

**Communicative pressure on caregivers leads to language input that supports  
children's word learning**

Ben Morris<sup>1</sup> & Daniel Yurovsky<sup>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 motivated to communicate with them. Does this communicative pressure on caregivers lead them to structure children’s linguistic input in a way that facilitates learning? As a case study, we focus on a canonical pedagogical behavior in caregivers: use of ostensive labeling to teach children new words. Using a multi-method approach, we show that this pedagogically supportive behavior can emerge from communicative pressure alone, without any explicit pedagogical goals. First, in a corpus study, we show that caregivers communicate by ostensive labeling precisely when children are likely to learn a new word. In an iterated reference game, we experimentally show that this strategy 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 communication model that includes planning, without any explicit teaching goal. Finally, in a series of simulations, we explore the language learning consequences of having a communicatively-motivated caregiver. We show that under many parameterizations, simple learning mechanisms interacting with a communicatively-motivated partner outperform more powerful learning mechanisms. This perspective offers a first step toward a unifying, formal account of both how linguistic input is structured and how children learn from it.

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

**Communicative pressure on caregivers leads to language input that supports  
children's 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 levels of language, including phonemic discrimination (Maye et al., 2002), word segmentation (Saffran, 2003), learning the meanings of both nouns (Smith & 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 have confirmed both the early availability of distributional learning mechanisms, and their potential utility across these diverse language phenomena (DeCasper & Fifer, 1980; DeCasper & Spence, 1986; Estes et al., 2007; Gomez & Gerken, 1999; Maye et al., 2002; Saffran et al., 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). In addition, computational models of statistical learning show that the rate of acquisition is highly sensitive to variation in environmental statistics.

For example, the Zipfian distribution of word frequencies and word meanings in the real world can lead to very long learning times for rare words (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—features likely typical of the naturalistic learning environment (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 children’s learning problem. Caregivers do adjust many aspects of their communication when conversing with young children, in both their language (infant directed speech, Snow, 1977) and actions (motionese, Brand et al., 2002)—and such modifications have been shown to yield learning benefits across a number of language phenomena. For example, in phoneme learning, infant-directed speech provides examples that seem to facilitate the acquisition of phonemic categories (Eaves Jr et al., 2016). In word segmentation tasks, infant-directed speech facilitates infant learning more than matched adult-directed speech (Thiessen et al., 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 et al., 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 process that generates structure in the language environment as a problem separate from explaining language learning. However, across a number of language phenomena, the language environment is not merely supportive, but

seems calibrated to children’s changing learning mechanisms (Yurovsky, 2018). For example, across development, caregivers engage in more multimodal naming of novel objects than familiar objects, and rely on this temporal synchrony between verbal labels and gesture most with young children (Gogate et al., 2000). The prevalence of such synchrony in child-directed speech parallels infant learning mechanisms: young infants appear to rely more on synchrony as a cue for word learning than older infants, and language input mirrors this developmental shift (Gogate et al., 2000). Beyond age-related changes, caregiver speech may also support learning through more local calibration to a child’s knowledge. Caregivers have been shown to provide more language about referents that are unknown to their child, and adapt their language in-the-moment to the knowledge their child displays during a referential communication game (Leung et al., 2021). The calibration of caregivers’ production to the child’s learning and knowledge suggests a co-evolution such that these processes should not be considered in isolation.

What then gives rise to the structure in early language input that mirrors children’s learning mechanisms? Because of widespread agreement that caregivers’ speech is not usually motivated by explicit pedagogical goals (Newport et al., 1977), the calibration of speech to learning mechanisms seems a happy accident; caregivers’ speech just happens to be calibrated to children’s learning needs. If caregivers’ speech was pedagogically-motivated, extant formal frameworks of teaching could be used to derive predictions and expectations (e.g., Shafto et al., 2014). This framework—in which caregivers choose what they say in order to maximize their children’s learning—has seen some success in explaining children’s developing phoneme discrimination (Eaves Jr et al., 2016). Even here, however, the statistical structure of phonemes in caregivers’ input also contains features that are not well explained by a pedagogical goal (McMurray et al., 2013).

Instead, the recent outpouring of work exploring optimal communication (the Rational Speech Act model, see Frank & Goodman, 2012) provides a different framework for understanding caregivers’ 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, most notably a range of pragmatic inferences (e.g., Frank & Goodman, 2012; Bohn & Frank, 2019; Goodman & Frank, 2016). 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 caregiver’s production: Both are driven by a pressure to communicate successfully.

Early, influential functionalist accounts of language learning focused on the importance of communicative goals (e.g., Brown, 1977). Recent work has demonstrated that many structural aspects of natural language may have arisen from pressure to communicate efficiently (Gibson et al., 2019). Our goal in this work is to formalize the intuitions in these functionalist accounts in a computational model of language input, 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 this in-the-moment pressure to communicate successfully (Yurovsky, 2018).

To examine the tuning hypothesis, we draw on evidence from naturalistic data, a reference game experiment, a formal model, and learning simulations. 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 caregiver communicative behavior in a longitudinal corpus of parent-child interaction in the home (Goldin-Meadow et al., 2014). We investigate the extent to which caregivers tune their ostensive labeling across their

child’s development to align to their child’s developing linguistic knowledge (Yurovsky et al., 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. Finally, we demonstrate potential benefits for the learner through a series of simulations to show that communicative pressure on caregivers’ speech facilitates learning. Under a variety of parameter settings, simple learners interacting with communicative partners outperform more complex statistical learners.

In sum, our goal in this work is to argue that the fundamental unit of analysis for understanding children’s language learning is not the child alone, but rather the caregiver-child dyad. To provide converging evidence for this claim, we use apply a multi-method approach to a case study of one pedagogically-supportive behavior. We use observational data to understand the contexts in which ostensive labeling happens, and then experimentally manipulate these contexts in a model system in order to demonstrate that the observed relationships are causal. We develop a formal model that quantitatively accounts for this experimental data, and then simulate the long-term learning outcomes of this model for children’s language learning. Together, these different methods show that rapid language learning can emerge from even highly-constrained child learners working together with communicative-motivated caregivers.

## Corpus Analysis

We first investigate parents' use of ostensive labeling in referential communication in a longitudinal corpus of parent-child interaction. Ostensive labeling—the behavior pointing to and labeling an object with its name—is a powerful source of information for word learning because it reduces the ambiguity about the possible meanings of the word the child is hearing. The word could of course still have many possible meanings—it could refer to a part of the object, or something about it's state for instance, but it is unlikely to refer to one of the other objects in the room (Quine, 1960). Prior work has shown that parents tend to use ostensive labeling when referring to objects with basic level words (e.g. “chair”, “dog”, Callanan, 1985, @ninio1980), and that ostensive labeling is a particularly powerful source of data for young children (e.g., Baldwin, 2000; Gogate et al., 2000). We take the ostensive labeling with 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 ostensive labeling in children's language environment at different ages. We find that this pedagogically-supportive form of input shows a key hallmark of adaptive tuning: caregivers using this information-rich cue more for young children and infrequent objects. Thus, parents' production of ostensive labeling is tuned to children's developing linguistic knowledge (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. We coded each of these communicative instances to identify 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 on caregivers' productions of ostensive labeling in the form of a multi-modal reference.

### *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 ten target children, five were girls, three were Black and two were Mixed-Race. Families spanned a broad range of incomes, with two 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 up for view) 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 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.

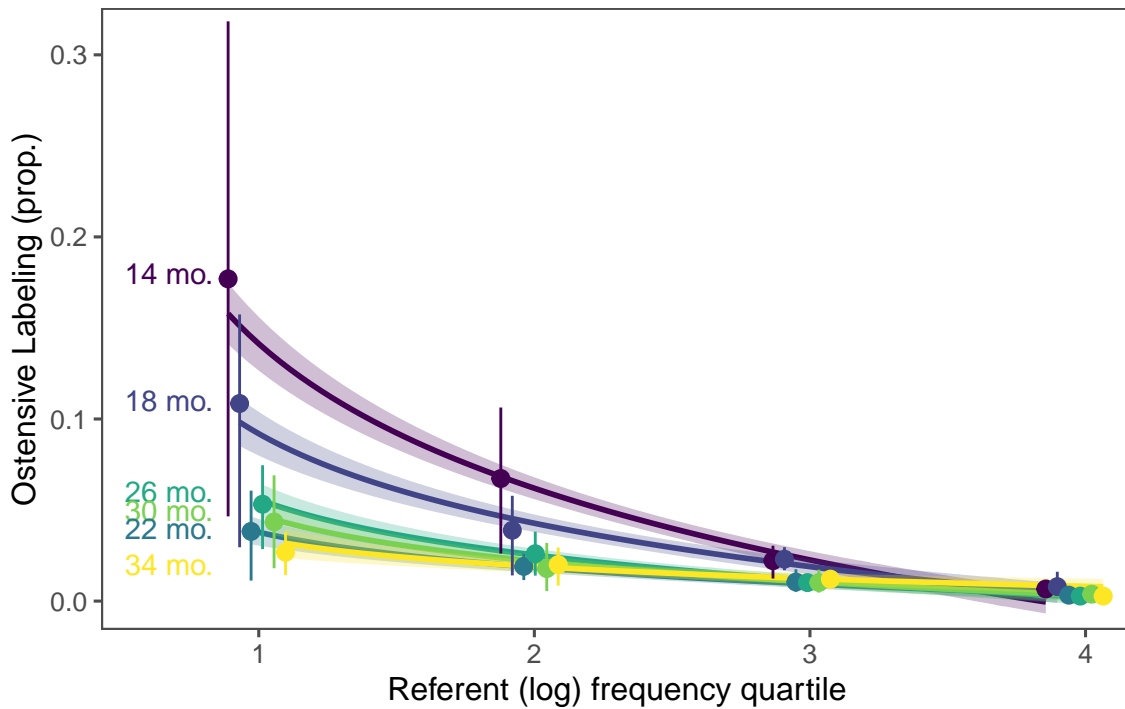
To ensure that each referent could potentially be referred to in gesture or speech, we focused on concrete nouns. We further wanted to ensure that the referents were physically present in the scene (and thus available for deictic gestures). Using the transcripts, a human rater judged whether the referent 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.

