

A communicative framework for early word learning

Benjamin C. Morris<sup>1</sup> & Daniel Yurovsky<sup>1,2</sup>

<sup>1</sup> University of Chicago

<sup>2</sup> Carnegie Mellon University

## Abstract

Children do not learn language from passive observation of the world, but from interaction with caregivers who want to communicate with them. These communicative exchanges are structured at multiple levels in ways that support language learning. We argue this pedagogically supportive structure can result from pressure to communicate successfully with a linguistically immature partner. We first characterize one kind of pedagogically supportive structure in a corpus analysis: caregivers provide more information-rich referential communication, using both gesture and speech to refer to a single object, when that object is rare and when their child is young. In an iterated reference game experiment on Mechanical Turk ( $n = 480$ ), we show how this behavior can arise from pressure to communicate successfully with a less knowledgeable partner. Then, we show that speaker behavior in our experiment can be explained by a rational planning model, without any explicit teaching goal. Lastly, in a series of simulations, we explore the language learning consequences of having a communicatively-motivated caregiver. In sum, this perspective offers the first steps toward a unifying, formal account of both the child's learning and the parents' production: Both are driven by a pressure to communicate successfully.

*Keywords:* language learning; communication; computational modeling; child-directed speech

Word count: X

## A communicative framework for early word learning

One of the most striking aspects of children's language learning is just how quickly they master the complex system of their natural language (Bloom, 2000). In just a few short years, children go from complete ignorance to conversational fluency in a way that is the envy of second-language learners attempting the same feat later in life (Newport, 1990). What accounts for this remarkable transition?

Distributional learning presents a unifying account of early language learning: Infants come to language acquisition with a powerful ability to learn the latent structure of language from the statistical properties of speech in their ambient environment (Saffran, 2003). Distributional learning mechanisms can be seen in accounts across language including phonemic discrimination (Maye, Werker, & Gerken, 2002), word segmentation (Saffran, 2003), learning the meanings of both nouns [Smith and Yu (2008) and verbs (Scott & Fischer, 2012), learning the meanings of words at multiple semantic levels (Xu & Tenenbaum, 2007), and perhaps even the grammatical categories to which a word belongs (Mintz, 2003). A number of experiments clearly demonstrate both the early availability of distributional learning mechanisms and their potential utility across these diverse language phenomena (DeCasper & Fifer, 1980; DeCasper & Spence, 1986; Estes, Evans, Alibali, & Saffran, 2007; Gomez & Gerken, 1999; Maye et al., 2002; Saffran, Aslin, & Newport, 1996; Smith & Yu, 2008; Xu & Tenenbaum, 2007).

However, there is reason to be suspicious about just how precocious statistical learning abilities are in early development. Although these abilities are available early, they are highly constrained by limits on other developing cognitive capacities. For example, infants' ability to track the co-occurrence information connecting words to their referents is constrained significantly by their developing memory and attention systems (Smith & Yu, 2013; Vlach & Johnson, 2013).

Computational models of these processes show that the rate of acquisition is highly sensitive to variation in environmental statistics (e.g., Vogt, 2012). Models of cross-situational learning have demonstrated that the Zipfian distribution of word frequencies and word meanings yields a learning problem that cross-situational learning alone cannot explain over a reasonable time frame (Vogt, 2012). Further, a great deal of empirical work demonstrates that cross-situational learning even in adults drops off rapidly when participants are asked to track more referents, and also when the number of intervening trials is increased (e.g., (Yurovsky & Frank, 2015)). Thus, precocious unsupervised statistical learning appears to fall short of a complete explanation for rapid early language learning.

Even relatively constrained statistical learning could be rescued, however, if caregivers structured their language in a way that simplified the learning problem and promoted learning. For example, in phoneme learning, infant-directed speech provides examples that seem to facilitate the acquisition of phonemic categories (Eaves Jr, Feldman, Griffiths, & Shafto, 2016). In word segmentation tasks, infant-directed speech facilitates infant learning more than matched adult-directed speech (Thiessen, Hill, & Saffran, 2005). In word learning scenarios, caregivers produce more speech during episodes of joint attention with young infants, which uniquely predicts later vocabulary (Tomasello & Farrar, 1986). Child-directed speech even seems to support learning at multiple levels in parallel— e.g., simultaneous speech segmentation and word learning (Yurovsky, 2012). For each of these language problems faced by the developing learner, caregiver speech exhibits structure that seems uniquely beneficial for learning.

Under distributional learning accounts, the existence of this kind of structure is a theory-external feature of the world that does not have an independently motivated explanation. Such accounts view the generative process of structure in the language environment as a problem separate from language learning. However, across a number of language phenomena, the language environment is not merely supportive, but seems

75 calibrated to children’s changing learning mechanisms (Daniel Yurovsky, 2018). For example,  
76 across development, caregivers engage in more multimodal naming of novel objects than  
77 familiar objects, and rely on this synchrony most with young children (Gogate, Bahrick, &  
78 Watson, 2000). The role of synchrony in child-directed speech parallels infant learning  
79 mechanisms: young infants appear to rely more on synchrony as a cue for word learning than  
80 older infants, and language input mirrors this developmental shift (Gogate et al., 2000).  
81 Beyond age-related changes, caregiver speech may also support learning through more local  
82 calibration to a child’s knowledge; caregivers have been shown to provide more language to  
83 refer to referents that are unknown to their child, and show sensitivity to the knowledge  
84 their child displays during a referential communication game (Leung, Tunkel, & Yurovsky,  
85 2019). The calibration of parents production to the child’s learning suggests a co-evolution  
86 such that these processes should not be considered in isolation.

87       What then gives rise to structure in early language input that mirrors child learning  
88 mechanisms? Because of widespread agreement that parental speech is not usually motivated  
89 by explicit pedagogical goals (Newport, Gleitman, & Gleitman, 1977), the calibration of  
90 speech to learning mechanisms seems a happy accident; parental speech just happens to be  
91 calibrated to children’s learning needs. Indeed, if parental speech was  
92 pedagogically-motivated, we would have a framework for deriving predictions and  
93 expectations (e.g., Shafto, Goodman, & Griffiths, 2014). Models of optimal teaching have  
94 been successfully generalized to phenomena as broad as phoneme discrimination (Eaves Jr et  
95 al., 2016) to active learning (Yang, Vong, Yu, & Shafto, 2019). These models take the goal  
96 to be to teach some concept to a learner and attempt to optimize that learner’s outcomes.  
97 While these optimal pedagogy accounts have proven impressively useful, such models are  
98 theoretically unsuited to explaining parent language production where there is widespread  
99 agreement that caregiver goals are not pedagogical (e.g., Newport et al., 1977).

100       Instead, the recent outpouring of work exploring optimal communication (the Rational

Speech Act model, see Frank and Goodman (2012)) provides another framework for understanding parent production. Under optimal communication accounts, speakers and listeners engage in recursive reasoning to produce and interpret speech cues by making inferences over one another's intentions (Frank & Goodman, 2012). These accounts have made room for advances in our understanding of a range of language phenomena previously uncaptured by formal modeling, notably a range of pragmatic inferences [e.g., Frank and Goodman (2012); other RSA papers]. In this work, we consider the communicative structure that emerges from an optimal communication system across a series of interactions where one partner has immature linguistic knowledge. This perspective offers the first steps toward a unifying account of both the child's learning and the parents' production: Both are driven by a pressure to communicate successfully (Brown, 1977).

Early, influential functionalist accounts of language learning focused on the importance of communicative goals (e.g., Brown, 1977). Our goal in this work is to formalize the intuitions in these accounts in a computational model, and to test this model against experimental data. We take as the caregiver's goal the desire to communicate with the child, not about language itself, but instead about the world in front of them. To succeed, the caregiver must produce the kinds of communicative signals that the child can understand and respond contingently, potentially leading caregivers to tune the complexity of their speech as a byproduct of in-the-moment pressure to communicate successfully (Daniel Yurovsky, 2018).

To examine this hypothesis, we focus on ostensive labeling (i.e. using both gesture and speech in the same referential expression) as a case-study phenomenon of information-rich structure in the language learning environment. We first analyze naturalistic parent communicative behavior in a longitudinal corpus of parent-child interaction in the home (Goldin-Meadow et al., 2014). We investigate the extent to which parents tune their ostensive labeling across their child's development to align to their child's developing linguistic knowledge (Yurovsky, Doyle, & Frank, 2016).

We then experimentally induce this form of structured language input in a simple model system: an iterated reference game in which two players earn points for communicating successfully with each other. Modeled after our corpus data, participants are asked to make choices about which communicative strategy to use (akin to modality choice). In an experiment on Mechanical Turk using this model system, we show that pedagogically-supportive input can arise from a pressure to communicate. We then show that participants' behavior in our game conforms to a model of communication as rational planning: People seek to maximize their communicative success while minimizing their communicative cost over expected future interactions. Lastly, we demonstrate potential benefits for the learner through a series of simulations to show that communicative pressure facilitates learning compared with various distributional learning accounts.

### Corpus Analysis

We first investigate parent referential communication in a longitudinal corpus of parent-child interaction. We analyze the production of multi-modal cues (i.e. using both gesture and speech) to refer to the same object, in the same instance. While many aspects of child-directed speech support learning, multi-modal cues (e.g., speaking while pointing or looking) are particularly powerful sources of data for young children (e.g., Baldwin, 2000; Gogate et al., 2000). We take multi-modal cues to be a case-study phenomenon of pedagogically supportive language input. While our account should hold for other language phenomena, by focusing on one phenomenon we attempt to specify the dynamics involved in the production of such input.

