

# Child-directed Listening: How Caregiver Inference Enables Children’s Early Verbal Communication

Stephan C. Meylan<sup>1,3</sup>, Ruthe Foushee<sup>2</sup>, Erika Bergelson<sup>3</sup>, and Roger P. Levy<sup>1</sup>

<sup>1</sup>Department of Brain and Cognitive Sciences, MIT ({smeylan, rplevy}@mit.edu)

<sup>2</sup>Department of Psychology, University of Chicago (foushee@uchicago.edu)

<sup>3</sup>Department of Psychology and Neuroscience, Duke University (elika.bergelson@duke.edu)

February 10, 2021

## Abstract

How do adults understand children’s speech? Children’s productions over the course of language development often bear little resemblance to typical adult pronunciations, yet caregivers nonetheless reliably recover meaning from them. Here, we employ a suite of Bayesian models of spoken word recognition to understand how adults overcome the noisiness of child language, showing that communicative success between children and adults relies heavily on adult inferential processes. By evaluating competing models on phonetically-annotated corpora, we show that adults’ recovered meanings are best predicted by prior expectations fitted specifically to the child language environment, rather than to typical adult-adult language. After quantifying the contribution of this “child-directed listening” over developmental time, we discuss the consequences for theories of language acquisition, as well as the implications for commonly-used methods for assessing children’s linguistic proficiency.

**Keywords:** language development, child-directed speech, noisy channel communication, spoken word recognition, Bayesian inference

## 1 Introduction

The past five decades have seen extensive research dedicated to characterizing how adults speak to infants and young children (Snow & Ferguson, 1977; Soderstrom, 2007), and to investigating the degree to which adults’ *child-directed speech* directly supports language learning (Golinkoff et al., 2015). By contrast, how caregivers understand the communicative acts of young children — *child-directed listening* (CDL) — has received far less attention. In this paper, we investigate how English-speaking adults interpret English-learning children’s verbal productions, making meaning out

of vocalizations that are often perceptually distant from targets in the adult language (e.g., /wid/ for *read*; see Table 1A).

This characterization of adults’ role in conversations with young learners dovetails with “noisy-channel” accounts of spoken language interpretation, which provide a framework for describing how listeners overcome imperfect acoustic information, verbal ambiguity, distractions, and speaker variability present in everyday conversation (Levy, 2008; Shannon, 1951; Gibson et al., 2013). To recover meanings from highly noisy input, adult listeners rely on their expectations about what speakers are likely to say, combined with the perceptual similarity between what the listener heard and guesses as to what the speaker might intend. We argue that child language represents a “noisier-than-usual” channel, where adults must use expectations fitted to the child language environment to recover meaning from child productions. That is, while hearing /wid/ might typically suggest *weed* or *wheat* as a speaker’s intended word (based solely on acoustic information), an adult caregiver might instead recover *read* as the intended word from a child speaker.<sup>1</sup>

In what follows, we seek evidence for the role of child-directed listening in language development. We present a computational framework to predict what adults are likely to recover from children’s imperfect speech, and compare it to what adults *actually* recovered. As a proxy for caregivers’ realtime interpretations, we use the orthographic annotations made by trained in-lab transcribers of spontaneous at-home child language recordings. This approach allows us to characterize the utility of adult listeners’ expectations, versus the acoustic/phonetic signal produced by the child. To capture the degree to which listening is truly *child*-directed (i.e., distinct from adult-directed listening), we compare the utility of expectations tuned on large-scale adult corpora, versus expectations tailored to reflect the child language environment.

