing a syntactically-complete phrase without an intonational phrase boundary. And when intonational phrase boundaries were embedded in multi-utterance turns, participants were tricked into pressing the “turn-end” button 29% of the time. Their results suggest that listeners actually do rely on prosodic cues to execute a response (see also de Ruiter et al. (2006):525). These experimental findings corroborate other corpus and experimental work promoting a combination of cues (lexicosyntactic, prosodic, and pragmatic) as key for accurate turn-end prediction (Duncan, 1972; Ford and Thompson, 1996; Hirvenkari et al., 2013).

In sum, adults accurately and spontaneously make predictions about upcoming turn structure. Their predictions rely on a sophisticated body of knowledge about linguistic structure, non-verbal signals, and social actions. Knowing this, we could expect that children’s acquisition of turn-taking skills is closely tied to their knowledge about language, gaze, gesture, and social cues. But children’s turn taking starts early in infancy, long before their first words or gestures emerge. So a primary role for lexicosyntactic cues doesn’t fit well with children’s pre-verbal turn taking.

1.3. Children’s turn prediction

1.3.1. Observational studies

The majority of work on children’s early turn taking has focused on observations of spontaneous interaction. Children’s first turn-like structures appear as early as two to three months in proto-conversation with their caregivers (Bruner, 1975, 1985). During proto-conversations, caregivers interact with their infants as if they were capable of making meaningful contributions; they take every look, vocalization, arm flail, and burp as “utterances” in the joint discourse (Bateson, 1975; Jaffe et al., 2001; Snow, 1977). Infants catch onto the structure of proto-conversations quickly. By three to four months they notice disturbances to the contingency of their caregivers’ response and, in reaction, change the rate and quality of their vocalizations (Bloom, 1988; Masataka, 1993). Infants at this age also notice changes to social contingency outside of turn structure. In the Still Face paradigm, caregivers interact with their infants and then suddenly halt, taking on a neutral expression with a sustained gaze. When faced with this sudden disappearance of social contingency, infants three months and older try a range of methods to reinitiate the

interaction, such as vocalization, reaching, and smiling before looking away or getting upset (Rochat et al., 1998; Toda and Fogel, 1993).

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

The protracted development of children’s timing may be attributable to their linguistic development: Taking turns on time is easier when the response is a simple vocalization rather than a linguistic utterance. Integrating language into the turn-taking system may be one major factor in children’s delayed responses (Casillas et al., In press). If response planning (i.e., language production) is the primary hurdle in children’s spontaneous turn taking, we should find evidence that children understand turn-taking behaviors before they are able to produce the behaviors themselves; this hypothesis has been recently explored in experimental settings.

1.3.2. Experimental studies

Children begin to develop specific expectations about conversational behavior before they begin to speak. Sometime between four and six months, children begin to attend differently to face-to-face and back-to-back conversation; six-month-olds follow conversational speakers more with their gaze when at least one speaker is looking at the other (Augusti et al., 2010). At ten months, infants expect people to look and talk at other people, and not to objects (Beier and Spelke, 2012). At twelve months infants expect to see responses to verbal (but not non-speech) utterances in face-to-face contexts (Thorgrímsson et al., 2015).

There are mixed results regarding when children begin to anticipate turn structure in conversation. One study found that 12-month-olds make more

predictive gaze shifts to a responder while watching human verbal conversation compared to conversation-like interactions with objects (Bakker et al., 2011), but another only found a similar effect at 36 months (von Hofsten et al., 2009). However, neither of these two studies had baselines to which the turn-relevant looking behavior could be compared. A baseline measurement is critical because there may be developmental differences in gaze shifting between conversational participants, even if the shifting is not related to turn structure. Such developmental differences could produce artifactual changes in measures of turn-contingent shifting.

Keitel and colleagues (2013) addressed the random baseline issue in their study of 6-, 12-, 24-, and 36-month-olds. They asked participants to watch short videos of conversation, and tracked their eye movements at points of speaker change. They found that children’s anticipatory gaze frequency was only greater than chance for 36-month-olds and adults. Their study was the first to focus on the role of linguistic processing in children’s turn predictions. They showed their participants two types of conversation videos: One normal and one with flattened pitch (i.e., with flattened intonation contours), finding that only 36-month-olds were affected by a lack of intonation contours. The adult control group made equal numbers of anticipatory looks in the videos, with and without intonation contours, consistent with prior adult findings (De Ruiter et al., 2006). Keitel and colleagues concluded that children’s ability to predict upcoming turn structure relies on their ability to comprehend the stimuli (emerging around 36 months), especially with respect to semantic access. They also suggest that intonation takes a secondary role in turn prediction, but only *after* children acquire more sophisticated, adult-like language comprehension systems (sometime after 36 months).

Although the Keitel et al. (2013) study constitutes a substantial advance over previous work, 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, during which switching is least likely to occur (Hirvenkari et al., 2013) and thereby maximizing their chances of finding a difference between gaze frequency at turn transitions and their baseline rate. A more conservative baseline would be to compare participants’ looking behavior at turn transitions to their looking behavior during randomly-selected windows of time throughout the stimulus. We follow this conservative approach in our work.

