

Developmental change in the speed of social attention and its relation to early word
learning

Daniel Yurovsky¹, Anna Wade², Allison M Kraus³, Grace W. Gengoux⁴, Antonio Y.
Hardan⁴, Michael C. Frank³

¹Department of Psychology, Carnegie Mellon University

²Department of Neurological Surgery, University of California, San Francisco

³Department of Psychology, Stanford University

⁴School of Medicine, Stanford University

Abstract

How do children learn words so rapidly? A powerful source of information about a new word's meaning is the set of social cues provided by its speaker (e.g. eye-gaze). Studies of children's use of social cues have tended to focus on the emergence of this ability in early infancy. We develop a paradigm to measure the development of this early-emerging ability and find slow, continuous changes in social attention over the first five years of life. Our findings suggest that the developing ability to allocate social attention may be a significant bottleneck on early word learning—in our paradigm, continuous changes in social attention are related to continuous changes in children's ability to learn new words. Further, we show that this bottleneck generalizes to children diagnosed with autism spectrum disorder, whose social information processing is atypical. These results suggest a potential route by which increases in social expertise can lead to changes in language learning ability, and more generally highlight the dependence of developmental outcomes not on just the existence of particular competencies, but on their proficient use in complex contexts.

Keywords: language acquisition, word learning, social cues, cognitive development

Developmental change in the speed of social attention and its relation to early word learning

Children's first years are a time of rapid change. One striking development is children's growing mastery of their native language: The typical child will go from saying her first word shortly before she turns one to producing 4000–5000 words by age five (Goulden, Nation, & Read, 1990). The pace and breadth of this transformation have led to a search for early-available, precocious mechanisms that support children's language acquisition. Perhaps because of the fundamentally discrete nature of the units of language itself, much of this search has focused specifically on pinpointing the earliest emergence of these mechanisms (Lidz, Waxman, & Freedman, 2003; Bulf, Johnson, & Valenza, 2011; Shukla, White, & Aslin, 2011). However, learning outside the laboratory is controlled not by early availability of these mechanisms, but by their proficient use in complex natural environments.

Here we take as a case study children's use of social information to infer the meanings of new words. Though even the earliest vocabularies contain words belonging to many grammatical categories, concrete nouns make up a large proportion (Bates et al., 1994). While acquiring a fully adult-like meaning for any of these nouns likely unfolds over multiple encounters, the very first problem a child faces when hearing a new word is referential uncertainty: Does the word refer to something in the current situation, and if so, what (Carey & Bartlett, 1978; Yu & Smith, 2007; Frank, Goodman, & Tenenbaum, 2009)? A powerful source of information for resolving this uncertainty is available in the social cues provided by the speaker; knowing where a speaker is looking is helpful for knowing what they are communicating about. Consequently, a large body of research documents and explores young infants' ability to track a speaker's social cues and use them to infer the target of her reference (e.g., Scaife & Bruner, 1975; Baldwin, 1993; Hollich, Hirsh-Pasek, & Golinkoff, 2000; Senju, Csibra, & Johnson, 2008).

Much of the work in this research program has focused on discovering the earliest point of infants' competence in using social information (e.g. Corkum & Moore, 1998; Brooks & Meltzoff, 2005; Csibra & Gergely, 2009). Although individual researchers often acknowledge the possibility of further developmental change, there has been relatively less emphasis on measuring this change. Thus, research that depends on this work often contains an implicit assumption that once the general ability to use social cues is demonstrated, children's proficiency with processing these cues is relatively high. This research strategy stands in contrast to work in domains like visual and motor development, or even spoken word recognition, in which researchers have sought to measure continuous improvement in children's abilities as they develop (Sokol, 1978; Banks, 1980; Forssberg, Eliasson, Kinoshita, Johansson, & Westling, 1991; Thelen, 1995; Fernald, Pinto, Swingley, Weinberg, & McRoberts, 1998). In these domains, continuous, quantitative changes can have as much impact on an infant's interaction with the world as qualitative changes (Adolph & Robinson, 2015). For example, the transition from crawling to walking leads to some advantages in children's ability to explore. However, infants' initial mobility after this transition is nowhere near as great as what they will achieve with another few months of practice (Adolph et al., 2012; Franchak, Kretch, & Adolph, 2018). The emergence of a behavior is only the beginning.

Our goal in this article is to follow this same developmental strategy—of measuring continuous improvement—to understand how children use social information to learn new words. Using social information to resolve referential uncertainty is a highly time-sensitive process of continuous re-allocation of attention between the speaker and the objects in the context. Indeed, recent work has shown that rapid gaze-following occurs relatively rarely in natural parent-child interactions, and is difficult even for older children and adults under some circumstances (Loomis, Kelly, Pusch, Bailenson, & Beall, 2008; Vida & Maurer, 2012; Yu & Smith, 2013). Thus, we hypothesize that competence in gaze following is not enough: the referential uncertainty problem in social interactions should remain a problem in

proportion to a children's developing ability to control their attention, process auditory and visual information, and hold this information in memory (Dempster, 1981; Kail, 1991; Gathercole, Pickering, Ambridge, & Wearing, 2004).

Previous studies have used lab-based exposures to investigate links between social attention and word learning (Hollich et al., 2000; Moore, Angelopoulos, & Bennett, 1999), but were designed primarily to discover the earliest points of emergence of competencies rather than measuring their change with high fidelity. Thus, to investigate this issue, we refined these previous paradigms, developing a video-based paradigm that allowed for the measurement of both social attention and novel word learning (c.f. Yurovsky & Frank, 2017). First, we constructed videos that used novel words in a series of naturalistic object-focused dialogues and monologues. These videos were designed to be sufficiently difficult in their structure that in-the-moment disambiguation would pose a significant challenge to young learners, yet sufficiently simple that the repeated co-occurrence of novel words and the objects they label could allow for successful word learning. Any experimental paradigm requires a compromise between different goals. In this case, because we were interested in using eye-tracking to measure social attention, a third-person, video-based paradigm was necessary; however, these elements of the paradigm might make it more difficult for younger children (Anderson, 2005). We return to these issues in the General Discussion.