## 2 Task and Modeling Setup

We focus here on the adult listener’s task of recovering meaning from noisy child productions. Specifically, we look at a large set of phonetically-transcribed productions (e.g., /aə wən də wɪd/ in Table 1A) from the Providence corpus (Demuth et al., 2006), and treat the challenge of inferring a word identity in context (here, an orthographic word like *read*) as a *masked word prediction* task (Devlin et al., 2019). To combine the contributions of caregiver expectations given the linguistic context with the specific sequence of phonemes produced by the child, we employ a Bayesian model of spoken word recognition in the vein of Norris & McQueen (2008), which assigns a probability to a candidate word identity  $w$  given corresponding perceptual input  $d$  in context  $c$ :

$$P(w|d, c) = \frac{P(d|w, c)P(w|c)}{\sum_{w' \in V} P(d|w', c)P(w'|c)} \quad (1)$$

This cashes out the intuition that the probability assigned to a candidate word  $w$  in spoken word recognition reflects the combination of (a) fit to perceptual data and (b) linguistic expectations. Fit to perceptual data is evaluated via a likelihood function,

<sup>1</sup>One intriguing deviation from the classic noisy-channel setup is that adults may “recover” messages when children do not intend to communicate anything at all (i.e., drawing from a noise distribution).

$P(d|w, c)$ , which reflects the probability that the word  $w$  would generate the observed data  $d$  in context  $c$ . Linguistic expectations are captured in the prior,  $P(w|c)$ , or the anticipated probability of the word in context  $c$ , absent any perceptual data. The denominator in Equation 1 reflects the summed strength of *all* competitor words  $w'$  in the candidate vocabulary  $V$ . Thus, the predictions derived from the model (a *posterior*) constitute a probability distribution over candidate words, with highly favored interpretations receiving more of the probability mass than disfavored ones.

Our principal goal is to find a model that best simulates how adults understand children. We discuss the likelihood and prior of the set of models under consideration in turn. All models used the same likelihood, derived from measures of pairwise string similarity a phonemic transcription of the child’s production and phonemic forms of all candidate words (translated into the International Phonetic Alphabet, IPA, via a dictionary of conventional English pronunciations). To illustrate, given the transcribed production /wid/, the likelihood term for the candidate word *weed* (citation phonetic form /wid/) will be higher than the likelihood term for the candidate word *read* (where the citation phonetic form /i:ɪd/ differs by one phoneme).

However, the inferential process sketched in Equation 3 foreshadows the inadequacy of the acoustic signal alone: if children often produce noisy, idiosyncratic phoneme sequences, the prior must do more “work.” The priors we evaluate take the form of probabilistic language models: computational models that return a probability distribution over word guesses, based on the surrounding linguistic context (Table 1C). When priors from each model are combined with the likelihood, they yield posterior distributions (Table 1D).

Here, we take advantage of a distinction within the transcripts of caregiver-child speech in the PhonBank database (Rose & MacWhinney, 2014), which allows us to evaluate competing models on two different dimensions. First, we evaluate models in their ability to reproduce the specific words recovered by annotators.<sup>2</sup> This analysis focuses specifically on what we term *communicative successes* (Table 1A) — instances where a phoneme sequence was not only phonemically transcribed (PhonBank %phon tier), but also received a *gloss*, or orthographic transcription. This allows us to assess the probability that each model assigns to the annotator-recovered word, with the best model being the one that assigns the highest average probability (alternatively, the lowest *surprisal*, or negative log probability<sup>3</sup>) to the glosses.

Second, we test whether models can predict when a child’s production will *not* receive a gloss (reflecting the annotator’s uncertainty as to the child’s intended word). This analysis relies on the communicative successes described above, as well as so-called *communicative failures* — instances where phoneme sequences are transcribed, but lack a gloss, due to difficulty in identifying the child’s intended word (Table 1B). In the absence of an annotator-recovered word, surprisal cannot be calculated. Instead, we measure the “peakedness” of the guesses regarding word identity by calculating the *information entropy* of the posterior distribution,

$$H(X) = - \sum_{i=1}^n P(x_i) \log P(x_i), \tag{2}$$

<sup>2</sup>We cannot know whether the word recovered by an annotator was the word intended by the child speaker.

<sup>3</sup>For statistics-oriented readers, this is the per-instance log-likelihood of the data under the model.

where  $P(x_i)$  is the probability of the  $i$ th candidate word. This provides us with a concise index of uncertainty: if posterior probability mass is centered on one or a few guesses for a given phoneme sequence, then entropy will be low; if the posterior is split across many candidate guesses, then entropy will be high. The best model under this analysis will be the one most able to discriminate failures from successes on the basis of entropy. We measure this with the *receiver operating characteristic*, or ROC, which measures the diagnostic ability of a classifier over the range of possible thresholds.

