

# Quantifying the impact of active choice in word learning

Shohei Hidaka (shhidaka@jaist.ac.jp) and Takuma Torii (tak.torii@jaist.ac.jp)

Japan Advanced Institute of Science and Technology  
1-1 Asahidai, Nomi, Ishikawa, Japan

George Kachergis (george.kachergis@gmail.com)

Dept. of Artificial Intelligence / Donders Institute, Radboud University  
Nijmegen, the Netherlands

## Abstract

Past theoretical studies on word learning sought to explain the speed of childrens word learning using models sampling from a uniform word frequency (WF) distribution. We consider more realistic nonuniform, long-tailed WF distributions (i.e., Zipfian or power-law). Our new mathematical analysis of a recently-proposed simple learning model suggests that the model is unable to account for word learning in feasible time for Zipfian distributions. Considering children do learn these distributions, we propose a type of self-directed learning where the learner can help construct the contexts from which they learn words. We show that active learners choosing optimal situations can learn words hundreds of times faster than learners given randomly-sampled situations. In agreement with past empirical studies, we find theoretical support for the idea that statistical structure in real-world situations—perhaps influenced by a self-directed learner, and/or by a teacher—is a potential remedy for learning words with Zipf-distributed frequency.

**Keywords:** cognitive models of language acquisition; cross-situational word learning; statistical learning

## Child word learning

One of the most prominent differences between human and nonhuman cognition is our language ability. Much research has been dedicated to understanding the human capability for language, with a great deal of discussion focused on the process of language acquisition. A central debate in this conversation considers whether acquisition is based on innate and language-specific mechanisms (Chomsky, 1965; Gleitman, 1990), or bootstrapped from domain-general mechanisms (Smith, 2000; Kachergis, 2012). From the former perspective, humans become competent language users—mastering a complex system of syntax to produce endless semantics—very rapidly, and with relatively little training.

Word learning has been treated as an indicator of language development, and has been compared with a number of other indicators of cognitive abilities, such as memory (Vlach & Johnson, 2013; Vlach & Sandhofer, 2012). Although there are multiple empirical estimates of the number of words that children acquire, many studies agree that child’s word learning is quite fast. Early word production starts when the child is 12 months old on average, and by 18 months children can produce 50 words and comprehend 100-150 (Hulit & Howard, 2002). By 18 years of age, it is estimated we know over 60,000

words (Bloom, 2000). Under the assumption that each child has 8 hours of word learning opportunity everyday, these estimates mean the child learns a new word every learning hour for 18 years of the life.

Given these empirical estimates of word learning, theoretical studies have attempted to account for the quantitative characteristics of word learning. The first question is: What combination of learning mechanisms and structure in the language environment allows children to learn at this rate? This question poses a good necessary condition for any account of child word learning, as it will need to address this quantitative aspect of word learning.

As a first-order approximation, child learning may be modeled as an independent sampling process in which each word is learned independently. To estimate the fastest possible learning rate, Blythe, Smith, and Smith (2010, 2016) proposed an idealized learning model to address acquiring a full lexicon in the long term: 60,000 words over 18 years. In their model, each word is learned with its first sample – known as *fast mapping* in the developmental literature, where it has found some empirical support (Gershkoff-Stowe & Hahn, 2007). Under the simplifying assumption that each word is independently learned via fast-mapping, and its word frequency is distributed uniformly, their mathematical analysis of the model showed that a cross-situational learner is sufficiently fast to learn all 60,000 words after experiencing a reasonably small number of spoken words.