Using this paradigm, we then measured children's eye-movements during viewing as an index of their online inferences about the current conversational referent, and then tested their retention of the words they learned via a series of forced-choice test trials. These data allow us to test two predictions: (1) although young infants show measurable competence in social information processing, this ability has a long developmental trajectory, and (2) this developing ability is linked to language learning: children who successfully attend to speakers' social cues are more likely to learn and retain the novel words they hear the speakers produce.

In Experiment 1, we measured developmental changes in social information processing and used these to predict word learning in typically developing children. In Experiment 2, we tested a sample of children diagnosed with autism spectrum disorder (ASD) in the same paradigm as Experiment 1. Social attention was related to this group's learning from these videos in the same way as for the typically developing children. In Experiment 3, we manipulated the timing of social information directly to test the causal role of fast social attention, again in typically developing children. In all three studies, we recruited children across a broad age range, giving us the ability to measure continuous changes in both social attention and novel word learning.

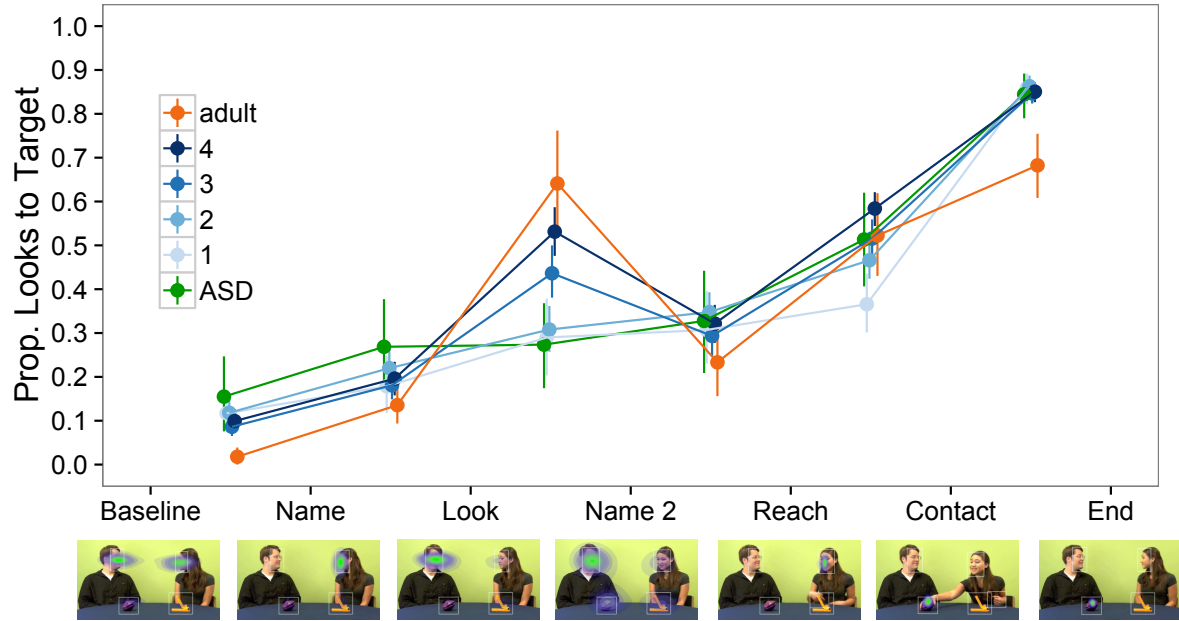


Figure 1. (top) Proportion of time looking to Target objects during the learning portion of Experiments 1 (typically-developing children and adults) and 2 (children with ASD), plotted by phase of the naming event. Points show means and error bars show 95% confidence interval across participants; points are offset on the horizontal to avoid overplotting. Colors indicate age in years for typically developing children. (bottom) Example frames from the first word learning dialogue in Experiment 1. Each image shows the regions of interest used for later analysis (white boxes) and a heat map of the distribution of all participants' points of gaze over time (brighter colors indicate more fixation; scale is constant across frames). See Supporting Information for a video representation.

Experiment 1: Social Attention and Word Learning

Experiment 1 was designed to estimate the developmental trajectory of children’s online social information processing in complex interactions, and to ask whether this processing is related to later word learning. The experiment consisted of two sections: Learning and Test. During the learning section, participants watched a series of dialogues and monologues in which actors introduced and discussed two novel objects and two familiar objects. Each video contained two naming sequences (pictured in Fig. 1, bottom); these sequences were broken into a series of phases during which the speaker named, looked at, named again, and reached for a particular object, allowing for the separate measurement of name-, gaze-, and reaching-related changes of attention.

We first measured the proportion of time that participants spent looking at the target referent during the naming sequences (Fig. 1, top). We analyzed these gaze trajectories by fitting a mixed-effects model predicting looking at the target from naming phase, age, and their interaction. Children increased their looking to the target object over the course of learning trials, with above-baseline looks to the target in all phases after the second naming. For all age groups, looking to the target toy reached nearly 100% by the end of these events—after the speaker had made contact with the toy ($\beta_{name2-reach} = .18, t = 3.73, p < .001$; $\beta_{reach-contact} = .13, t = 2.56, p < .01$; $\beta_{contact-end} = .75, t = 15.8, p < .001$). The effect of age was not significant, but there was a significant interaction between age and phase—older children looked significantly more in the phases after the speaker’s initial look and initial reach ($\beta_{look-name2} = .11, t = 7.38, p < .001$; $\beta_{reach-contact} = .08, t = 5.72, p < .001$). While the youngest children were only occasionally able to follow the speaker’s social cues, older children and adults were much more consistent. Thus, in our paradigm, the ability to process social information quickly and reliably improves markedly over the course of the first five years, and even further into adulthood.