To ensure our transcript-based coding of referent presence 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 available in the open access dataset for further analysis.

## **Results**

Corpus data were analyzed using a mixed effects regression to predict parents' use of ostensive labeling 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 use ostensive labeling less with older children ( $\beta = -0.84, p < .001$ ) and less for more frequent referents ( $\beta = -0.09, p = .045$ ). In addition, the interaction between the two was significant, indicating that for parents ostensively label more for younger children when referents are infrequent ( $\beta = 0.18, p = .001$ ). Thus, in these data, we evidence that parents provide more pedagogically-supportive input about rarer things in the world for their younger children (Figure 1).



**Figure 1**

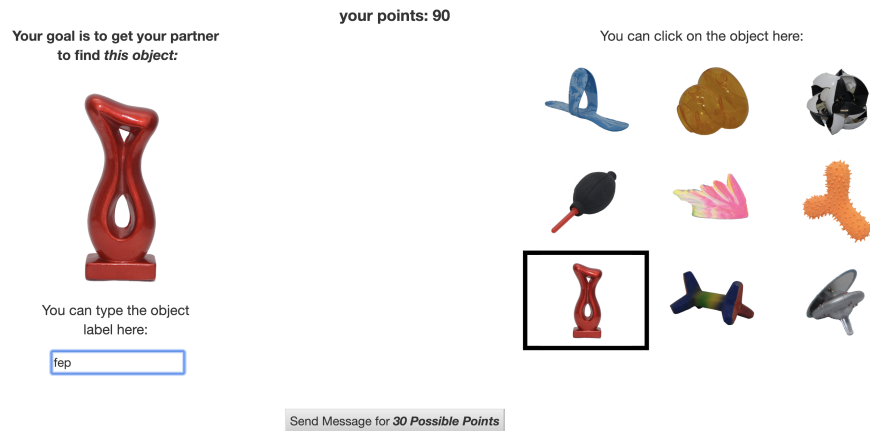
*Parents' rate of ostensive labelling via multi-modal reference for referents of varying frequency. Log frequencies were binned into quartiles for visualization. Points show empirical frequencies along with 95% confidence intervals computed by non-parametric bootstrapping. Lines show estimates from a logisitic regression along with ribbons indicating one standard error. Parents used ostensive labeling more for younger children and infrequent referents.*

## Discussion

Caregivers are not indiscriminate in their use of ostensive labeling; 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 caregiver alignment: parents are sensitive to their child’s linguistic knowledge and adjust their communication accordingly (Yurovsky et al., 2016).

This caregiver alignment is straightforwardly predicted from a dyadic perspective that privileges communicative goals, yet wholly missing from learner-centric accounts of language learning (Yurovsky, 2018). While some early accounts of language learning suggested a important role for caregiver tuning of linguistic input (e.g., Snow, 1972), these accounts fell out of favor in light of influential investigations that did not find tuning of syntactic structure (Newport et al., 1977). However, this work, along with other recent papers, shows that caregivers may be well calibrated to their children’s semantic knowledge and this tuning may be quite relevant for language problems like word learning where communicative success is central (e.g., Yurovsky et al., 2016, @leung2021).

While language input that is tuned to the child’s linguistic competence could undoubtedly aid in language learning, the presence of such input does not necessarily imply pedagogical goals. These data could potentially be explained by a simpler, less altruistic goal: to communicate successfully. Because we do not have access to parents’ goals, we these data alone cannot distinguish between these competing accounts not can they determine whether word frequency and children’s linguistic knowledge are *causally* related to the use of ostensive labeling. To demonstrate that communicative pressure alone is sufficient to explain these patterns in ostensive labeling, we next developed an experimental paradigm to study these and other factors in a controllable model system.

**Figure 2**

*Screenshot showing the participant view during gameplay.*

## Reference Game Experiment

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 nine 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 2). Although we told participants that their partner would be played by another human, all partners were actually pre-programmed bots.

Participants could choose to refer either using the novel labels they had been exposed to, or they could use a deictic gesture (i.e. point) to indicate the referent to their partner. The point 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.

In choosing a communicative signal, participants should try to maximize their likelihood of success while minimizing their communicative cost (see e.g. Frank &

Goodman, 2012). This cost should be related to the energy required to produce the signal. Though it need not be the case for all gestures and contexts, our task compares simple lexical labeling and unambiguous deictic gestures, which are likely are slower and more effortful to produce (e.g., see Yurovsky et al., 2018). We thus assumed that pointing should be more costly than speaking. Nonetheless, because we do not have a way of estimating the costs of pointing and speaking from the observational data, we manipulated them experimentally across conditions to understand how they impact peoples' behavior.

Critically, participants were told that they would play this game repeatedly with their partner. In these repeated interactions, participants are able to learn about an interlocutor and potentially influence their learning. Thus, there is a third type of signal participants could send: using both pointing and speech within a single trial to effectively teach the listener through ostensive labeling. This strategy necessitates making inferences about their partner's knowledge state, so we induced knowledge asymmetries between the participant and their partner. 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 participant pay a high cost upfront. However, rational communicators may understand that if one is accounting for future trials, paying the cost upfront to teach their partner allows them to use a less costly message strategy on subsequent trials (namely, speech). Manipulating the partner's 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.

While our reference game setting has limited ecological validity, this setup allows us to explicitly manipulate the crucial features of the communicative setting (e.g., communicative cost, strategy, and partner knowledge). In this controlled task, we can look for the emergence of structure that parallels the naturalistic input described in our corpus

evidence, while also experimentally testing for possible drivers of such structure. This experimental setup further allows us to straightforwardly test and compare key predictions using a formal model to explain participant behavior.

## 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.

### *Transparency and Openness*

This sample size, experimental design, and analysis plan were pre-registered at <https://osf.io/63qdg>. Sample size was determined based on prior pilot experiments. All data, analysis code, and research materials are available at <https://osf.io/d9gkw/>. Data were analyzed using R version 4.1.1 (R Core Team, 2021). Models were estimated using version 1.1-27.1 of the `lme4` package (Bates et al., 2015).

### *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 (accurately reporting their partners knowledge prior to gameplay) 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, three of them two times, and three of them 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, participants saw the target object in addition to an array of all nine objects. Participants had the option of either directly selecting the target object from the array (pointing)—a higher cost, but unambiguous cue—or typing a label for the object (speech)—a lower cost cue contingent on their partner’s knowledge. After sending the message, participants were shown which object their partner 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 three between-subjects conditions, participants were told that their partner had either no experience with the lexicon, had the same experience as them, or had twice their experience. As a manipulation check, participants were then asked to report their partner’s level of exposure, and were corrected if they answered incorrectly. Note that in order for participants to account for future interactions during the game, participants were explicitly told that they would be asked to refer to each object three times during the game.

Partners were programmed with starting knowledge states initialized according to the partner knowledge condition. Partners with no exposure began the game with knowledge of 0 object-label pairs. Partners with the same exposure as the participant



began with knowledge of five object-label pairs (three high-frequency, one mid-frequency, one low-frequency), based on the learning we observed from participants in a pilot experiment. Lastly, partners with twice as much exposure as the participant began with knowledge of all nine object-label pairs.

To simulate knowledgeable behavior, when the participant typed an object label, the partner 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 partner’s vocabulary. Partners then selected the candidate with the smallest edit distance (e.g., if a participant typed the message “tomi”, the programmed partner would select the referent corresponding to “toma”, provided toma was found in its vocabulary). If the participant’s message was more than two edits away from all of the words in the partner’s vocabulary, the partner selected an object whose label they did not know. If the participant clicked on an object (pointing), the partner was programmed to always select that referent.

Participants could win up to 100 points per trial if their partner correctly selected the target referent based on their message. If the partner failed to identify the target object, participants received no points. We manipulated the relative utility of the speech cue between subjects across two conditions: Higher Speech Efficiency and Lower Speech Efficiency. In the *Higher Speech Efficiency* condition, participants received 30 points for gesturing and 100 points for labeling, and thus speech had very little cost relative to pointing and participants should be highly incentivized to speak. In the *Lower Speech Efficiency* condition, participants received 50 points for gesturing and 80 points for labeling, and thus gesturing is still costly relative to speech, but the difference between them is smaller lowering the incentive to speak.

Participants were told about a third type of possible message: using both pointing and speech within a single trial to effectively teach their partner an object-label mapping. This action directly mirrors the ostensive labeling behavior parents produced in the corpus

data—it yields an information-rich, pedagogically-supportive learning moment. In order to produce this teaching behavior, participants had to pay the cost of producing both cues (i.e. both pointing 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). Partners were programmed to integrate new taught words into their knowledge of the lexicon, and check those taught labels on subsequent trials when evaluating participants’ messages.

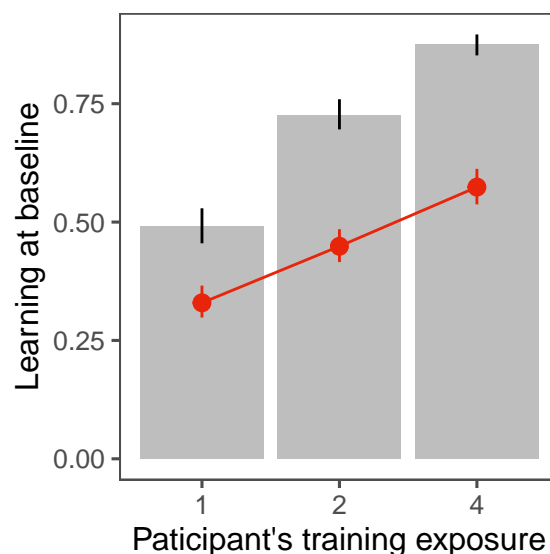
Crossing our 2 between-subjects manipulations yielded 6 conditions (2 utility manipulations: Higher Speech Efficiency and Lower Speech Efficiency; and 3 levels of partner’s exposure: None, Same, Twice), with 80 participants in each condition. We expected that participants would be sensitive to our utility manipulation, such that rates of labeling and teaching would be higher in the Higher Speech Efficiency conditions than the other conditions. Crucially, 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 participant knew more (None manipulation) compared with when there was equal knowledge (Same manipulation) or when the partner was more familiar with the language (Twice manipulation).

