# Learning Object Names in Real Time with Little Data

Jon Scott Stevens (jonsteve@ling.upenn.edu)

Department of Linguistics, 255 S. 36th Street Philadelphia, PA 19104 USA

## Abstract

#### We present an online learning model of early cross-situational word learning which maps words to objects from context with relatively sparse input. The model operates by rewarding and penalizing probabilities of possible word-to-object mappings based on real-time observation, and using those probabilities to determine a lexicon. We integrate prosodic and gestural cues and allow the learner to evaluate lexical entries. These enrichments allow efficient learning with minimal computational effort, producing results comparable to that of more complex models.

Keywords: Cross-situational word learning; online learning

## Introduction

The problem of how children learn the meanings of their first words, a problem for philosophers at least since the time of Augustine, has become an object of scrutiny in psychology and computational cognitive science. On one hand, experimental research shows us that young children can use a variety of cues (Bloom, 2000) to learn meanings from context with only a few exposures (Carey, 1978); on the other hand, computational modeling work underlines the difficulty of the process, requiring either complex statistical algorithms (Yu and Ballard, 2007; Frank, Goodman, and Tenenbaum, 2009) or large amounts of data (Fazly, Alishahi, and Stevenson, 2008) to achieve adequate learning. The goal of this paper is to simplify the computational problem of early word learning by integrating empirically motivated cues into a simple statistical model that learns object names in real time with speed and precision.

Following Yu and Ballard (2007), we integrate prosodic and gestural cues into a statistical learning algorithm for object names (which comprise the bulk of an early child's vocabulary in many languages, including English). Unlike other models, we give the learner access to a lexicon that adaptively changes as new observations are processed. This allows the learner to check hypothesized word meanings against new input and to enforce a preference for one-to-one mappings between words and objects. These enrichments have roots in experimental research and allow us to construct a simple, effective, and principled online learning model based on rewarding and penalizing probabilities (following Yang, 2002) associated with semantic hypotheses. We believe that this foundation of minimal complexity and empirical motivation produces a more psychologically plausible model.

Below we briefly outline some recent experimental and computational work in this area before presenting the details of our model and discussing its advantages. Children are able to learn words from context, often taking entire sentences as input and breaking that input down to create word-to-meaning mappings. In addition to this task, which is far from trivial, children must filter out an infinite number of erroneous but logically possible hypotheses of word meaning, as Quine (1960) famously noticed. The quantum leap between considering each member of an infinite set (an impossible task) and considering each member of a finite set, no matter how large, is the basis for the claim that word learning relies fundamentally on innately given hypothesis space constraints. The question, then, is not whether learning is constrained, but how it is constrained. We take computational models to be tests of purported answers to this question. As such, models should reflect the representations and, more loosely, the mechanisms present in human learners.

**Previous Work** 

### **Experimental Work**

We make use of three principles of early word learning that have emerged from experimental research: (1) mutual exclusivity, (2) the availability of gestural and prosodic cues, and (3) the apparent ability of learners to evaluate hypothesized word meanings against new data.

Markman (1992) and others have proposed that word learning is guided by a mutual exclusivity assumption, a default assumption that objects have only one name. There is independent experimental evidence (Ichinco, Frank, and Saxe, 2009) that suggests that children disprefer many-to-one word-toobject mappings, and such a preference improves the performance of a simple learning model, excluding would-be distractors from the semantic hypothesis space when those distractors already have a name in the learner's lexicon. Markman's view is that mutual exclusivity acts in concert with other default assumptions to extract a finite hypothesis space from Quine's infamous infinity.

Not only must the learner's hypothesis space be made finite, but it must interact with the learning mechanism in a way that produces quick results. Since Carey (1978) it has been noted that children learn words with impressive speed, often after only a few exposures. To achieve this end, we hold that word learning is guided not only by constraints like mutual exclusivity, but also by principles of salience and knowledge. This view allows the learning algorithm itself to be quite simple.