## Theoretical approach

Blythe *et al.*’s theoretical estimate has been treated as a theoretical implication that shows learning via independent fast-mapping of words is efficient enough to be a model of child word learning. In this study, we reinspect this theoretical implication by introducing a more realistic word frequency distribution. Intuitively, the difficulty with quickly learning a realistically long-tailed distribution is that a learner can expect to wait a very long time to hear the many low-frequency words in a power-law distribution such as is found in natural copora (Zipf, 1949). Our mathematical analysis implies that the learning rate of the independent fast mapping is quite sensitive to the word frequency distribution. More importantly, even fast mapping—the most efficient learn-

ing, requiring only a single sample, can be too slow to learn 60,000 words in 18 years, if word frequency follows Zipf’s law. Thus, our analysis implies that the independent fast-mapping model cannot be an account for child word learning, if there are many words sampled less frequently. This mathematical implication leads to an empirical test of whether the word distribution in the child-directed speech is uniform or non-uniform such as a Zipf distribution. Thus, in the second study, we analyzed the CHILDES corpora for word frequency distribution in child-directed speech.

Given this result of the mathematical analysis, we explore an extension of the word learning mechanism by additionally assuming that the word learning is more *active* than that is supposed to be traditionally. Typically, as analyzed in the past studies above (Blythe et al., 2010, 2016), the learning is supposed passive – the learner has no choice but observing samples words and objects from a given probability distribution. This is certainly oversimplified, as actual child word-learners choose when, where and from whom they would like to learn words. Thus, our second analysis estimates the impact of a form of active choice of situations in word learning. Our analysis shows that active learning is likely to have a sufficiently beneficial impact to make word learning fast enough to happen on a realistic timescale.

## Independent fast-mapping learning

### Uniformly distributed word frequency

Blythe, Smith, & Smith (2010) proposed a mathematical model of word learning, which has a closed-form expression under a certain simplification. In their recent study, Blythe, Smith, & Smith (2016) analyzed essentially the same model, although slightly modified for analytic convenience. Here we briefly introduce the most recent form (2016) of their model.

Blythe et al. originally consider cross-situational word learning. Suppose there are  $W$  words and  $O$  objects in the hypothetical world. Further the numbers of words and objects are equal,  $W = O$ , in their cross-situational learning scheme, and every object has its name and no objects have two names. Namely, there are  $W$  correct pairs of words and objects. Without loss of generality, denote the  $W$  pairs by  $1, 2, \dots, W$ , and suppose the  $k^{\text{th}}$  object is paired with the  $k^{\text{th}}$  word.

Given these pairs are unknown, a word learner is to infer correct pairs by experiencing a series of *episodes*. In each episode, the learner is exposed to  $M \leq W$  words and  $M$  objects, without any explicit information about which word is paired with which object. With one episode with  $M \geq 2$  objects and words, the learner cannot tell which of  $M$  words should be associated to which of  $M$  objects.

The simplest model in a series of extended ones is called the fast-mapping learning model. In the litera-

ture of language development, it is well-known that children as young as three years old can quickly generalize a novel name to a novel object even on the first occurrence (Mervis & Bertrand, 1994). Due to this one-shot nature of their word learning, it is called fast mapping (Carey & Bartlett, 1978). Capturing this empirically-supported constraint, the fast-mapping learner in the model is supposed to learn a new pair of word and object only with the first experience of it. The fast-mapping learner is equivalent to the cross-situational learner if there is one correct word-object pair present in each episode ( $M = 1$ ).

As fast-mapping learning is the most efficient scheme (at least for independent word learning), it gives a good baseline estimate of the number of samples to learn all the words in a given list. Blythe et al. (2010) model a fast-mapping learner acquiring words independently drawn from a uniform distribution of  $W$  words given in each episode. As every episode has one word with probability  $1/W$ , this is equivalent to the so-called Coupon Collector’s Problem (Blom, Holst, & Sandell, 1994). In this problem, the expected time  $T$  to finish sampling all the words is

$$E[T] = \sum_{i=1}^W E[t_i].$$

where  $t_i$  is the time to sample a  $i^{\text{th}}$  new word given  $(i-1)$  words being learned. Thus,

$$E[T] = W \sum_{i=1}^W 1/i \approx W \log W. \quad (1)$$

Setting the number of words  $W = 60,000$ , which is an empirical estimate of the number of words 18 years old knows on average,  $T = 660,126$ . This estimate is comparable with the “reasonable” number of samples justified by Blythe et al. (2010) which individual children can be exposed to for their 18 years of lives.