## Results

In each trial, participants could choose one of 3 communicative strategies: pointing, speech, or teaching. We expected participants to flexibly use communicative strategies in response to their relative utilities, their partner’s knowledge of the lexicon, and participants’ own lexical knowledge. To test our predictions about each communicative behavior (pointing, speech, and teaching), we fit separate logistic mixed effects models for each behavior, reported below. It should be noted that these three behaviors are mutually exhaustive. First, we report how well participants learned our novel lexicon during training.

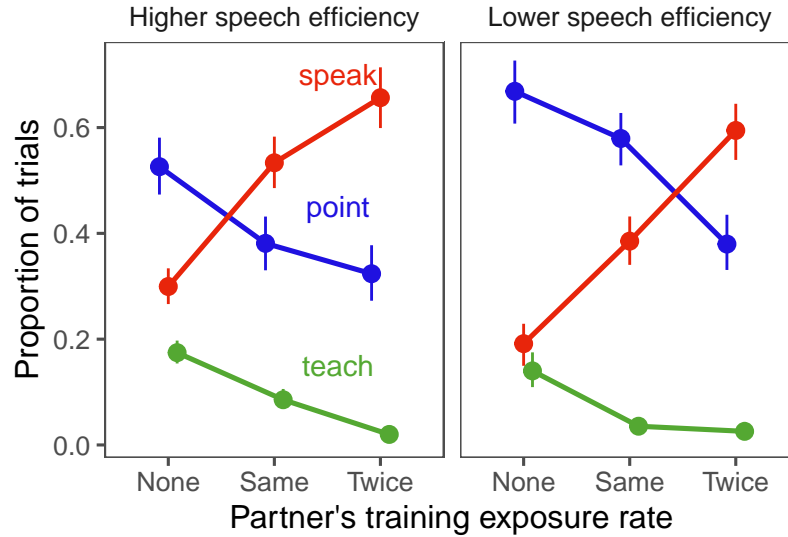
406 *Learning*

407 As an initial check of our exposure manipulation, we first fit a logistic regression  
 408 predicting accuracy at test from a fixed effect of exposure rate and random intercepts and  
 409 slopes of exposure rate by participant as well as random intercepts by item. We found a  
 410 reliable effect of exposure rate, indicating that participants were better able to learn items  
 411 that appeared more frequently in training ( $\beta = 1.08$ ,  $p < .001$ , see Figure 3). On average,  
 412 participants knew approximately 6 of the 9 words in the lexicon ( $M(sd) = 6.28 (2.26)$ ). As  
 413 a simple check of baseline performance across groups, an analysis of variance confirmed  
 414 that learning did not differ systematically across participants by partner's exposure, utility  
 415 manipulation, or their interaction ( $ps > 0.05$ ).



**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 during gameplay in which participants used the learned labels, excluding teaching behaviors. Error bars show 95% confidence intervals computed by non-parametric bootstrapping.*

**Figure 4**

*Participants' communicative method choice as a function of exposure and the utility manipulation. Error bars indicate 95% confidence intervals computed by non-parameteric bootstrapping*

### Pointing

When will participants rely on pointing? Pointing 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 the utility scheme is relatively biased toward pointing (i.e., the Lower Speech Efficiency condition). To test these predictions, we ran a mixed effects logistic regression to predict whether participants chose to point 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 (further details can be found in the supplemental materials).

Consistent with our predictions, exposure rate during training was a significant negative predictor of pointing during the game, such that participants were less likely to rely on pointing for well trained (and thus well learned) objects ( $\beta = -0.50$ ,  $p < .001$ ).

Additionally, participants were significantly more likely to point in the Lower Speech Efficiency condition where pointing is relatively less costly, compared with the Higher Speech Efficiency condition ( $\beta = 1.20, p < .001$ ; see Figure 4). We also found a significant negative effect of partner’s knowledge, such that participants pointed more for partners with less knowledge of the lexicon ( $\beta = -0.81, p < .001$ ).

Participants pointing behavior in this game is not just reflecting their own knowledge or the general efficiency of pointing, but is crucially modulated by their beliefs about their partner’s knowledge. These patterns mirror previous corpus analyses demonstrating parents’ use of pointing in naturalistic parental communicative behaviors (see Yurovsky et al., 2018). Note that these effects cannot be explained by solely participants’ knowledge; all patterns above hold when looking *only* at words known by the participant at pretest ( $ps < 0.01$ ).

### *Speech*

When will participants rely on 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), and when utility scheme is relatively biased toward speech (i.e., the Higher Speech Efficiency condition). To test these predictions, we ran a mixed effects logistic regression to predict whether participants 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 (further details can be found in the supplemental materials).

Consistent with our predictions, speech seemed to largely trade off with gesture. Exposure rate during training was a significant positive predictor of speaking during the game, such that participants were more likely to use speech for well trained (and thus well learned) objects ( $\beta = 0.35, p < .001$ ). Additionally, participants were significantly less likely

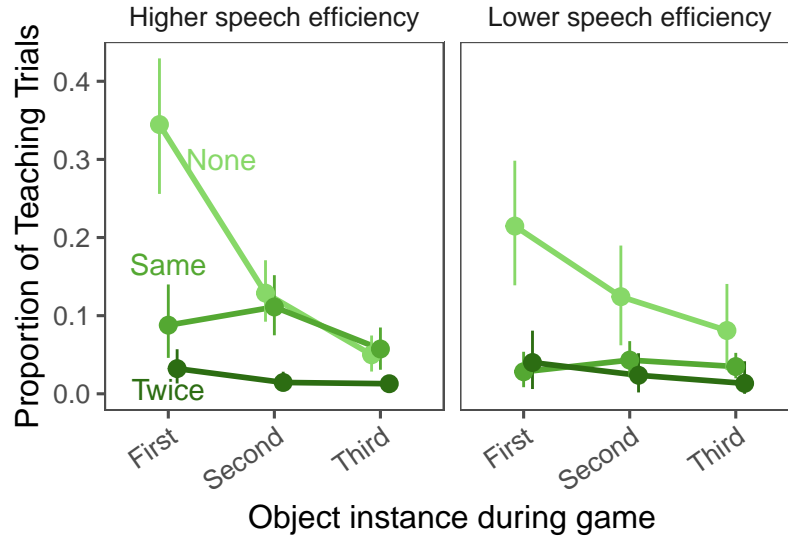
to speak in the Lower Speech Efficiency condition where speech is relatively more costly, compared with the Higher Speech Efficiency condition ( $\beta = -0.87$ ,  $p = .001$ ). Crucially, we also found a significant positive effect of partner's knowledge, such that participants used speech more for partners with more knowledge of the lexicon ( $\beta = 1.95$ ,  $p < .001$ ).

Unlike for gesture, there was a significant effect of object instance in the game (i.e., first, second, or third trial with this target object) on the rate of speaking, such that later trials were 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 participant 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 participants to trade off between gesture and speech sensibly. Next, we turn to the emergence of teaching behavior.

When will participants rely on teaching? 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), and when utility scheme is relatively biased toward speech (i.e., the Higher Speech Efficiency condition). In this utility scheme, it is more valuable to pay the cost of teaching early because of the increased benefit of using speech later. To test these predictions, we ran a mixed effects logistic regression to predict whether participants 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 (further details can be found in the supplemental materials).

**Figure 5**

*Rates of teaching across the six utility and partner knowledge conditions as a function of how many times the current target referent object had previously been the target. Error bars show 95% confidence intervals computed by non-parametric bootstrapping.*

Consistent with our predictions, rates of teaching were higher for more highly 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 Lower Speech Efficiency condition where speech is relatively more costly, compared with the Higher Speech Efficiency condition ( $\beta = -0.96, p = .001$ ). Crucially, participants' teaching behavior depends on their beliefs about their partner's knowledge; we 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 (while rates of speaking increase, as we saw above). Consistent with this prediction, participants were significantly less likely to teach on the later appearances of the target object ( $\beta = -1.09, p < .001$ ).

## Discussion

As predicted, the results of this experiment 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. Participants choice of communicative behavior depends not just on the participants knowledge (i.e. exposure rate) or the communicative cost (i.e. efficiency condition), but crucially on the participant's beliefs about their partner's knowledge.

For speaking and teaching, we also see effects of object instance where teaching is in effective strategy early on because it enables speaking on subsequent trials. In our experiment, participants are told that they will see each object 3 times throughout the game. The utility of teaching relies on the use of cheaper strategies on subsequent trials, so it is crucial that participants know they will be interacting with their partner repeatedly.

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 trends 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 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 been successfully applied to a variety of linguistic phenomena 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 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 that 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, 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.

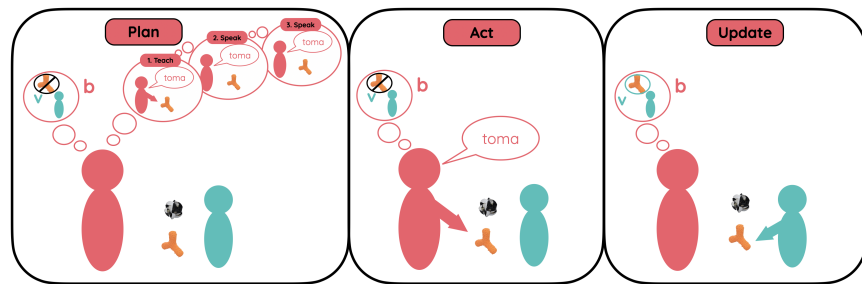
Alternatively, pedagogically-supportive input could emerge from an explicit pedagogical goal. Shafto et al. (2014) have developed an framework of rational pedagogy built on the same recursive reasoning principles as in the Rational Speech Act Framework: Teachers aim to teach a concept by choosing a set of examples that would maximize learning for students who reason about the teachers choices as attempting to maximize their learning. Rafferty et al. (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 caregivers' 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 caregivers' 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. The Utilities ( $U$ ) are given directly by the rules of this game. Because communication is a repeated game, people should take actions that maximize their Expected Utility ( $EU$ ) 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 et al., 1998).

Optimal planning in a Partially Observable Markov Decision Process involves a cycle of three phases: (1) Plan, (2) Act, and (3) Update beliefs. On each trial of the referential game, the model first makes a plan—reasoning about which action it should take on this trial and on subsequent trials to come. To help build intuition for the model, we first describe how it might operate in a series of example trials. We then give a formal description of the model’s learning and behavior.