After watching these naming events, children’s learning of the novel words was tested using the looking-while-listening procedure (Fernald et al., 1998; Fernald, Zangl, Portillo, &

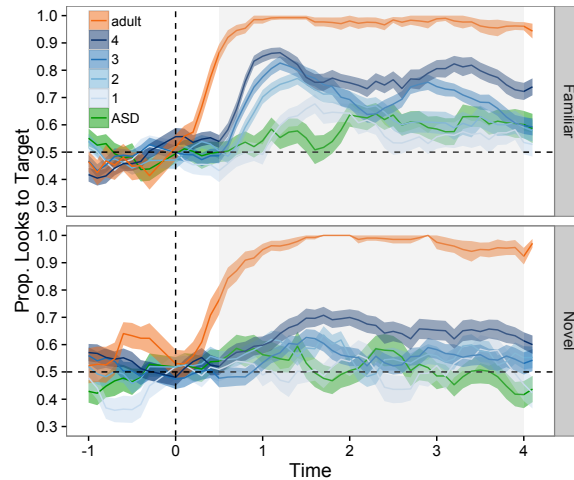


Figure 2. Children’s and adults’ looks to the Target and Competitor objects over the course of Test trials for typically developing children and adults (Experiment 1) and children with ASD (Experiment 2). Lines show age-group means, and shaded regions show standard errors. The light gray rectangle shows the window over which looking proportions were computed for subsequent statistical analyses.

Marchman, 2008). On Test trials, children saw two toys on the screen and heard a voice asking them to find the target toy. On some trials, the target was a one of the Novel toys from the Learning trials (e.g. “fep”). On other trials, the target was a familiar object whose label is typically in the comprehension vocabularies of young children (e.g. “dog”). These Familiar trials allowed us to measure children’s language comprehension more generally (Fig. 2). Children in all age groups successfully looked at Familiar referents at above chance levels (smallest $\mu_{1\text{-year}} = .57$, $t(28) = 4.05$, $p < .001$), and all but the one-year-olds reliably learned the Novel words (smallest $\mu_{2\text{-year}} = .56$, $t(51) = 3.96$, $p < .001$). A linear mixed-effects model showed that children performed better on Test trials over development ($\beta_{age} = .07$, $t = 7.89$, $p < .001$), and that the effect of age was greater for Familiar than Novel trials ($\beta_{age*novel} = -.03$, $t = -2.47$, $p < .05$, Fig. 2)

Thus, older children learned more from the same naming events. These children also more quickly followed the speakers’ social gaze and reaches, and more quickly processed

Familiar words. All of these factors were independently, significantly correlated with success on Test trials ($r_{age}(196) = .29, p < .001$; $r_{familiar}(196) = .20, p < .01$; $r_{look-name2}(189) = .34, p < .001$; $r_{reach-contact}(191) = .17, p < .05$). Which factor was most responsible for improvements in learning? To answer this question, we fit a linear regression, predicting learning from age, familiar word processing, and looking to the target in the two relevant windows—after the look, and after the reach. Only looking to the target referent following the initial look reached significance ($\beta_{look-name2} = .14, t = 3.26, p < .01$), although age was marginal ($\beta_{age} = .02, z = 1.60, p = .11$). Because these predictors were all correlated, we also fit this same model after first residualizing out the effect of age on learning. In this model, gaze-following still remained significant ($\beta_{look-name2} = .11, t = 2.95, p < .01$).

Children who did not follow the speaker’s social cues to determine the toy that was the topic of conversation did not learn that it was associated with the novel word in the discourse. Thus, age-related improvements in social attention, whatever their cause, were strongly related to how much children learned from naming events. If these results generalize from our paradigm to children’s real-world environment, this finding would suggest that the early emergence of competence in following social cues is no guarantee of its proficient use in complex contexts.

Experiment 2: Social Attention in Children with Autism

The naming events in Experiment 1 contained information sufficient to infer the meanings of the novel words at two timescales. First, children could learn the meanings within the interactions by following the informative social cues. Second, children could learn these meanings across the interactions by using co-occurrence statistics between the words and objects (Smith & Yu, 2008). In our data, typically developing children’s social information processing predicted considerable variance in their word learning, suggesting

that they may have learned largely via the social cues. Do children with atypical trajectories of social development rely more on co-occurrence statistics instead?

Children on the autism spectrum are one population whose strategy might be expected to differ. Deficits in social information processing in children with autism spectrum disorders have profound consequences for language learning (Baron-Cohen, Baldwin, & Crowson, 1997; Leekam, Hunnisett, & Moore, 1998). One possibility is that because of these deficits, if children with autism learn in our paradigm, it would be due to cross-situational statistical learning. Another possibility however is that their social information processing impairments lie on a continuum with other, less extreme changes in social information processing that are still impactful for language learning (Brooks & Meltzoff, 2005). On this second account, their social information processing should be related to their learning outcomes, just as in typically-developing children.

To address this question, we tested a group of 40 2–8-year-old children diagnosed with autism spectrum disorder (ASD). These children did indeed process social information less efficiently—following the speaker’s gaze no better than the youngest children in our sample (Fig. 1). They also learned the Novel words, and attended to the referents of Familiar words at approximately the same levels (Fig. 2). However, while age was uncorrelated with Novel word learning for the children with ASD ($r(28) = .11, p = .57$), following the speaker’s social gaze was highly correlated with Novel word learning ($r(21) = .58, p < .01$), just as in the typically developing sample (Fig. 3). We again fit a linear mixed-effects regression, predicting Novel word learning from age, Familiar word processing, and looking to the Target in the two relevant windows—after the look, and after the reach. This model showed significant effects of Familiar word processing ($\beta_{\text{familiar}} = .51, t = 2.48, p < .05$) and gaze following ($\beta_{\text{look-name2}} = .25, t = 3.06, p < .01$).