In the third analysis, we quantify how much the estimate of word identity changes as a function of 1) conditioning on context (using a fitted prior) 2) conditioning on data (the posterior when using a uniform prior), or 3) conditioning on *both* context *and* data (the posteriors reflecting the fitted priors). As a baseline for comparison, we start with a uniform prior, where all words in the vocabulary are equiprobable. We then measure the per-word average *information gain*, or Kullback-Liebler divergence, between that uniform prior distribution and each of the distributions identified above. Information gain can be interpreted as a measure of *entropy reduction*, corresponding to the difference between the uniform prior and the somewhat more peaked estimates of word identity under the fitted priors, and the (usually) yet more peaked estimates under the posteriors. If the models are using the perceptual signal to identify words, then the prior information gain will be small in comparison to the posterior information gain. If, by contrast, caregivers are relying heavily on their prior expectations, then the prior information gain will be larger with respect to the posterior information gain.

A further question is how these measures of information gain will track with developmental time. We expect prior information gain to *increase* over developmental time: as the child says more words in the surrounding context, the priors can better constrain guesses for the masked words (placing more mass on a smaller set of words, reflected in lower entropy). At the same time, as children’s productions approximate conventional pronunciations, we expect to see an increase in posterior information gain. It remains to be seen how these two quantities will interact.

### 3 Methods and Model Details

We test several language models in their ability to predict adult caregivers’ interpretations of children’s linguistic and proto-linguistic vocalizations in the Providence corpus (Demuth et al., 2006). Utterances and phonological transcripts with both phonemic and orthographic transcription were retrieved through *chldes-db 2020.1*. (Sanchez et al., 2019).

#### 3.1 Selecting Communicative Successes and Failures

We selected as communicative successes all tokens produced by children in the intersection of four criteria: (1) possessing monosyllabic IPA forms (motivated below) (2) possessing *no* unintelligible (CHILDES code *xxx*) or phonology-only (*yyy*) tokens in the same utterance (3) whose gloss is extant as a token in BERT (motivated below) (4) whose gloss is included in the Carnegie Mellon Pronunciation Dictionary (henceforth CMU dictionary). Communicative failures had to meet the first criterion, but must

Table 1: Examples of communicative success and failure, with samples from highest-ranked prior and posterior candidates.

PhonBank transcript	A. Communicative Success <sup>†</sup>		B. Communicative Failure <sup>†</sup>	
		MOT	this is	MOT
	MOT	you want mamma let’s see	CHI	no
⊗phon	CHI	/ ɑə wʌn də <span style="border: 1px solid black; padding: 0 2px;">wid*</span> /	CHI	/ ju məɪk yoʊ <span style="border: 1px solid black; padding: 0 2px;">fɛt*</span> /
gloss		<i>I want to &lt;read&gt;</i>		<i>you make your &lt;unintelligible&gt; → yyy</i>
	MOT	okay that’s fine	MOT	can I make one?
	MOT	okay mommy’s gonna pick out a book	MOT	no

Language Model	C. Best PRIOR Guesses for <span style="border: 1px solid black; padding: 0 2px;">wid / fɛt</span>	
CDL+CONTEXT <sup>‡</sup>	see (.86) look (.03) go (.02) play (.01)	own (.74) house (.01) shapes (.01) friends (.01)
BERT+CONTEXT <sup>‡</sup>	<i>read</i> (.49) see (.28) play (.04) know (.04)	own (.25) choice (.24) point (.04) bed (.03) call (.03)
CHILDES-1GRAM	I (.04) a (.03) the (.03) yeah (.03)	I (.04) a (.03) the (.03) yeah (.03) it (.02)