**Figure 6**

*A schematic of the plan-act-update cycle of communication as a partially observable Markov decision process (POMDP). The red figure corresponds to caregivers in the world, or participants in the experiment. The blue figure corresponds to children or experimental partners. The caregiver is using their beliefs ( $b$ ) about their child’s true vocabulary ( $v$ ) to plan their next communicative act. Teaching by ostensive labeling is costly in the moment, but has utility because it enables a plan to speak in the future.*

At the start of each trial, the model makes a plan consider the actions it will take on this trial and future trials. If the model speaks—producing the label for the target object, two outcomes are possible. First, its partner could select the correct referent. This would give maximal utility on the current trial, and also cause the model to know that its partner knows the right label for this object. In that case, the model could speak again on future trials and expect to succeed and be rewarded. On the other hand, its partner could select the wrong referent. This could happen either because the model itself does not know the right label, or because its partner does not. In that case, the model would get low utility on this round, and would likely teach or point for this referent on subsequent trials. Alternatively the model could point. This would give some utility on the current trial, and would leave the model in the same state of uncertainty about whether its partner knows the correct label for the referent on subsequent trials. Finally the model could teach. This would lead to very little utility on the current trial, but would cause the model to know that it can speak to refer to this referent on future trials (as in Figure 6). Reasoning forward about however many trials are left to play for this referent, the model makes a plan with the appropriate number of steps. At the start of the game, the model knows it will play three times for each referent so it might make a plan like  $\{speak, speak, speak\}$

After formulating this plan, the model will take the first action in the plan sequence (e.g. *speak*). It will then observe its partner’s behavior. In this case, suppose that its partner selects the incorrect referent. The model will then update its beliefs—its partner must not know the correct label for the target object. The next time it needs to communicate about the same object, it will be very unlikely to plan to speak, even though its previous plan was to do so. This is because the model’s belief about the world has changed and now it will be more likely to  $\{point, point\}$  or to  $\{teach, speak\}$ .

We next formally specify each step in the cycle and finally define how people form initial beliefs about their partner’s language. All code for implementing the model is available on the Open Science Foundation project page associated with this paper.

638 ***Plan***

639 When people plan, they compute the expected utility of each possible action ( $a$ ) by  
 640 combining the expected utility of that action now with the Discounted Expected Utility  
 641 they will get in all future actions. The amount of discounting ( $\gamma$ ) reflects how much people  
 642 care about success now compared to success in the future. Because utilities depend on the  
 643 communicative partner’s vocabulary, people should integrate over all possible vocabularies  
 644 in proportion to the probability that their belief assigns to that vocabulary ( $\mathbb{E}_{v \sim b}$ ).

$$EU[a|b] = \mathbb{E}_{v \sim b} (U(a|v) + \gamma \mathbb{E}_{v', o', a'} (EU[a'|b'])) \quad (1)$$

645 ***Act***

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

$$P(a|b) \propto \alpha e^{EU[a|b]} \quad (2)$$

652 ***Update beliefs***

653 After taking an action, people observe ( $o$ ) their partner’s choice—sometimes they  
 654 correctly select the intended object, and sometimes they do not. People then update their  
 655 beliefs about the partner’s vocabulary based on this observation. For simplicity, we assume  
 656 that people think their partner should always select the correct target if they point to it, or  
 657 if they teach, and similarly should always select the correct target if they produce its label

and the label is in their partner’s vocabulary. Otherwise, they assume that their partner will select the wrong object. People could of course have more complex inferential rules, e.g., assuming that if their partner does know a word they will choose among the set of objects whose labels they do not know (mutual exclusivity, Markman & Wachtel, 1988). Empirically, however, our simple model appears to accord well with people’s behavior.

$$b'(v') \propto P(o|v', a) \sum_{v \in V} P(v'|v, a) b(v) \quad (3)$$

The critical feature of a repeated communication game is that people can change their partner’s vocabulary. In teaching, people pay the cost of both talking and pointing together, but can leverage their partner’s new knowledge on future trials. Note here that teaching has an upfront cost and the only benefit to be gained comes from using less costly 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.

$$P(v'|v, a) = \begin{cases} 1 & \text{if } v_w \in v \& v' \mid a = \text{point+talk} \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

### *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 WebPPL—a programming language designed for specifying probabilistic models (Goodman & Stuhlmüller, 2014). We began with the POMDP specification developed by Evans et al. (2017). 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 (i.e., *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 its rationality ( $\alpha$ ). In order to determine the values of these parameters that best characterize human participants, we used empirical Bayesian inference to estimate the posterior means of both. Using posterior mean estimates rather than the maximum likelihood estimates naturally penalizes models for their ability to predict patterns of data that were not observed, applying a kind of Bayesian Occam’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 hyper-parameters forming the priors for individual participants’ parameters. This hierarchical estimation process achieves the same partial pooling as as subject-level random effects in mixed-effects models, giving estimates of the group-level parameters (Gelman & Hill, 2006). Details of the estimation

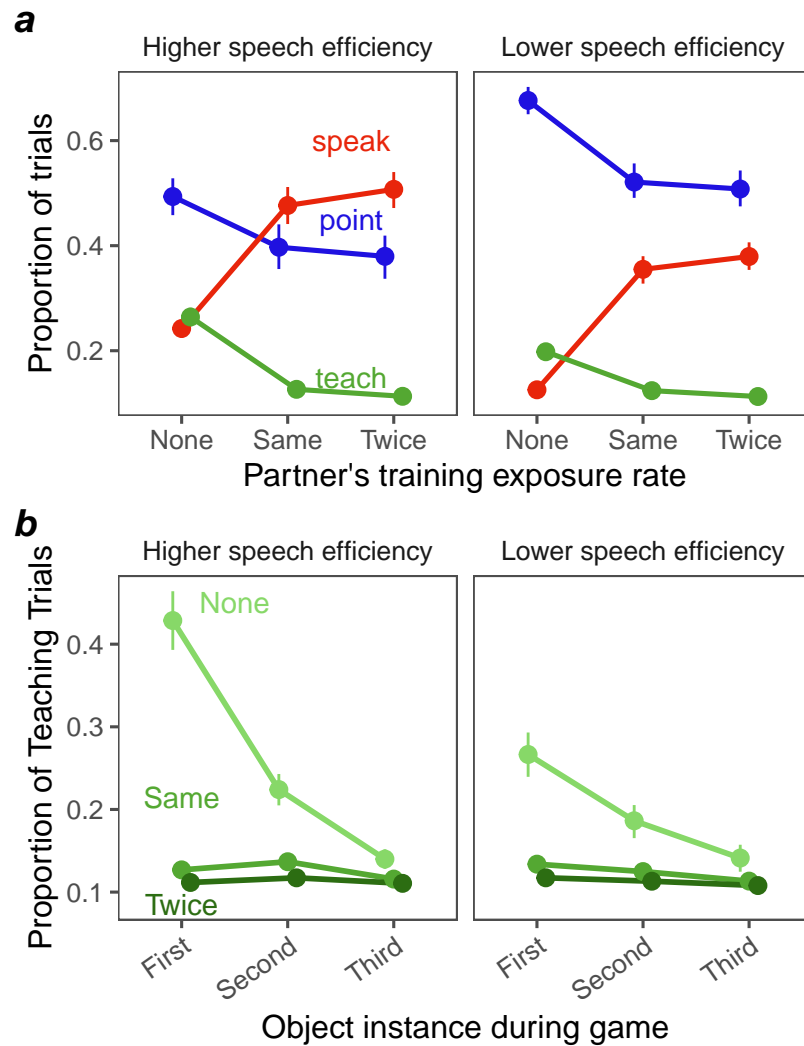


procedure can be found in the Supplemental Materials.

## 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, 2014). We estimated the posterior mean rationality ( $\alpha$ ) to be 1.33 with a 95% credible interval of [1.24, 1.42]. We did not have strong expectations for the value of the discounting parameter ( $\gamma$ ), but estimated it to be 0.42 [0.39, 0.44], suggesting that on average participants weighed the next occurrence of a referent as slightly less than half as important as the current occurrence.

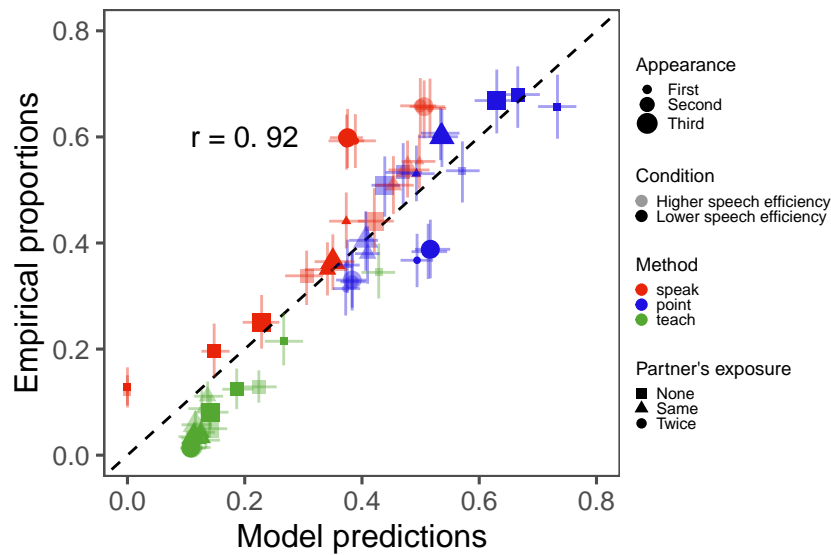
To derive predictions from the model, we ran 100 simulations of the model’s choices participant-by-participant and trial-by-trial using our posterior estimates of the hyper-parameters  $\alpha$  and  $\gamma$ . Because we did not use our participant-level parameter estimates, this underestimates the correlations between model predictions and empirical data (as it ignores variability across participants). Instead, it reflects the model’s best predictions about the results of a replication of our experiment, where individual participants’ parameters will not be known apriori. Figure 7a shows the predictions from the model in analogous format to the empirical data in Figure 4. The model correctly captures the qualitative trends in participants’ behavior: It speaks more and points less in the Higher speech efficiency condition. Figure 7b shows the model’s predicted teaching behavior in detail in an analogous format to the empirical data in Figure 5. The model again captures the qualitative trends apparent in participants’ behavior. The model teaches less knowledgeable partners, especially those who it believes have no language knowledge at all. The model teaches more when speech is relatively more efficient, and thus the future utility of teach a partner is higher. And finally the model teaches most on the first occurrence of each object, and becomes less likely to teach on future occurrences when (1) partners should be more likely to know object labels, and (2) the expected future

**Figure 7**

(a) Model prediction choice of communicative method choice as a function of exposure and the utility manipulation. (b) Model predicted probability of teaching by Partner's language knowledge and exposure rate.

rewards of teaching are smaller.