In this analysis of naturalistic communication, we examine the prevalence of multi-modal cues in children's language environment, to demonstrate that it is a viable, pedagogically supportive form of input. Beyond being a prevalent form of communication, multi-modal reference may be especially pedagogically supportive if usage patterns reflect adaptive linguistic tuning, with caregivers using this information-rich cue more for young

children and infrequent objects. The amount of multi-modal reference should be sensitive to the child’s age, such that caregivers will be more likely to provide richer communicative information when their child is younger (and has less linguistic knowledge) than as she gets older (Yurovsky et al., 2016).

## Methods

We used data from the Language Development Project– a large-scale, longitudinal corpus of naturalistic parent child-interaction in the home (Goldin-Meadow et al., 2014). The Language Development Project corpus contains transcription of all speech and communicative gestures produced by children and their caregivers over the course of the 90-minute home recordings. An independent coder analyzed each of these communicative instances and identified each time a concrete noun was referenced using speech, gesture, or both in the same referential expression (so called ostensive labeling). In these analyses, we focus only caregiver’s productions of ostensive labeling.

**Participants.** The Language Development Project aimed to recruit a sample of families who are representative of the Chicago community in socio-economic and racial diversity (Goldin-Meadow et al., 2014). These data are drawn from a subsample of 10 families from the larger corpus. Our subsample contains data taken in the home every 4-months from when the child was 14-months-old until they were 34-months-old, resulting in 6 timepoints (missing one family at the 30-month timepoint). Recordings were 90 minute sessions, and participants were given no instructions.

Of the 10 target children, 5 were girls, 3 were Black and 2 were Mixed-Race. Families spanned a broad range of incomes, with 2 families earning \$15,000 to \$34,999 and 1 family earning greater than \$100,000. The median family income was \$50,000 to \$74,999.

**Procedure.** From the extant transcription and gesture coding, we specifically coded all concrete noun referents produced in either the spoken or gestural modality (or both). Spoken reference was coded only when a specific noun form was used (e.g., “ball”), to



exclude pronouns and anaphoric usages (e.g., “it”). Gesture reference was coded only for deictic gestures (e.g., pointing to or holding an object) to minimize ambiguity in determining the intended referent. In order to fairly compare rates of communication across modalities, we need to examine concepts that can be referred to in either gesture or speech (or both) with similar ease. Because abstract entities are difficult to gesture about using deictic gestures, we coded only on references to concrete nouns.

**Reliability.** To establish the reliability of the referent coding, 25% of the transcripts were double-coded. Inter-rater reliability was sufficiently high (Cohen’s  $\kappa = 0.76$ ). Disagreements in coding decisions were discussed and resolved by hand.

To ensure that our each referent could potentially be referred to in gesture or speech, we focused on concrete nouns. We further wanted to ensure that the referents were physically present in the scene (and thus accessible to deictic gestures). Using the transcripts, a human rater judged whether the referent was likely to be present, primarily relying on discourse context (e.g., a referent was coded as present if the deictic gesture is used or used at another timepoint for the reference, or if the utterance included demonstratives such as “This is an X”). A full description of the coding criteria can be found in the Supporting Materials. **MAKE SURE WE MAKE THIS.**

To ensure our transcript-based coding of presentness was sufficiently accurate, a subset of the transcripts (5%) were directly compared to corresponding video data observation. Reliability across the video data and the transcript coding was sufficiently high ( $\kappa = 0.72$ ). Based on transcript coding of all the referential communication about concrete nouns, 90% of the references were judged to be about referents that were likely present. All references are included in our dataset for further analysis.

## Results

These corpus data were analyzed using a mixed effects regression to predict parent use of multi-modal reference for a given referent. The model included fixed effects of age in months, frequency of the referent, and the interaction between the two. The model included a random intercept and random slope of frequency by subject and a random intercept for each unique referent. Frequency and age were both log-scaled and then centered both because age and frequency tend to have log-linear effects and to help with model convergence. The model showed that parents teach less to older children ( $\beta = -0.78$ ,  $t = -7.88$ ,  $p < .001$ ), marginally less for more frequent targets ( $\beta = -0.08$ ,  $t = -1.81$ ,  $p = .071$ ), and that parents teach their younger children more often for equally frequent referents ( $\beta = 0.18$ ,  $t = 3.25$ ,  $p = .001$ ). Thus, in these data, we see early evidence that parents are providing richer, structured input about rarer things in the world for their younger children (Figure \ref{fig:corpus-plot}).

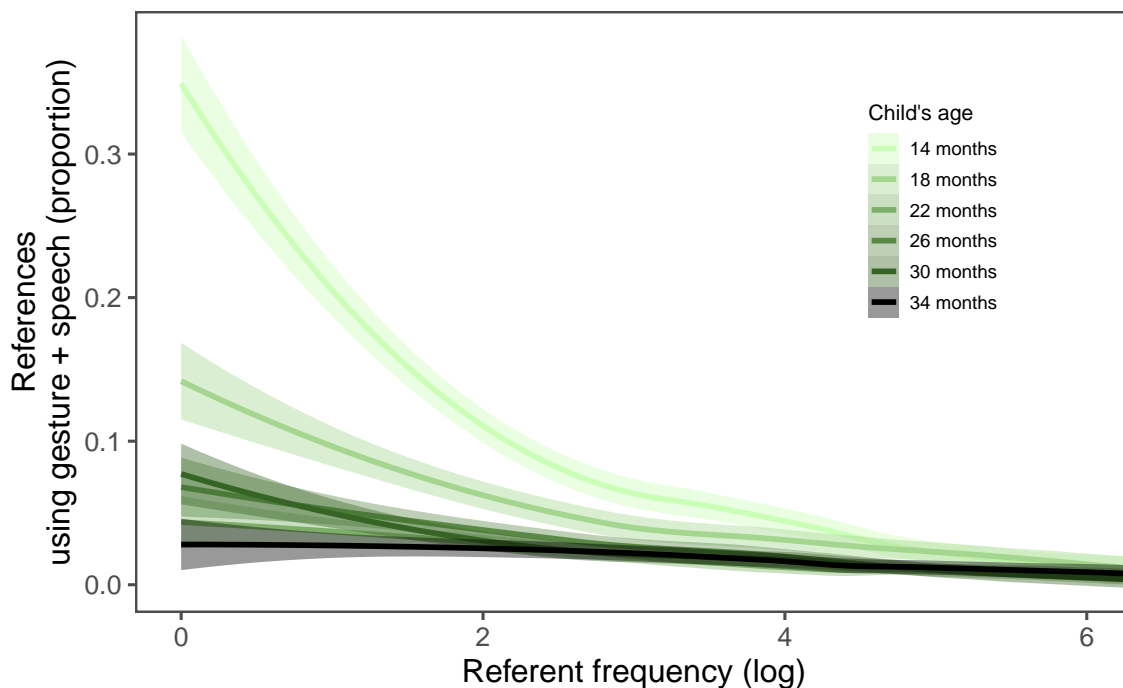


Figure 1. Proportion of parent multi-modal referential talk across development. The log of a referent's frequency is given on the x-axis, with less frequent items closer to zero.

## Discussion

Caregivers are not indiscriminate in their use of multi-modal reference; in these data, they provided more of this support when their child was younger and when discussing less familiar objects. These longitudinal corpus findings are consistent with an account of parental alignment: parents are sensitive to their child's linguistic knowledge and adjust their communication accordingly (Yurovsky et al., 2016). Ostensive labeling is perhaps the most explicit form of pedagogical support, so we chose to focus on it for our first case study. We argue that these data could be explained by a simple, potentially-selfish pressure: to communicate successfully. The influence of communicative pressure is difficult to draw in naturalistic data, so we developed a paradigm to try to experimentally induce richly-structured, aligned input from a pressure to communicate in the moment.

## Experimental Framework

To study the emergence of pedagogically supportive input from communicative pressure, we developed a simple reference game in which participants would be motivated to communicate successfully. After giving people varying amounts of training on novel names for 9 novel objects, we asked them to play a communicative game in which they were given one of the objects as their referential goal, and they were rewarded if their partner successfully selected this referent from among the set of competitors (Figure ??).

Participants could choose to refer either using the novel labels they had been exposed to, or they could use a deictic gesture to indicate the referent to their partner. The gesture was unambiguous, and thus would always succeed. However, in order for language to be effective, the participant and their partner would have to know the correct novel label for the referent.

Across conditions, we manipulated the relative costs of these two communicative methods (gesture and speech), as we did not have a direct way of assessing these costs in our

naturalistic data, and they likely vary across communicative contexts. In all cases, we assumed that gesture was more costly than speech. Though this need not be the case for all gestures and contexts, our framework compares simple lexical labeling and unambiguous deictic gestures, which likely are more costly and slower to produce (see Yurovsky et al., 2018). We set the relative costs by explicitly implementing strategy utility, assigning point values to each communicative method.

If people are motivated to communicate successfully, their choice of referential modality should reflect the tradeoff between the cost of producing the communicative signal with the likelihood that the communication would succeed. We thus predicted that peoples' choice of referential modality would reflect this tradeoff: People should be more likely to use language if they have had more exposures to the novel object's correct label, and they should be more likely to use language as gesture becomes relatively more costly.