### Non-uniformly distributed word frequency

Here we extend this analysis on the fast-mapping learner to the case with word frequency distributed non-uniformly. Our extended analysis will reveal that the estimate based on Equation (1) by Blythe et al. is quite “optimistic”, as an estimate with non-uniform word distribution is larger than that in general.

Here let us derive the number of episodes  $T$  that, for  $0 \leq \epsilon \leq 1$ , the  $(1 - \epsilon)$  of children learned all the  $W$  words listed. Suppose a set of  $W$  words in which each word  $1, \dots, W$  is drawn from the distribution  $p = (p_1, \dots, p_W)$ . The proportion of children who finished learning all the words is  $(1 - \epsilon)$  for  $0 < \epsilon < 1$  requires the number of episodes  $T$ , which is the root of

$$\prod_{i=1}^W (1 - (1 - p_i)^T) = 1 - \epsilon. \quad (2)$$

The left hand side of (2) is the probability that every word is present at least once in the  $T$  episodes.

Write

$$f_{W,\epsilon}(x) := \frac{\log(1 - (1 - \epsilon)^{1/W})}{\log(1 - x)}.$$

For the uniform distribution,  $p_i = 1/W$  for every  $i = 1, \dots, W$ , the root of (2) is given by

$$T = f_{W,\epsilon}(1/W). \quad (3)$$

This  $T$  is the number of episodes with which the proportion of children finished learning all the words is  $(1 - \epsilon)$ . Setting  $\epsilon = 1/2$  in (3), we obtain the median of  $T$ ,  $f_{W,1/2}(1/W)$ , that is comparable with the mean of  $T$  in (1).

Unlike (3) for the uniform distribution, the root  $T$  of Equation (2) in general is not closed-form. Thus, let us consider the upper and lower bound for the root instead of the rigorous form of it. For the general word distribution  $p = (p_1, \dots, p_W)$ , the intermediate value theorem states that there exists a unique constant  $c$  holding  $\min p \leq c \leq \max p$ , with which the root of (2) is expressed as

$$T = f_{W,\epsilon}(c).$$

Equivalently, we have inequality

$$f_{W,\epsilon}(\max p) \leq T \leq f_{W,\epsilon}(\min p).$$

As we are interested in the worst possible estimate of  $T$ , this inequality states that the upper bound  $T_+ := f_{W,\epsilon}(\min p)$  of  $T$  is characterized with the probability to sample the least frequent word  $\min p$ .

This extended mathematical analysis implies that the uniform distribution  $q = (1/W, \dots, 1/W)$  of words gives the minimal possible upper bound  $T_+$  among any frequency distribution of  $W$  words, as any distribution  $\min p$  of  $W$  words holds  $\min p \leq \min q$ . Therefore, the expectation of  $T$  in the form of (1) with the uniform distributed words is the most optimistic, which may underestimate the number of episodes required for learning with a realistic word distribution.

For example, let us consider an alternative case that the  $W$  word follows the Zipf distribution  $p = (1^{-1}/H_W, 2^{-1}/H_W, \dots, W^{-1}/H_W)$ , where  $H_W$  is the harmonic number  $H_W = \sum_{i=1}^W i^{-1}$ . In this case, the minimal probability is  $\min p \approx 1.44 \times 10^{-6}$ , and the upper bound  $T_+$  is  $1.08 \times 10^7$  for  $\epsilon = 0.01$ . This estimate means that learning of Zipf-distributed words requires 16.4 times as many samples as learning of uniformly-distributed words. That means that 206 independent episodes exposed to a word learner every hour (or three episodes every minute), assuming 8 hours of learning everyday of 18 years of life. This estimate cannot possibly be considered “reasonable” with respect to ordinary life of children in any culture.

