

Distributional Language Learning: Mechanisms and Models of Category Formation

Richard N. Aslin^a and Elissa L. Newport^b

^aUniversity of Rochester and ^bGeorgetown University

In the past 15 years, a substantial body of evidence has confirmed that a powerful distributional learning mechanism is present in infants, children, adults and (at least to some degree) in nonhuman animals as well. The present article briefly reviews this literature and then examines some of the fundamental questions that must be addressed for any distributional learning mechanism to operate effectively within the linguistic domain. In particular, how does a naive learner determine the number of categories that are present in a corpus of linguistic input and what distributional cues enable the learner to assign individual lexical items to those categories? Contrary to the hypothesis that distributional learning and category (or rule) learning are separate mechanisms, the present article argues that these two seemingly different processes—acquiring specific structure from linguistic input and generalizing beyond that input to novel exemplars—actually represent a single mechanism. Evidence in support of this single-mechanism hypothesis comes from a series of artificial grammar-learning studies that not only demonstrate that adults can learn grammatical categories from distributional information alone, but that the specific patterning of distributional information among attested utterances in the learning corpus enables adults to generalize to novel utterances or to restrict generalization when unattested utterances are consistently absent from the learning corpus. Finally, a computational model of distributional learning that accounts

This article is a summary of a talk presented at the Workshop on Infant Language Development held in San Sebastian, Spain on June 22, 2013. The workshop was organized by Monika Molnar, Arty Samuel, Manuel Carreiras, Kepa Paz-Alonso, and Melissa Baese-Berk, with support from the Basque Center on Cognition, Brain, and Language. Much of the research reviewed in this article was conducted in collaboration with Patricia Reeder and Ting Qian, with support from grants from NIH (HD-037082, HD-067250).

Correspondence concerning this article should be addressed to Richard N. Aslin, Department of Brain and Cognitive Sciences, Meliora Hall, River Campus University of Rochester, Rochester, NY 14627. E-mail: aslin@cvs.rochester.edu

for the presence or absence of generalization is reviewed and the implications of this model for linguistic-category learning are summarized.

Keywords statistical learning; rule learning; generalization; grammatical categories; infancy; child language

Introduction

A seminal event in the history of the study of child language acquisition was the publication of Eimas, Siqueland, Jusczyk, and Vigorito (1971)—which concluded that “the means by which the categorical perception of speech, that is, perception in a linguistic mode, is accomplished may well be part of the biological make up of the organism” (p. 306)—is now more nuanced than originally conceived. But that does not detract from the impact of their findings. Eimas et al. set the stage for literally hundreds of experiments that documented the remarkably sophisticated language-processing skills of preverbal infants and in turn raised questions about how those skills arise. What was initially viewed as evidence that linguistic experience plays a relatively minor role in language acquisition is now interpreted as evidence for a species-general set of constraints on how complex acoustic signals are discriminated and categorized. The role of experience comes later in infancy, beginning at 4–6 months for vowels and then 6–10 months for consonants. But there can be little doubt, even before any infant data on speech perception were gathered, that experience must play a substantial role in language processing. Cross-language differences in the surface properties of linguistic systems and imperfect acquisition of any system later in life document the power of exposure in shaping a native-language linguistic system.

Foundations of Statistical Learning

Further fueling the importance of early experience on language acquisition was the finding 25 years after Eimas et al. (1971) by Saffran, Aslin, and Newport (1996). Saffran et al. (1996) showed that infants can use the distributional properties of a corpus composed of an uninterrupted stream of syllables to extract information about the statistical coherence of samples drawn from that corpus. The stream was composed of 12 consonant-vowel (CV)-syllables arranged into four trisyllabic strings, with each syllable occurring in only a single triplet. Thus, the transitional probabilities from syllable-1 to syllable-2 and from syllable-2 to syllable-3 were 1.0, but the transitional probability

after each triplet to the first syllable of the next triplet was 0.33. Importantly, other cues to the grouping of syllables present in natural language input, such as syllable lengthening or variations in pitch, were eliminated to determine whether distributional cues alone were sufficient for 8-month-old infants to parse the continuous stream into its underlying components. To be clear, this was a critical test of the *sufficiency* of statistical cues and not a test of whether these cues are the sole determinant for parsing streams of speech into their underlying auditory word forms. As noted by Saffran et al., “Although experience with speech in the real world is unlikely to be as concentrated as it was in these studies, infants in more natural settings presumably benefit from other types of cues correlated with statistical information” (p. 1928).