Under our conception of the process, word learning is guided both by word stress and by gestures, with greater weight being given to semantic hypotheses that map stressed words to gesturally indicated objects. Together we call these two cues "salience cues", reflecting their function of highlighting particularly important words and objects and making them salient to the learner. Without these crucial components, the data is simply too noisy for a simple learner to navigate. But these cues are independently justified. It is well known that babies are attentive to eye gaze and gestures. By nine months, they are capable of joint attention (Baldwin, 1991; Bloom, 2000), even responding to the emotional reactions of others. In short, humans seem to be programmed to pay attention to the actions of other humans from an early age. Thus, a gesture can serve as an "attentional magnet" for a young word learner.

If gesture serves to draw attention within the visual field, then patterns of prosodic prominence can be thought of as auditory gesture. Since the prosodic peaks of natural language have audible acoustic correlates (which are exaggerated in infant-directed speech), and since babies are known to be sensitive to these correlates (Soderstrom, Seidl, Nelson, and Jusczyk, 2003; Thiessen, Hill, and Saffran, 2005), we can posit that phonological phenomena such as word stress can be brought to bear on the question of how young learners figure out which words in an utterance are meant to refer. Indeed, prosodic information has been shown to be a good guide to word segmentation (Yang, 2004), an ability that must precede word learning.

Finally, recent work suggests that word learning involves a form of hypothesis evaluation, whereby learners will guess at a word's meaning and then, as further utterances of that word are processed, search the object space for evidence supporting their guess. Medina, Trueswell, Snedeker, and Gleitman (2009) assess mechanisms of cross-situational learning in adults using the human simulation paradigm (Gillette, Gleitman, Gleitman, and Lederer, 1999), a method whereby subjects are given video vignettes of naming events with the audio track removed and a single nonsense word uttered in place of some real word. Subjects were asked to give their best guesses as to the meaning of the nonsense words uttered in the vignettes. The vignettes were divided into "high informative" (HI) and "low informative" (LI) vignettes. The HI vignettes were those which were guessed correctly a majority of the time in isolation (determined in a separate experiment), and everything else was coded as a LI vignette.

Interestingly, subjects who saw a HI vignette followed by four LI vignettes were more likely to guess word meanings correctly at the end of the experiment than subjects who saw the same five vignettes in a different order. The authors hypothesize that early low informative instances handicap the learner, because rather than using the high informativity of later instances to make correct guesses, learners instead waste their time checking and rejecting the erroneous guesses they made previously. Subsequent eye-tracking studies show similar effects (Medina, Hafri, Trueswell, and Gleitman, 2010). Subjects behave as if they are choosing a hypothesized meaning for a novel item, and then verifying or falsifying that meaning as new data is received. This process of hypothesis evaluation opposes the traditional view of cross-situational word learning as a process of associating words with sets of multiple co-present objects. The computational model presented here reflects these developments; we show that it is helpful for the learner to be able to evaluate the semantic hypotheses contained in their lexicon against new data.

## **Previous Models**

Beginning with Siskind (2000), computational modeling has been a valuable tool for investigating the early word learning process. Various approaches have been taken, including Bayesian (Niyogi, 2002; Xu and Tenenbaum, 2007; Frank et al., 2009) and machine translation (Yu and Ballard, 2007; Fazly et al., 2008) approaches.

Yu and Ballard's (2007) work is particularly interesting for our purposes because it demonstrates the positive effects that prosodic and gestural cues can have on model performance. A machine translation algorithm (Brown et al. 1990) serves as a purely statistical core which is expanded by external social factors. The authors code corpus data for both prosodic peaks and indication by gesture or eye gaze. The words that represent peaks on an utterance's pitch track are given more weight than the other words in the utterance, and objects that are judged to be indicated in the visual field are given an analogous boost. We use a similar coding method, but our model differs from that of Yu and Ballard in a crucial way: it operates in real time. Yu and Ballard's is a batch learning model, which has a complexity disadvantage. Firstly, batch learning requires all tokens to be stored in memory, whereas online learning only requires types to be stored. Secondly, a real time implementation of a batch learning model would necessitate constant recalculation over all observed stimuli; as a result, the run time of such an algorithm will increase with the square of the number of observed stimuli, a sharper increase than that of an equivalent online model.