To estimate the quantitative fit between model predictions and empirical data, we compute the Pearson correlation between the model’s probability of using each action and participants’ probability of using that same action as a function of appearance, condition, and partner’s exposure. Across experimental manipulations, the model’s predictions were highly correlated with participant behavior ( $r = 0.92$  [0.86, 0.95],  $t(52) = 16.67$ ,  $p < .001$ ; Figure 8).



**Figure 8**

*Fit between model predictions and empirical data.*

Finally, we compare this model to two simpler alternative models: (1) A no-cost model in which people are indifferent to the costs of communication, and (2) a myopic model in which people do not plan for future interactions, and instead only care about the utility of their communicative choices on the immediate communicative event. We estimated parameters for these two simpler models using the same procedure as the full model: We first fit individual participant-level parameters and then estimated the posterior mean parameters for the population of participants. To compare these reduced models to our full model, we computed the log likelihood of observing the experimental data if

participants behaved according to each of the three models. These likelihoods combine both the probability of observing the empirical the data under the model, and the probability of the model parameters under the model priors. This prior probability implements a kind of Bayesian Occam’s razor, penalizing the two models which involve planning and thus fit a discounting parameter (full, no-cost) relative to the no-planning model which has only a rationality parameter. For the full model, the average likelihood across 100 runs of the model was -15,771.93. By comparison, the likelihoods for the no-cost model and myopic model were -21,016.96 and -17,257.91 respectively. Thus, the probability of observing the empirical data was thousands of times more likely under the full model than either of the simpler alternatives.

## Discussion

In both qualitative and quantitative analyses, participants’ behavior in our communication task was well explained by a model of communication as rational planning under uncertainty. The key intuition formalized by this model is that the value of a communicative acts derives from (1) the immediate effect on resolving the current communicative need, and (2) the potential benefit of the act for communication with this conversational partner in the future. Crucially, this model is able to predict a putatively altruistic behavior—teaching by ostensive labeling—without any altruistic goals at all. Because ostensive labeling can increase the efficiency of future communication, it can be beneficial even under a purely self-interested utility function. What’s more, the model correctly predicts the circumstances under which participants will engage in teaching behavior: early interactions with linguistically naive communicative partners in circumstances where language is a relatively efficient communicative modality.

Importantly, this model does not rule out the possibility that participants in our experiment—and more broadly people in the real world—may teach because of other more altruistic mechanisms or pressures. The model simply shows that appealing to such

mechanisms is not necessary to explain the ostensive labeling observed in caregivers' conversations with their children, and by extension other behaviors that may at first blush appear to be pedagogically motivated. By the same logic, the model predicts that there should be other pedagogically supportive behaviors in the interactions between caregivers and their children, and likely in the interactions between any two communicative partners who have some expectation that they will communicate again in the future. This framework thus provides a potential explanation for the occurrence of these behaviors and a framework for understanding their impact on language learning.

Of course, not all potentially pedagogically-supportive behaviors will yield an immediate or future communicative benefit. For instance, correcting children's syntactic errors could be helpful for their language development, but unless it resolves a communicative ambiguity, it will have little impact on communicative success. Our framework would predict that these behaviors should be rare, and indeed such behaviors appear to be generally absent in children's input (Marcus, 1993). We return this issue at greater length in the General Discussion. Before turning to that, however, we first consider the consequences of this model of communication for children's language. In the next section, we use simulation methods to ask how caregivers' communicative motivation may impact their children's learning, and how this impact changes as a function of the complexity of the world and the efficacy of children's learning mechanisms.

### Consequences for Learning

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

speech and pointing, learning rate, etc.), this modulation is quantitatively well predicted by a rational model of planning under uncertainty about a partner’s vocabulary.

In this final section, we take up the consequences of communicatively-motivated linguistic input for a child learning language. To do this, we adapt a framework used by Blythe et al. (2010) to estimate the learning times for an idealized child learning language under different models of the child and their caregiver. We derive estimates by simulating exposure to successive communicative events, and measuring the probability that successful learning happens after each event. The question of how different caregiver goals during caregiver-child impact children’s learning can then be formalized as a question of how much more quickly learning happens under one simulation of child and caregiver model than another.

We consider children’s word learning in three different simulated language environments:

1. *Teaching* - The caregiver’s goal in each interaction is to maximize their child’s learning (by teaching on every trial). This goal is equivalent to a model in which the goal is to maximize communicative success without minimizing communicative cost. In this case, we model the child’s learning with the same simple binomial model we proposed for participants in the experiment. After every teaching event, if the child has not yet learned the label, they do so with a probability defined by their learning rate.
2. *Communication* - The caregiver’s goal in each interaction with their child is to maximize their communicative success while minimizing their communicative cost. This the model described in the Model section above. In this case, we model the child’s learning with the same simple binomial model whenever the caregiver produces a teaching event. If the caregiver instead points or speaks, the child cannot learn at all. We make this assumption for simplicity—a fairer simulation would also

include the child’s learning from co-occurrence statistics of speaking events as in the model above. However, simultaneously simulating the child’s learning from ambiguous input and the caregivers’ reasoning about their child’s likelihood of having learned from ambiguous input is computationally challenging. This model thus underestimates the rate at which children would learn from communication and should be considered a lower bound.

3. *Talking* - The caregiver’s goal in each interaction is to refer to their intended referent so that a knowledgeable listener would understand them, without accounting for the child’s language knowledge. This goal is equivalent to minimizing communicative cost without maximizing communicative success. Because the caregiver never teaches under this model, the child must learn from the co-occurrence statistics of the words they hear and objects they see. That is, the child needs to solve the cross-situational learning problem (Yu & Smith, 2007). In this case, the child’s learning is affected not only by their own ability to remember words and track statistics, but also by the ambiguity of the learning environment (i.e. how many objects are around when they hear a word).

Formalizing these models allows us to ask three questions: (1) What is the lower bound on time to learn if caregivers are motivated to teach and always engage in ostensive labeling? (2) If caregivers have a less altruistic goal—communication—how much longer would it take to learn? (3) If the child instead had to rely on learning from co-occurrence statistics, how powerful a statistical learner would a child have to be in order to match the rate of learning of a very simple learner in the context of goal-motivated caregivers?

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 et al., 2009; Yu, 2008; Yurovsky et al., 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 et al., 2013). We assume independence primarily for pragmatic reasons here—it makes the simulations significantly more tractable (although it is also 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 et al., 2013). Thus, the ordering across conditions should remain unchanged. However, the magnitude of the difference across conditions could potentially increase or decrease.

## Method

In each of the sections below, we describe the join models of caregivers' communication and children's learning that predict learning times under each of the three models of caregivers' goals.

### *Teaching.*

Because the teaching caregiver is indifferent to communicative cost, they rely on ostensive labeling (pointing + speaking) for 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 assume that the learner does not forget, 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 successes 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 Model section. 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 stand-in for their child’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 a priori 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 three future events were indistinguishable from models that planned over four or more, so we simulated three steps of recursion<sup>1</sup>. Finally, to increase the speed

---

<sup>1</sup> It is an interesting 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 et al.,

of the simulations we re-implemented them in the R programming language. All other aspects of the model were identical.

In our simulations, we varied the children’s learning rate ( $p$ ) from .1 to 1 in steps of .1 as in the Teaching simulation, caregivers’ future-weighting ( $\gamma$ ) from .1 to 1 in steps of .1, the caregivers’ rationality ( $\alpha$ ) from .5 to 3 in steps of .5, and considered three values each of the cost of speaking ( $S = (0, 10, 20)$ ) and pointing ( $P = (50, 60, 70)$ ). The utility of communicating successfully was always 100.

### *Talking.*

When caregivers are producing a label for one of the objects in the environment, but the child does not know which one, learning requires tracking how often that words occurs with each potential object. This learning problem has been studied extensively in the language acquisition literature under the guise of “cross-situational learning” (Yu & Smith, 2007). Models of cross-situational learning have taken a variety of forms in order to instantiate different theoretical positions about the mechanisms involved in learning, most centrally whether learning is hypothesis-driven or associative (see Yu & Smith, 2012 for a review). In our analyses, we do not attempt to distinguish among these classes of models as has been done in other learning-time simulations (e.g. Smith et al., 2011). Because of its natural alignment with the learning models we use in the Teaching and Communication simulations, we implemented a simple positive hypothesis testing model. Our choice to focus on hypothesis testing rather than other learning frameworks is purely a pragmatic choice—the learning parameter  $p$  in this model 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.

In this model, learners begin with no hypotheses and add new ones to their store as

---

2015). This future work is outside the scope of the current project.

they encounter data. Upon first encountering a word and a set of objects, the model encodes up to  $h$  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 Talking 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 event, we sampled the correct target as well as  $(C-1)$  competitors from a total set of  $M$  objects. We then simulated 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  $M = (8, 16, 32, 64, 128)$  total objects,  $C = (1, 2, 4, 8)$  objects per trial,  $h = (1, 2, 3, 4)$  concurrent hypotheses, as the child’s learning rate  $p$  varied from .1 to 1 in increments of .1.

## Results

In order to understand how learning rates vary with model parameters, we first discuss the dependence of each of the three tested models on its parameters, and then

944 discuss relationships between the models. For clarity of exposition, we analyze the number  
945 of events required for 75% of simulated learners to acquire the target word, and plot a  
946 representative subset of parameter values.