	D. Best POSTERIOR Guesses for <span style="border: 1px solid black; padding: 0 2px;">wid / fɛt</span>	
CDL+CONTEXT <sup>‡</sup>	see (.967) watch (.012) <i>read</i> (.005) look (.001)	own (.59) feet (.27) foot (.02) food (.01) hat (0.01)
BERT+CONTEXT <sup>‡</sup>	<i>read</i> (.61) see (.35) watch (.01) hear (.01)	bet (.31) own (.24) cut (.06) shot (.04) bed (.03)
CHILDES-1GRAM	we (.34) need (.11) and (.06) would (.04)	it (.15) that (.11) fit (.06) what (.06) feet (.05)

\* masked phoneme sequence    † MOT= Mother, CHI= Child    ‡ Model considers +/- 20 utterances of surrounding context.

have received the gloss of *yyy* (with no other *yyy* or *xxx* in the same utterance). Under these definitions, an utterance could contain several communicative successes, but at most one failure.

### 3.2 Candidate Vocabulary

The inventory of candidate words considered by each model was the intersection of (1) words in the CMU dictionary with one or two syllables and (2) tokens present in BERT (motivated below) (3) tokens that appeared 3 or more times in CHILDES (to limit to words that might reasonably be said in this context). This means that while only one-syllable phoneme sequences were analyzed, two-syllable words were also considered as possible candidate interpretations. The final inventory of candidates,  $V$ , included 7,904 words. We reconcile IPA formats following a procedure detailed in our code.

### 3.3 Priors: Language Models

For each communicative success and failure, we retrieve prior probabilities over candidate words using a suite of probabilistic language models. As a “best” prior architecture, we use BERT (Devlin et al., 2019), which has demonstrated extremely competitive performance for single-word completion tasks, including spoken word recognition (Salazar et al., 2020). By virtue of its attentional mechanisms, BERT is able to effectively model long distance dependencies (Jawahar et al., 2019), and capture speech register and discourse-level information. We compute the probabilities for the masked word  $P(w)$  from BERT, using a language modeling head with the `transformers` library (Wolf et al., 2020). For each masked phoneme sequence, we take the real-valued vector of predictions corresponding to the model’s vocabulary, extract the activations corresponding to the candidate words, and compute the softmax to yield a vector of probabilities over the candidate words (Table 1).

We test an “off-the-shelf” model of BERT trained on large quantities of (principally adult-directed) language scraped from the internet, predicting the word from the immediate utterance only (BERT+ONEUTT). We additionally test the predictions of a BERT model meant to best capture adult expectations about children’s utterances. To do this we “fine-tune” the above model on adult and child CHILDES utterance glosses, excluding PhonBank. In fine tuning, a new model is initialized with an “off-the-shelf” model, then the weights in the model are updated to best predict masks inserted into a new training set — in this case, the lines of 80% of CHILDES transcripts (20% were held out for model validation). This fine-tuned model (CDL+ONEUTT) should be expected to be more representative of adult linguistic expectations in understanding child speech than the off-the-shelf model for three reasons. First, it should assign higher probability to words that are common in speech to and from children. Second, it should assign higher probability to non-sentence fragments, which are ubiquitous in conversational speech but somewhat less prevalent in adult-directed written language. Third, it may prove capable of developing an expectation for the dyadic, back-and-forth structure of scenes typically captured in transcripts.

In addition to fine-tuning the model, we manipulate whether prior estimates reflect access to the larger discourse context as captured by the transcript before and after a phoneme sequence. In that these models are meant to be representative of *caregiver* expectations, these models condition the prediction of the masked token on what the caregiver and child *both* say, both before and after the masked token. We create priors parallel to those above by feeding the models 20 utterances preceding and following each mask (CDL+CONTEXT and BERT+CONTEXT).

BERT has its own vocabulary, which imposes limitations on the vocabulary in the analysis. Standard implementations of BERT split longer words into “word pieces”, or most common repeated sub-sequences. In English, this often yields morphological segmentation (e.g., *fishing* → *fish ##ing*), but the process is highly noisy. For the purposes of predicting a masked word, BERT predicts only one word or word piece. We limit the vocabulary to word-initial word pieces like *fish*, and exclude continuations like *##ing* from consideration. This also motivates the choice to predict monosyllabic phoneme sequences, in that the model does not allow us to predict multiple words (which might be contained in *yyy*).