Thus, social information processing is a bottleneck on word learning for children with autism spectrum disorder, just as it is for typically developing children. Although as a group these children’s looking behavior during learning trials was different from the looking

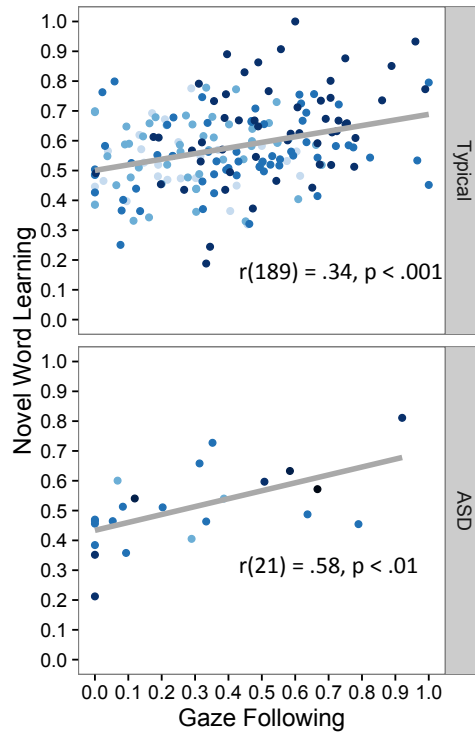


Figure 3. Correlations between gaze-following on Learning trials and accuracy on Novel Test trials for typically developing children in Experiment 1 and children with autism spectrum disorder (ASD) in Experiment 2. Darker colored points indicate older children.

behavior of their typically developing peers, children with ASD who successfully learned the Novel words did so by following social cues.

Experiment 3: Varying Demands on Social Attention

The naming events in Experiments 1 and 2 provided children with a number of informative cues to the referents of the novel words: social gaze, speaker's manual interaction, and also cross-situational co-occurrence statistics. Our analyses showed that fast processing of social gaze was a strong predictor of children's ultimate word learning. In Experiment 3, we tested this prediction directly by manipulating the accessibility of the social cue.

In Experiment 3, the two Novel words were introduced in two different kinds of naming events. In Extended Hold events, the speaker made contact with and manipulated the Target toy for the duration of the naming event, providing an extended cue to the target of her referential intention. In contrast, in Brief Look events, the speaker only provided punctate gaze information, looking to the target of her referential utterance briefly after naming it and then looking forward towards the camera for the remainder of the naming event. On the basis of the results from Experiments 1 and 2, we predicted that Extended Hold trials would show less developmental differentiation in social attention and be easier to learn from. In contrast, we predicted that following the Brief Look would require rapid reallocation of social attention, and thus that older children would succeed in following the speaker's gaze while younger children failed. Further, we predicted that this difference in gaze-following would produce down-stream differences in learning.

As predicted, on Extended Hold trials, children in all age groups spent a large proportion time looking at the Target object ($\mu_{1-year} = .58$, $\mu_{2-year} = .62$, $\mu_{3-year} = .65$, $\mu_{4-year} = .64$), and age was weakly correlated with looking behavior ($r(288) = .14$, $p = < .05$). In contrast, children spent a much lower proportion of Brief Look events looking at the Target, and this proportion increased across much more across development ($\mu_{1-year} = .10$, $\mu_{2-year} = .12$, $\mu_{3-year} = .21$, $\mu_{4-year} = .25$; $r(287) = .45$, $p < .001$). To test whether these correlations were different, we fit a mixed-effects model predicting looking at the target on these learning trials from cue type (Brief vs. Extended), age, and their interaction. The interaction term was significant, indicating that age predicted more of the variance in looking on Brief Look trials ($\beta_{BriefLook*age} = .03$, $t = 3.64$, $p < .001$).

As before, children in all age groups showed evidence of processing Familiar words at above chance levels (smallest $\mu_{1-year} = .59$, $t(63) = 5.04$, $p < .001$). Children two-years-old and older showed evidence of learning from the Extended Hold trials (smallest $\mu_{2-year} = .61$, $t(66) = 5.70$, $p < .001$), but only the 3- and 4-year-olds learned from the

Brief Look trials (smallest $\mu_{3-year} = .58$, $t(67) = 4.25$, $p < .001$). To confirm these analyses, we fit a linear mixed effects model predicting looking at the correct referent on Test trials from children's age and the trial type. Children improved significantly over development ($\beta_{age} = .05$, $t = 13.83$, $p < .001$). Children's Test trial performance also varied significantly across trial types. Relative to novel words encountered in Extended Hold events, children performed better on Familiar words ($\beta_{Familiar} = .11$, $t = 8.85$, $p < .001$) and worse on Novel words encountered in Brief Look events ($\beta_{BriefLook} = -.08$, $t = -6.74$, $p < .001$).

These results confirm our first predictions: Children improved in following social cues across developments, and the size of this improvement was larger when the the cue required rapid re-allocation of attention. Children also learned less from these hard-to-follow Brief Looks. We next provide two convergent analyses that confirm our final prediction: Individual differences in looking behavior on learning trials predicted individual differences in learning.

First, we fit a mixed effects model predicting children's looking on test trials from their age, their Familiar word processing, and their looking on learning trials. This model showed significant effects of all three predictors, confirming that individual differences in social information processing predicted individual differences in learning ($\beta_{age} = .05$, $t = 5.08$, $p < .001$; $\beta_{Familiar} = .17$, $t = 2.41$, $p < .05$, $\beta_{learning} = .15$, $t = 6.23$, $p < .001$). This model was not improved by adding social cue type, and predicted significantly more variance than a model in which we included cue type instead of individual learning scores ($\chi^2 = .75$, $p < .001$). Thus, experimentally manipulating the social cue induced individual differences in children's social attention, and these individual differences predicted individual differences in downstream learning.