## Sensitivity to non-uniformity of word frequency distribution

To analyze the sensitivity to the non-uniformity, here we analyze the Zipf distribution with different exponent parameters. Denote the Zipf distribution with the exponent parameter  $a \geq 0$  by  $p = (1^{-a}/H_{W,a}, 2^{-a}/H_{W,a}, \dots, W^{-a}/H_{W,a})$  where  $H_{W,a}$  is the generalized harmonic number  $H_{W,a} = \sum_{i=1}^W i^{-a}$ . It is reduced to the uniform distribution by  $a = 0$ . The larger the exponent  $a$  is, the minimal probability  $\min p$  is smaller. Thus, here we analyze the upper bound  $T_+$  as a function of the exponent parameter  $a$ .

Write  $T_+ = f_{N,\epsilon}(\min p)$ , which gives a reasonable estimate of the upper bound of the root  $T$  of (2). As a function of the exponent  $a$ , we have

$$\frac{\partial \log T_+}{\partial a} = \frac{p_{\min} \left( \frac{\partial H_{N,a}}{\partial a} / H_{N,a} + \log W \right)}{(1 - p_{\min}) \log(1 - p_{\min})}$$

and further we have

$$\frac{\partial^2 \log T_+}{\partial a^2} \geq 0.$$

This implies the  $T_+$  is a super-exponential monotone function of the exponent  $a$ . It is also numerically confirmed in Figure 1, in which the numbers of episodes are shown as functions of the exponent for  $W = 10000, 60000$ . In this plot,  $a = 0$  shows an estimate for the uniform distribution, and  $a = 1$  shows that of the standard Zipfian distribution. It is striking that even the fastest learning such as fast mapping can be quite slow (exponentially as a function of  $a$ ) for with distributions with some item with a very small probability.

## Empirical dataset

Given theoretical implication in the previous study, let us analyze an empirical word distribution, which children typically are exposed to. It is difficult to exactly count “episodes” or “pairs of word and object” in a real dataset, due to its ambiguity of definition and it is also up to children’s subjective perspective. Here, as a proxy of them, we counted the word frequency based on child-directed speech in the CHILDES corpus (MacWhinney & Snow, 1990). Figure 2 shows a representative word distribution of 51,446 words aggregated over 4,163 transcripts of all the corpora in CHILDES retrieved in December 2007. The minimal word probability was  $1.089 \times 10^{-7}$ , which gives the upper bound  $T_+ = f_{51446,0.01}(1.089 \times 10^{-7}) = 1.420 \times 10^8$  or the median estimate  $T_+ = f_{51446,0.5}(1.089 \times 10^{-7}) = 1.030 \times 10^8$ . These estimates of required samples, an order of magnitude larger than the optimistic theoretical estimate, suggest that it is difficult to learn these empirical words with this Zipfian-like frequency distribution.

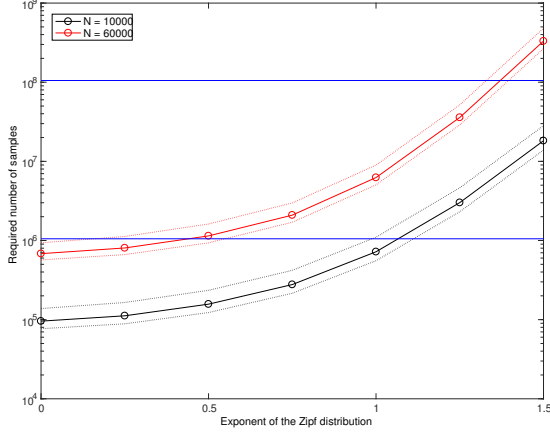


Figure 1: For  $\epsilon = 0.01, 0.5, 0.99$  (broken and solid lines),  $N = 10000, 60000$  and the exponent  $a = 0, 0.25, \dots, 1.5$ , the required number of samples  $M$  for a generalized Zipf distribution  $p_k = k^{-a} / \sum_{k=1}^N k^{-a}$  is numerically calculated by the root of Equation (2).