In addition to the BERT models, we also test two simpler priors. The first is a simple smoothed unigram model estimated from counts in CHILDES. This model, CHILDES 1-GRAM, assigns probability to all word types proportional to their counts in the same CHILDES dataset used in the CDL models, above. To account for unseen data, we add a small pseudocount (.001) smoothing to all counts before computing probabilities. The second is the UNIFORMPRIOR model, which assigns equal probability to all words ( $1/|V|$ , where  $|V|$  is the number of candidates). This provides the comparison case of a maximally uninformative prior.

### 3.4 Likelihood

For the likelihood,  $P(d|w)$ , we use a transformation of string edit distance between the phoneme sequence produced by the child and all candidate words. Specifically, we use exponentiated negative edit distance (Levy, 2008):

$$P(d|w) \propto e^{-\beta \times \text{dist}(d':w', d)} \quad (3)$$

where *dist* is the Levenshtein distance (minimal number of deletions, insertions and substitutions) between citation form  $d'$  for candidate word  $w'$ , designated here ( $d' : w'$ ), and the observed transcription ( $d$ ). For the results presented here, we grid sample  $\beta$  values between 1 and 6 by 0.1 increments, and take the value that assigns the highest posterior probability to a sample of 1000 communicative successes across models ( $\beta = 3.2$ ). This treatment of edit distance does not take into account phoneme similarity, *i.e.*, that certain phonemes are much more perceptually similar. We propose another more sophisticated likelihood function that captures this in the Discussion.

All model training and analysis code, as well as the fine-tuned model can be accessed at [https://osf.io/v7c3e/?view\\_only=176bb0f538af424da59007c53eff7e05](https://osf.io/v7c3e/?view_only=176bb0f538af424da59007c53eff7e05).

## 4 Results

### 4.1 Predicting Adult Recoveries

A comparison of Bayesian speech recognition models reflecting different priors reveals that the CDL+CONTEXT prior assigns the lowest average surprisal (highest average probability) to the recovered word gloss in the transcript. As Table 2 reveals, BERT models making use of context perform better than those that do not. CHILDES-tuned BERT models outperform the respective off-the-shelf BERT models. All BERT models outperform the CHILDES 1GRAM model, and all models with fitted priors assign significantly higher probability to the recovered glosses than the UNIFORMPRIOR model. These results mean that the model that is (1) fine-tuned to the child environment and (2) uses the surrounding utterance context is best able to predict the recoveries made by adults.

We next investigate how the prior probabilities in the previous analysis combine with likelihoods to predict word identity. That is, how do the adults’ prior expectations support inference when children’s productions are more or less adult-like? Comparing average surprisal across edit distances (Figure 1) reveals that models using BERT-based priors assign massively higher probability to word identities posited by annotators. For child productions that are 2 phonemes away from the citation form ( $x = 2$  in Fig. 1), CDL+CONTEXT assigns on average a probability of .24 ( $2^{-1 \times \text{surprisal}}$ ) to the correct gloss. This compares favorably to .12 under BERT+CONTEXT, .08 under CDL+ONEUTT, .03 under BERT+ONEUTT, .006 under the CHILDES 1GRAM, and .002 under UNIFORMPRIOR. CDL+CONTEXT assigns uniformly higher probability (lower surprisal) to the correct word identity, particularly when the phonetic form is more dissimilar (3 or more edits). This means that priors support recognition more when the perceptual input is noisier.

### 4.2 Predicting communicative failures

A separate question is which model best predicts whether a particular phoneme sequence will be a communicative success or failure. We address this by testing how

Table 2: Average prior surprisal on communicative successes from the Providence corpus (lower is better). The difference in average probability assigned to the actual gloss is  $2^{\text{diff}}$ , where diff is the difference between two model scores. \*Paired  $t$ -tests confirm sig. differences between models,  $p < 10^{-5}$ .