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

Despite these minor methodological issues, the Keitel et al. (2013) study still demonstrates intriguing age-based differences in children’s ability to predict upcoming turn structure, and the results suggest that both semantic and intonational development *do* play a role in children’s looking patterns. Our current work thus takes this paradigm as our starting point.²

1.3.3. Prosodic development

The roles of prosody and lexicosyntax in children’s turn predictions are currently unknown, but children understand more about prosody than lexicosyntax early in life. Children begin to acquire prosody in the womb, and can distinguish their native language’s rhythm type from others (e.g., syllable-timed vs. stress-timed) 2–5 days after birth (Mehler et al., 1988; Moon et al., 1993; Nazzi and Ramus, 2003). Beginning between four and five months, infants 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 (Jusczyk et al., 1995; Soderstrom et al., 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; Jusczyk et al., 1993; Morgan and Saffran, 1995). In comparison, children show only a limited lexical inventory at six months, and begin to recognize function words just before their first birthdays, with syntactic categorization beginning around 14 months (Bergelson

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

and Swingley, 2013; Shi and Melancon, 2010). Two-month-olds also notice changes in word order, but this ability appears to rely on prosodic cueing (Mandel et al., 1996). Generally speaking then, our current knowledge about children’s linguistic development points to a possible early advantage for prosody in children’s turn-taking predictions.

1.4. The Current Study

We report here on the role of linguistic processing in children’s predictions about upcoming turn structure. We focus in particular on how children use prosodic and lexicosyntactic information to make their predictions. Prior work has focused mainly on lexicosyntax and intonation, and not on prosody proper (De Ruiter et al., 2006; Keitel et al., 2013, but see Torreira et al., 2015), even though infants seem to acquire the basic rhythmic properties of the prosodic signal first (Mehler et al., 1988; Moon et al., 1993; Nazzi and Ramus, 2003).

In two eye-tracking experiments, we measured children’s anticipatory gaze to upcoming responders while controlling for the amount of lexicosyntactic and prosodic information available. In Experiment 1, English-speaking participants viewed video clips of naturalistic conversation from several different languages. We used multiple languages to control for the presence of lexicosyntactic information while keeping prosodic and non-linguistic cues intact. The results showed minimal differences between the predictive looking behavior of preschoolers and adults. In Experiment 2, we created artificial (puppet) visual scenes, enabling us to separately control lexicosyntactic and prosodic cues in the conversational stimuli. In this more controlled paradigm, we found 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 consistently looked faster to responders after hearing questions, compared to non-questions. Both prosodic and lexicosyntactic information played a role in children’s predictions about turn structure, but the two information sources were used differently at different ages. Our findings overall support an account in which predictive processes for turn taking in conversation are present early, but their integration with linguistic information takes substantial practice.

2. Experiment 1

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

2.1. Methods

2.1.1. Participants

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

2.1.2. Materials

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

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



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

topics (favorite foods, entertainment, hometown layout, and pets). Then we asked talkers to choose a new topic—one relevant to young children (e.g., riding a bike, eating breakfast)—and to improvise a dialogue on that topic. We asked them to speak as if they were on a children’s television show in order to elicit child-directed speech toward each other. We recorded until the talkers achieved at least 30 seconds of uninterrupted discourse with enthusiastic, child-directed speech. Most talker pairs took less than five minutes to complete the task, usually by agreeing on a rough script at the start. We encouraged talkers to ask at least a few questions to each other during the improvisation. The resulting conversations were therefore not entirely spontaneous, but were as close as possible while still remaining child-oriented in topic, prosodic pattern, and lexicosyntactic construction.³

After recording, we combined the audio and video files by hand, and cropped each recording to the 30-second interval with the most turn activity. Because we recorded the conversations in stereo, the male and female voices came out of separate speakers during video playback. This gave each voice in the videos a localized source (from the left or right loudspeaker). We coded each turn transition in the videos for language condition (English vs. non-English), inter-turn gap duration (in milliseconds), and speech act (question

³All of the non-English talkers were fluent in English as a second language, and some fluently spoke a third or more language. We chose male-female pairs as a natural way of creating contrast between the two talker voices. See an example run-through of the videos here: www.youtube.com/channel/UCGEZQcM9t8Zfjqqi_B1Q5Sw.

vs. non-question). The non-English stimuli were coded for speech act from a monolingual English-speaker’s perspective, i.e., which turns “sound like” questions, and which don’t: we asked five native American English speakers to listen to the audio signal for each turn and judge whether it sounded like a question. We then coded turns with at least 80% “yes” responses as questions.

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

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

2.1.3. Procedure

Participants sat in front of an SMI 120Hz corneal reflection eye-tracker mounted beneath a large flatscreen display. The display and eye-tracker were

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.

secured to a table with an ergonomic arm that allowed the experimenter to position the whole apparatus at a comfortable height, approximately 60 cm from the viewer. We placed stereo speakers on the table, to the left and right of the display.