Critically, participants were told that they will play this game repeatedly with their partner. In these repeated interactions, participants are then able to learn about an interlocutor and potentially influence their learning. Thus, there is a third type of message: using both gesture and speech within a single trial to effectively teach the listener an object-label mapping. This strategy necessitates making inferences about the listener's knowledge state, so we induced knowledge asymmetries between speaker and listener. To do so, we manipulated how much training they thought their partner had received. Our communicative game was designed to reward in-the-moment communication, and thus teaching required the speaker pay a high cost upfront. However, rational communicators may understand that if one is accounting for future trials, paying the cost upfront to teach the listener allows a speaker to use a less costly message strategy on subsequent trials (namely, speech). Manipulating the listener knowledge and the utility of communicative strategies, we aimed to experimentally determine the circumstances under which richly-structured input emerges, without an explicit pedagogical goal.

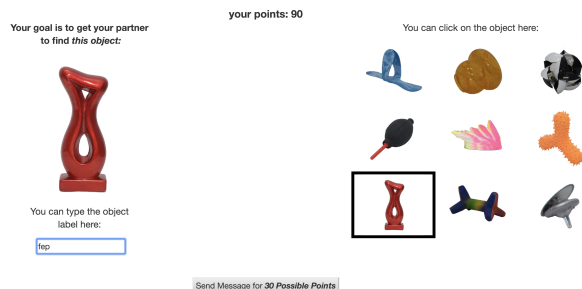


Figure 2. (#fig:exp\_screenshot)Screenshot of speaker view during gameplay.

## Method

In this experiment, participants were recruited to play our reference game via Amazon Mechanical Turk, an online platform that allows workers to complete surveys and short tasks for payment. In this study, all participants were placed in the role of speaker and listener responses were programmed.

**Participants.** 480 participants were recruited through Amazon Mechanical Turk and received \$1 for their participation. Data from 51 participants were excluded from subsequent analysis for failing the critical manipulation check and a further 28 for producing pseudo-English labels (e.g., “pricklyyone”). The analyses reported here exclude the data from those participants, but all analyses were also conducted without excluding any participants and all patterns hold ( $ps < 0.05$ ).

**Design and Procedure.** Participants were told they would be introduced to novel object-label pairs and then asked to play a communication game with a partner wherein they would have to refer to a particular target object. Participants were exposed to nine novel objects, each with a randomly assigned pseudo-word label. We manipulated the exposure rate within-subjects: during training participants saw three of the nine object-label mappings four times, two times, or just one time, yielding a total of 21 training trials. Participants were then given a simple recall task to establish their knowledge of the novel lexicon (pretest).

During gameplay, speakers saw the target object in addition to an array of all six

objects. Speakers had the option of either directly selecting the target object from the array (deictic gesture)—a higher cost cue but without ambiguity—or typing a label for the object (speech)—a lower cost cue but contingent on the listener’s knowledge. After sending the message, speakers are shown which object the listener selected.

We also manipulated participants’ expectations about their partner’s knowledge to explore the role of knowledge asymmetries. Prior to beginning the game, participants were told how much exposure their partner had to the lexicon. Across 3 between subjects conditions, participants were told that their partner had either no experience with the lexicon, had the same experience as the speaker, or had twice the experience of the speaker. As a manipulation check, participants were then asked to report their partner’s level of exposure, and were corrected if they answer incorrectly. Participants were then told that they would be asked to discuss each object three times during the game.

Listeners were programmed with starting knowledge states initialized according to the partner knowledge condition. Listeners with no exposure began the game with knowledge of 0 object-label pairs. Listeners with the same exposure of the speaker began with knowledge of five object-label pairs (3 high frequency, 1 mid frequency, 1 low frequency), based average retention rates found previously. Lastly, the listener with twice as much exposure as the speaker began with knowledge of all nine object-label pairs.

To simulate knowledgeable listener behavior when the speaker typed an object label, the listener was programmed to consult their own knowledge. Messages were evaluated by taking the Levenshtein distance (LD) between the typed label and each possible label in the listener’s vocabulary. Listeners then selected the candidate with the smallest edit distance (e.g., if a speaker entered the message “tomi”, the programmed listener would select the referent corresponding to “toma”, provided toma was found in its vocabulary). If the speaker message had an LD greater than two with each of the words in the listener’s vocabulary, the listener selected an unknown object. If the speaker clicked on object (gesture message), the

listener was programmed to simply make the same selection.

Speakers could win up to 100 points per trial if the listener correctly selected the target referent based on their message. If the listener failed to identify the target object, the speaker received no points. We manipulated the relative utility of the speech cue between-subjects across two conditions: low relative cost (“Low Relative Cost”) and higher relative cost (“Higher Relative Cost”). In the “Low Relative Cost” condition, speakers received 30 points for gesturing and 100 points for labeling, and thus speech had very little cost relative to gesture and participants should be highly incentivized to speak. In the “Higher Relative Cost” condition speakers received 50 points for gesturing and 80 points for labeling, and thus gesturing is still costly relative to speech but much less so and participants should be less incentivized to speak.

Participants were told about a third type of possible message using both gesture and speech within a single trial to effectively teach the listener an object-label mapping. This action directly mirrors the multi-modal reference behavior from our corpus data– it presents the listener with an information-rich, potentially pedagogical learning moment. In order to produce this teaching behavior, speakers had to pay the cost of producing both cues (i.e. both gesture and speech). Note that, in all utility conditions, teaching yielded participants 30 points (compared with the much more beneficial strategy of speaking which yielded 100 points or 80 points across our two utility manipulations). Listeners were programmed to integrate new taught words into their knowledge of the lexicon, and check those taught labels on subsequent trials when evaluating speaker messages.

Crossing our 2 between-subjects manipulations yielded 6 conditions (2 utility manipulations: “Low Relative Cost” and “Higher Relative Cost”; and 3 levels of partner’s exposure: None, Same, Double), with 80 participants in each condition. We expected to find results that mirrored our corpus findings such that rates of teaching would be higher when there was an asymmetry in knowledge where the speaker knew more (None manipulation)

compared with when there was equal knowledge (Same manipulation) or when the listener was more familiar with the language (Double manipulation). We expected that participants would also be sensitive to our utility manipulation, such that rates of labeling and teaching would be higher in the “Low Relative Cost” conditions than the other conditions.

## Results

In each trial, participants are able to choose one of 3 communicative strategies: gesture, speech, or teaching. We primarily expect flexible trade-off between the use of each strategy given their relative utilities, participant’s knowledge of the lexicon, and the listener’s knowledge of the lexicon. To test our predictions about each communicative behavior (gesture, speech, and teaching), we conducted separate logistic mixed effects models for each behavior, reported below. It should be noted that these three behaviors are mutually exhaustive. First, we establish how well participants learned our novel lexicon during training.

**Learning.** As an initial check of our exposure manipulation, we first conducted a logistic regression predicting accuracy at test from a fixed effect of exposure rate and random intercepts and slopes of exposure rate by participant as well as random intercepts by item. We found a reliable effect of exposure rate, indicating that participants were better able to learn items that appear more frequently in training ( $\beta = 1.08$ ,  $p < .001$ ). On average, participants knew at least 6 of the 9 words in the lexicon ( $M(sd) = 6.28 (2.26)$ ). An analysis of variance confirmed that learning did not differ systematically across participants by partner’s exposure, utility manipulation, or their interaction ( $ps > 0.05$ ).

**Gesture.** When should we expect participants to rely on gesture? Gesturing has the highest utility for words you failed to learn during training, words you think your partner is unlikely to know (i.e., for lower partner knowledge conditions), and when utility scheme is relatively biased toward gesturing (i.e., the “Higher Relative Cost” condition). To test these predictions, we ran a mixed effects logistic regression to predict whether speakers chose to



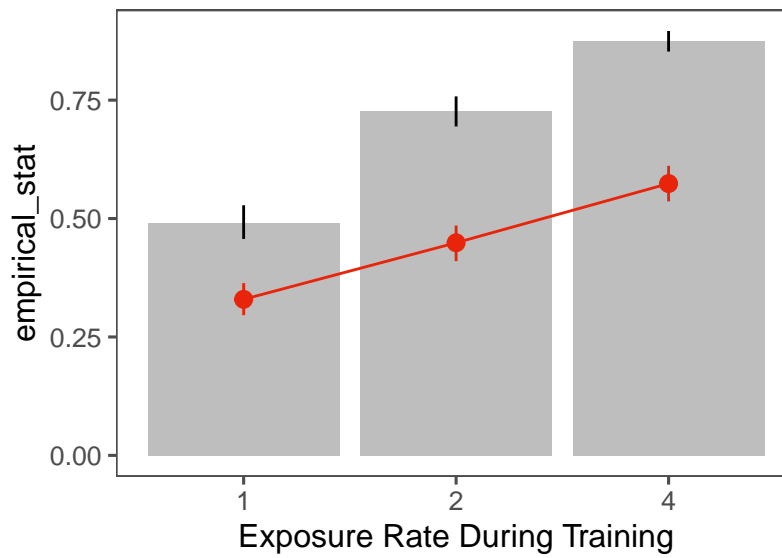


Figure 3. Participants' performance on the baseline recall task for the lexicon, as function of amount of exposure during training (grey bars). The red line shows the proportion of trials in the game in which participants used the learned labels.