Model	Avg. Prior Surprisal* (bits)
CDL+CONTEXT	3.17
BERT+CONTEXT	4.59
CDL+ONEUTT	5.28
BERT+ONEUTT	7.09
CHILDES 1GRAM	8.80
UNIFORMPRIOR	12.95



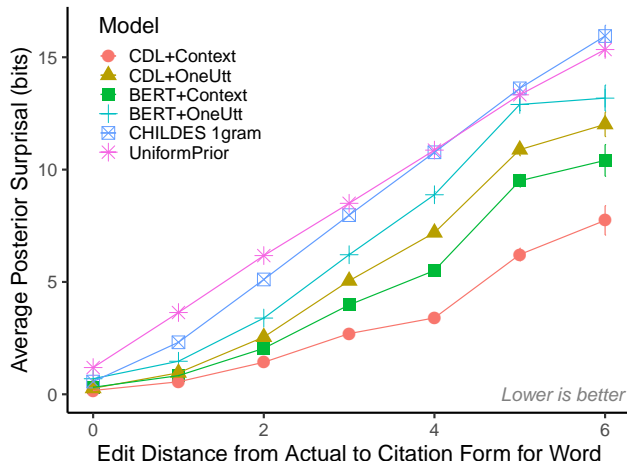


Figure 1: Posterior surprisal (negative log probability) of the recovered meaning for communicative successes. Error bars indicate standard error of the mean.

well posterior entropy under the models can predict communicative failures. As with the first analysis, the CDL+CONTEXT model provides the best trade-off between the prevalence of true positives and false positives (Figure 2). As both UNIFORMPRIOR and CHILDES 1GRAM models assign constant entropy to phoneme sequences (prior probabilities of candidates do not change as a function of context), their posterior entropy *only* reflects the contribution of the perceptual data. This analysis provides converging evidence that a model that is tuned specifically to child language and uses the surrounding utterance context — the one that best instantiates child-directed listening — is best able to replicate adult inferences.

### 4.3 Quantifying prior vs. posterior information

Finally, we quantify the information gain over time in conditioning on context (the fitted priors), conditioning on data (the posterior under the UNIFORMPRIOR model), and conditioning on both (the posteriors corresponding to the fitted priors). This analysis shows a larger shift in the probability distribution over candidates (greater information gain) going from the uniform prior to the CDL+CONTEXT prior compared to going from the uniform prior to its corresponding posterior (red line vs. green line in panel 1 of Figure 3). That is, the prior under the CDL+CONTEXT model contributes *more* information (better constrains guesses to word identity) than perceptual information alone. Contrary to our predictions, we find that the information gain for the prior is relatively constant over time for the CHILDES-fitted models. This suggests that child-directed listening can helpfully constrain adult listeners’ interpretations of children’s earliest verbal productions. As expected, children’s improving articulatory abilities result in an increase of all models’ posteriors over developmental time, as the likelihood function shared across models is able to contribute more and more to the task of interpretation.

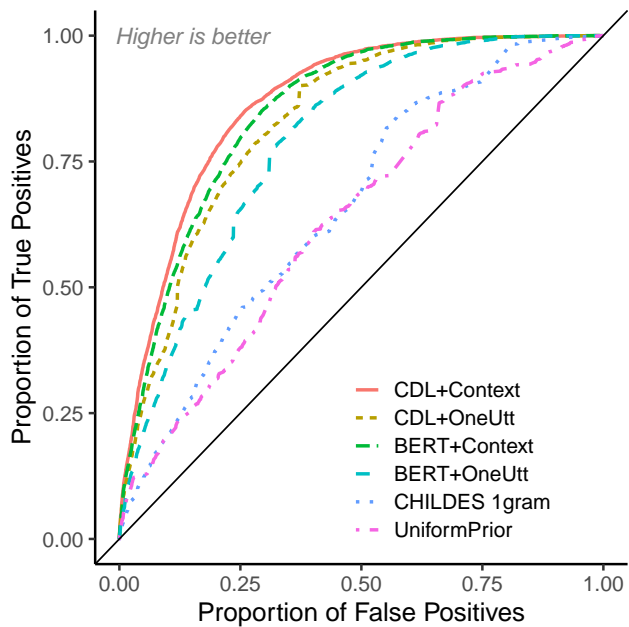


Figure 2: Classification performance in predicting communicative failures, as measured by the ROC of posterior entropy. The solid line with slope = 1 indicates chance. The area above this line indicates better classification performance.