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

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

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

2.1.4. Data preparation and coding

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

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

As mentioned, most anticipatory switches happen in the inter-turn gap, but we also allowed anticipatory gaze switches that occurred in the final syllables of the prior turn. Early switches are consistent with the distribution of responses in explicit turn-boundary prediction tasks. For example, in a button press task, adult participants anticipate turn ends approximately 200 msec in advance of the turn’s end, and anticipatory responses to pitch-flattened stimuli come even earlier (De Ruiter et al., 2006). We therefore allowed switches to occur as early as 200 msec before the end of the prior turn. For very early and very late switches, our requirement for 100 msec of fixation on each speaker would sometimes extend outside of the transition window boundaries (200 msec before and after the inter-turn gap). The maximally

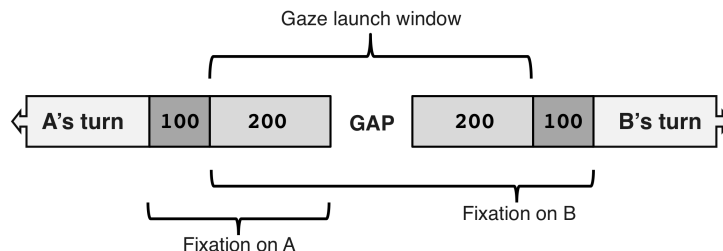


Figure 2: Schematic summary of criteria for anticipatory gaze shifts from speaker A to speaker B during a turn transition.

available fixation window was 100 msec before and after the earliest and latest possible switch point (300 msec before and after the inter-turn gap). We did not count switches made during the fixation window as anticipatory. We *did* count switches made during the inter-turn gap. The period of time from the beginning of the possible fixation window on the prior speaker to the end of the possible fixation window on the responder was our total analysis window (300 msec + the inter-turn gap + 300 msec).

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

2.2. Results and discussion

Participants looked at the screen most of the time during video playback (81% and 91% on average for children and adults, respectively). They primarily kept their eyes on the person who was currently speaking in both English and non-English videos: They gazed at the current speaker between 38% and 63% of the time, looking back at the addressee between 15% and 20% of the time (Table 1). Even three-year-olds looked more at the current speaker than anything else, whether the videos were in a language they could understand or not. Children looked at the current speaker less than

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.

adults did during the non-English videos. Despite this, their looks to the addressee did not increase substantially in the non-English videos, indicating that their looks away were probably related to boredom rather than confusion about ongoing turn structure. Overall, participants’ pattern of gaze to current speakers indicated that they performed basic turn tracking during the videos, regardless of language.

2.2.1. Statistical models

We identified anticipatory gaze switches for all 36 usable turn transitions, based on the criteria outlined in Section 2.1.4, and analyzed them for effects of language, transition type, and age with two mixed-effects logistic regressions (Bates et al., 2014; R Core Team, 2014). We built one model each for children and adults. We modeled children and adults separately because effects of age are only pertinent to the children’s data. The child model included condition (English vs. non-English)⁵, transition type (question vs. non-question), age (3, 4, 5), and duration of the inter-turn gap (seconds,

⁵Because each non-English language was represented by a single stimulus, we cannot treat individual languages as factors. Gaze behavior might be best for non-native languages that have the most structural overlap with participants’ native language: English speakers can make predictions about the strength of upcoming Swedish prosodic boundaries nearly as well as Swedish speakers do, but Chinese speakers are at a disadvantage in the same task (Carlson et al., 2005). We would need multiple items from each of the languages to check for similarity effects of specific linguistic features.

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 also 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 ($\beta=3.4$, $SE=1.229$, $t=2.77$, $p<.01$) with additional effects of a language condition-transition type interaction ($\beta=2.68$, $SE=1.35$, $t=1.99$, $p<.05$) and an age-language condition interaction ($\beta=0.52$, $SE=0.212$, $t=2.46$, $p<.05$). There were no significant effects of age or transition type alone ($\beta=-0.18$, $SE=0.175$, $t=-1.04$, $p=.3$ and $\beta=-1.10$, $SE=0.865$, $t=-1.27$, $p=.2$, respectively).

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

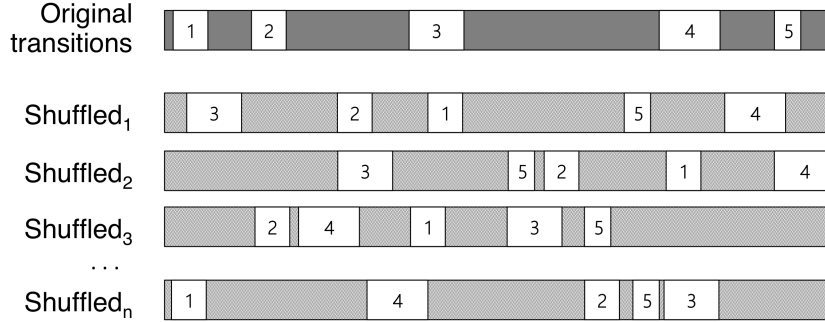


Figure 3: Example of shuffling for five turn transition analysis windows. The windows were ± 300 msec around the inter-turn gap.