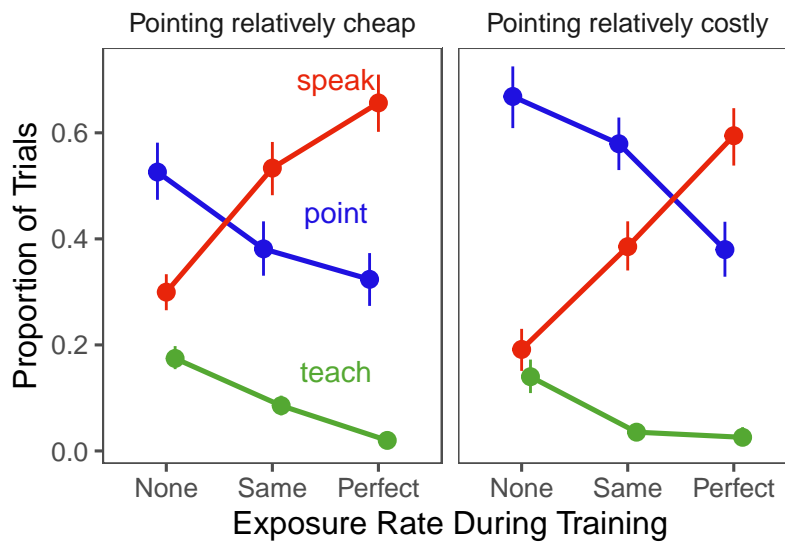


Figure 4. Speaker communicative method choice as a function of exposure and the utility manipulation.

gesture during a given trial as a function of the target object's exposure rate during training, object instance in the game (first, second, or third), utility manipulation, and partner manipulation. Random effects terms for subject and object were included in the model.

Consistent with our predictions, exposure rate during training was a significant negative predictor of gesturing during the game (see Figure 4), such that participants were less likely to rely on gesture for well trained (and thus well learned) objects ( $\beta = -0.50$ ,  $p < .001$ ). Additionally, participants were significantly more likely to gesture in the Higher Relative Cost condition where gesture is relatively less costly, compared with the Low Relative Cost condition ( $\beta = 1.20$ ,  $p < .001$ ) (see Figure 4). We also found a significant negative effect of partner's knowledge, such that participants used gesture more for partners with less knowledge of the lexicon ( $\beta = -0.81$ ,  $p < .001$ ).

Note that these effects cannot be explained by solely speaker knowledge; all patterns above hold when looking *only* at words known by the speaker at pretest ( $ps < 0.01$ ). Further, these patterns directly mirror previous corpus analyses demonstrating adult's use of gesture in naturalistic parental communicative behaviors, and parents likely have lexical knowledge of even even the least frequent referent (see Daniel Yurovsky, 2018).

**Speech.** When should we expect participants to use speech? Speech has the highest utility for words you learned during training, words you think your partner is likely to know (i.e., for higher partner knowledge conditions), when utility scheme is relatively biased toward speech (i.e., the “Low Relative Cost” condition). To test these predictions, we ran a mixed effects logistic regression to predict whether speakers chose to speak during a given trial as a function of the target object's exposure rate during training, object instance in the game (first, second, or third), utility manipulation, and partner manipulation. Random effects terms for subjects and object were included in the model.

Consistent with our predictions, speech seemed to largely tradeoff with gesture. Exposure rate during training was a significant positive predictor of speaking during the

game, such that participants were more likely to utilize speech for well trained (and thus well learned) objects ( $\beta = 0.35, p < .001$ ). Additionally, participants were significantly less likely to speak in the High Relative Cost condition where speech is relatively more costly, compared with the Low Relative Cost condition ( $\beta = -0.87, p .001$ ). We also found a significant positive effect of partner's knowledge, such that participants used speech more for partners with more knowledge of the lexicon ( $\beta = 1.95, p < .001$ ). Unlike for gesture, there is a significant effect of object instance in the game (i.e., whether this is the first, second, or third trial with this target object) on the rate of speaking, such that later trials are more likely to elicit speech ( $\beta = 0.72, p < .001$ ). This effect of order likely stems from a trade-off with the effects we see in teaching (described below); after a speaker teaches a word on the first or second trial, the utility of speech is much higher on subsequent trials.

**Emergence of Teaching.** Thus far, we have focused on relatively straightforward scenarios to demonstrate that a pressure to communicate successfully in the moment can lead speakers to trade-off between gesture and speech sensibly. Next, we turn to the emergence of teaching behavior.

When should we expect participants to teach? Teaching has the highest utility for words you learned during training, words you think your partner is unlikely to know (i.e., for lower partner knowledge conditions), when utility scheme is relatively biased toward speech (i.e., the “Low Relative Cost” condition). To test these predictions, we ran a mixed effects logistic regression to predict whether speakers chose to teach during a given trial as a function of the target object's exposure rate during training, object instance in the game (first, second, or third), utility manipulation, and partner manipulation. Random effects terms for subjects and object were included in the model.

Consistent with our predictions, rates of teaching were higher for better trained words, less knowledgeable partners, and when speech had the highest utility. Exposure rate during training was a significant positive predictor of teaching during the game, such that

participants were more likely to teach for well trained (and thus well learned) objects ( $\beta = 0.14, p .001$ ). While costly in the moment, teaching can be a beneficial strategy in our reference game because it subsequently allows for lower cost strategy (i.e. speaking), thus when speaking has a lower cost, participants should be more incentivized to teach. Indeed, participants were significantly less likely to teach in the High Relative Cost condition where speech is relatively more costly, compared with the Low Relative Cost condition ( $\beta = -0.96, p .001$ ). We also found a significant negative effect of partner's knowledge, such that participants taught more with partners that had less knowledge of the lexicon ( $\beta = -2.23, p < .001$ ). There was also a significant effect of object instance in the game (i.e., whether this is the first, second, or third trial with this target object) on the rate of teaching. The planned utility of teaching comes from using another, cheaper strategy (speech) on later trials, thus the expected utility of teaching should decrease when there are fewer subsequent trials for that object, predicting that teaching rates should drop dramatically across trials for a given object. Participants were significantly less likely to teach on the later appearances of the target object ( $\beta = -1.09, p < .001$ ).

## Discussion

As predicted, the data from our paradigm corroborate our findings from the corpus analysis, demonstrating that pedagogically supportive behavior emerges despite the initial cost when there is an asymmetry in knowledge and when speech is less costly than other modes of communication. While this paradigm has stripped away much of the interactive environment of the naturalistic corpus data, it provides important proof of concept that the structured and tuned language input we see in those data could arise from a pressure to communicate. The paradigm's clear, quantitative predictions also allow us to build a formal model to predict our empirical results.

The results from this experiment are qualitatively consistent with a model in which participants make their communicative choices to maximize their expected utility from the

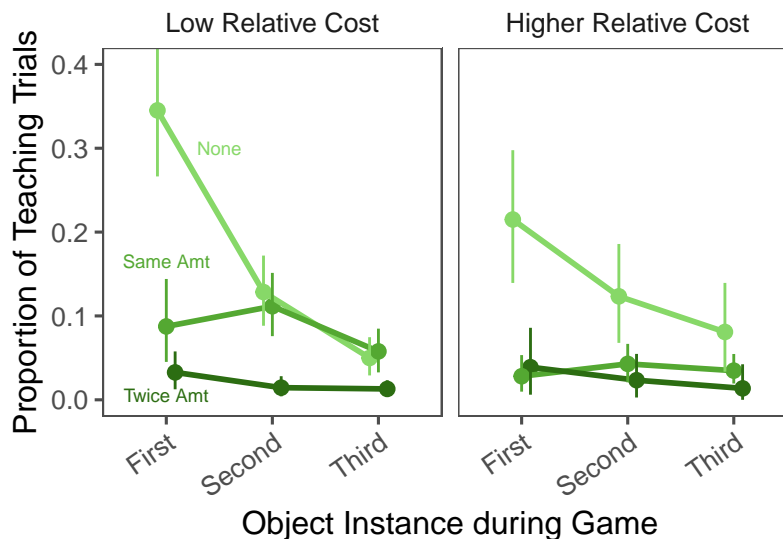


Figure 5. Rates of teaching across the 6 conditions, plotted by how many times an object had been the target object.

reference game. We next formalize this model to determine if these results are predicted quantitatively as well.

### Model: Communication as planning

In order to model when people should speak, point, or teach, we begin from the problem of what goal people are trying to solve (Marr, 1982). Following a long history of work in philosophy of language, we take the goal of communication to be causing an action in the world by transmitting some piece of information to one's conversational partner (e.g. Wittgenstein, 1953; Austin, 1975). If people are near-optimal communicators, they should choose communicative signals that maximize the probability of being understood while minimizing the cost of producing the signal (Clark, 1996; Grice, 1975). In the special case of reference, solving this problem amounts to producing the least costly signal that correctly specifies one's intended target referent in such a way that one's conversational partner can select it from the set of alternative referents.

Recently, Frank and Goodman (2012) developed the Rational Speech Act framework—a

formal instantiation of these ideas. In this model, speakers choose from a set of potential referential expressions in accordance to a utility function that maximizes the probability that a listener will correctly infer their intended meaning while minimizing the number of words produced. This framework has found successful application in a variety of linguistic applications such as scalar implicature, conventional pact formation, and production and interpretation of hyperbole (Goodman & Frank, 2016; see also related work from Franke, 2013). These models leverage recursive reasoning—speakers reasoning about listeners who are reasoning about speakers—in order to capture cases in which the literal meaning and the intended meaning of sentences diverge.

To date, this framework has been applied primarily in cases where both communicative partners share the same linguistic repertoire, and thus communicators know their probability of communicating successfully having chosen a particular signal. This is a reasonable assumption for pairs of adults in contexts with shared common ground. But what if partners do not share the same linguistic repertoire, and in fact do not know the places where their knowledge diverges? In this case, communicators must solve two problems jointly: (1) Figure out what their communicative partner knows, and (2) produce the best communicative signal they can given their estimates of their partner’s knowledge. If communicative partners interact repeatedly, these problems become deeply intertwined: Communicators can learn about each-other’s knowledge by observing whether their attempts to communicate succeed. For instance, if a communicator produces a word that they believe identifies their intended referent, but their partner fails to select that referent, the communicator can infer that their partner must not share their understanding of that word. They might then choose not to use language to refer to this object in the future, but choose to point to it instead.