In fact, the design of Saffran et al. (1996) did not provide definitive evidence of using transitional probabilities. This is because each statistically coherent triplet occurred three times more often than each less coherent triplet that spanned a word boundary and served as the part-word test. That is, the joint probability of words (Syl-1, Syl-2, Syl-3) was higher than the joint probability of Syl-3, Syl-1, Syl-2). To unconfound transitional probability from joint probability, Aslin, Saffran, and Newport (1998) varied the frequency of occurrence of the statistically coherent triplets. By doubling the frequency of two of the four words, and using the part-words created when these two high-frequency words abutted each other, the frequency of occurrence of the lower frequency words and the higher frequency part-words was exactly equated. Yet, despite identical joint probabilities, the transitional probabilities within the tested words and part-words remained different. Words continued to have Syl-1 to Syl-2 and Syl-2 to Syl-3 transitional probabilities of 1.0, but part-words had transitional probabilities of 0.5 and 1.0. Although this difference was subtle and joint probabilities were equated, 8-month-old infants showed the same pattern of listening times on the posttest as in Saffran et al. These findings confirm that infants can use transitional probabilities even in the absence of other statistical cues that typically cooccur in natural languages.

To be clear, we are not now claiming, nor have we ever claimed, that the sole or even the primary distributional cue for word segmentation is syllable transitional probabilities. Raw frequency is undoubtedly a more robust source of distributional information and much easier to compute than any conditionalized statistic. Our reason for emphasizing the importance of transitional probabilities is that they are not subject to the errors of prediction that arise from lower-order statistics. Many high frequency sequences (e.g., “How are you?”) are not single units; the best test of their structure is the ability of the parts to appear in other sequences as well as together (e.g., “How is John?,” “Are they running?”). This

variation in sequencing, as contrasted with the greater stability of units that are actually parts of a single word, is what transitional probability measures. We will take up this important point about predictiveness with regard to grammatical categories in a subsequent section. But for now, let's complete this brief review of statistical learning for parsing a stream of speech into its underlying words.

Constraints on Statistical Learning

One could ask, following on the Saffran et al. (1996) findings, whether this ability to segment words from fluent speech is a specialization for linguistic materials. Two lines of research provided strong evidence against this language-specific hypothesis. First, Saffran, Johnson, Newport, and Aslin (1999) substituted 12 tones from a single musical octave for the 12 syllables used in Saffran et al. (1996) while keeping the statistical structure identical. The results from 8-month-olds were the same as with speech syllables: infants discriminated statistically coherent tone-triplets from slightly less coherent tone-triplets (analogous to words and part-words). Beyond the auditory modality, Kirkham, Slemmer, and Johnson (2002) showed that 2-, 5-, and 8-month-olds, after viewing a repeating sequence of eight temporally paired visual shapes, discriminated statistically coherent shape-pairs from randomly ordered shape-pairs. In further studies of the visual modality, where the statistical coherence between shapes was spatial rather than temporal, 8-month-olds also extracted the statistically coherent shape-pairs (Fiser & Aslin, 2002; Wu, Gopnik, Richardson, & Kirkham, 2010).

The second line of research suggesting that statistical learning is not a language-specific ability comes from studies of nonhuman species. Toro and Trobalon (2005) showed that rats could parse streams of human speech based on their statistical coherence, and many other studies of animals (e.g. Gentner, Fenn, Margoliash, & Nusbaum, 2006) have shown sensitivity to the temporally ordered statistics of auditory stimuli. Thus, the powerful statistical learning "engine" used to parse temporally ordered and spatially arranged elements across a corpus of input appears, at least to some degree, to be modality, domain, and species general.

Of course, only humans acquire natural languages, so one can ask what factors limit language acquisition. Are humans uniquely endowed with an additional special mechanism that is designed to acquire language (or a more powerful version of this statistical-learning engine than the one employed in other domains and species), or are there other constraints that are not language specific, but nevertheless prevent other species from acquiring (and using)

a linguistic system? It is important to note that this one simple transitional probability mechanism would not, in any case, be adequate to acquire the complex structure of human languages (cf. Chomsky, 1957; Saffran et al., 1996). Perhaps the more complex statistical computations called into play for acquiring other aspects of linguistic structure are unique to human language acquisition. There are also undoubtedly nonlinguistic factors (e.g., social pressures or skills for communicating complex propositions) that may differentiate humans from nonhuman species.