## Active choice of situations

### Formulation

The implication of the mathematical analysis above, which suggested that even fast-mapping may not be efficient enough for non-uniformly distributed words, raises a controversy between past theoretical analyses and empirical findings of quantitative aspects of word learning.

Here, we explore a possibility to reconcile the discrepancy between theory and empirical findings, by considering a further relaxation of past theoretical assumptions about children’s word learning. In the conventional theoretical framework, the learner is assumed to be *passive*, having no choice but to observe and learn from a given context: a randomly-sampled set of objects, of which a (random) subset are labeled with words. This assumption of a passive learner simplifies the theory, but surely underestimates real learners, who have some choice about which contexts they experience. Here, we consider a type of active learner who is able to choose from which situation/context he or she learns words.

Suppose that there are  $N$  word-object pairs and  $M$  situations, and that the conditional probability to observe the  $i^{\text{th}}$  word-object pair is  $p_{ij}$  given the  $j^{\text{th}}$  situation. Thus, the active learner has a choice of the situation out of the given  $M$  situations from which he or she learns the word-object pairs. Suppose that the active learner chooses the  $j^{\text{th}}$  situation with probability  $q_j$ . Let us denote the  $N \times M$  matrix of the conditional probability by  $P = \{p_{ij}\}_{ij}$  and the  $N \times 1$  vector of the choice probability by  $q = (q_1, q_2, \dots, q_M)^T$ . With this notation, the marginal probability of word-object pairs is given by the vector  $Pq \in \mathbb{R}^N$ . According to our mathematical anal-

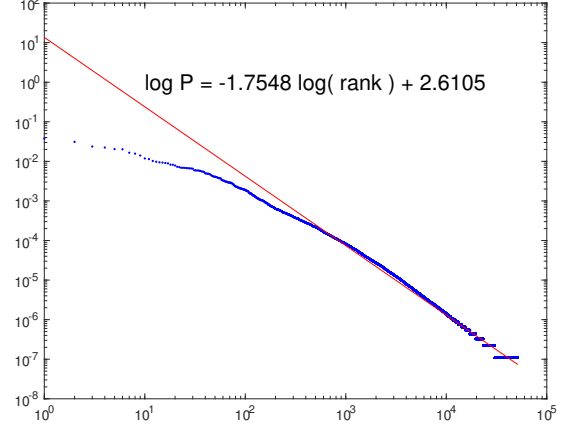


Figure 2: Word frequency in a corpus aggregated from the CHILDES transcripts.

ysis in the previous section, the minimal probability of objects decides the number of samples required to complete the word learning, the best choice for the active learner is given by the choice probability

$$\hat{q} = \arg \max_q \min(Pq).$$

This minimal probability,  $\min(P\hat{q})$ , gives the theoretical upper bound for the minimal number of samples  $f_{W,\epsilon}(\min(P\hat{q}))$ , as  $P$  is not known before empirical learning, and the active learner also needs to estimate  $P$  from the sample. For a given matrix  $P$ , the optimal  $\hat{q}$  can be computed by the iterated linear programming algorithm (See also Appendix for the detail).

As a baseline for the passive learner, we consider the average  $\min(Pq)$  with the uniform distribution over the vector  $q$ , whose lower bound is given by the Jensen’s inequality

$$\int_{q \in \mathbb{S}^N} \min(Pq) (N-1)! dq \geq \min(P\mathbf{1}_N/N),$$

where the integral is taken over the  $N-1$  dimensional unit simplex  $q \in \mathbb{S}^N$ . For a sufficiently small  $x \ll 1$  and  $y \ll 1$ ,  $f_{W,\epsilon}(x)/f_{W,\epsilon}(y) \approx y/x$ . Thus, the rate  $R = \min(P\hat{q})/\min(P\mathbf{1}_N/N)$  gives a good estimate for the rate of efficiency  $R$ , by which the active learning with the optimal probability  $\hat{q}$   $R$  times faster than the passive learning with a fixed probability  $q$ .