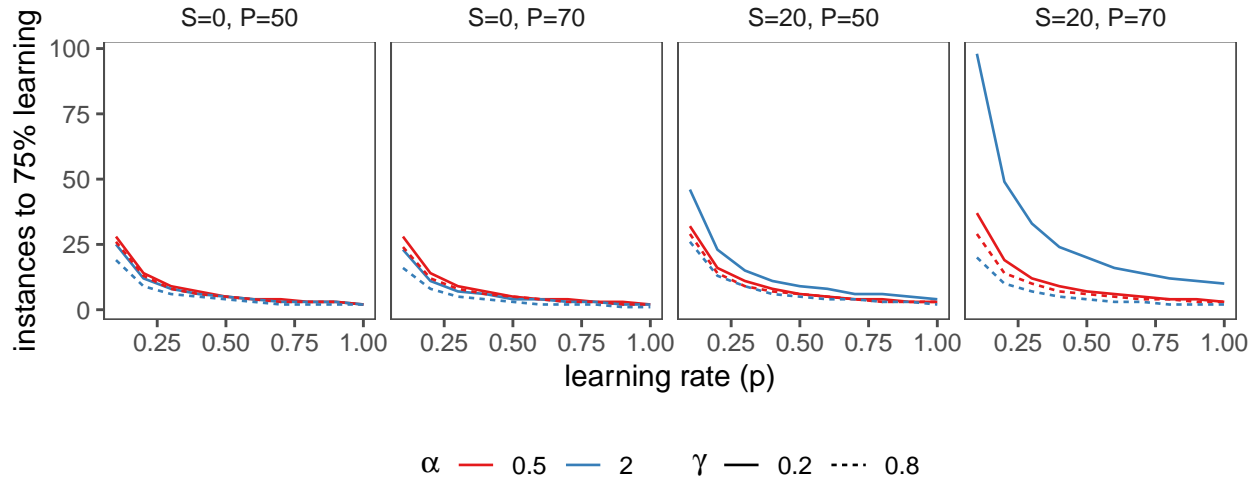
947 In addition the results reported here, we have made the full set of simulated results  
948 available in an interactive web application at [dyurovsky.shinyapps.io/ref-sims](http://dyurovsky.shinyapps.io/ref-sims). We  
949 encourage readers to fully explore the relationships among the models beyond the summary  
950 we provide.

### 951 *Teaching.*

952 Because the Teaching model behaves identically on each trial regardless of the  
953 learner, the rate of learning under this model depends entirely on the learner's learning  
954 rate  $p$ . If the learning rate was high (e.g. .8), more than 75% of learners acquired the word  
955 after a single learning instance. If the learning rate was medium, closer to the range we  
956 estimated for adult learners (.6), more than 75% of learners acquired the word after only 2  
957 instances. Finally, if the learning rate was very low (.2), the same threshold was reached  
958 after 7 instances. Thus, the model is predictably sensitive to learning rate, but even very  
959 slow learners are expected to acquire words after a small number of communicative events.

### 960 *Communication*

961 The Communication model's behavior depends on parameters of both the child  
962 learner and the caregiver communicator. In general, parameters of both participants had  
963 predictable effects on learning: Children learned faster when they had higher learning  
964 rates, when caregivers were more rational, and when caregivers gave greater weight to the  
965 future. Further, the effects of caregivers' parameters were more pronounced at the lowest  
966 learning rates. However, as the cost of speaking increased relative to pointing, the effects of  
967 caregivers' parameters changed. In particular, highly rational caregivers who heavily  
968 discounted the future lead to significantly slower learning. At these parameter settings, the

**Figure 9**

*Number of exposures required for 75% of children to learn a word under the Communication model as parameters vary. Color shows rationality ( $\alpha$ ), Linetype shows future weighting ( $\gamma$ ), facets indicate the the cost of speaking ( $S$ ) and pointing ( $P$ ). The middle two facets corresponds to Higher Speech Efficiency and Lower Speech efficiency conditions of the experiment.*

caregiver becomes very likely to point on any given trial in order to maximize the local utility at the expense of discounted future utility gained from teaching. In addition, as the cost of both modalities increases, the utility of communicating successfully (here defined as 100 points) becomes less motivating. Thus, caregivers become less discriminating among their communicative choices. Figure 9 shows the number of trials required for 75% of learners to acquire a word as a function of parameters in the Communication model.

### ***Talking.***

Finally, when caregivers spoke on each trial and children had to learn from cross-situational statistics, learning was controlled by the the child's learning rate, the number of hypotheses the child could entertain, the number of objects per event, and to a small extent the total vocabulary size. In general, children learned faster when they had a higher learning rate, and could entertain more hypotheses. Learning was also predictably slower when there were more objects on each event and thus ambiguity was higher. Finally,

as the total vocabulary size increased, the rate of learning increased slightly, as it does with human cross-situational learners (Yu & Smith, 2007). This counter-intuitive outcome occurs because the rate of spurious co-occurrences, in which the target word consistently co-occurs with an object that is not its referent, decreases as the set of potential foils expands. The effect of context size ( $C$ ) and number of hypotheses can be seen along with the learning rates of the other two models in Figure 10.

## Comparing the Models

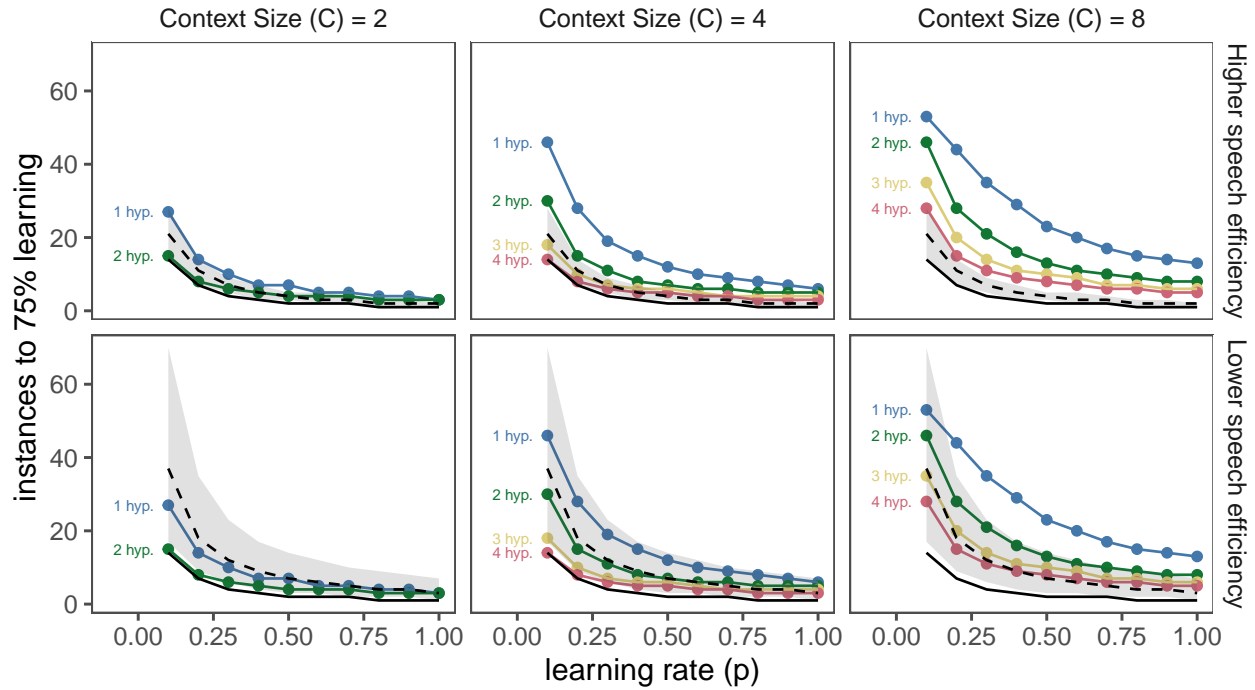
Because the real-world parameters appropriate for each model are difficult to determine, we consider the relationship between the models over the range of their possible parameters. Figure 10 shows the time for 75% of learners to acquire a word in each of the three models. Across all possible child learning rates ( $p$ ), the Teaching model lead to the fastest learning as expected. We can treat this model as a lower bound how quickly learning could possibly happen.

For the Communication model, we considered the range of all possible rates of learning that could unfold as the parameters of both child and caregiver varied. The range was substantial. If caregivers weigh the future near equally to the present, and are highly rational, the child's resultant rate of learning is nearly identical to the rate of learning under the Teaching model: Children required 1.07 times as many learning instances under the Communication model as the Teaching model when averaging over all child learning rates. In contrast, if the caregiver weighs the future much less than the present, and is relatively irrational about maximizing utility, the rate of learning can be quite slow—in the worst case requiring children to have 24.30 as many learning instances as under the Teaching model. Despite this worst case scenario, if caregivers' parameters are close to the ones we estimated in our experiment, Communication would require only 1.75 as many instances as Teaching if speech is high efficiency relative to pointing, and 3.12 as many instances if speech is lower efficiency.

For the Talking model, we also observed a wide range of learning times as a function of both the ambiguity of the learning environment and the number of simultaneous hypotheses that the child can maintain. When the environment was unambiguous—only 2 objects were present at a time—and the child could encode both, learning under Talking took only 2.03 times as many instances as Teaching. In contrast, if ambiguity was high, and learners could only track a single hypothesis, learning was significantly slower under Talking than Teaching, (requiring 10.05 times as many instances).

Comparing Communication and Talking to each-other, we find that that Talking can lead to faster learning under some parameter settings. In particular, if events are low in ambiguity, or children can maintain a very large number of hypotheses about the meaning of a word relative the number of objects in each event, children can learn rapidly even if caregivers are just Talking. This learning can be faster than simpler child models learning from highly myopic or relatively irrational caregivers Communicating, especially if speech is high-cost. At medium levels of ambiguity, Communication and Talking are similar and their ordering depends on other parameters. At high levels of ambiguity Communication is the clear winner.

Together, these results suggest that if the set of possible candidate referents is small, even simple cross-situational learners can cope just fine even if their caregiver is just Talking; they learn roughly two to three times more slowly than if their caregiver was Teaching them. However, if the set of possible referents is four, or, eight, or even more on average, cross-situational learners need to have very high bandwidth or their rates of learning will be an order of magnitude slower than if their caregiver were Teaching them. In these cases, even the simplest learner—who can encode a single hypothesis about the meaning of a word and gets no information from co-occurrence statistics—can learn quite rapidly if they are learning from a caregiver that Communicates with them.