Critically, communicators can also change each-other’s knowledge. When a communicator both points to an object and produces a linguistic label, they are in effect teaching their partner the word that they use to refer to this object. While this this behavior

is costly in the moment, and no more referentially effective than pointing alone, it can lead to more efficient communication in the future—instead of pointing to this referent forever more, communicators can now use the linguistic label they both know they share. This behavior naturally emerges from a conception of communication as planning: Communicators’ goal is to choose a communicative signal today that will lead to efficient communication not just in the present moment, but in future communications as well. If they are likely to need to refer to this object frequently, it is worth it to be inefficient in this one exchange in order to be more efficient future. In this way, pedagogically supportive behavior can emerge naturally from a model with no separate pedagogical goal. In the following section, we present a formal instantiation of this intuitive description of communication as planning and show that it accounts for the behavior we observed in our experiments.

Alternatively, pedagogically-supportive input could emerge from an explicit pedagogical goal. Shafto et al. (2014) have developed an framework of rational pedagogy built on the same recursive reasoning principles as in the Rational Speech Act Framework: Teachers aim to teach a concept by choosing a set of examples that would maximize learning for students who reason about the teachers choices as attempting to maximize their learning. Rafferty, Brunskill, Griffiths, and Shafto (2016) et al. expanded this framework to sequential teaching, in which teachers use students in order to infer what they have learned and choose the subsequent example. In this case, teaching can be seen as a kind of planning where teachers should choose a series of examples that will maximize students learning but can change plans if an example they thought would be too hard turns out too easy—or vice-versa. In the case of our reference game, this model is indistinguishable from a communicator who seeks to maximize communicative success but is indifferent to communicative cost. A cost-indifferent model makes poor predictions about parents’ behavior in our corpus, and also adults’ behavior in our experiments, but we return to it in the subsequent section to consider how differences in parents’ goals and differences in children’s learning contribute to changes in the rate of language acquisition.

## Formal Model

We take as inspiration the idea that communication is a kind of action—e.g. talking is a speech act (Austin, 1975). Consequently, we can understand the choice of *which communicative act* a speaker should take as a question of which act would maximize their utility: achieving successful communication while minimizing their cost (Frank & Goodman, 2012). In this game, speakers can take three actions: talking, pointing, or teaching. In this reference game, these Utilities ( $U$ ) are given directly by the rules. Because communication is a repeated game, people should take actions that maximize their Expected Utility ( $EU$ ) over not just for the current round, but for all future communicative acts with the same conversational partner. We can think of communication, then as a case of recursive planning. However, people do not have perfect knowledge of each-other’s vocabularies ( $v$ ). Instead, they only have uncertain beliefs ( $b$ ) about these vocabularies that combine their expectations about what kinds of words people with as much linguistic experience as their partner are likely to know with their observations of their partner’s behavior in past communicative interactions. This makes communication a kind of planning under uncertainty well modeled as a Partially Observable Markov Decision Process (POMDP, Kaelbling, Littman, & Cassandra, 1998).

Optimal planning in a Partially Observable Markov Decision Process involves a cycle of three phases: (1) Plan, (2) Act, and (3) Update beliefs. We describe those in turn and finally define how people form initial beliefs about their partner’s language now.

**Plan.** When people plan, they compute the expected utility of each possible action ( $a$ ) by combining the expected utility of that action now with the Discounted Expected Utility they will get in all future actions. The amount of discounting ( $\gamma$ ) reflects how much people care about success now compared to success in the future. Because utilities depend on the communicative partner’s vocabulary, people should integrate over all possible vocabularies in proportion to the probability that their belief assigns to that vocabulary



534  $(\mathbb{E}_{v \sim b})$ .

$$EU[a|b] = \mathbb{E}_{v \sim b} (U(a|v) + \gamma \mathbb{E}_{v', o', a'} (EU[a'|b']))$$

535 **Act.** Next, people take an action as a function of its expected utility. Following other  
 536 models in the Rational Speech Act framework, we use the Luce Choice Axiom, in which each  
 537 choice is taken in probability proportional to its exponentiated utility (Frank & Goodman,  
 538 2012; Luce, 1959). This choice rule has a single parameter  $\alpha$  that controls the noise in this  
 539 choice—as  $\alpha$  approaches 0, choice is random and as  $\alpha$  approaches infinity, choice is optimal.

$$P(a|b) \propto \alpha e^{EU[a|b]}$$

540 **Update beliefs.** After taking an action, people observe ( $o$ ) their partner's  
 541 choice—sometimes they correctly select the intended object, and sometimes they do not.  
 542 People then update their beliefs about the partner's vocabulary based on this observation.  
 543 For simplicity, we assume that people think their partner should always select the correct  
 544 target if they point to it, or if they teach, and similarly should always select the correct  
 545 target if they produce its label and the label is in their partner's vocabulary. Otherwise, they  
 546 assume that their partner will select the wrong object. People could of course have more  
 547 complex inferential rules, e.g. assuming that if their partner does know a word they will  
 548 choose among the set of objects whose labels they do not know (mutual exclusivity,  
 549 Markman & Wachtel, 1988). Empirically, however, our simple model appears to accord well  
 550 with people's behavior.

$$b'(v') \propto P(o|v', a) \sum_{v \in V} P(v'|v, a) b(v)$$

551 The critical feature of a repeated communication game is that people can change their  
 552 partner's vocabulary. In teaching, people pay the cost of both talking and pointing together,  
 553 but can leverage their partner's new knowledge on future trials. Note here that teaching has  
 554 an upfront cost and the only benefit to be gained comes from using less costly  
 555 communication modes later. There is no pedagogical goal—the model treats speakers as

selfish agents aiming to maximize their own utilities by communicating successfully. We assume for simplicity that teaching is always successful in this very short game, that communicative partners do not forget words once they have learned them, and that no learning happens by inference from mutual exclusivity.

$$P(v'|v, a) = \begin{cases} 1 & \text{if } v_w \in v \& v' \mid a = \text{point+talk} \\ 0 & \text{otherwise} \end{cases}$$

**Initial Beliefs.** The final detail is to specify how people estimate their partner’s learning rate ( $p$ ) and initial vocabulary ( $v$ ). We propose that people begin by estimating their own learning rate by reasoning about the words they learned at the start of the task: Their learning rate ( $p$ ) is the rate that maximizes the probability of them having learned their initial vocabularies from the trials they observed. People can then expect their partner to have a similar  $p$  (per the “like me” hypothesis, Meltzoff, 2005). Having an estimate of their partner’s  $p$ , they can estimate their vocabulary by simulating their learning from the amount of prior exposure to language their partner had before the game. In our experiments, we explicitly manipulated this expectation by telling participants how much exposure their partner had relative to their own exposure.

## Method

We implemented the planning model using the WebPPL— a programming language designed for specifying probabilistic models (Goodman & Stuhlmüller, 2014). To derive predictions from the model, we exposed it to the same trial-by-trial stimuli as the participants in our experiment, and used the probabilistic equations defined above to determine the likelihood of choosing each behavior (e.g. “speak”, “point”, or “teach”) on every trial. Separate predictions were made for each trial for each participant on the basis of all of the information available to each participant at that point in time (e.g. how many words they had learned, their partner’s observed behavior previously, etc).

The model's behavior is contingent on two parameters—discounting ( $\gamma$ ), and it's rationality ( $\alpha$ ). In order to determine the values of these parameters that best characterize human participants, we used Bayesian inference to estimate the posterior means of both. Using estimates rather than the maximum likelihood estimates naturally penalizes the models for their ability to predict patterns of data that were not observed, applying a kind of Bayesian Ockham's razor (MacKay, 1992). Because of we found substantial variability in the best parameter estimates across individual participants, we estimated parameters hierarchically, with group-level parameters forming the priors for individual participants' parameters. This hierarchical estimation process achieves the same partial pooling as subject-level random effects in mixed-effects models, giving estimates of the group-level parameters (Gelman & Hill, 2006). Details of the estimation procedure can be found in the Supplemental Materials.

## Model Results

In line with previous work on rational speech act models, and decision making, we expected rationality ( $\alpha$ ) to be around 1 or 2 (Frank & Goodman, 2012, p. @frank2014). We estimated the posterior mean rationality ( $\alpha$ ) to be 0.44 with 95% credible intervals of [1.10, 1.26]. We did not have strong expectations for the value of the discounting parameter ( $\gamma$ ), but estimated it to be 0.41 [0.41, 0.47], suggesting that on average participants weighed the next occurrence of a referent as slightly less than half as important as the current occurrence.

## Consequences for Learning

In the model and experiments above, we asked whether the pressure to communicate successfully with a linguistically-naïve partner would lead to pedagogically supportive input. These results confirmed its' sufficiency: As long as linguistic communication is less costly than deictic gesture, speakers should be motivated to teach in order to reduce future communicative costs. Further, the strength of this motivation is modulated by predictable factors (speaker's linguistic knowledge, listener's linguistic knowledge, relative cost of speech

and gesture, learning rate, etc.), and the strength of this modulation is well predicted by a rational model of planning under uncertainty about listener's vocabulary.