One of the most powerful recent models is another batch learning model, the Bayesian model of Frank et al. (2009). Using Bayesian inference, this model assigns a posterior probability score to individual lexicons given a corpus of data. MCMC stochastic search is used to find the lexicon with the highest score; no claims are made about how human learners do this. The scoring algorithm considers all possible intended sets of referents for a given scene. For example, if two objects, a pig and a horse, are visible to the learner during a particular utterance, four possible intentions must be considered: the speaker could be talking about the horse, the pig, both, or neither. Each possible intention yields some probability value, and those values are added together to obtain the contribution of that utterance to a lexicon's overall score.

Although the lack of explicitly given clues about speaker intent is perceived as an advantage, there is no indication that this reflects the behavior of human learners. Furthermore, considering all possible intents adds considerable complexity to the model in that the lexicon scoring algorithm becomes exponentially more demanding the more cluttered the room is. Since values are computed over the power set of visible objects, a naming event involving n candidate objects will

contribute  $2^n$  calculations to the scoring process. This is not too problematic with relatively clean data, but one can easily imagine a naturalistic learning environment with 30 distinct objects in the visual field, which would require over a billion calculations just to score one lexicon.

The authors claim that Bayesian inference explains mutual exclusivity. However, it is a choice by the modelers to make the likelihood term of their probability calculation dependent on the conditional probability P(word|object), rather than P(object|word). Thus, mutual exclusivity is built into the inference mechanism, not explained by it. In the absence of a deep explanation, we treat ME as an external cue rather than an architectural fact.

Fazly et al. (2008) present a more computationally plausible incremental model, but rather than focusing on object names as other models do, their model learns rich conceptual structures and as a result necessitates larger amounts of data to converge on correct meanings. Where their model requires as many as 20,000 utterance-situation pairs for accurate learning, our model learns with precision after fewer than 500, with some words being learned after fewer than six exposures. This reflects young children's famous ability to learn effectively from sparse input via fast-mapping.

## **Model Overview**

Word learning is mediated by a probability matrix with word types on the vertical axis and object types on the horizontal axis, illustrated in Figure 1.

The semantic hypothesis space for a potential object name is both open-ended and contingent on observation. This means that:

- A word-to-object mapping gets a value if and only if the word and the object have co-occurred.
- New words and objects can be introduced into the matrix at any time.

For example, the words *can*, *read*, and *books* are never uttered in the presence of the object coded 'EYES' in our evaluation corpus. Therefore, mappings from these words to 'EYES' have no value in Figure 1. This has the effect of reducing the size of a word's hypothesis space and preventing completely unfounded mappings from receiving a positive value when other mappings are penalized.

Novel words are mapped to 'NULL' with probability 1, with co-occurring objects receiving a value of 0. The 'NULL' mapping corresponds to the hypothesis that a word does not refer to an object. We take this to be the learner's default assumption. New objects are introduced into an old word's hypothesis space with a probability value of  $\frac{1}{n}$ , where *n* is the new size of that word's hypothesis space. The rest of the probability vector is normalized to accommodate the addition. This gives new semantic hypotheses a fair shot at lexicon inclusion.