**Figure 10**

Comparing the number of exposures required for 75% of children to learn a word under all three models as parameters vary. Columns show variation in context size ( $C$ ), a parameter of the Talking model. Rows show the two variations in the costs of Speech and Pointing for the Communication model used in our experiments. In each facet, the solid black line shows learning under the Teaching model, the light gray region shows an envelope of learning times corresponding to all variations in Communication model parameters, and the black dotted line shows learning time under the Communication model with parameters equal to the empirical estimates from experiments. Colored lines show learning times under the Talking model with varying numbers of hypotheses. Because there was little effect of the total number of objects ( $M$ ) in the Talking model, all panels show results for 128 objects. Note that Communication model parameters vary across rows, while Talking model parameters vary across columns.



## Discussion

Most of the language that children hear from their caregivers is unlikely to be designed to teach them language. However, the language that caregivers direct to them *is* designed to communicate successfully. Here we consider the learning consequences of these differences in design. How different are the learning consequences of language designed for teaching, language designed for communication, and ambient language not designed for the child at all?

If input is not communicatively motivated, the rate of learning depends entirely on what the learner brings to the table. In line with prior analyses of cross-situational learning, we find that learning can be quite rapid if environments are low in ambiguity or the learner has very high bandwidth for storing candidate hypotheses (Smith et al., 2011; Yu & Smith, 2012). However, the child’s environment is neither guaranteed to be unambiguous nor are young children likely to have high bandwidth for statistical information (Medina et al., 2011; Vlach & Johnson, 2013; Woodard et al., 2016). In fact, when the set of candidate referents is small, it is quite likely to be small in part because caregivers have designed the context to support communication (Tomasello & Farrar, 1986).

Learning from communication consistently outperforms learning from ambient language for all but the most precocious learners. If we take learning from teaching as an upperbound, we see that the rate of learning from communication is almost as fast under many possible parameter settings we explored. On average, across all possible parameter values, learning from communication is only 2.5 times slower than learning from teaching. Further, in this model, the learner gets no information from co-occurrence statistics at all. Combining learning from communication with low-bandwidth cross-situational learning could bring the expected rate of learning down to very close to learning from teaching (MacDonald et al., 2017).

While these simulations cannot directly speak to empirical differences in children’s

language learning, they provide an important proof-of-concept predicting that such learning differences should arise. We thus might make significant progress on understanding how children learn language so quickly not just by studying children, but also by understanding how caregivers design the language they produce in order to support successful communication (Leung et al., 2021).

## General Discussion

Across naturalistic corpus data, experimental data, and model predictions and simulation, we see evidence that pressure to communicate successfully with a linguistically immature partner could fundamentally structure caregiver production and shape child learning. 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 caregiver 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 simple learners whose caregivers want to communicate with them out-learn more powerful statistical learners whose caregivers do not have a communicative goal.

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, and many more (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. 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; our account’s emphasis on calibration to such learning mechanisms suggests the role of ordered or incremental input from the environment may be crucial. Our findings support the idea that rapid language learning may be facilitated by the *combination* of the learner’s limited statistical learning skills combined with communicatively (but not pedagogically) motivated caregiver input. Such results emphasize the importance of a dyadic learning approach, whereby considering the joint contributions of learner and caregiver can yield new insights (Yurovsky, 2018).

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. 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. 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 et al., 2018). Frameworks of visual learning must incorporate both the

relevant learning abilities and this motivated, contingent structure in the environment.

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 adapting 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 must explain caregiver production as well as learning itself. Recent work has shown that many aspects of the the structure of natural languages, both in semantics and in syntax, can be explained by evolutionary optimization of language to be more efficient for communication (see Gibson et al., 2019 for a review). Further, people talking with other native language speakers make in-the-moment word choices that are well predicted by models of optimal communication (e.g., Mahowald et al., 2013). Here we extend these ideas to asymmetric communications between speakers and less knowledgeable listeners, and show that they can predict how caregivers modify their communicative acts when talking with their children.

Explaining caregiver 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, caregiver 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 et al., 2016; although c.f. Shneidman & Goldin-Meadow, 2012). Our argument is that the rate and ultimate attainment of language learners will vary substantially as a function of caregiver modification, and that describing the cause of this variability is a necessary feature of models of language learning.

Our account aims to explain caregivers' 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. Some such phenomena will be easily accounted for; aspects of language that shape communicative efficiency should shift in predictable patterns across development. For example, the exaggerated pitch contours seen in infant-directed speech serve to draw infants' attention and facilitate phoneme learning. These language modifications are well-explained by our proposed framework, though incorporating them will likely require altering aspects of our account and decisions about which alterations are most appropriate. In the example of exaggerated pitch, one could expand the definition of communicative success to include the goal of maintaining attention, or 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 alterations to the extant structure of the framework.

### *Limitations*

Our account is formulated primarily around referential communication and future work must address its viability in other aspects of language learning. Of course, not all aspects of language should be calibrated to the child's language development. Communication goals and learning goals are not always aligned, and thus aspects of language that minimally affect communicative efficiency are unlikely to show a pattern of calibration. Indeed, early explorations of possible calibration in child-directed speech focused on syntax and found little evidence of calibration (Newport et al., 1977). Furthermore, some aspects of caregiver production are unrepresented in our framework, such as aspects of production driven by speaker-side constraints.

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 (Hoff, 2003). This is to be expected to the extent that caregiver-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. Indeed, child-directed speech itself varies cross-linguistically, both in its features (Fernald et al., 1989) and quantity (e.g., Shneidman & Goldin-Meadow, 2012)—although, there is some evidence that child-directed speech predicts learning even in cultures where it 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 learning for a given learning problem, 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 caregiver 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 on early functional accounts of language learning, our perspective considers the caregiver-child dyad as the fundamental unit of analysis and emphasizes the

1191 importance of communicative success in shaping language input and language learning. We  
1192 have developed an initial formal account for jointly considering caregiver productions and  
1193 child language learning within the same system.

1194       Such an account helps to explain caregivers' increased production of ostensive  
1195 labeling in naturalistic communication for infrequent referents and when children are very  
1196 young. In an experimental context, pressure to with a less knowledgeable partner results in  
1197 analogous behavior from participants in an iterated reference game. A rational  
1198 communication model that includes planning can explain these behaviors without an  
1199 explicit teaching goal, and such a model produces learning outcomes that outperform all  
1200 but the best cross-situational learners in simulations. In sum, this work demonstrates that  
1201 the pressure to communicate successfully across knowledge asymmetries may help create a  
1202 learning environment that fosters language learning. Rapid language learning in early  
1203 childhood may be best explained by considering that the child's powerful, though not  
1204 precocious, learning mechanisms are met with a language environment specifically designed  
1205 to communicate successfully with them.

### 1206                                   **Acknowledgement**

1207       The authors are grateful to Madeline Meyers for her work coding referential  
1208 communication in the corpus data, and to Mike Frank and Chen Yu for their feedback on  
1209 the manuscript. This research was funded by James S. McDonnell Foundation Scholar  
1210 Award in Understanding Human Cognition #220020506 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.
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48.  
<https://doi.org/10.18637/jss.v067.i01>
- 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.
- Bohn, M., & Frank, M. C. (2019). The pervasive role of pragmatics in early language. *Annual Review of Developmental Psychology*, 1, 223–249.
- Brand, R. J., Baldwin, D. A., & Ashburn, L. A. (2002). Evidence for “motionese”: Modifications in mothers’ infant-directed action. *Developmental Science*, 5(1), 72–83.
- Brown, R. (1977). Introduction. In C. E. Snow & C. A. Ferguson (Eds.), *Talking to children: Language input and interaction*. Cambridge, MA.: MIT Press.
- Callanan, M. A. (1985). How parents label objects for young children: The role of input in the acquisition of category hierarchies. *Child Development*, 508–523.
- Chomsky, N. (1965). *Aspects of the theory of syntax*. MA: MIT Press.
- Clark, H. H. (1996). Using language. In *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.



- 1235 DeCasper, A. J., & Spence, M. J. (1986). Prenatal maternal speech influences newborns'  
1236 perception of speech sounds. *Infant Behavior and Development*, 9(2), 133–150.
- 1237 Eaves Jr, B. S., Feldman, N. H., Griffiths, T. L., & Shafto, P. (2016). Infant-directed  
1238 speech is consistent with teaching. *Psychological Review*, 123(6), 758.
- 1239 Elman, J. L. (1993). Learning and development in neural networks: The importance of  
1240 starting small. *Cognition*, 48(1), 71–99.
- 1241 Estes, K. G., Evans, J. L., Alibali, M. W., & Saffran, J. R. (2007). Can infants map  
1242 meaning to newly segmented words? Statistical segmentation and word learning.  
1243 *Psychological Science*, 18(3), 254–260.
- 1244 Evans, O., Stuhlmüller, A., Salvatier, J., & Filan, D. (2017). *Modeling Agents with*  
1245 *Probabilistic Programs*. <http://agentmodels.org>.
- 1246 Fausey, C. M., Jayaraman, S., & Smith, L. B. (2016). From faces to hands: Changing  
1247 visual input in the first two years. *Cognition*, 152, 101–107.
- 1248 Fernald, A., Taeschner, T., Dunn, J., Papousek, M., Boysson-Bardies, B. de, & Fukui, I.  
1249 (1989). A cross-language study of prosodic modifications in mothers' and fathers'  
1250 speech to preverbal infants. *Journal of Child Language*, 16(3), 477–501.
- 1251 Foushee, R., Griffiths, T. L., & Srinivasan, M. (2016). *Lexical complexity of child-directed*  
1252 *and overheard speech: Implications for learning*.
- 1253 Frank, M. C., & Goodman, N. D. (2012). Predicting pragmatic reasoning in language  
1254 games. *Science*, 336, 998–998.
- 1255 Frank, M. C., & Goodman, N. D. (2014). Inferring word meanings by assuming that  
1256 speakers are informative. *Cognitive Psychology*, 75, 80–96.
- 1257 Frank, M. C., Goodman, N., & Tenenbaum, J. (2009). Using speakers' referential intentions  
1258 to model early cross-situational word learning. *Psychological Science*, 20, 578–585.
- 1259 Franke, M. (2013). Game theoretic pragmatics. *Philosophy Compass*, 8(3), 269–284.