### Empirical evaluation

To evaluate the potential impact of the active learning, we study the SUN database (Xiao, Hays, Ehinger, Oliva, & Torralba, 2010) as an empirical object distributions in an collection of real-life scenes. The SUN database (retrieved on September 25th in 2016) has  $N = 3,458$  objects and  $M = 1,111$  scenes in it. This data is supposed

to give the  $N \times M$  matrix  $P$  in which each column is the conditional probability of the objects given each scene. If the scene choice probability is the uniform distribution  $q = \mathbf{1}_N/N$ , the  $\min(Pq)$  was  $8.30 \times 10^{-9}$ . Meanwhile, with the optimal  $\hat{q}$ , the  $\min(P\hat{q})$  was  $1.95 \times 10^{-6}$ , which implies the active learning was approximately  $\min(P\mathbf{1}_N/N)/\min(P\hat{q}) = 235.3$  times faster than the baseline passive learning. The marginal probability distributions of objects for the baseline and optimal  $q$  are shown in Figure 3. The difference between the two marginal distributions is visible at their tails – the tail for the uniform  $q$  decreases like an exponential function, but that for the optimal  $\hat{q}$  decreases as a power function (linear in the double log plot). This empirical evaluation suggests that the active learning of interest can boost the fast mapping a few orders more efficiently.

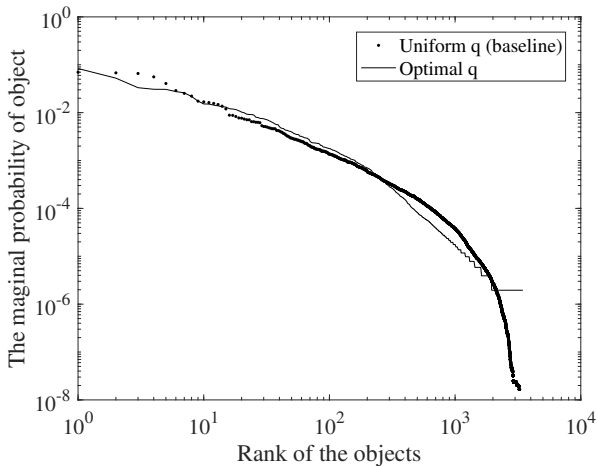


Figure 3: The marginal probability of objects for the optimal  $\hat{q}$  (line) and its baseline (dots).

### Online active learning

The quantification of the efficiency of active learning is based on the optimal  $\hat{q}$  with the knowledge of  $P$ . This gives an optimistic estimate for the active learner, as the matrix  $P$  is not fully known in reality. Here we performed a Monte Carlo simulation to quantify the efficiency of an *online* active learner who gradually updates knowledge in the matrix  $P$  and estimates  $q$  on the basis of the sample estimate of  $P$ . If this online active learner is comparable with the optimal active learner with  $\hat{q}$ , we can treat the performance analysis on the optimal active learner above (a few orders more efficient) as holding for the online active learner. For this purpose, we generated a  $N \times M$  matrix  $P$  with  $N = 1000, M = 100$ , which has the elements in each column are Zipfian probabilities  $P_{\pi(i)} \propto i^{-a}$  with the random coefficients  $a \in [1, 1.5]$ , where  $\pi : \{1, \dots, N\} \mapsto \{1, \dots, N\}$  is a random permutation. The online active learner has the uniform