In this final section, we take up the consequences of communicatively-motivated teaching for the listener. To do this, we adapt a framework used by Blythe, Smith, and Smith (2010) and colleagues to estimate the learning times for an idealized child learning language under a variety of models of both the child and their parent. We come to these estimates by simulating exposure to successive communicative events, and measuring the probability that successful learning happens after each event. The question of how different models of the parent impact the learner can then be formalized as a question of how much more quickly learning happens in the context of one model than another.

We consider three parent models:

1. *Teacher* - under this model, we take the parents' goal to be maximizing the child's linguistic development. Each communicative event in this model consists of an ostensive labelling event (Note: this model is equivalent to a *Communicator* that ignores communicative cost).
2. *Communicator* - under this model, we take the parents' goal to be maximizing communicative success while minimizing communicative cost. This is the model we explored in the previous section.
3. *Indifferent* - under this model, the parent produces a linguistic label in each communicative event regardless of the child's vocabulary state. (Note: this model is equivalent to a *Communicator* who ignores communicative success).

## SOME STUFF ABOUT CROSS SITUATIONAL LEARNING

One important point to note is that we are modeling the learning of a single word rather than the entirety of a multi-word lexicon (as in Blythe et al., 2010). Although

learning times for each word could be independent, an important feature of many models of word learning is that they are not (Frank, Goodman, & Tenenbaum, 2009; Yu, 2008; Yurovsky, Fricker, Yu, & Smith, 2014; although c.f. McMurray, 2007). Indeed, positive synergies across words are predicted by the majority of models and the impact of these synergies can be quite large under some assumptions about the frequency with which different words are encountered (Reisenauer, Smith, & Blythe, 2013). We assume independence primarily for pragmatic reasons here—it makes the simulations significantly more tractable (although it is what our experimental participants appear to assume about learners). Nonetheless, it is an important issue for future consideration. Of course, synergies that support learning under a cross-situational scheme must also support learning from communicators and teachers (Frank et al., 2009; Markman & Wachtel, 1988; Yurovsky, Yu, & Smith, 2013). Thus, the ordering across conditions should remain unchanged. However, the magnitude of the difference across teacher conditions could potentially increase or decrease.

## Method

**Teaching.** Because the teaching model is indifferent to communicative cost, it engages in ostensive an ostensive labeling (pointing + speaking) on each communicative event. Consequently, learning on each trial occurs with a probability that depends entirely on the learner’s learning rate ( $P_k = p$ ). Because we do not allow forgetting, the probability that a learner has failed to successfully learn after  $n$  trials is equal to the probability that they have failed to learn on each of  $n$  successive independent trials (The probability of zero success on  $n$  trials of a Binomial random variable with parameter  $p$ ). The probability of learning after  $n$  trials is thus:

$$P_k(n) = 1 - (1 - p)^n$$

The expected probability of learning after  $n$  trials was thus defined analytically and

required no simulation. For comparison to the other models, we computed  $P_k$  for values of  $p$  that ranged from .1 to 1 in increments of .1.

**Communication.** To test learner under the communication model, we implemented the same model described in the paper above. However, because our interest was in understanding the relationship between parameter values and learning outcomes rather than inferring the parameters that best describe people’s behavior, we made a few simplifying assumptions to allow many runs of the model to complete in a more practical amount of time. First, in the full model above, speakers begin by inferring their own learning parameters ( $P_s$ ) from their observations of their own learning, and subsequently use their maximum likelihood estimate as a standin for their listener’s learning parameter ( $P_l$ ). Because this estimate will converge to the true value in expectation, we omit these steps and simply stipulate that the speaker correctly estimates the listener’s learning parameter.

Second, unless the speaker knows apriori how many times they will need to refer to a particular referent, the planning process is an infinite recursion. However, each future step in the plan is less impactful than the previous step (because of exponential discounting), this infinite process is in practice well approximated by a relatively small number of recursive steps. In our explorations we found that predictions made from models which planned over 3 future events were indistinguishable from models that planned over four or more, so we simulated 3 steps of recursion<sup>1</sup>. Finally, to increase the speed of the simulations we re-implemented them in the R programming language. All other aspects of the model were identical.

**Hypothesis Testing.** The literature on cross-situational learning is rich with a variety of models that could broadly be considered to be “hypothesis testers.” In an eliminative hypothesis testing model, the learner begins with all possible mappings between

---

<sup>1</sup> It is an intersting empirical question to determine how the level of depth to which that people plan in this and similar games (see e.g. bounded rationality in Simon, 1991; resource-rationality in Griffiths, Lieder, & Goodman, 2015). This future work is outside the scope of the current project.

words and objects and prunes potential mappings when they are inconsistent with the data according to some principle. A maximal version of this model relies on the principle that every time a word is heard its referent must be present, and thus prunes any word-object mappings that do not appear on the current trial. This model converges when only one hypothesis remains and is probably the fastest learner when its assumed principle is a correct assumption (Smith, Smith, & Blythe, 2011).

A positive hypothesis tester begins with no hypotheses, and on each trial stores one or more hypotheses that are consistent with the data, or alternatively strengthens one or more hypotheses that it has already stored that are consistent with the new data. A number of such models have appeared in the literature, with different assumptions about (1) how many hypotheses a learner can store, (2) existing hypotheses are strengthened, (3) how existing hypotheses are pruned, and (4) when the model converges (Siskind, 1996; Smith et al., 2011; Stevens, Gleitman, Trueswell, & Yang, 2017; Trueswell, Medina, Hafri, & Gleitman, 2013; Yu & Smith, 2012).

Finally, Bayesian models have been proposed that leverage some of the strengths of both of these different kinds of model, both increasing their confidence in hypotheses consistent with the data on a given learning event and decreasing their confidence in hypotheses inconsistent with the event (Frank et al., 2009).

Because of its more natural alignment with the learning models we use Teaching and Communication simulations, we implemented a positive hypothesis testing model<sup>2</sup>. In this model, learners begin with no hypotheses and add new ones to their store as they encounter data. Upon first encountering a word and a set of objects, the model encodes up to  $h$

---

<sup>2</sup> Our choice to focus on hypothesis testing rather than other learning frameworks is purely a pragmatic choice—the learning parameter  $p$  in this models maps cleanly onto the learning parameter in our other models. We encourage other researchers to adapt the code we have provided to estimate the long-term learning for other models.

hypothesized word-object pairs each with probability  $p$ . On subsequent trials, the model checks whether any of the existing hypotheses are consistent with the current data, and prunes any that are not. If no current hypotheses are consistent, it adds up to  $h$  new hypotheses each with probability  $p$ . The model has converged when it has pruned all but the one correct hypothesis for the meaning of a word. This model is most similar to the Propose but Verify model proposed in Trueswell et al. (2013), with the exception that it allows for multiple hypotheses. Because of the data generating process, storing prior disconfirmed hypotheses (as in Stevens et al., 2017), or incrementing hypotheses consistent with some but not all of the data (as in Yu & Smith, 2012) has no impact on learner and so we do not implement it here. We note also that, as described in Yu and Smith (2012), hypothesis testing models can mimic the behavior of associative learning models given the right parameter settings (Townsend, 1990).

In contrast to the Teaching and Communication simulations, the behavior of the Hypothesis Testing model depends on which particular non-target objects are present on each naming event. We thus began each simulation by generating a corpus of 100 naming events, on each sampling the correct target as well as  $(C-1)$  competitors from a total set of  $M$  objects. We then simulated a hypothesis tester learning over this set of events as described above, and recorded the first trial on which the learner converged (having only the single correct hypothesized mapping between the target word and target object). We repeated this process 1000 times for each simulated combination of  $M = (16, 32, 64, 128)$  total objects,  $C = (1, 2, 4, 8)$  objects per trial,  $h = (1, 2, 3, 4)$  concurrent hypotheses, as the learning rate  $p$  varied from .1 to 1 in increments of .1.

## General Discussion

Across naturalistic corpus data, experimental data, and model predictions, we see evidence that pressure to communicate successfully with a linguistically immature partner could fundamentally structure parent production. In our experiment, we showed that people



tune their communicative choices to varying cost and reward structures, and also critically to their partner’s linguistic knowledge—providing richer cues when partners are unlikely to know the language and many more rounds remain. These data are consistent with the patterns shown in our corpus analysis of parent referential communication and demonstrate that such pedagogically supportive input could arise from a motivation to maximize communicative success while minimizing communicative cost—no additional motivation to teach is necessary. In simulation, we demonstrate that such structure could have profound implications for child language learning, simplifying the learning problem posed by most distributional accounts of language learning.

Accounts of language learning often aim to explain its striking speed in light of the sheer complexity of the language learning problem itself. Many such accounts argue that simple (associative) learning mechanisms alone seem insufficient to explain the rapid growth of language skills and appeal instead to additional explanatory factors, such as the so-called language acquisition device, working memory limitations, word learning biases, etc. (e.g., Chomsky, 1965; Goldowsky & Newport, 1993; Markman, 1990). While some have argued for the simplifying role of language distributions (e.g., McMurray, 2007), these accounts largely focus on learner-internal explanations. For example, Elman (1993) simulates language learning under two possible explanations to intractability of the language learning problem: one environmental, and one internal. He first demonstrates that learning is significantly improved if the language input data is given incrementally, rather than all-at-once (Elman, 1993). He then demonstrates that similar benefits can arise from learning under limited working memory, consistent with the “less-is-more” proposal (Elman, 1993; Goldowsky & Newport, 1993). Elman dismisses the first account arguing that ordered input is implausible, while shifts in cognitive maturation are well-documented in the learner (Elman, 1993); however, our account’s emphasis on changing calibration to such learning mechanisms suggests the role of ordered or incremental input from the environment may be crucial.