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

Children

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.96146	0.84901	-1.132	0.257446
Age	-0.18268	0.17507	-1.043	0.296725
LgGroupNE	-3.29347	0.96045	-3.429	0.000606
TypeS	-1.10129	0.86494	-1.273	0.202925
Duration	3.40169	1.22826	2.770	0.005614
Age:LgGroupNE	0.52065	0.21190	2.457	0.014008
Age:TypeS	-0.01628	0.19437	-0.084	0.933232
LgGroupNE:TypeS	2.68166	1.35016	1.986	0.047013
Age:LgGroupNE:TypeS	-0.45632	0.30163	-1.513	0.130315

Adults

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.1966	0.6942	-0.283	0.776988
LgGroupNE	-0.8812	0.9602	-0.918	0.358754
TypeS	-4.4953	1.3139	-3.421	0.000623
Duration	-1.1227	1.9880	-0.565	0.572238
LgGroupNE:TypeS	3.2972	1.6101	2.048	0.040581
LgGroupNE:Duration	1.3626	3.0077	0.453	0.650527
TypeS:Duration	10.5107	3.3459	3.141	0.001682
LgGroupNE:TypeS:Duration	-6.3156	4.4926	-1.406	0.159790

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

2.2.2. Random baseline comparison

We estimated the probability that these patterns were the result of random looking by running the same regression models on participants’ real eye-tracking data, only this time calculating their anticipatory gaze switches with respect to randomly permuted turn transitions (Figure 3). This involved: (1) randomizing the order and temporal placement of the analysis windows (Figure 2) for each stimulus, thereby randomly redistributing the windows of analysis across the eye-tracking signal for each stimulus, then (2) re-running each participant’s eye tracking data through switch identification (described in 2.1.4) with the randomly permuted analysis windows, and finally (3) entering the resulting anticipatory gazes into the same statis-

tical model we used with the original data (above). Properties of each turn transition (e.g., prior speaker identity, transition type, duration, etc.) stayed constant, even though its onset and offset time within the eye-tracking signal was shuffled.

In effect, this procedure de-links participants' gaze data from the turn structure in the original stimulus, thereby allowing us to compare turn-related (original) and non-turn-related (randomly permuted) looking behavior from the same data. The random permutations represent average anticipatory gaze rate over all possible starting points; a random baseline. If we anticipatory gaze rates turned out to be the same for the original and randomly permuted versions of the data, we would have no evidence for turn-relevant gaze switching. Further, by running both data sets through the identical statistical models, we can estimate how likely it is that predictor effects in the original data (e.g., the effect of language condition) arose from random looking.

We completed this random baseline procedure on 5,000 permutations of the original turn transition analysis windows and then compared the beta estimates from the original data models (above) to the distribution of beta estimates for the models of the 5,000 randomly permuted datasets. We estimated whether each effect in our original statistical models differed from chance by calculating how many random beta estimates were smaller than the original beta estimate for each predictor, using the absolute value of all beta estimates for a two-tailed test. For example, children's original "language condition" beta estimate was $|-3.29|$, which is greater than 96.9% of $|\text{beta estimates}|$ from the models of randomly-permuted data. This leads us to conclude that the effect of language condition in the original model is unlikely to be the result of random gaze shifting. We excluded the output of random-permutation models that did not converge to remove unreliable β estimates from our percentile calculations below (22.5% and 24.4% of models for children and adults, respectively).

Most of the significant predictors from models of the original, turn-related data can not be explained by random looking. The children's data showed strong evidence of differentiation between the random and original beta estimates for the effect of language condition ($>96.9\%$ of random beta estimates) and the age-language condition interaction ($>93.9\%$), with somewhat weaker support for the language condition-transition type interaction ($>85.1\%$), and no support for the effect of gap duration ($>73\%$; Figure 4a). The adult data showed moderate evidence for the effect of transition type ($>88.6\%$), weak