choice probability  $q_1 = \mathbf{1}_N/N$ . For  $k^{\text{th}}$  batch of 1000 steps, the online learner samples the objects according to the probability  $Pq_k$ , and constructs the sample probability matrix  $\hat{P}_k$  according to the sample frequency. After the  $k^{\text{th}}$  sampling step, the online learner estimates  $q_k := \arg \max_q \min(\hat{P}_k q)$ . In each run of this procedure, we repeat up to  $100 \times 1000$  samples, and obtain one sample for the number of required samples to finish learning all the 1000 objects. With 100 runs, we obtain the Monte Carlo estimate of the online learner shown in Figure 4. Figure 4 shows the sample probability distribution of the number of required samples in the Monte Carlo simulation (circles: histogram, line: smoothed estimate), and its comparable median estimate for the optimal learner (green vertical line) and the passive learner with the uniform  $q$  (red vertical line). This simulation result shows that the online learner is as fast as the optimal learner, and is likely to be faster than the passive learner.

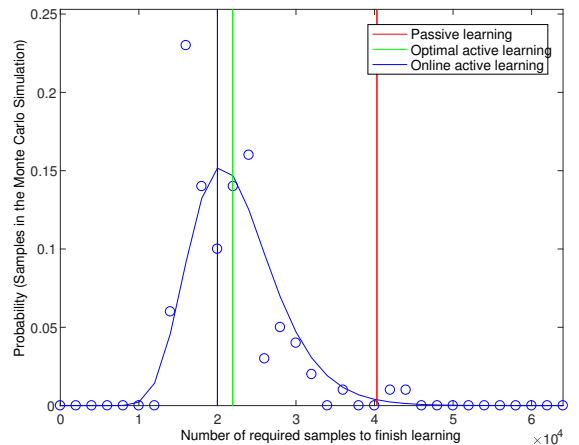


Figure 4: The probability distribution of the number of required samples to finish learning for passive (red), optimal (green), and online active learner (blue).

### Discussion

This study has provided mathematical analyses of quantitative aspects of word learning that provide key constraints which any theoretical account for children’s word learning should satisfy. We reinspected the past theoretical claim by Blythe et al. (2010) that learning via independent fast mapping was efficient enough to account for the average number of words known by 18-year-olds. Our new analysis extends their analysis to fast mapping with non-uniform word frequency distributions, and shows that even learning via fast mapping is not efficient enough to learn words whose distribution has rarely sampled words—including the Zipf (i.e. power-law) distribution, which describes empirical word frequency distributions from natural language.

Given that this new analysis implies learning would be too slow under realistic distributions, we consider a more efficient learning scheme, in which the learner can choose preferred situations from which words are learned. This type of active control over situations or contexts seems natural with respect to general observations of children’s behavior, and has been shown to benefit adult word learners (Kachergis, Yu, & Shiffrin, 2013), but has not been subjected to theoretical analysis as far as we know. We quantify and evaluate the effect of this type of self-directed learning in word learning. As the least probable word in the distribution determines learning efficiency, we analyzed the active choice for the situations maximizing this key parameter. Analyzing an empirical dataset of the words given situations, we estimate that active learning is over two hundred times more efficient in learning time than passive learning. This result suggests that active choice in word learning can resolve the issue that naturalistic non-uniform word distributions greatly slows passive fast mapping.

Our analyses in this paper utilized one of the simplest learning schemes, fast mapping, in order to highlight the effects of varied word frequency distributions, and of active learning. However, we expect the analytic techniques we employed would also allow analysis of other learning algorithms, including many proposed variants of cross-situational learning. In future work, we will report similar analyses for learning schemes with perhaps greater cognitive plausibility. On this path towards ever more realistic assumptions about the language environment and learners’ ability to shape it, we expect to make progress toward a general theoretical framework spanning many proposed word learning schemes.

### Acknowledgment

This study is supported by the JSPS KAKENHI Grant-in-Aid for Young Scientists JP 16H05860.