This matrix provides us with a way to add and track probabilities of word-to-object mappings. Learning proceeds by

	BOOK	BIRD	RATTLE	FACE	EYES	NULL
look	0.45	0.00	0.01	0.38	NA	0.16
we	0.00	0.01	0.00	0.01	NA	0.98
can	0.00	0.01	0.00	0.01	NA	0.98
read	0.01	0.00	0.01	0.00	NA	0.98
books	0.23	0.00	0.00	0.36	NA	0.41
david	0.36	0.00	0.00	0.23	NA	0.41

Figure 1: A partial probability matrix for words and objects

updating these probabilities. We use Bush and Mosteller's (1951) Linear Reward-Penalty (LR-P) scheme, which was first applied to linguistic learning by Yang (2002). Below are the LR-P functions for rewarding and penalizing the probability of a hypothesis.

Table 1: Linear Reward-Penalty functions for a hypothesis *h*.

REWARD( <i>h</i> )	$p(h) = p(h) + \gamma(1 - p(h))$ where $\gamma$ is some constant between 0 and 1 For all $h' \neq h$ : $p(h') = p(h') * (1 - \gamma)$
PENALIZE( <i>h</i> )	$p(h) = p(h) * (1 - \gamma)$ For all $h' \neq h$ : $p(h') = \frac{\gamma}{n-1} + p(h') * (1 - \gamma)$ where <i>n</i> is the number of hypotheses being considered

The learning coefficient  $\gamma$  determines the severity of rewards and penalties. The final version of our model uses variable  $\gamma$  values to represent the privileged status of salient words and objects. Using these functions we update probabilities on the fly, and we use the results to update the learner's current lexicon of word-object pairs by including all and only those pairs whose probability values exceed a given threshold. This threshold (set to 0.65 in our simulations) serves to transform the probabilities into a discrete set of mappings that the learner can evaluate.

We implement different versions of the model to test the effect of each ability we give the learner. We use as our baseline a simple nested loop which rewards, in random order, all candidate objects for all words in each utterance (we call this process "multiple-candidate rewarding"). This is essentially a real-time equivalent of simple association frequency. We then add the hypothesis evaluation component by treating words that are in the current lexicon differently than other words. If a word is already mapped to an object, then the probability associated with that mapping is rewarded or penalized depending on whether that object is in the present situation (i.e. depending on whether the learner's hypothesis is consistent with current observation). In this case, no other candidates are rewarded. In all models, mutual exclusivity

For each observation, consisting of an utterance $U$ and a randomly-ordered set of possible object referents $O$ :
For each word $w$ in $U$ :
<ol> <li>If w is novel, assign probability 1 to w → NULL</li> <li>Else, add new objects to w's hypothesis space.</li> <li>If w is in the current lexicon:</li> <li>⇒ If w's hypothesized meaning m is an element of O, reward(w → m).</li> <li>⇒ Else, penalize (w → m).</li> <li>If w is not in the current lexicon:</li> <li>⇒ For each o in O:</li> <li>⇒ If o is not in the current lexicon, reward(w → o).</li> </ol>
Update the current lexicon.

Figure 2: An online cross-situational learning algorithm [The arrow  $(\rightarrow)$  in the algorithm should be read "maps to".]

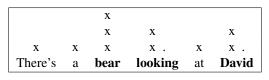


Figure 3: Stress on a prosodic grid

is enforced by exempting objects that already have names from multiple-candidate rewarding. The algorithm with both multiple-candidate rewarding and single-hypothesis evaluation is outlined in Figure 2.

The final component of our model is the integration of the salience cues. Objects in our video corpus were coded for gesture. An object was considered to be indicated by gesture during an utterance if it any point it was both (1) judged to be in the baby's field of vision, and (2) pointed to or held up in front of the baby. Eye gaze, being less obvious in the videos and therefore more prone to errors, was not coded.

Words were coded for prosodic accent. Utterances were given prosodic grid structures like the one in Figure 3, representing peaks in stress. Any word that received stress above the lexical level was coded as a stressed word. In typical adult speech the acoustic correlates of stress are subtle, and thus coding in this way is prone to subjectivity. However, this problem is ameliorated here, at least in part, by the exaggerated pronunciations utilized in the child-directed speech in the evaluation corpus.