Whatever the domain and species status of these mechanisms turns out to be, it is clear that statistical learning must operate with some rather severe constraints in order to be tractable. This necessity of constraints comes from the so-called *computational explosion* problem that emerges when one considers how many statistics could be computed from any reasonably complex set of inputs. In the simple case of Saffran et al. (1996), with only 12 CV syllables, there are still a large number of things that could be computed: the frequencies of each syllable, each syllable pair (bigram), and all higher-order N-grams; the forward and backward transitional probabilities of all adjacent syllables; the joint and transitional probabilities of all nonadjacent syllables, and so on. With a larger and much more complicated inventory of elements, the number of computations would quickly become intractable. What are the limitations or constraints that might make this type of statistical mechanism focused enough to be helpful in word segmentation (and not become an intractable search for the “right” statistics)?

Since the publication of Saffran et al. (1996) and Aslin et al. (1998), a large number of useful constraints on learners’ computations have been identified. These include *attentional constraints* (Toro, Sinnett, & Soto-Faraco, 2005; Turk-Browne, Junge, & Scholl, 2005), *preferences for certain types of units* as the basic elements of computation (Newport & Aslin, 2004; Bonatti, Peña, Nespor, & Mehler, 2005), *Gestalt perceptual principles* that favor relations among certain types of elements (Creel, Newport, & Aslin, 2004; Endress & Bonatti, 2007; Shukla, Nespor, & Mehler, 2007), *language-specific prosodic groupings* (Johnson & Jusczyk, 2001; Thiessen & Saffran, 2003), and *gaze or action cues* that direct learners’ attention between some units or relations over others (Baldwin, Andersson, Saffran, & Meyer, 2008; Yu, Ballard, & Aslin, 2005). Each of these constraints directs learning toward linguistically useful relations that occur in natural languages of the world. Though some have been demonstrated only in adults rather than infants, they suggest a variety of constraints on statistical learning that could, in principle, render it tractable as a mechanism for word segmentation and domain-general parsing of temporally ordered stimuli.

Another concern about laboratory-based studies of statistical learning is whether they could plausibly scale-up to natural corpora to which infants are actually exposed in the real world. We cannot provide a definitive answer to that question at present, but some evidence from computational analyses of parental speech directed to infants and young children and larger scale studies of adults lend support to the plausibility of this scaling-up. Frank, Tenenbaum, and Gibson (2013) had adults listen on their iPod over a 2-week period to approximately 1,000 nonsense words embedded in short utterances. On a posttest one day after their last exposure, subjects were able to reliably discriminate syllable strings of variable lengths that had been statistically more coherent compared to those that were less coherent in terms of the underlying statistics of the nonsense words in the corpus. Swingley (2005) analyzed a corpus of child-directed speech and computed several statistical metrics, including transitional probabilities and mutual information, under various assumptions (e.g., using a minimum frequency criterion). The same general pattern of higher statistical coherence within words and lower coherence across word boundaries was revealed, mirroring the statistical structures in the lab-based studies of Saffran et al. (1996) and Aslin et al. (1998), though of course with considerably more variability.

Another question one can ask about the word segmentation studies is whether they represent an early stage of prelexical development, after which true word learning commences. This staged view of language acquisition is quite common in the field and was tested in a study by Graf Estes, Evans, Alibali, and Saffran (2007). They had 17-month-olds listen to a stream of speech syllables as in Saffran et al. (1996), and then the infants were introduced to a word-referent mapping task. In the word-referent mapping task, infants could associate either an isolated presentation of a word or an isolated presentation of a part-word with a novel visual object. At issue was whether a statistically coherent auditory word form that had been extracted from a continuous stream of speech in the previous task, would be a better token for subsequent word learning in a referential context. The answer was yes: infants learned to associate the statistically coherent word with the referent, but they did not associate the less statistically coherent part-word with the referent.

Although the staged model of language acquisition must be true in the limit (e.g., there must be at least one putative auditory word form before any mapping with a referent can occur), it does not necessarily imply that a large inventory of segmented words must enter the lexicon before systematic mapping to referents begins. In fact, McMurray (2007) has provided a compelling model of early word learning that suggests a highly parallel process of