Finally, we asked whether this relationship varied between the Extended Hold and the Brief Look. Individual differences in children's ability to follow the Extended Hold were marginally correlated with their learning ($r(255) = .11$, $p = .08$). In contrast individual differences in children's success in following the Brief Look trials predicted individual

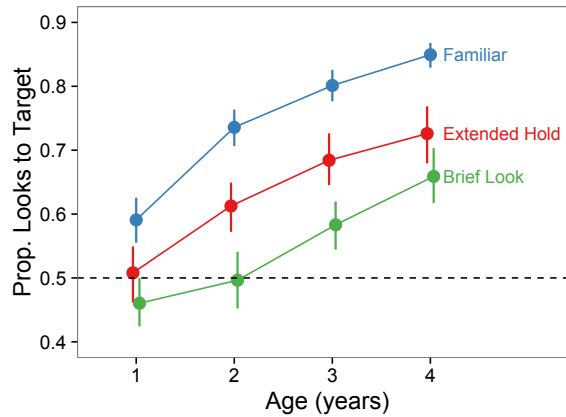


Figure 4. Test trial performance in Experiment 3 for the Familiar words, as well as for Novel words that appeared in Brief Look and in Extended Hold learning events. Error bars show 95% confidence intervals computed by non-parametric bootstrap; points are offset on the horizontal to avoid overplotting.

differences in learning from Brief Looks ($r(256) = .32, p < .001$). To test whether these correlations were different, we fit a mixed-effects model predicting test accuracy from cue type (Brief vs. Extended), proportion of time looking at the target during learning trials, and their interaction. We found a significant interaction, indicating that variability in children’s social attention were more consequential when demands on social information processing were greater; when the speaker provided only a Brief Look ($\beta_{ExtendedHold*learning} = -.22, t = -2.22, p < .05$).

Together, these results demonstrate the powerful role of fast social information processing in early word learning. All of the children were able to follow the speaker’s social cues when they Extended over the entire naming event, and this supported their learning. In contrast, only the older children were able to follow the Brief Look, and only they were consequently able to learn from it.

Discussion

The early emergence of children's linguistic and communicative capacities is extraordinary. Infants show evidence of following social cues and understanding the communicative function of language by 6 months of age (Senju et al., 2008; Vouloumanos, Martin, & Onishi, 2014). By the same early age, infants can learn about the structure of their languages by tracking the distributional properties of the speech they hear, and even appear to have at least nascent meanings for some words (Thiessen & Saffran, 2003; Bergelson & Swingley, 2012). A large body of research in developmental psychology has taken the early emergence of these abilities as evidence for expertise, implicitly or explicitly endorsing the idea that children learn language quickly because they are expert social and distributional information processors. But early competence is not enough to produce rapid learning. Children must also be able to *perform* these abilities rapidly and robustly in complex real-world settings.

We built on previous work to develop a paradigm intended to measure this type of performance. To do so, we used eye-tracking during children's viewing of a video display. In our paradigm, Social attention this had a long developmental trajectory, improving dramatically over the first five years of life. While one-year-olds were able to follow social gaze occasionally in our experiments, this level of performance did not translate into novel word learning. And even four-year-olds, whose performance was much better, still did not shift their social attention as flexibly as adults.

This relationship between individual performance in social attention and novel word learning also characterized the behavior of our sample of children with autism spectrum disorder. Autism is a complex and heterogeneous disorder, affecting different children to different degrees, and in different ways. Nonetheless, one of the core deficits appears to be a disorder of social attention. This deficit was manifest in our data—children with autism spectrum disorder performed substantially lower than typically developing children in following social cues and learning new words. But critically, variability in social attention

among these children was related to variability in learning just as with our typically developing sample. We take this as evidence for the generality of the importance of social information to word learning: Successful learners were successful in the same way.

These experiments are a first step towards measuring continuous developmental change in social attention and its relation to word learning, and the generality of our results must be tested by future work.

Although we endeavored to recruit a large and developmentally diverse set of participants, the primary data in Experiments 1 and 2 are correlational in nature. Experiment 3 was designed to causally manipulate social information, providing stronger evidence for the correlation observed in Experiments 1 and 2.

In addition, to minimize variability in stimulus presentation across children, we used a set of fixed video stimuli. Because videos are non-contingent in their timing, they may be difficult, especially for younger children. Thus the use of videos might lead to underestimates of social attention and learning (Anderson, 2005), especially for the youngest children in the sample. However, we note that some work shows that the “video deficit” in learning is ameliorated when children observe reciprocal interactions like the dialogues in Experiment 1 (O’Doherty et al., 2011).

Finally, to minimize measurement noise, in our experiments we artificially fixed the learning environment for each child and made the child a static, third-person observer. In contrast, in the natural context of learning, input varies considerably across children and across development and is observed from a first-person perspective (Yu & Smith, 2013). Further, the types of cues to social attention may vary as well (Franchak et al., 2018; Kretch, Franchak, & Adolph, 2014). We do not know how changes in speed of social attention would affect learning in a more natural context.

Although these experiments provide evidence for pronounced improvement in children’s social attention, they leave open the question of what is responsible for these changes. One possibility is that these changes in social attention are caused by

domain-general developmental changes in attentional control more broadly (Rueda & Rothbart, 2005; Smith, 2013) or even general changes in speed of processing (Kail, 1991). Alternatively, these children could be refining their domain-specific representations of the visual and temporal structure of conversations, producing better predictions about where speakers will look and reach next (Acheson & MacDonald, 2009; Krogh-Jespersen, Liberman, & Woodward, 2015). In either case, we propose that these changes might play a role in explaining the rapidly accelerating pace of children's language learning and that future work should investigate this question. Indeed, while much has been made of the changes in word learning that occur in the first and second years, the rate at which children learn words continues to increase over the third, fourth, and fifth years (Bloom, 2000).

If infants' social information processing performance is so poor, why do they learn words so rapidly? Although substantial work will be required to flesh out this connection, one intriguing possibility is that parents may tune their visual and linguistic input to their children's developing social information processing skills (Snow, 1972; Gogate, Bahrick, & Watson, 2000; Brand, Baldwin, & Ashburn, 2002). That is, caregivers may not use Brief Looks to indicate their referents for young children, but instead use Extended Holds (Experiment 3). Perhaps rapid early word learning is not a result of early expertise per se, but instead emerges from the interaction between children's developing processing skills and their caregivers' coordinated linguistic and social input.