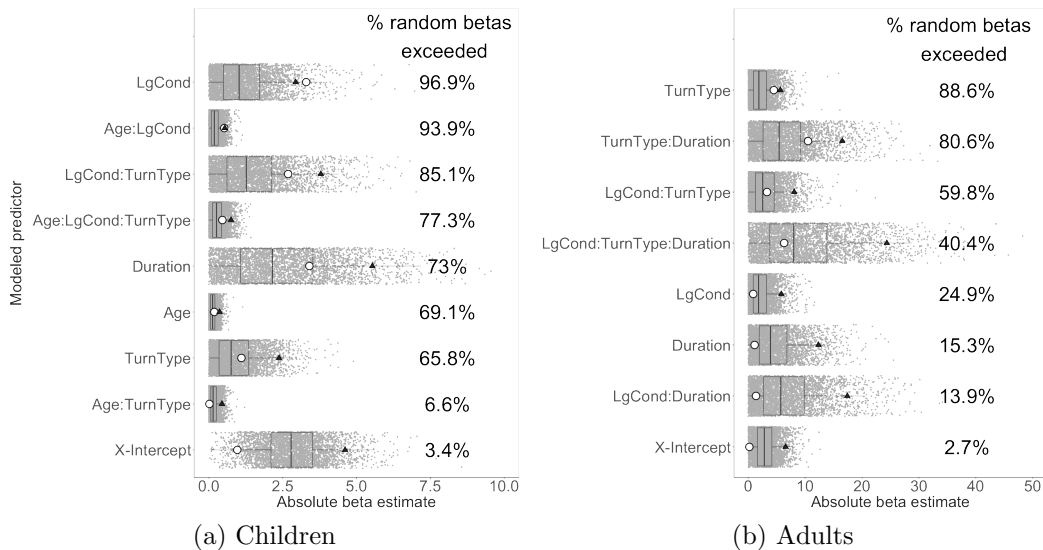


Figure 4: Random-permutation and original $|\text{beta estimates}|$ for predictors of children and adults’ anticipatory gaze rates. Gray dots = random model estimates, White dots = original model estimates, Triangles = 95th percentile for each beta estimate distribution.

evidence for the interaction of transition type and duration ($>80.6\%$), and no evidence for the interaction of language condition and transition type ($>59.8\%$; Figure 4b).

2.2.3. Summary

3. Experiment 2

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

Non-verbal cues in Experiment 1 (e.g., gaze and gesture) could have

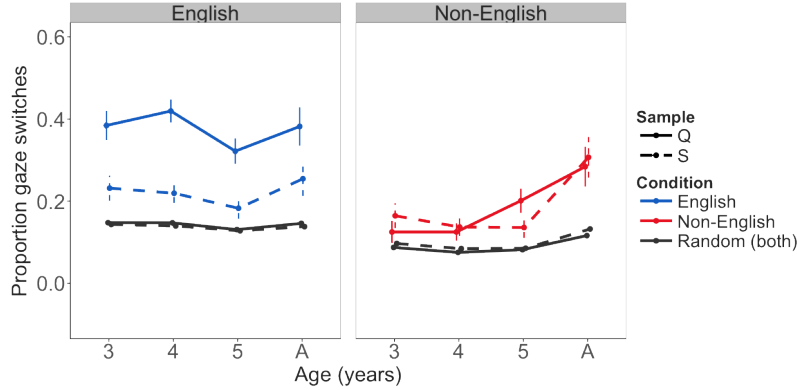


Figure 5: Anticipatory gaze rates across language condition and transition type for the real (red and blue) and randomly permuted (gray) data.

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

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

3.1. Methods

3.1.1. Participants

We recruited 27 undergraduate adults and 129 children between ages 1;0–6;11 to participate in our experiment. We recruited our child participants from the Children’s Discovery Museum in San Jose, California, targeting approximately 20 children for each of the six 1-year age groups (range=20–23).

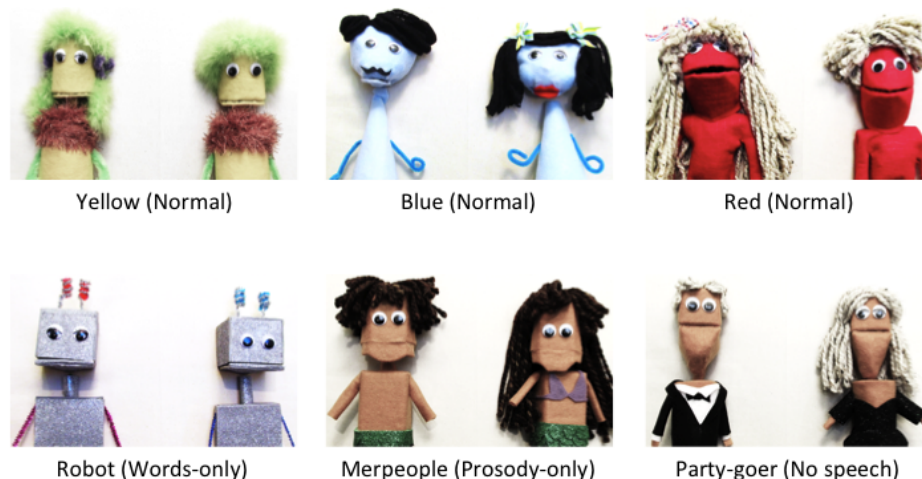


Figure 6: The six puppet pairs (and associated audio conditions). Each pair was linked to three distinct conversations from the same condition across the three experiment versions.

All participants were native English speakers, though some parents ($N=27$) reported that their child heard a second (and sometimes third) language at home. None of the adult participants reported fluency in a second language. We ran Experiment 2 at a local children’s museum because it gave us access to children with a more diverse range of ages.