simultaneous word-referent mapping in the very early stage of lexical development. To further examine this hypothesis of parallel processes of word-segmentation and referent mapping, Shukla, White, and Aslin (2011) conducted a study of 6-month-olds that combined these two processes in the same experimental session. Infants viewed a multiobject visual display within which a single object (e.g., the red ball) underwent motion to attract the infant's gaze. As the visual display was being viewed, the auditory stimuli were being presented, consisting of a family of short sentences composed of nonsense syllables. The syllables were structured across sentences so that some syllables were statistically coherent, and others were less coherent, just as in the words and part-words from Saffran et al. (1996). But the sentences had the normal prosodic patterns of natural language, with the language-universal property of intonational phrases (falling pitch and word-final lengthening). For one group of infants, the end of a statistically coherent target word was aligned with the end of an intonational phrase; for a second group of infants the statistical and prosodic information was misaligned (i.e., the word boundary was *within* an intonational phrase, as in "The pretty ba, by ate her cereal"). After exposure, infants were tested with repetitions of isolated words and part-words while the visual display showed multiple objects, none of which were moving. The dependent measure was how long the infants looked at the target object (i.e., the red ball) in the presence of word and part-word labels. Only the infants in the exposure condition where the statistical and prosodic information was aligned looked reliably longer at the target object. These results provide evidence that (1) infants at 6 months of age can segment auditory word forms from sentences, (2) these infants are biased to segment words based on statistical information when that information is packaged within an intonational phrase, and (3) while segmentation is happening, infants are also treating the putative auditory word form as a label for the object that was most salient during the exposure phase.

As impressive as these multiple language-processing skills are in 6-month-olds, they also highlight a potentially larger point about how we design experiments. There is a long tradition in experimental psychology of holding the myriad of independent variables constant and only allowing one of these variables to be manipulated at a time. That is precisely the design strategy used by Saffran et al. (1996). But this is not the way that these many factors covary in the natural environment. By creating such artificial designs, we may be creating experimental contexts in which infants are confused by the *absence* of the normal correlation among factors, thereby reducing their performance. This conjecture suggests that we should expand our range of experimental designs

to avoid overly constrained contexts and thereby reveal how infants are affected by covarying factors within a normative operating range.

Beyond Word Segmentation

A surprising amount of attention has been directed to the word-segmentation problem since the Saffran et al. (1996) article appeared (Google Scholar indicates that it has been cited over 2,500 times). But advances in understanding language acquisition will require moving beyond the simple mechanism proposed there—potentially adequate for word segmentation, but certainly not powerful enough for learning other more complex aspects of language. One notable effort in moving beyond the simple Saffran et al. mechanism was the study by Marcus, Vijayan, BandiRao, and Vishton (1999). In contrast to the continuous streams of speech used by Saffran et al., Marcus et al. presented nonsense syllables in short strings. While every string was composed of only three syllables and the corpus contained an entire inventory of only eight syllables, the crucial focus was that the materials combined to form a *pattern*: such as AAB or ABB. During the exposure phase, 7-month-olds heard multiple examples of three-syllable strings, such as *leledi*, *wewewi*, and *dededi*, all following an AAB pattern. Then in a posttest infants heard two types of strings: AAB strings composed of an entirely new set of syllables, or ABA or ABB strings composed of these same new syllables. The results showed that infants listened longest to the novel pattern (ABA or ABB) and not as long to the familiar pattern (AAB), even though both were composed of unfamiliar syllables. Importantly, infants could not solve this pattern-matching problem merely by computing joint or transitional probabilities among the syllables they had heard, because none of the test strings had any familiar syllables (i.e., all transitional probabilities were zero). Rather, infants must have extracted some common and more abstract rule, such as AAB, to distinguish between the test strings.

Marcus et al. (1999) argued that statistical learning operated only at the level of surface statistics, whereas rule learning operated at a deeper level involving abstract patterns. This exciting finding drew attention to the fact that statistics alone cannot provide a complete description of language learning. It also resurrected the foundational debate between Chomsky (1959) and Skinner (1957), although now in a slightly more modern form regarding higher-order statistics and how learners might solve the poverty-of-the-stimulus problem. What allows a learner to induce a rule based on sparse evidence? This is old

territory in the animal learning literature that was couched in the terminology of “gradients of generalization.” Why does a pigeon who has been trained to peck an orange key to receive a food pellet continue pecking a purple key but not a green key? One could argue that generalization is based on sensory similarity, but what aspect of the AAB pattern is “sensory”? Is the mechanism of generalization from *leledi* to *gagabu* merely the result of encoding the strings as sharing an initial repetition and not encoding the identity of the syllables themselves?

This type of rule learning is not unique to language materials. Saffran, Pollak, Seibel, and Shkolnik (2007) and Johnson et al. (2009) both showed that AAB rule learning is present for visual materials, although Marcus, Fernandes, and Johnson (2007) showed that such rules are more readily learned when instantiated in speech materials than in nonspeech materials. AAB rule learning has also been demonstrated in rats (de la Mora & Toro, 2013; Murphy, Mondragon, & Murphy, 2008). There is thus ample evidence that the generalization seen in these simple rule learning paradigms is modality, domain, and species general, just as with simple transitional probability-type statistical learning. Moreover, the detection of repetition-based rules is a robust mechanism present even in newborns (Gervain, Macagno, Cogoï, Peña, & Mehler, 2008).