This account is consonant with work in other areas of development, such as recent demonstrations that the infant’s visual learning environment has surprising consistency and incrementality, which could be a powerful tool for visual learning. Notably, research using head mounted cameras has found that infant’s visual perspective privileges certain scenes and that these scenes change across development (Fausey, Jayaraman, & Smith, 2016). In early infancy, the child’s egocentric visual environment is dominated by faces, but shifts across infancy to become more hand and hand-object oriented in later infancy (Fausey et al., 2016). This observed shift in environmental statistics mirrors learning problems solved by infants at those ages, namely face recognition and object-related goal attribution respectively (Fausey et al., 2016). These changing environmental statistics have clear implications for learning and demonstrate that the environment itself is a key element to be captured by formal efforts to evaluate statistical learning (Smith, Jayaraman, Clerkin, & Yu, 2018). Frameworks of visual learning must incorporate both the relevant learning abilities and this motivated, contingent structure in the environment (Smith et al., 2018).

By analogy, the work we have presented here aims to draw a similar argument for the language environment, which is also demonstrably beneficial for learning and changes across development. In the case of language, the contingencies between learner and environment are even clearer than visual learning. Functional pressures to communicate and be understood make successful caregiver speech highly dependent on the learner. Any structure in the language environment that is continually suited to changing learning mechanisms must come in large part from caregivers themselves. Thus, a comprehensive account of language learning that can successfully grapple with the infant curriculum (Smith et al., 2018) must explain parent production, as well as learning itself. In this work, we have taken first steps toward providing such an account.

Explaining parental modification is a necessary condition for building a complete theory of language learning, but modification is certainly not a sufficient condition for

language learning. No matter how calibrated the language input, non-human primates are unable to acquire language. Indeed, parental modification need not even be a necessary condition for language learning. Young children are able to learn novel words from (unmodified) overheard speech between adults ((Foushee, Griffiths, & Srinivasan, 2016), although there is reason to think that overheard sources may have limited impact on language learning broadly (e.g., Shneidman & Goldin-Meadow, 2012). Our argument is that the rate and ultimate attainment of language learners will vary substantially as a function of parental modification, and that describing the cause of this variability is a necessary feature of models of language learning.

**Generalizability and Limitations.** Our account aims to think about parent production and child learning in the same system, putting these processes into explicit dialogue. While we have focused on ostensive labeling as a case-study phenomenon, our account should reasonably extend to the changing structure found in other aspects of child-directed speech— though see below for important limitations to this extension. Some such phenomena will be easily accounted for: aspects of language that shape communicative efficiency should shift in predictable patterns across development.

While these language phenomena can be captured by our proposed framework, incorporating them will likely require altering aspects of our account and decisions about which alterations are most appropriate. For example, the exaggerated pitch contours seen in infant-directed speech could be explained by our account if we expand the definition of communicative success to include the goal of maintaining attention. Alternatively, one could likely accomplish the same goal by altering the cost structure to penalize loss of engagement. Thus, while this account should generalize to other modifications found in child-directed speech, such generalizations will likely require non-trivial alterations to the extant structure of the framework.

Of course, not all aspects of language should be calibrated to the child’s language

development. Our account also provides an initial framework for explaining aspects of communication that would not be modified in child-directed speech: namely, aspects of communication that minimally effect communicative efficiency. In other words, communication goals and learning goals are not always aligned. For example, young children sometimes overregularize past and plural forms, producing incorrect forms such as “runned” or “foots” (rather than the irregular verb “ran” or irregular plural “feet”; Marcus et al., 1992). Mastering the proper tense endings (i.e. the learning goal) might be aided by feedback from parent; however, adults rarely provide explicit corrective feedback for these errors (Marcus, 1993). This is perhaps because incorrect grammatical forms nonetheless successfully communicate their intended meaning, and thus do not prevent the successful completion of the communicative goal of language (Chouinard & Clark, 2003). The degree of alignment between communication and learning goals should predict the extent to which a linguistic phenomenon is modified in child-directed speech. Fully establishing the degree to which modification is expected for a given language phenomena will likely require working through a number of limitations in the generalizability of the framework as it stands.

Some aspects of parent production are likely entirely unrepresented in our framework, such as aspects of production driven by speaker-side constraints. Furthermore, our account is formulated primarily around concrete noun learning and future work must address its viability in other language learning problems. We chose to focus on ostensive labeling as a case-study phenomenon because it is an undeniably information-rich cue for young language learners, however ostensive labeling varies substantially across socio-economic, linguistic, and cultural groups (???). This is to be expected to the extent that parent-child interaction is driven by different goals (or goals given different weights) across these populations—variability in goals could give rise to variability in the degree of modification. Nonetheless, the generalizability of our account across populations remains unknown. Indeed, child-directed speech itself varies cross-linguistically, both in its features (citation) and quantity (e.g., Shneidman & Goldin-Meadow, 2012). There is some evidence that CDS

predicts learning even in cultures where CDS is qualitatively different and less prevalent than in American samples (Shneidman & Goldin-Meadow, 2012). Future work is needed to establish the generalizability of our account beyond the western samples studied here.

We see this account as building on established, crucial statistical learning skills—distributional information writ large and (unmodified) language data from overheard speech are undoubtedly helpful for some learning problems (e.g., phoneme learning). There is likely large variability in the extent to which statistical learning skills drive the learning for a given learning problem. The current framework is limited by its inability to account for such differences across learning problems, which could derive from domain or cultural differences. Understanding generalizability of this sort and the limits of statistical learning will likely require a full account spanning both parent production and child learning.

A full account that explains variability in modification across aspects of language will rely on a fully specified model of optimal communication. Such a model will allow us to determine both which structures are predictably unmodified, and which structures must be modified for other reasons. Nonetheless, this work is an important first step in validating the hypothesis that language input that is structured to support language learning could arise from a single unifying goal: The desire to communicate effectively.

## Conclusion

Building of early functional account of language learning (e.g., Brown, 1977), our account emphasizes the importance of communicative success in shaping language input and language learning. We have developed an initial formal framework for jointly considering parent productions and child language learning within the same system. We showed that such an account helps to explain parents' naturalistic communicative behavior and participant behavior in an iterated reference game. Formalized model predictions explain these behaviors without an explicit teaching goal, and show demonstrable effects on learning

854 in model simulations. In sum, this work

855 **Acknowledgement**

856 The authors are grateful to XX and YY for their thoughtful feedback on this  
857 manuscript. The authors are grateful to Madeline Meyers for her work coding referential  
858 communication in the corpus data. This research was supported by a James S McDonnell  
859 Foundation Scholars Award to DY.

## References

- Austin, J. L. (1975). *How to do things with words* (Vol. 88). Oxford university press.
- Baldwin, D. (2000). Interpersonal understanding fuels knowledge acquisition. *Current Directions in Psychological Science*, 9, 40–45.
- Bloom, P. (2000). *How children learn the meanings of words*. MIT press: Cambridge, MA.
- Blythe, R. A., Smith, K., & Smith, A. D. M. (2010). Learning times for large lexicons through cross-situational learning. *Cognitive Science*, 34, 620–642.
- Brown, R. (1977). Introduction. In C. E. Snow & C. A. Ferguson (Eds.), *Talking to children: Language input and interaction*. Cambridge, MA.: MIT Press.
- Chomsky, N. (1965). *Aspects of the theory of syntax*. MA: MIT Press.
- Chouinard, M. M., & Clark, E. V. (2003). Adult reformulations of child errors as negative evidence. *Journal of Child Language*, 30, 637–669.
- Clark, H. H. (1996). *Using language*. *Journal of Linguistics* (pp. 167–222). Cambridge Univ Press.
- DeCasper, A. J., & Fifer, W. P. (1980). Of human bonding: Newborns prefer their mothers' voices. *Science*, 208(4448), 1174–1176.
- DeCasper, A. J., & Spence, M. J. (1986). Prenatal maternal speech influences newborns' perception of speech sounds. *Infant Behavior and Development*, 9(2), 133–150.
- Eaves Jr, B. S., Feldman, N. H., Griffiths, T. L., & Shafto, P. (2016). Infant-directed speech is consistent with teaching. *Psychological Review*, 123(6), 758.
- Elman, J. L. (1993). Learning and development in neural networks: The importance of

starting small. *Cognition*, 48(1), 71–99.

Estes, K. G., Evans, J. L., Alibali, M. W., & Saffran, J. R. (2007). Can infants map meaning to newly segmented words? Statistical segmentation and word learning. *Psychological Science*, 18(3), 254–260.

Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016). From faces to hands: Changing visual input in the first two years. *Cognition*, 152, 101–107.

Foushee, R., Griffiths, T. L., & Srinivasan, M. (2016). Lexical complexity of child-directed and overheard speech: Implications for learning. In.

Frank, M. C., & Goodman, N. D. (2012). Predicting pragmatic reasoning in language games. *Science*, 336, 998–998.

Frank, M. C., & Goodman, N. D. (2014). Inferring word meanings by assuming that speakers are informative. *Cognitive Psychology*, 75, 80–96.

Frank, M. C., Goodman, N., & Tenenbaum, J. (2009). Using speakers' referential intentions to model early cross-situational word learning. *Psychological Science*, 20, 578–585.

Franke, M. (2013). Game theoretic pragmatics. *Philosophy Compass*, 8(3), 269–284.