Method

Disclosures

We report how we determined our sample size, all data exclusions (if any), all manipulations, and all measures in the study.

Participants

All sample sizes were a result of data collection across a defined period of time (e.g., summer session). No data collection decisions were dependent on data analysis. Data from typically-developing children were collected at the San Jose Children's Discovery Museum, where parents and their children were invited to participate in an experiment investigating children's early word learning after providing informed consent. In Experiment 1, we collected both demographic and eye-tracking data from 349 children in the target age range, of whom 113 were excluded from the final sample for one or more of the following reasons: unacceptable eye-tracker calibration or data ($N = 49$), parent-reported atypical development ($N = 30$), and less than 75% parent-reported exposure to English ($N = 41$). Our final sample included 238 children (42 1-year-olds (21 girls); 66 2-year-olds (32 girls); 74 3-year-olds (36 girls); and 56 4-year-olds (30 girls)). In Experiment 3, we collected demographic and eye-tracking data from 425 children in the target age range, of whom 200 were excluded from the final sample for one or more of the following reasons: unacceptable eye-tracker calibration or data ($N = 118$), parent-reported atypical development ($N = 36$), and less than 75% parent-reported exposure to English ($N = 46$). Our final sample included 225 children, ages 1—5 (61 1-year-olds (25 girls); 57 2-year-olds (30 girls); 51 3-year-olds (25 girls); and 50 4-year-olds (24 girls)).

Data from children with ASD were collected at Stanford University. Children were primarily recruited through the Autism and Developmental Disorders Research Registry, and by flyers posted in the Autism and Developmental Disorders Clinic. Children with a diagnostic history of ASD underwent a comprehensive diagnostic evaluation to determine the accuracy of the previous diagnosis based on DSM-5 criteria, which was confirmed with research diagnostic methods. These diagnostic methods included the ADI-R (Lord, Rutter, & Le Couteur, 1994; Le Couter et al., 1989) and the Autism Diagnostic Observation Schedule—Generic (ADOS-G) (Lord, Rutter, DiLavore, & Risi, 1999; Lord, Risi, Lambrecht, & Cook Jr, 2000). Exclusion criteria included: 1) a genetic, metabolic, or

infectious etiology for ASD on the basis of medical history, neurological history, and available laboratory testing for inborn errors of metabolism and chromosomal analysis; and 2) a DSM-5 diagnosis of any severe mental disorder (e.g. schizophrenia or bipolar disorder). We collected demographic and eye-tracking data from 51 1–7-year old-children, of whom 10 were excluded for unacceptable eye-tracker calibration. The final sample comprised 41 children ($M_{age} = 4.22$ -years, range = (2.24–7.97-years), 4 girls).

Adult participants in Experiment 1 were 17 Stanford undergraduates who participated in exchange for course credit. Informed consent was obtained from parents of all children, and from all adult participants, before experiments began.

Stimulus and Design

After an initial eye-tracker calibration phase, children and adults in both experiments watched a ~ 6 minute video. Each video presented two kinds of trials. Learning trials consisted of videos in which speakers seated at a table with two toys provided social cues and labeled one of the toys. Test trials showed pictures of two objects on a black background while a voice asked the participant to look at one of them (as in Fernald et al., 1998).

In addition, each video included a re-calibration stimulus in which a small, brightly colored object move around the screen. These phases of the videos were used to correct the calibrations estimated at the beginning of the experiment (see Frank, Vul, & Saxe, 2012). Finally, videos contained a small number of filler trials consisting of engaging pictures or videos designed to maintain participants’ attention.

Learning trials varied across the Experiments. In Experiment 1 and 2, these learning trials consisted of two dialogues in which two speakers sat at a table together, and one referred to each of the two novel toys whose names were taught in the experiment (“toma” and “fep”). The other two learning trials were monologues in which one of the speakers sat

at a table alone, with one of the novel toys and one familiar toy and referred to each in turn. Each naming event consisted of six events: first a naming phrase, then a look at the object accompanied by a comment, a second naming, a reach for the object, and finally a demonstration of the object's function accompanied by a third naming. Each novel object was named nine times in total over the course of the video. Participants watched all of the learning trials before they began the test trials. Eight of these test trials tested familiar objects (e.g. dog/car or lamp/carrot); the other eight paired the two novel objects. Naming phrases were of the form "Look at the [car/fep]! Do you see it?" and were spoken by the actors in the video.

In Experiment 3, we simplified the learning trials. All naming sequences were monologues and both toys on the table were novel. Four of the naming sequences were Extended Hold trials, in which the speaker reached for and interacted with one of the two toys while describing its function and producing its name three times. The other four were Brief Look trials in which the speaker produced the same kind of description but indicated the target toy only with a brief look after first producing its label. Brief Look and Extended Hold trials used a distinct but consistent set of two toys, and the same toy was consistently either the target or the competitor on each trial type. Because children in Experiment 1 sometimes lost interest in these videos before reaching the test trials, in Experiment 3 test trials were interspersed with learning trials so that at least some test data could be acquired from each child. In total, Experiment 3 contained 20 test trials. Eight of these test familiar objects as in Experiments 1 and 2. Eight paired the named objects from the learning trials against their foils from the learning trials, four for the object named in Brief Look trials, and four for the object named in Extended Hold trials. The remaining four test trials paired the two named objects against each other, with children being asked to find each two times.

Data Analysis

In all experiments, raw gaze data were transformed before statistical analysis. First to ensure appropriate precision in area-of-interest analyses, infants' calibrations were corrected and verified via robust regression (described in Frank et al., 2012), and calibration corrections were assessed by two independent coders ($\kappa_{Exp1} = .8$, $\kappa_{Exp2} = .87$, $\kappa_{Exp3} = .77$). Children whose calibrations could not be verified and corrected were excluded from further analyses.