3.1.2. Materials

We created 18 short videos of improvised, child-friendly conversation (Figure 6). To eliminate non-verbal cues to turn transition and to control the types of linguistic information available in the stimuli we (1) recorded improvised conversations, (2) phonetically manipulated those recordings to limit the availability of prosodic and lexical information, and (3) recorded a new set of videos that featured puppets as talkers, using the manipulated audio as the puppets’ speech.

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

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

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

The *words only* clips were manipulated to have robot-like speech: We flattened the intonation contours to each talker’s average pitch (F0) and we reset the duration of every nucleus and coda to each talker’s average nucleus and coda duration.⁷ We made duration and pitch manipulations using PSOLA resynthesis in Praat (Boersma and Weenink, 2012). Thus, the *words only* versions of the audio clips had no pitch or durational cues to upcoming turn boundaries, but did have intact lexicosyntactic cues (and residual phonetic correlates of prosody, like intensity).

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

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

Finally, because low-pass filtering removes significant acoustic energy, the *prosody only* clips were much quieter than the other three conditions. Our last step was to downscale the intensity of videos from the three other conditions to match the volume of the *prosody only* clips. We referred to the

⁷We excluded hyper-lengthened words like [wau:] ‘wooooow!’. These were rare in the clips.

conditions as “normal”, “robot”, “mermaid”, and “birthday party” speech when interacting with participants.

Video recordings. We created puppet video recordings to match the manipulated 20-second audio clips. The puppets were minimally expressive; the experimenter could only control the opening and closing of their mouths; their head, eyes, arms, and body stayed still. Puppets were positioned looking forward to eliminate shared gaze as a cue to turn structure (Thorgrímsson et al., 2015). We took care to match the puppets’ mouth movements to the syllable onsets as closely as possible, and avoided any mouth movement before the onset of a turn. We then added the manipulated audio clips to the puppet video recordings by hand.

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

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

There were 79 turn transitions for analysis, after excluding transitions longer than 550 msec and shorter than 0 msec. The remaining turn transitions were distributed evenly across transition types (questions N=47 and non-questions N=32) and conditions, keeping in mind that there were three *normal* videos and only one limited cue video for each experiment version: *Normal* (N=36), *words only* (N=13), *prosody only* (N=12), and no

⁸See the videos here: www.youtube.com/channel/UCGEZQcM9t8Zfjqqi_B1Q5Sw.

speech (N=18). Inter-turn gaps for questions (mean=358, median=405) were longer than those for non-questions (mean=296, median=288) on average, but gap duration was overall comparable across conditions: *Normal* (mean=333, median=337), *words only* (mean=345, median=383), *prosody only* (mean=320, median=336), and *no words* (mean=324, median=328). The longer gaps for question transitions could give them an advantage because our anticipatory measure includes shifts initiated during the gap between turns (Figure 2).

3.2. Procedure

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

Data preparation, coding, and random baseline analysis. We coded each turn transition for its linguistic condition (*normal*, *words only*, *prosody only*, and *no speech*) and transition type (question/non-question)⁹, identified anticipatory gaze switches to the upcoming speaker, and estimated a random baseline of anticipatory switches using the methods from Experiment 1.

3.3. Results and discussion

3.3.1. Summary

4. General Discussion

Acknowledgements

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

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).

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

References

- Augusti, E.M., Melinder, A., Gredebäck, G., 2010. Look who’s talking: Pre-verbal infants’ perception of face-to-face and back- to-back social interactions. *Developmental Psychology* 1, 161.
- Bakker, M., Kochukhova, O., von Hofsten, C., 2011. Development of social perception: A conversation study of 6-, 12-and 36-month-old children. *Infant Behavior and Development* 34, 363–370.
- Barr, D.J., Levy, R., Scheepers, C., Tily, H.J., 2013. Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language* 68, 255–278.
- Bates, D., Maechler, M., Bolker, B., Walker, S., 2014. lme4: Linear mixed-effects models using Eigen and S4. URL: <https://github.com/lme4/lme4/http://lme4.r-forge.r-project.org/>. [Computer program] R package version 1.1-7.
- Bateson, M.C., 1975. Mother-infant exchanges: The epigenesis of conversational interaction. *Annals of the New York Academy of Sciences* 263, 101–113.
- Beier, J.S., Spelke, E.S., 2012. Infants’ developing understanding of social gaze. *Child Development* 83, 486–496.
- Bergelson, E., Swingle, D., 2013. The acquisition of abstract words by young infants. *Cognition* 127, 391–397.
- Bloom, K., 1988. Quality of adult vocalizations affects the quality of infant vocalizations. *Journal of Child Language* 15, 469–480.
- Boersma, P., Weenink, D., 2012. Praat: doing phonetics by computer. URL: <http://www.praat.org>. [Computer program] Version 5.3.16.