Information about stress and gesture is used to determine the value of the learning coefficient  $\gamma$  for each rewarding or penalizing event. We give the model three parameters:

- $\gamma_H$  is the learning coefficient used when rewarding or penalizing a mapping that is already in the lexicon (hypothesis evaluation).
- $\gamma_M$  is the default learning coefficient used when rewarding possible mappings that are not already in the lexicon (multiple-candidate rewarding).

• *b* determines how much weight is given to hypotheses that map stressed words to gesturally indicated objects during multiple-candidate rewarding.

For words already in the lexicon, single hypotheses are rewarded or penalized with  $\gamma = \gamma_H$ . For words not in the lexicon, multiple possible mappings are rewarded with a different gamma value; mappings between stressed words and gesturally indicated objects are rewarded with  $\gamma = \gamma_M * b$ , while other mappings are rewarded with  $\gamma = \gamma_M * (1 - b)$ . The best performance is achieved when  $\gamma_H$  and  $\gamma_M$  are relatively high (0.4 and 0.36, respectively), and when most of the weight is given to salient mappings (b = 0.98).

To restate, the learner rewards and penalizes more drastically when checking their current lexicon against the world than when making multiple associations, and when the learner is making multiple associations, more weight is given to hypotheses that map stressed words to gesturally indicated objects.

To illustrate, consider the utterance in Figure 3. Assume, as shown in Fig. 4, that there are five visible objects accompanying this utterance, and only one of them is indicated by gesture (the mother is pointing to the bear and ignoring the other objects). Upon hearing this utterance, the learner possesses a lexicon of one entry: the word "david" maps erroneously to the object 'MIRROR'.

These data will be processed incrementally by the learner in the following way:

- 1. Since *there's* is not in the lexicon, it undergoes multiplecandidate rewarding rather than single hypothesis evaluation. Since it is not stressed, all present object meanings are rewarded using the coefficient  $\gamma_M * (1-b)$ .
- 2. The unstressed article *a* undergoes the same process as *there's*.
- 3. The lexicon does not have a mapping for *bear*, so it undergoes multiple-candidate rewarding, but since *bear* is stressed, the learning coefficient can vary. The gesturally indicated object referent 'BEAR' is rewarded with the higher coefficient  $\gamma_M * b$ , while the other non-indicated objects are rewarded with  $\gamma_M * (1-b)$ .
- 4. The stressed verb *looking* undergoes the same process as *bear*.
- 5. The unstressed preposition *at* behaves like *there's* and *a*.
- 6. Since *david* has a mapping in the learner's current lexicon, only that mapping is considered. In this case, *david* maps to 'MIRROR', and the object 'MIRROR' is not present in the current scene, so the learner's hypothesis is penalized. If the penalty lowers the probability value below the given threshold, then *david* → 'MIRROR' is kicked out of the lexicon.

uttered:	{there's, a, bear, looking, at, david}
stressed:	{bear, looking, david}
visible:	{BOOK, BIRD, RATTLE, BEAR, BOTTLE}
indicated:	{BEAR}
lexicon:	$\{ david \rightarrow MIRROR \}$

Figure 4: Example stimulus and accompanying lexicon

## **Performance and Comparisons**

All models were run on hand codings of two videos of mother-child interaction from the Rollins corpus (CHILDES, MacWhinney, 2000). Together the videos consist of 496 utterance-situation pairs (about 20 minutes of video). Performance was evaluated by aggregating the precision and recall against a gold standard over 100 simulations<sup>1</sup>, and taking the harmonic mean of the average precision and recall to produce an F-score. Model performance is detailed in Table 2. Three online models were tested: the baseline model, which does not utilize hypothesis evaluation, and two versions of the model given in Figure 2, one with a fixed  $\gamma$  value, and one which uses stress and gesture to determine  $\gamma$ . These models are compared to two implementations of Frank et al.'s (2009) Bayesian model: a direct implementation and a variant that only computes over stressed words and indicated objects.<sup>2</sup>

T 11 0	36 1 1	C	•
Table 7.	Model	nortormanco	comparison
1 a D R 2.	WIUUUI	performance	comparison.