Analyses were performed use an Area of Interest (AOI) approach. On learning trials, AOIs were hand-coded frame-by-frame for the speakers' faces and for the two on-screen objects. On test trials, these AOIs corresponded to the screen positions of the two alternatives. To use standard statistical analyses, we transformed the timecourse data to looking proportions within relevant windows. On test trials, the start of this window was set to 500ms after the point of disambiguation—the onset of the target label. The end of the window was set to the length of the shortest test trial for accurate averaging. This was 4s for Experiments 1 and 2, and 4.5s for Experiment 3. These windows were chosen to give children sufficient time to process the label, and to maximize signal (see Fig. 2).

In Experiments 1 and 2, Learning trials contained six distinct phases (Fig. 1): a baseline period (2s), name 1 to look ($M = 1.7s$), look to name 2 ($M = 2s$), name 2 to initiation of reach ($M = 4.8s$), reach to point of contact ($M = .8s$), and after contact with the toy ($M = 1s$). Proportion of looks to the target AOI were computed in each of these windows. In Experiment 3, learning trials were designed to separate demands on social attention across rather than within trials. In Experiment 3 we computed proportion of looking of the entirety of learning trials.

Due to occasional bouts of inattention, eye-gaze data were not available for all children during all portions of the experiment. To correct for statistical errors introduced by averaging over small windows, data from individual learning and test trials were

excluded from analysis if less than 50% of the window contained eye-tracking data (regardless of where children were looking). Second, if more than 50% of the trials for a given participant were excluded in this manner, all of the remaining trials were dropped as well.

Acknowledgments

We gratefully acknowledge the parents and children who participated in this research and the staff of the San Jose Children's Discovery Museum for their collaboration on this line of research. This work was supported by NIH NRSA F32HD075577 to DY, grants from the John Merck Scholars program, the Stanford Child Health Research Initiative to MCF, and a grant from the Mosbacher Family Fund for Autism Research at Stanford University.

Contributions

Contributed to conception and design: DY, AW, MF

Contributed to acquisition of data: AW, AMK, DY

Contributed to analysis and interpretation of data: DY, MF

Drafted and/or revised the article: DY, MF

Approved the submitted version for publication: DY, AW, AMK, GWG, AYH, MF

Competing Interests

No authors have competing interests with the publication of this manuscript.

Data Accessibility Statement

All data and code are available through the Open Science Framework (<https://osf.io/kjr98/>) and a public GitHub repository (<http://github.com/dyurovsky/refword>).

References

- Acheson, D. J., & MacDonald, M. C. (2009). Verbal working memory and language production: Common approaches to the serial ordering of verbal information. *Psychol Bull*, 135, 50–68.
- Adolph, K. E., Cole, W. G., Komati, M., Garciaguirre, J. S., Badaly, D., Lingeman, J. M., ... Sotsky, R. B. (2012). How Do You Learn to Walk? Thousands of Steps and Dozens of Falls per Day. *Psychol Sci*, 23, 1387–1394.
- Adolph, K. E., & Robinson, S. R. (2015). Motor Development. In *Handbook of child psychology and developmental science vol. cognitive processes* (pp. 114–157). Hoboken, NJ, USA: John Wiley & Sons, Inc.
- Anderson, D. R. (2005). Television and Very Young Children. *Am Behav Sci*, 48, 505–522.
- Baldwin, D. A. (1993). Early referential understanding: Infants' ability to recognize referential acts for what they are. *Dev Psychol*, 29, 832–843.
- Banks, M. S. (1980). The development of visual accommodation during early infancy. *Child Dev*, 646–666.
- Baron-Cohen, S., Baldwin, D., & Crowson, M. (1997). Do children with autism use eye-direction to infer linguistic reference. *Child Dev*, 68, 48–57.
- Bates, E., Marchman, V., Thal, D., Fenson, L., Dale, P., Reznick, J. S., ... Hartung, J. (1994). Developmental and stylistic variation in the composition of early vocabulary. *J Child Lang*, 21, 85–123.
- Bergelson, E., & Swingle, D. (2012). At 6–9 months, human infants know the meanings of many common nouns. *P Natl Acad Sci USA*, 109, 3253–3258.
- Bloom, P. (2000). *How children learn the meanings of words*. MIT press Cambridge, MA.
- Brand, R. J., Baldwin, D. A., & Ashburn, L. A. (2002). Evidence for 'motionese': modifications in mothers' infant-directed action. *Developmental Sci*, 5, 72–83.
- Brooks, R., & Meltzoff, A. N. (2005). The development of gaze following and its relation to language. *Developmental Sci*, 8, 535–543.

- Bulf, H., Johnson, S. P., & Valenza, E. (2011). Visual statistical learning in the newborn infant. *Cogn*, 121(1), 127–132.
- Carey, S., & Bartlett, E. (1978). Acquiring a single new word. *Pap R Child*, 15, 17–29.
- Corkum, V., & Moore, C. (1998). The origins of joint visual attention in infants. *Dev Psychol*, 34, 28.
- Csibra, G., & Gergely, G. (2009). Natural pedagogy. *Trends Cogn Sci*, 13, 148–153.
- Dempster, F. N. (1981). Memory span: Sources of individual and developmental differences. *Psychol Bull*, 89, 63.
- Fernald, A., Pinto, J. P., Swingle, D., Weinberg, A., & McRoberts, G. W. (1998). Rapid gains in speed of verbal processing by infants in the 2nd year. *Psychol Sci*, 9, 228–231.
- Fernald, A., Zangl, R., Portillo, A. L., & Marchman, V. A. (2008). Looking while listening: Using eye movements to monitor spoken language. *Developmental psycholinguistics: On-line methods in children's language processing*, 113–132.
- Forssberg, H., Eliasson, A., Kinoshita, H., Johansson, R., & Westling, G. (1991). Development of human precision grip i: basic coordination of force. *Exp Brain Res*, 85, 451–457.
- Franchak, J. M., Kretch, K. S., & Adolph, K. E. (2018). See and be seen: Infant–caregiver social looking during locomotor free play. *Developmental science*, 21(4), e12626.
- Frank, M. C., Goodman, N., & Tenenbaum, J. (2009). Using speakers' referential intentions to model early cross-situational word learning. *Psychol Sci*, 20, 578–585.
- Frank, M. C., Vul, E., & Saxe, R. (2012). Measuring the development of social attention using free-viewing. *Infancy*, 17, 355–375.
- Gathercole, S. E., Pickering, S. J., Ambridge, B., & Wearing, H. (2004). The structure of working memory from 4 to 15 years of age. *Developmental psychology*, 40, 177–190.
- Gogate, L. J., Bahrick, L. E., & Watson, J. D. (2000). A study of multimodal motherese: The role of temporal synchrony between verbal labels and gestures. *Child Dev*, 71,

878–894.