- Bögels, S., Magyari, L., Levinson, S.C., 2015. Neural signatures of response planning occur midway through an incoming question in conversation. *Scientific Reports* 5.
- Bruner, J., 1985. Child's talk: Learning to use language. *Child Language Teaching and Therapy* 1, 111–114.
- Bruner, J.S., 1975. The ontogenesis of speech acts. *Journal of Child Language* 2, 1–19.
- Carlson, R., Hirschberg, J., Swerts, M., 2005. Cues to upcoming swedish prosodic boundaries: Subjective judgment studies and acoustic correlates. *Speech Communication* 46, 326–333.
- Casillas, M., Bobb, S.C., Clark, E.V., In press. Turn taking, timing, and planning in early language acquisition. *Journal of Child Language* .
- Casillas, M., Frank, M.C., 2012. Cues to turn boundary prediction in adults and preschoolers. *Proceedings of SemDial* .
- Casillas, M., Frank, M.C., 2013. The development of predictive processes in children's discourse understanding, in: *Proceedings of the 35th Annual Meeting of the Cognitive Science Society*.
- De Ruiter, J.P., Mitterer, H., Enfield, N.J., 2006. Projecting the end of a speaker's turn: A cognitive cornerstone of conversation. *Language* 82, 515–535.
- De Vos, C., Torreira, F., Levinson, S.C., 2015. Turn-timing in signed conversations: coordinating stroke-to-stroke turn boundaries. *Frontiers in Psychology* 6.
- Dingemanse, M., Torreira, F., Enfield, N., 2013. Is “Huh?” a universal word? Conversational infrastructure and the convergent evolution of linguistic items. *PloS one* 8, e78273.
- Duncan, S., 1972. Some signals and rules for taking speaking turns in conversations. *Journal of Personality and Social Psychology* 23, 283.
- Ervin-Tripp, S., 1979. Children's verbal turn-taking, in: Ochs, E., Schieffelin, B.B. (Eds.), *Developmental Pragmatics*. Academic Press, New York, pp. 391–414.

- Ford, C.E., Thompson, S.A., 1996. Interactional units in conversation: Syntactic, intonational, and pragmatic resources for the management of turns. *Studies in Interactional Sociolinguistics* 13, 134–184.
- Garvey, C., 1984. *Children’s Talk*. volume 21. Harvard University Press.
- Gísladóttir, R., Chwilla, D., Levinson, S.C., 2015. Conversation electrified: ERP correlates of speech act recognition in underspecified utterances. *PloS one* 10, e0120068.
- Griffin, Z.M., Bock, K., 2000. What the eyes say about speaking. *Psychological science* 11, 274–279.
- Hedberg, N., Sosa, J.M., Görgülü, E., Mamani, M., 2010. The prosody and meaning of Wh-questions in American English, in: *Speech Prosody 2010–Fifth International Conference*.
- Hilbrink, E., Gattis, M., Levinson, S.C., 2015. Early developmental changes in the timing of turn-taking: A longitudinal study of mother-infant interaction. *Frontiers in Psychology* 6.
- Hirvenkari, L., Ruusuvuori, J., Saarinen, V.M., Kivioja, M., Peräkylä, A., Hari, R., 2013. Influence of turn-taking in a two-person conversation on the gaze of a viewer. *PloS one* 8, e71569.
- von Hofsten, C., Uhlig, H., Adell, M., Kochukhova, O., 2009. How children with autism look at events. *Research in Autism Spectrum Disorders* 3, 556–569.
- Jaffe, J., Beebe, B., Feldstein, S., Crown, C.L., Jasnow, M.D., Rochat, P., Stern, D.N., 2001. Rhythms of dialogue in infancy: Coordinated timing in development. *Monographs of the Society for Research in Child Development*. JSTOR.
- Johnson, E.K., Jusczyk, P.W., 2001. Word segmentation by 8-month-olds: When speech cues count more than statistics. *Journal of Memory and Language* 44, 548–567.
- Jusczyk, P.W., Cutler, A., Redanz, N.J., 1993. Infants’ preference for the predominant stress patterns of English words. *Child Development* 64, 675–687.