Gelman, A., & Hill, J. (2006). *Data analysis using regression and multilevel/hierarchical models*. Cambridge university press.

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 Development*, 71(4), 878–894.

Goldin-Meadow, S., Levine, S. C., Hedges, L. V., Huttenlocher, J., Raudenbush, S. W., & Small, S. L. (2014). New evidence about language and cognitive development based



on a longitudinal study: Hypotheses for intervention. *American Psychologist*, 69(6), 588–599.

Goldowsky, B. N., & Newport, E. L. (1993). Limitations on the acquisition of morphology: The less is more hypothesis. In *The proceedings of the twenty-fourth annual child language research forum* (p. 124). Center for the Study of Language (CSLI).

Gomez, R. L., & Gerken, L. (1999). Artificial grammar learning by 1-year-olds leads to specific and abstract knowledge. *Cognition*, 70(2), 109–135.

Goodman, N. D., & Frank, M. C. (2016). Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Sciences*, 20(11), 818–829.

Goodman, N. D., & Stuhlmüller, A. (2014). The Design and Implementation of Probabilistic Programming Languages. <http://dippl.org>.

Grice, H. P. (1975). Logic and conversation. In P. Cole & J. L. Morgan (Eds.), *Syntax and semantics: Vol. 3: Speech acts* (pp. 41–58). New York, NY: Academic Press.

Griffiths, T. L., Lieder, F., & Goodman, N. D. (2015). Rational use of cognitive resources: Levels of analysis between the computational and the algorithmic. *Topics in Cognitive Science*, 7(2), 217–229.

Kaelbling, L. P., Littman, M. L., & Cassandra, A. R. (1998). Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101, 99–134.

Leung, A., Tunkel, A., & Yurovsky, D. (2019). Parents calibrate speech to their children’s vocabulary knowledge. In *CogSci* (pp. 651–656).

Luce, R. D. (1959). Individual choice behavior.

MacKay, D. J. (1992). Bayesian interpolation. *Neural Computation*, 4(3), 415–447.

- 925 Marcus, G. F. (1993). Negative evidence in language acquisition. *Cognition*, 46(1), 53–85.
- 926 Marcus, G. F., Pinker, S., Ullman, M., Hollander, M., Rosen, T. J., Xu, F., & Clahsen, H.  
927 (1992). Overregularization in language acquisition. *Monographs of the Society for*  
928 *Research in Child Development*, i–178.
- 929 Markman, E. M. (1990). Constraints children place on word meanings. *Cognitive Science*,  
930 14(1), 57–77.
- 931 Markman, E. M., & Wachtel, G. F. (1988). Children’s use of mutual exclusivity to constrain  
932 the meanings of words. *Cognitive Psychology*, 20(2), 121–157.
- 933 Marr, D. (1982). *Vision: A Computational Investigation into the Human Representation and*  
934 *Processing of Visual Information*. New York, NY: W. H. Freeman.
- 935 Maye, J., Werker, J. F., & Gerken, L. (2002). Infant sensitivity to distributional information  
936 can affect phonetic discrimination. *Cognition*, 82(3), B101–B111.
- 937 McMurray, B. (2007). Defusing the childhood vocabulary explosion. *Science*, 317(5838),  
938 631–631.
- 939 Meltzoff, A. N. (2005). Imitation and other minds: The “like me” hypothesis. *Perspectives*  
940 *on Imitation: From Neuroscience to Social Science*, 2, 55–77.
- 941 Mintz, T. H. (2003). Frequent frames as a cue for grammatical categories in child directed  
942 speech. *Cognition*, 90(1), 91–117.
- 943 Newport, E. L. (1990). Maturational constraints on language learning. *Cognitive Science*,  
944 14(1), 11–28.
- 945 Newport, E. L., Gleitman, H., & Gleitman, L. R. (1977). Mother, I’d rather do it myself:  
946 Some effects and non-effects of maternal speech style. In C. A. Ferguson (Ed.),

*Talking to children language input and interaction* (pp. 109–149). Cambridge University Press.

Rafferty, A. N., Brunskill, E., Griffiths, T. L., & Shafto, P. (2016). Faster teaching via pomdp planning. *Cognitive Science*, 40(6), 1290–1332.

Reisenauer, R., Smith, K., & Blythe, R. A. (2013). Stochastic dynamics of lexicon learning in an uncertain and nonuniform world. *Physical Review Letters*, 110(25), 258701.

Saffran, J. R. (2003). Statistical language learning: Mechanisms and constraints. *Current Directions in Psychological Science*, 12(4), 110–114.

Saffran, J. R., Aslin, R. N., & Newport, E. L. (1996). Statistical learning by 8-month-old infants. *Science*, 274, 1926–1928.

Scott, R. M., & Fischer, C. (2012). 2.5-year-olds use cross-situational consistency to learn verbs under referential uncertainty. *Cognition*, 122, 163–180.

Shafto, P., Goodman, N. D., & Griffiths, T. L. (2014). A rational account of pedagogical reasoning: Teaching by, and learning from, examples. *Cognitive Psychology*, 71, 55–89.

Shneidman, L. A., & Goldin-Meadow, S. (2012). Language input and acquisition in a mayan village: How important is directed speech? *Developmental Science*, 15(5), 659–673.

Simon, H. A. (1991). Bounded rationality and organizational learning. *Organization Science*, 2(1), 125–134.

Siskind, J. M. (1996). A computational study of cross-situational techniques for learning word-to-meaning mappings. *Cognition*, 61, 39–91.

Smith, K., Smith, A. D., & Blythe, R. A. (2011). Cross-situational learning: An

experimental study of word-learning mechanisms. *Cognitive Science*, 35(3), 480–498.

Smith, L. B., Jayaraman, S., Clerkin, E., & Yu, C. (2018). The developing infant creates a curriculum for statistical learning. *Trends in Cognitive Sciences*, 22(4), 325–336.

Smith, L. B., & Yu, C. (2008). Infants rapidly learn word-referent mappings via cross-situational statistics. *Cognition*, 106, 1558–1568.

Smith, L. B., & Yu, C. (2013). Visual attention is not enough: Individual differences in statistical word-referent learning in infants. *Language Learning and Development*, 9, 25–49.

Stevens, J. S., Gleitman, L. R., Trueswell, J. C., & Yang, C. (2017). The pursuit of word meanings. *Cognitive Science*, 41, 638–676.

Thiessen, E., Hill, E. A., & Saffran, J. R. (2005). Infant-directed speech facilitates word segmentation. *Infancy*, 7, 53–71.

Tomasello, M., & Farrar, M. J. (1986). Joint attention and early language. *Child Development*, 57, 1454–1463.

Townsend, J. T. (1990). Serial vs. Parallel processing: Sometimes they look like tweedledum and tweedledee but they can (and should) be distinguished. *Psychological Science*, 1(1), 46–54.

Trueswell, J. C., Medina, T. N., Hafri, A., & Gleitman, L. R. (2013). Propose but verify: Fast mapping meets cross-situational word learning. *Cognitive Psychology*, 66(1), 126–156.

Vlach, H. A., & Johnson, S. P. (2013). Memory constraints on infants' cross-situational statistical learning. *Cognition*, 127(3), 375–382.

- 991 Vogt, P. (2012). Exploring the robustness of cross-situational learning under zipfian  
992 distributions. *Cognitive Science*, 36(4), 726–739.
- 993 Wittgenstein, L. (1953). *Philosophical investigations*. Oxford, UK: Basil Blackwell.
- 994 Xu, F., & Tenenbaum, J. B. (2007). Word learning as Bayesian inference. *Psychological*  
995 *Review*, 114(2), 245–272.
- 996 Yang, S. C.-H., Vong, W. K., Yu, Y., & Shafto, P. (2019). A unifying computational  
997 framework for teaching and active learning. *Topics in Cognitive Science*, 11(2),  
998 316–337.
- 999 Yu, C. (2008). A statistical associative account of vocabulary growth in early word learning.  
1000 *Language Learning and Development*, 4(1), 32–62.
- 1001 Yu, C., & Smith, L. B. (2012). Modeling cross-situational word-referent learning: Prior  
1002 questions. *Psychological Review*, 119, 21–39.
- 1003 Yurovsky, D. (2012). Statistical speech segmentation and word learning in parallel:  
1004 scaffolding from child-directed speech. *Frontiers in Psychology*, 374.
- 1005 Yurovsky, D. (2018). A communicative approach to early word learning. *New Ideas in*  
1006 *Psychology*, 50, 73–79.
- 1007 Yurovsky, D., Doyle, G., & Frank, M. C. (2016). Linguistic input is tuned to children’s  
1008 developmental level. In *Proceedings of the annual meeting of the cognitive science*  
1009 *society* (pp. 2093–2098).
- 1010 Yurovsky, D., & Frank, M. C. (2015). An integrative account of constraints on  
1011 cross-situational learning. *Cognition*, 145, 53–62.
- 1012 Yurovsky, D., Fricker, D. C., Yu, C., & Smith, L. B. (2014). The role of partial knowledge in

- 1013            statistical word learning. *Psychonomic Bulletin & Review*, 21, 1–22.
- 1014   Yurovsky, D., Meyers, M., Burke, N., & Goldin-Meadow, S. (2018). Children gesture when  
1015            speech is slow to come. In J. Z. Chuck Kalish Martina Rau & T. Rogers (Eds.),  
1016            *CogSci 2018* (pp. 2765–2770).
- 1017   Yurovsky, D., Yu, C., & Smith, L. B. (2013). Competitive processes in cross-situational word  
1018            learning. *Cognitive Science*, 37, 891–921.