- Goulden, R., Nation, P., & Read, J. (1990). How large can a receptive vocabulary be? *App Linguist*, 11, 341–363.
- Hollich, G. J., Hirsh-Pasek, K., & Golinkoff, R. M. (2000). Breaking the Language Barrier: An Emergentist Coalition Model for the Origins of Word Learning. *Monogr Soc Res Child*.
- Kail, R. (1991). Developmental change in speed of processing during childhood and adolescence. *Psychol Bull*, 109, 490.
- Kretch, K. S., Franchak, J. M., & Adolph, K. E. (2014). Crawling and walking infants see the world differently. *Child development*, 85(4), 1503–1518.
- Krogh-Jespersen, S., Liberman, Z., & Woodward, A. L. (2015). Think fast! The relationship between goal prediction speed and social competence in infants. *Developmental Sci*, 18, 815–823.
- Le Couter, A., Rutter, M., Lord, C., Rios, P., Robertson, H., Holdgrafer, M., & McLennan, J. (1989). Autism diagnostic interview: A standardized investigator-based instrument. *J Autism Dev Disord*, 19, 363–387.
- Leekam, S. R., Hunnisett, E., & Moore, C. (1998). Targets and cues: Gaze-following in children with autism. *J Child Psychol Psyc*, 39, 951–962.
- Lidz, J., Waxman, S., & Freedman, J. (2003). What infants know about syntax but couldn't have learned: experimental evidence for syntactic structure at 18 months. *Cogn*, 89(3), 295–303.
- Loomis, J. M., Kelly, J. W., Pusch, M., Bailenson, J. N., & Beall, A. C. (2008). Psychophysics of perceiving eye-gaze and head direction with peripheral vision: Implications for the dynamics of eye-gaze behavior. *Perception*, 37, 1443–1457.
- Lord, C., Risi, S., Lambrecht, L., & Cook Jr, E. H. (2000). The Autism Diagnostic Observation Schedule—Generic: A standard measure of social and communication deficits associated with the spectrum of autism. *J Autism Dev Disord*, 205–223.

- Lord, C., Rutter, M., DiLavore, P. C., & Risis, S. (1999). *Autism diagnostic observation schedule: Manual*. Los Angeles, CA: Western Psychological Services.
- Lord, C., Rutter, M., & Le Couteur, A. (1994). Autism Diagnostic Interview-Revised: A revised version of a diagnostic interview for caregivers of individuals with possible pervasive developmental disorders. *J Autism Dev Disord*, *24*, 659–685.
- Moore, C., Angelopoulos, M., & Bennett, P. (1999). Word learning in the context of referential and salience cues. *Developmental psychology*, *35*(1), 60.
- O'Doherty, K., Troseth, G. L., Shimpi, P. M., Goldenberg, E., Akhtar, N., & Saylor, M. M. (2011). Third-party social interaction and word learning from video. *Child Dev*, *82*, 902–915.
- Rueda, M. R., & Rothbart, M. K. (2005). Training, maturation, and genetic influences on the development of executive attention. *P Natl Acad Sci USA*, *102*, 14931–14936.
- Scaife, M., & Bruner, J. S. (1975). The capacity for joint visual attention in the infant. *Nature*.
- Senju, A., Csibra, G., & Johnson, M. H. (2008). Understanding the referential nature of looking: infants' preference for object-directed gaze. *Cognition*, *108*, 303–319.
- Shukla, M., White, K. S., & Aslin, R. N. (2011). Prosody guides the rapid mapping of auditory word forms onto visual objects in 6-mo-old infants. *P Natl Acad Sci USA*, *108*(15), 6038–6043.
- Smith, L. B. (2013). It's all connected: Pathways in visual object recognition and early noun learning. *Am Psychol*, *68*, 618–629.
- Smith, L. B., & Yu, C. (2008). Infants rapidly learn word-referent mappings via cross-situational statistics. *Cogn*, *106*, 1558–1568.
- Snow, C. E. (1972). Mothers' speech to children learning language. *Child Dev*, *43*, 549–565.
- Sokol, S. (1978). Measurement of infant visual acuity from pattern reversal evoked potentials. *Vision Res*, *18*, 33–39.

- Thelen, E. (1995). Motor development: A new synthesis. *Am Psychol*, 50, 79.
- Thiessen, E. D., & Saffran, J. R. (2003). When cues collide: Use of stress and statistical cues to word boundaries by 7- to 9-month-old infants. *Dev Psychol*, 39, 706–716.
- Vida, M. D., & Maurer, D. (2012). The development of fine-grained sensitivity to eye contact after 6years of age. *J Exp Child Psychol*, 112, 243–256.
- Vouloumanos, A., Martin, A., & Onishi, K. H. (2014). Do 6-month-olds understand that speech can communicate? *Developmental Sci*, 17, 872–879.
- Yu, C., & Smith, L. B. (2007). Rapid word learning under uncertainty via cross-situational statistics. *Psychol Sci*, 18, 414–420.
- Yu, C., & Smith, L. B. (2013). Joint attention without gaze following: Human infants and their parents coordinate visual attention to objects through eye-hand coordination. *PLOS ONE*, 8, e79659.
- Yurovsky, D., & Frank, M. C. (2017). Beyond naïve cue combination: Salience and social cues in early word learning. *Developmental science*, 20(2), e12349.