- Jusczyk, P.W., Hohne, E., Mandel, D., Strange, W., 1995. Picking up regularities in the sound structure of the native language. *Speech perception and linguistic experience: Theoretical and methodological issues in cross-language speech research* , 91–119.
- Kamide, Y., Altmann, G., Haywood, S.L., 2003. The time-course of prediction in incremental sentence processing: Evidence from anticipatory eye movements. *Journal of Memory and Language* 49, 133–156.
- Keitel, A., Prinz, W., Friederici, A.D., Hofsten, C.v., Daum, M.M., 2013. Perception of conversations: The importance of semantics and intonation in childrens development. *Journal of Experimental Child Psychology* 116, 264–277.
- Lemasson, A., Glas, L., Barbu, S., Lacroix, A., Guilloux, M., Remeuf, K., Koda, H., 2011. Youngsters do not pay attention to conversational rules: is this so for nonhuman primates? *Nature Scientific Reports* 1.
- Levelt, W.J., 1989. *Speaking: From intention to articulation*. MIT press.
- Levinson, S.C., 2013. Action formation and ascriptions, in: Stivers, T., Sidnell, J. (Eds.), *The Handbook of Conversation Analysis*. Wiley-Blackwell, Malden, MA, pp. 103–130.
- Magyari, L., Bastiaansen, M.C.M., De Ruiter, J.P., Levinson, S.C., 2014. Early anticipation lies behind the speed of response in conversation. *Journal of Cognitive Neuroscience* 26, 2530–2539.
- Magyari, L., De Ruiter, J.P., 2012. Prediction of turn-ends based on anticipation of upcoming words. *Frontiers in Psychology* 3:376, 1–9.
- Mandel, D.R., Kemler Nelson, D.G., Jusczyk, P.W., 1996. Infants remember the order of words in a spoken sentence. *Cognitive Development* 11, 181–196.
- Masataka, N., 1993. Effects of contingent and noncontingent maternal stimulation on the vocal behaviour of three-to four-month-old Japanese infants. *Journal of Child Language* 20, 303–312.
- Mehler, J., Jusczyk, P., Lambertz, G., Halsted, N., Bertoncini, J., Amiel-Tison, C., 1988. A precursor of language acquisition in young infants. *Cognition* 29, 143–178.

- Moon, C., Cooper, R.P., Fifer, W.P., 1993. Two-day-olds prefer their native language. *Infant Behavior and Development* 16, 495–500.
- Morgan, J.L., Saffran, J.R., 1995. Emerging integration of sequential and suprasegmental information in preverbal speech segmentation. *Child Development* 66, 911–936.
- Nazzi, T., Ramus, F., 2003. Perception and acquisition of linguistic rhythm by infants. *Speech Communication* 41, 233–243.
- R Core Team, 2014. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing. Vienna, Austria. URL: <http://www.R-project.org>. [Computer program] Version 3.1.1.
- Ratner, N., Bruner, J., 1978. Games, social exchange and the acquisition of language. *Journal of Child Language* 5, 391–401.
- Rochat, P., Neisser, U., Marian, V., 1998. Are young infants sensitive to interpersonal contingency? *Infant Behavior and Development* 21, 355–366.
- Ross, H.S., Lollis, S.P., 1987. Communication within infant social games. *Developmental Psychology* 23, 241.
- Rossano, F., Brown, P., Levinson, S.C., 2009. Gaze, questioning and culture, in: Sidnell, J. (Ed.), *Conversation Analysis: Comparative Perspectives*. Cambridge University Press, Cambridge, pp. 187–249.
- Sacks, H., Schegloff, E.A., Jefferson, G., 1974. A simplest systematics for the organization of turn-taking for conversation. *Language* 50, 696–735.
- Schegloff, E.A., 2007. *Sequence organization in interaction: Volume 1: A primer in conversation analysis*. Cambridge University Press.
- Shi, R., Melancon, A., 2010. Syntactic categorization in French-learning infants. *Infancy* 15, 517–533.
- Snow, C.E., 1977. The development of conversation between mothers and babies. *Journal of Child Language* 4, 1–22.

- Soderstrom, M., Seidl, A., Kemler Nelson, D.G., Jusczyk, P.W., 2003. The prosodic bootstrapping of phrases: Evidence from prelinguistic infants. *Journal of Memory and Language* 49, 249–267.
- Stivers, T., Enfield, N.J., Brown, P., Englert, C., Hayashi, M., Heinemann, T., Hoymann, G., Rossano, F., De Ruiter, J.P., Yoon, K.E., et al., 2009. Universals and cultural variation in turn-taking in conversation. *Proceedings of the National Academy of Sciences* 106, 10587–10592.
- Stivers, T., Rossano, F., 2010. Mobilizing response. *Research on Language and Social Interaction* 43, 3–31.
- Takahashi, D.Y., Narayanan, D.Z., Ghazanfar, A.A., 2013. Coupled oscillator dynamics of vocal turn-taking in monkeys. *Current Biology* 23, 2162–2168.
- Thorgrímsson, G., Fawcett, C., Liszkowski, U., 2015. 1- and 2-year-olds’ expectations about third-party communicative actions. *Infant Behavior and Development* 39, 53–66.
- Tice (Casillas), M., Henetz, T., 2011. Turn-boundary projection: Looking ahead, in: *Proceedings of the 33rd Annual Meeting of the Cognitive Science Society*.
- Toda, S., Fogel, A., 1993. Infant response to the still-face situation at 3 and 6 months. *Developmental Psychology* 29, 532.
- Torreira, F., Bögels, S., Levinson, S.C., 2015. Intonational phrasing is necessary for turn-taking in spoken interaction. *Journal of Phonetics* 52, 46–57.
- Wittenburg, P., Brugman, H., Russel, A., Klassmann, A., Sloetjes, H., 2006. Elan: a professional framework for multimodality research, in: *Proceedings of LREC*.