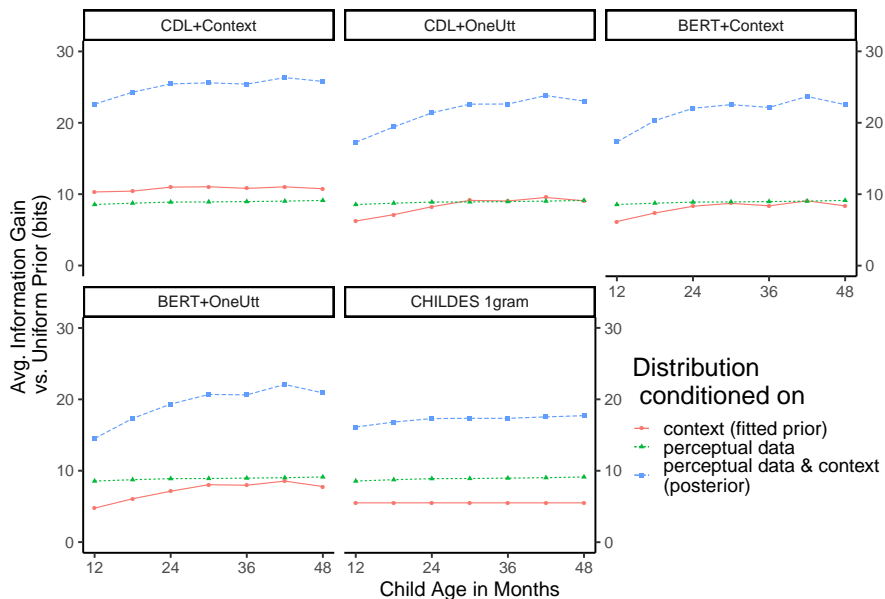


Figure 3: Average information gain from conditioning word prediction on context only (red, corresponding to the prior), perceptual data only (green), and context and perceptual data (blue, corresponding to the posterior) relative to a uniform prior.

## 5 Discussion

Language development is often characterized in terms of an increasing facility with processes on the side of the learner: developing motor planning, recognizing regularities of linguistic structure at different levels, and relating structure to entities and communicative contexts in the world. The current work suggests that early verbal communication depends not only on these well-studied developmental processes, but also on cognitive processes in the minds of adult caregivers.

We note two limitations with the current work before discussing its implications. First, the simple measure of edit distance does not capture the perceptual confusability of phonemes: *bug* and *rug* are equally good candidates for *pug*. One potential elaboration would be to use a *weighted* edit distance measure that takes into account the perceptual confusability of the phonemes. For example using a *probabilistic finite state string transducer* would allow assigning edits different “costs” according to experimentally-obtained confusion probabilities, *e.g.*, (Cutler et al., 2004).

Second, we make the simplifying assumption that inferences made by adult annotators in the lab are representative of the inferences made by adult caregivers in the moment, communicating in real time with children. While the inferential capacities of annotators are likely substantially *less* than those of adult caregivers (who have access to the non-linguistic context, as well as significantly more shared history with the child), research assistants may well be a decent proxy for adult listeners, due to their training as transcribers and exposure to child language. Potential differences in the inferential capacities of caregivers relative to other adult “listeners” should be tested experimentally.

These results additionally call attention to the interpretation of common methods in child language research. For example, vocabulary production measures on the Communicative Development Inventories (Fenson et al., 2007), have been historically interpreted as an index of children’s vocabulary and articulatory maturity. However, the current work suggests that successful communication – adult recognition of a word as a conventional form — relies additionally on adult inferential processes. Indeed the measure of a word’s “babiness,” a significant predictor of the order of children’s reported vocabulary production, may reflect the degree to which a word is more likely in child-directed speech compared to adult-directed speech.