Model type	Precision	Recall	F-score
Bayesian (FGT 09)	0.36	0.29	0.32
Bayesian (FGT 09)			
+ stress and gesture	0.72	0.38	0.52
Real-time updating	0.24	0.06	0.10
Real-time w/ evaluation	0.36	0.06	0.10
Real-time w/ evaluation			
+ stress and gesture	0.92	0.32	0.48

We see that adding prosodic and gestural information is a boost to both types of models; however, the cues have a more drastic effect on the real-time model. Once the cues are integrated, the F-scores for both types of models are comparable. Though the Bayesian model achieves a slightly higher F-score, the real-time model has a decided advantage in precision, with almost no erroneous mappings remaining in the lexicon. This is a desirable result because as learning continues beyond 20 minutes of interaction, the absence of misleading lexical entries will make for a more efficient process. The

Word	Object	Word	Object
book	book	piggies	pig
bear	bear	hat	hat
bunny	bunny	moocow	cow
kittycat	cat	meow	cat
sheep	sheep	bigbird	bird
bird	duck	ring	ring

Figure 5: Most frequent output lexicon

majority of simulations using this model produce the lexicon seen in Figure 5.

Performance is comparable to the Bayesian model of Frank et al. (2009), and our online learning model represents a computational simplification. Beal and Roberts (2009) argue for the importance of complexity analysis in computational cognitive science. A cognitive model should operate within known limits of human computational power, and complexity analysis is necessary to evaluate how realistic a model could be. Beal and Roberts show the Bayesian model of Xu and Tenenbaum (2007) to be quite costly from this perspective. Frank et al.'s model is even more costly. As mentioned above, it is problematic to sum probabilities for all possible intention sets for each situation. If the number of objects seen at one time has some upper bound N, then the upper bound asymptotic complexity will be  $O(2^N)$ ; the time it takes to process one situation will grow exponentially with the number of visible objects. This is not a problem for relatively clean data like the videos from the Rollins corpus, where the number of visible objects does not typically exceed 6 or 7, but an especially cluttered room may force the learner to make billions of calculations to score one lexicon against one interaction. This problem does not arise in our model. Furthermore, in contrast to batch learning models, our model necessitates only one pass through the input data.

Finally, the model presented here holds the promise of further unification with experimental research. Experiments like those described by Medina et al. (2009, 2010) may prove to be valuable both as a testing ground and as a source of refinement for research of this type, whose goal is to incorporate observable human behaviors into a psychologically plausible computational learning model.

## Conclusion

We have presented a model of object name learning that relies on gestural and prosodic cues and utilizes both singlecandidate and multiple-candidate probability updating mechanisms. The model operates in real time, making only one pass through a corpus and updating a lexicon after each successive utterance-situation pair. Performance is close to that of a comparable Bayesian model. The simplicity and success of the model suggests two things: (1) having access to word stress and gestural information makes word learning considerably easier, and (2) the ability to test beliefs about individual words makes learning more efficient. The next step

<sup>&</sup>lt;sup>1</sup>Multiple simulations account for slight variations in output caused by randomizing the order in which multiple candidates are rewarded.

 $<sup>^{2}</sup>$ We used our own hand-coding of the same videos that were used by Frank et al. For the Bayesian implementations, the authors' original code was used, strongly suggesting that the discrepancy between the performance reported here and the performance reported in Frank et al. (2009) is due to differences in the coding of the data.

in this line of research is to link up this computational approach even closer with experimental findings, and it is our hope that in doing so we may contribute to the growing pool of knowledge about how children learn the meanings of their first words.