- 1260 Gelman, A., & Hill, J. (2006). *Data analysis using regression and multilevel/hierarchical*  
1261 *models*. Cambridge university press.
- 1262 Gibson, E., Futrell, R., Piantadosi, S. P., Dautriche, I., Mahowald, K., Bergen, L., & Levy,  
1263 R. (2019). How efficiency shapes human language. *Trends in Cognitive Sciences*,  
1264 *23*(5), 389–407.
- 1265 Gogate, L. J., Bahrick, L. E., & Watson, J. D. (2000). A study of multimodal motherese:  
1266 The role of temporal synchrony between verbal labels and gestures. *Child*  
1267 *Development*, *71*(4), 878–894.
- 1268 Goldin-Meadow, S., Levine, S. C., Hedges, L. V., Huttenlocher, J., Raudenbush, S. W., &  
1269 Small, S. L. (2014). New evidence about language and cognitive development based  
1270 on a longitudinal study: Hypotheses for intervention. *American Psychologist*, *69*(6),  
1271 588–599.
- 1272 Goldowsky, B. N., & Newport, E. L. (1993). Limitations on the acquisition of morphology:  
1273 The less is more hypothesis. *The Proceedings of the Twenty-Fourth Annual Child*  
1274 *Language Research Forum*, 124.
- 1275 Gomez, R. L., & Gerken, L. (1999). Artificial grammar learning by 1-year-olds leads to  
1276 specific and abstract knowledge. *Cognition*, *70*(2), 109–135.
- 1277 Goodman, N. D., & Frank, M. C. (2016). Pragmatic language interpretation as  
1278 probabilistic inference. *Trends in Cognitive Sciences*, *20*(11), 818–829.
- 1279 Goodman, N. D., & Stuhlmüller, A. (2014). *The Design and Implementation of*  
1280 *Probabilistic Programming Languages*. <http://dippl.org>.
- 1281 Grice, H. P. (1975). Logic and conversation. In P. Cole & J. L. Morgan (Eds.), *Syntax and*  
1282 *semantics: Vol. 3: Speech acts* (pp. 41–58). New York, NY: Academic Press.
- 1283 Griffiths, T. L., Lieder, F., & Goodman, N. D. (2015). Rational use of cognitive resources:  
1284 Levels of analysis between the computational and the algorithmic. *Topics in*

*Cognitive Science*, 7(2), 217–229.

Hoff, E. (2003). The specificity of environmental influence: Socioeconomic status affects early vocabulary development via maternal speech. *Child Development*, 74(5), 1368–1378.

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. (2021). Parents fine-tune their speech to children’s vocabulary knowledge. *Psychological Science*, 32(7), 975–984.

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

MacDonald, K., Yurovsky, D., & Frank, M. C. (2017). Social cues modulate the representations underlying cross-situational learning. *Cognitive Psychology*, 94, 67–84.

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

Mahowald, K., Fedorenko, E., Piantadosi, S. T., & Gibson, E. (2013). Info/information theory: Speakers choose shorter words in predictive contexts. *Cognition*, 126(2), 313–318.

Marcus, G. F. (1993). Negative evidence in language acquisition. *Cognition*, 46(1), 53–85.

Markman, E. M. (1990). Constraints children place on word meanings. *Cognitive Science*, 14(1), 57–77.

Markman, E. M., & Wachtel, G. F. (1988). Children’s use of mutual exclusivity to constrain the meanings of words. *Cognitive Psychology*, 20(2), 121–157.

Marr, D. (1982). *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. New York, NY: W. H. Freeman.

Maye, J., Werker, J. F., & Gerken, L. (2002). Infant sensitivity to distributional information can affect phonetic discrimination. *Cognition*, 82(3), B101–B111.

- 1310 McMurray, B. (2007). Defusing the childhood vocabulary explosion. *Science*, 317(5838),  
1311 631–631.
- 1312 McMurray, B., Kovack-Lesh, K. A., Goodwin, D., & McEchron, W. (2013). Infant directed  
1313 speech and the development of speech perception: Enhancing development or an  
1314 unintended consequence? *Cognition*, 129(2), 362–378.
- 1315 Medina, T. N., Snedeker, J., Trueswell, J. C., & Gleitman, L. R. (2011). How words can  
1316 and cannot be learned by observation. *Proceedings of the National Academy of*  
1317 *Sciences*, 108(22), 9014–9019.
- 1318 Meltzoff, A. N. (2005). Imitation and other minds: The “like me” hypothesis. *Perspectives*  
1319 *on Imitation: From Neuroscience to Social Science*, 2, 55–77.
- 1320 Mintz, T. H. (2003). Frequent frames as a cue for grammatical categories in child directed  
1321 speech. *Cognition*, 90(1), 91–117.
- 1322 Newport, E. L. (1990). Maturation constraints on language learning. *Cognitive Science*,  
1323 14(1), 11–28.
- 1324 Newport, E. L., Gleitman, H., & Gleitman, L. R. (1977). Mother, I’d rather do it myself:  
1325 Some effects and non-effects of maternal speech style. In C. A. Ferguson (Ed.),  
1326 *Talking to children language input and interaction* (pp. 109–149). Cambridge  
1327 University Press.
- 1328 Ninio, A. (1980). Ostensive definition in vocabulary teaching. *Journal of Child Language*,  
1329 7(3), 565–573.
- 1330 Quine, W. V. O. (1960). Word and object. In *Cambridge, Mass.* Cambridge, Mass.: MIT  
1331 Press.
- 1332 Rafferty, A. N., Brunskill, E., Griffiths, T. L., & Shafto, P. (2016). Faster teaching via  
1333 pomdp planning. *Cognitive Science*, 40(6), 1290–1332.
- 1334 R Core Team. (2021). *R: A language and environment for statistical computing*. R

Foundation for Statistical Computing. <https://www.R-project.org/>

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.

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,

1360 25–49.

1361 Snow, C. E. (1972). Mothers' speech to children learning language. *Child Development*, 43,  
1362 549–565.

1363 Snow, C. E. (1977). Mothers' speech research: From input to interaction. *Talking to*  
1364 *Children: Language Input and Acquisition*, 3149.

1365 Stevens, J. S., Gleitman, L. R., Trueswell, J. C., & Yang, C. (2017). The pursuit of word  
1366 meanings. *Cognitive Science*, 41, 638–676.

1367 Thiessen, E., Hill, E. A., & Saffran, J. R. (2005). Infant-directed speech facilitates word  
1368 segmentation. *Infancy*, 7, 53–71.

1369 Tomasello, M., & Farrar, M. J. (1986). Joint attention and early language. *Child*  
1370 *Development*, 57, 1454–1463.

1371 Townsend, J. T. (1990). Serial vs. Parallel processing: Sometimes they look like  
1372 tweedledum and tweedledee but they can (and should) be distinguished.  
1373 *Psychological Science*, 1(1), 46–54.

1374 Trueswell, J. C., Medina, T. N., Hafri, A., & Gleitman, L. R. (2013). Propose but verify:  
1375 Fast mapping meets cross-situational word learning. *Cognitive Psychology*, 66(1),  
1376 126–156.

1377 Vlach, H. A., & Johnson, S. P. (2013). Memory constraints on infants' cross-situational  
1378 statistical learning. *Cognition*, 127(3), 375–382.

1379 Vogt, P. (2012). Exploring the robustness of cross-situational learning under zipfian  
1380 distributions. *Cognitive Science*, 36(4), 726–739.

1381 Wittgenstein, L. (1953). *Philosophical investigations*. Oxford, UK: Basil Blackwell.

1382 Woodard, K., Gleitman, L. R., & Trueswell, J. C. (2016). Two-and three-year-olds track a  
1383 single meaning during word learning: Evidence for propose-but-verify. *Language*  
1384 *Learning and Development*, 12(3), 252–261.

- 1385 Xu, F., & Tenenbaum, J. B. (2007). Word learning as Bayesian inference. *Psychological*  
1386 *Review*, 114(2), 245–272.
- 1387 Yu, C. (2008). A statistical associative account of vocabulary growth in early word  
1388 learning. *Language Learning and Development*, 4(1), 32–62.
- 1389 Yu, C., & Smith, L. B. (2007). Rapid word learning under uncertainty via cross-situational  
1390 statistics. *Psychological Science*, 18, 414–420.
- 1391 Yu, C., & Smith, L. B. (2012). Modeling cross-situational word-referent learning: Prior  
1392 questions. *Psychological Review*, 119, 21–39.
- 1393 Yurovsky, D. (2018). A communicative approach to early word learning. *New Ideas in*  
1394 *Psychology*, 50, 73–79.
- 1395 Yurovsky, D., Doyle, G., & Frank, M. C. (2016). Linguistic input is tuned to children’s  
1396 developmental level. *Proceedings of the Annual Meeting of the Cognitive Science*  
1397 *Society*, 2093–2098.
- 1398 Yurovsky, D., & Frank, M. C. (2015). An integrative account of constraints on  
1399 cross-situational learning. *Cognition*, 145, 53–62.
- 1400 Yurovsky, D., Fricker, D. C., Yu, C., & Smith, L. B. (2014). The role of partial knowledge  
1401 in statistical word learning. *Psychonomic Bulletin & Review*, 21, 1–22.
- 1402 Yurovsky, D., Meyers, M., Burke, N., & Goldin-Meadow, S. (2018). Children gesture when  
1403 speech is slow to come. In J. Z. Chuck Kalish Martina Rau & T. Rogers (Eds.),  
1404 *Proceedings of the 40th annual meeting of the cognitive science society* (pp.  
1405 2765–2770).
- 1406 Yurovsky, D., Yu, C., & Smith, L. B. (2012). Statistical speech segmentation and word  
1407 learning in parallel: Scaffolding from child-directed speech. *Frontiers in Psychology*,  
1408 3, 374.
- 1409 Yurovsky, D., Yu, C., & Smith, L. B. (2013). Competitive processes in cross-situational

1410

word learning. *Cognitive Science*, 37, 891–921.