Furthermore, our data invite a reconstrual of the nature of feedback in early language development. For example, if we assume that successful communication is itself reinforcing, child-directed *listening* might provide feedback to the child learner even in the absence of child-directed *speech*: a caregiver who interprets a child’s production of “uh” to mean “up” may not *say* anything in response to the child’s production, but provides feedback by effecting change on the part of the child when they pick the child up. This, in turn, leads to new puzzles: if adult caregivers can help many deficient communicative acts succeed, what presses children to get better?

Finally, we speculate regarding the role that child-directed listening might contribute to the emergence of language, both on evolutionary timescales and cases of rapid language emergence like Nicaraguan Sign Language. The current work suggests that successful recovery of meaning from child speech acts reflect not only the inductive biases, linguistic knowledge, and articulatory maturity of speakers, but also the

inferential biases of listeners.

## 6 Conclusion

We present a suite of Bayesian models of spoken word recognition to characterize the process of *child-directed listening*, or how adult caregivers find meaning in the noisy and often non-conventional speech productions of young children. We find that priors capitalizing on recent neural architectures — when trained specifically on child speech samples, and taking advantage of the greater linguistic context to make predictions — are best able to simulate adult inferential processes when interpreting noisy child speech. This research paves the way for understanding how children learn to employ language as goal-seeking agents in the presence of others.

## References

- Cutler, A., Weber, A., Smits, R., & Cooper, N. (2004, Dec). Patterns of English phoneme confusions by native and non-native listeners. *J Acoust Soc Am*, 116(6), 3668–3678.
- Demuth, K., Culbertson, J., & Alter, J. (2006). Word-minimality, epenthesis and coda licensing in the early acquisition of English. *Lang Speech*, 49(2), 137–174.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1* (pp. 4171–4186). Association for Computational Linguistics.
- Fenson, L., et al. (2007). *Macarthur-Bates communicative development inventories*. Paul H. Brookes Publishing Company Baltimore, MD.
- Gibson, E., Bergen, L., & Piantadosi, S. (2013). Rational integration of noisy evidence and prior semantic expectations in sentence interpretation. *Proceedings Natl. Acad. Sci. U.S.A.*, 110(20), 8051–8056.
- Golinkoff, R. M., Can, D. D., Soderstrom, M., & Hirsh-Pasek, K. (2015). (Baby)Talk to Me: The Social Context of Infant-Directed Speech and Its Effects on Early Language Acquisition. *Current Directions in Psychological Science*, 24(5), 339–344.
- Jawahar, G., Sagot, B., & Seddah, D. (2019, July). What does BERT learn about the structure of language? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (pp. 3651–3657). Association for Computational Linguistics.
- Levy, R. (2008). A noisy-channel model of human sentence comprehension under uncertain input. In *Proceedings of the 2008 conference on Empirical Methods in Natural Language Processing* (pp. 234–243).
- Norris, D., & McQueen, J. M. (2008, Apr). Shortlist B: a Bayesian model of continuous speech recognition. *Psychol Rev*, 115(2), 357–395.
- Rose, Y., & MacWhinney, B. (2014). The PhonBank project..
- Salazar, J., Liang, D., Nguyen, T. Q., & Kirchoff, K. (2020, July). Masked language model scoring. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics* (pp. 2699–2712). Online: Association for Computational Linguistics.
- Sanchez, A., Meylan, S., Braginsky, M., MacDonald, K., Yurovsky, D., & Frank, M. (2019). childes-db: A flexible and reproducible interface to the child language data exchange system. *Behavior Research Methods*, 51(4), 1928–1941.
- Shannon, C. (1951). Prediction and Entropy of Printed English. *Bell Systems Technical Journal*, 30, 50–64.
- Snow, C. E., & Ferguson, C. A. (1977). *Talking to children*. Cambridge University Press.
- Soderstrom, M. (2007). Beyond babytalk: Re-evaluating the nature and content of speech input to preverbal infants. *Developmental Review*, 27(4), 501–532.
- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., ... Rush, A. M. (2020, October). Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations* (pp. 38–45). Online: Association for Computational Linguistics.