### References

- Blom, G., Holst, L., & Sandell, D. (1994). Problems and snapshots from the world of probability. In (p. 85-87). New York, NY: Springer-Verlag New York.
- Bloom, P. (2000). *How children learn the meaning of words*. Cambridge, MA: MIT Press.
- Blythe, R. A., Smith, A. D. M., & Smith, K. (2016). Word learning under infinite uncertainty. *Cognition*, 151, 18–27.
- Blythe, R. A., Smith, K., & Smith, A. D. M. (2010, January). Learning Times for Large Lexicons Through Cross-Situational Learning. *Cognitive Science*, 34(4), 620–642.
- Carey, S., & Bartlett, E. (1978). Acquiring a single new word. *Papers and Report on Child Language Development*, 15, 17–29.
- Chomsky, N. (1965). *Aspects of the theory of syntax*. MIT Press.
- Gershkoff-Stowe, L., & Hahn, E. R. (2007). Fast mapping skills in the developing lexicon. *Journal of Speech, Language, and Hearing Research*, 50, 682–697.
- Gleitman, L. (1990). The structural sources of word meaning. *Language Acquisition*, 1, 3–55.
- Hulit, L., & Howard, M. R. (2002). *Born to talk*. Toronto: Allyn and Bacon.
- Kachergis, G. (2012). Learning nouns with domain-general associative learning mechanisms. In N. Miyake, D. Peebles, & R. P. Cooper (Eds.), *Proceedings of the 34th annual conference of the cognitive science society* (p. 533-538). Austin, TX: Cognitive Science Society.
- Kachergis, G., Yu, C., & Shiffrin, R. M. (2013). Actively learning object names across ambiguous situations. *Topics in Cognitive Science*.
- MacWhinney, B., & Snow, C. (1990). The child language data exchange system: An update. *Journal of Child Language*, 17(02), 457–472.
- Smith, L. B. (2000). How to learn words: An associative crane. In R. Golinkoff & K. Hirsh-Pasek (Eds.), *Breaking the word learning barrier* (pp. 51–80). Oxford: Oxford University Press.
- Vlach, H. A., & Johnson, S. P. (2013). Memory constraints on infants’ cross-situational statistical learning. *Cognition*, 127, 375–382.
- Vlach, H. A., & Sandhofer, C. M. (2012). Fast mapping across time: Memory processes support children’s retention of learned words. *Frontiers in Developmental Psychology*, 3(46), 1–8.
- Xiao, J., Hays, J., Ehinger, K., Oliva, A., & Torralba, A. (2010). Sun database: Large-scale scene recognition from abbey to zoo. In *Ieee conference on computer vision and pattern recognition*.
- Zipf, G. (1949). *Human behavior and the principle of least effort*. Cambridge, MA: Addison-Wesley.

### Appendix: Iterated linear programming

For a  $N \times M$  matrix  $P$ , write its  $i^{\text{th}}$  row by  $P_i$ . Let  $I = \{1, 2, \dots, N\}$  be the set of all indices. At the initial step, define

$$K_0 := \emptyset, C_0 := \mathbb{S}^M, q_0 := e_1,$$

where  $e_1 := (1, 0, \dots, 0)^T \in \mathbb{R}^N$ . Then for  $0 < n \leq N$ , define

$$\begin{aligned} k_n &:= \arg \min_{k \in I \setminus K_{n-1}} P_k q_{n-1}, K_n := K_{n-1} \cup \{k_n\}, \\ C_n &:= \left\{ q \in C_0 \mid \bigwedge_{k \in K_n} (P_k - P_{k_n})q \leq 0 \right\}, \\ q_n &:= \arg \max_{q \in C_n} P_{k_n} q, \end{aligned}$$

until  $n = m$  such that  $\min_{k \in K_m} P_k q_m \leq \min_{k \in I} P_k q_m$ . The algorithm stops the iterative procedure by outputting  $q := q_m$ .