## Acknowledgments

Thanks to Charles Yang, John Trueswell, and Tamara Medina for their help with this project, and thanks to the anonymous reviewers for their valuable comments.

### References

- Baldwin, D. (1991). Infants' contribution to the achievement of joint reference. *Child Development*, 62.
- Beal, J., & Roberts, J. (2009). Enhancing methodological rigor for computational cognitive science: Computational complexity. *Cognitive Science Conference*.
- Bloom, P. (2000). *How children learn the meanings of words*. Cambridge, MA: MIT Press.
- Brown, P., Cocke, J., Pietra, S. D., Pietra, V. D., Jelinek, F., Lafferty, J., et al. (1990). A statistical approach to machine translation. *Computational Linguistics*, *16*(2), 79-85.
- Bush, R., & Mosteller, F. (1951). A mathematical model for simple learning. *Psychological Review*, 68, 313-323.
- Carey, S. (1978). The child as word learner. In M. Halle, J. Bresnan, & G. Miller (Eds.), *Linguistic theory and psychological reality*. Cambridge, MA: MIT Press.
- Fazly, A., Alishahi, A., & Stevenson, S. (2008). A probabilistic incremental model of word learning in the presence of referential uncertainty. *Proceedings of the 30th Annual Conference of the Cognitive Science Society.*
- Frank, M., Goodman, N., & Tenenbaum, J. (2009). Using speakers' referential intentions to model early crosssituational word learning. *Psychological Science*, 20(5).
- Gillette, J., Gleitman, H., Gleitman, L., & Lederer, A. (1999). Human simulations of vocabulary learning. *Cognition*, 73, 135-176.
- Ichinco, D., Frank, M., & Saxe, R. (2009). Cross-situational word learning respects mutual exclusivity. *Proceedings of* the 31st Annual Meeting of the Cognitive Science Society.
- MacWhinney, B. (2000). *The childes project: Tools for analyzing talk* (3rd ed., Vol. 2). Mahwah, NJ: Erlbaum.
- Markman, E. (1992). Constraints on word learning: Speculations about their nature, origin, and domain specificity. In M. Gunnar & M. Maratsos (Eds.), *Modularity and constraints on language and cognition: The minnesota symposium on child psychology.* Mahwah, NJ: Erlbaum.
- Medina, T., Hafri, A., Trueswell, J., & Gleitman, L. (2010). Propose but verify: Fast-mapping meets cross-situational word learning. *Boston University Conference on Language Development*.
- Medina, T., Trueswell, J., Snedeker, J., & Gleitman, L. (2009). Rapid word learning under realistic learning conditions. LSA Annual Meeting, San Francisco, CA.

Niyogi, S. (2002). Bayesian learning at the syntax-semantics

interface. Proceedings of the 24th Annual Conference of the Cognitive Science Society, 697-702.

- Quine, W. (1960). *Word and object*. Cambridge, MA: MIT Press.
- Siskind, J. (2000). Learning word-to-meaning mappings. In P. Broeder & J. Murre (Eds.), *Models of language acquisition*. Oxford: Oxford University Press.
- Soderstrom, M., Seidl, A., Nelson, D., & Jusczyk, P. (2003). The prosodic bootstrapping of phrases: Evidence from prelinguistic infants. *Journal of Memory and Language*, 49, 249-267.
- Thiessen, E., Hill, E., & Saffran, J. (2005). Infant directed speech facilitates word segmentation. *Infancy*, *7*, 49-67.
- Xu, F., & Tenenbaum, J. (2007). Word learning as bayesian inference. *Psychological Review*, 114(2).
- Yang, C. (2002). Knowledge and learning in natural language. Oxford: Oxford University Press.
- Yang, C. (2004). Universal grammar, statistics, or both? *TRENDS in Cognitive Sciences*, 8(10), 451-456.
- Yu, C., & Ballard, D. (2007). A unified model of early word learning: Integrating statistical and social cues. *Neurocomputing*, 70, 2149-2165.