Gerken (2006) provided some important insights into the nature of AAB rule learning in language materials by conducting a follow-up experiment that presented two groups of infants with different subsets of the three-syllable strings used by Marcus et al. (1999). One group heard 4 of the 16 strings, but each AAB string was entirely unique: *leledi*, *wiwije*, *jijili*, *dedewe*. The other group heard a different subset of four strings: *leledi*, *wiwidi*, *jijidi*, *dededi*—which also followed the AAB pattern, but all ended in the same syllable *di*. Notice that in both groups, the four strings conform to the AAB rule. But the second group has a less variable B syllable. Thus, for this second group, an alternative “rule” is: AA + ends in *di*. Infants in the first group performed on the posttest like infants in Marcus et al.—they detected the familiarity of the AAB rule and listened longer to the ABB rule violation. But infants in the latter group did not generalize to novel strings if they failed to end in *di*. That is, infants in the latter group acted as if the rule was more narrowly defined as “AA + ends in *di*” and not the broader “AAB.”

An important issue raised by the Gerken (2006) results is whether Marcus et al. (1999) should be interpreted as engaging a rule-learning mechanism fundamentally different from statistical learning, or rather that every learning task involves a gradient of generalization based on the most likely abstraction

indicated by the composition of the exposure corpus. Aslin and Newport (2012) have argued that the evidence across many subsequent studies suggests a single statistical learning mechanism with a gradient of generalization.

Beyond Repetition-Based Rule Learning

A crucial next step in understanding how a distributional-learning mechanism could be utilized in natural languages begins with defining what types of rules are needed. The essential idea behind the Marcus et al. (1999) study and all of its follow-ups, as well as many studies of natural language learning, is that words form grammatical categories, such as noun, verb, and adjective, and the basic syntax of every language is framed in terms of the order and phrasal groupings of these categories of words. The language learner is therefore confronted with the tasks of (a) discovering how many grammatical categories there are in the natural language spoken by the infant's parents and (b) correctly assigning words to the appropriate category. In the artificial language designs used by Marcus et al. and many others, categories are defined by the repetition of identical words. But in natural languages this is not how categories are defined; they are defined by their functional roles in the underlying grammar. A correlate of such roles is the patterning of adjacent and nonadjacent words. For example, determiners such as "a" and "the" always precede nouns in English, although there may be intervening adjectives, as in "the blue car" (but not "blue car the"). Mintz, Newport, and Bever (2002) and Mintz (2002, 2003) have proposed that a first step in inferring grammatical categories is keeping track of which words come before and after each word, thereby building up an inventory of cooccurrence statistics for the *relative* position of a word and seeing how these statistics pattern with those of other words in the language. Crucially, one should not rely solely on the *absolute* position of a word in an utterance (e.g., "the" in utterance-initial position, "blue" in second position, "car" in third position), because absolute position varies in natural languages and is irrelevant to correct category assignment.

In a series of experiments with adults, Reeder, Newport, and Aslin (2013) asked whether such distributional information alone—the patterning of words that occur before and after each target word—was sufficient to enable learners of an artificial grammar to solve the two crucial tasks of inducing the number of categories and determining which words are assigned to each. In addition, the experiments were designed to assess the gradient of generalization alluded to earlier in the rules versus statistics debate. That is, are learners who simply

listen to a corpus of speech able to acquire rudimentary information about grammatical categories and, when tested with novel utterances, can they appropriately generalize the category rules to predict the correct contexts for words they have not yet heard in that context? Crucially, how do learners adjust their tendency to generalize when the evidence in the corpus is sparse? Do they withhold generalization when there are consistent gaps, or do they generalize, implicitly assuming that the gaps are due merely to small samples drawn from the input?

The designs of the Reeder et al. (2013) experiments were all based on a core paradigm in which adults listened to a set of utterances that conformed to a (Q)AXB(R) grammar. Each of these letters refers to a category of words, with optional categories indicated by parentheses. A set of nonsense words, some consisting of single syllables (e.g., *glim*, *zub*) and others bisyllabic (e.g., *fluggit*, *klidam*), was assigned to each category, with two words in the Q and R categories and three words in the A, X, and B categories. Thus, the corpus consisted of three-, four-, and five-word utterances. The optional Q and R words were included to prevent learners from simply relying on absolute position in the utterance. Even in this small language, the total number of possible utterances (given the number of words per category) was $2 \times 3 \times 3 \times 3 \times 2 + 2 \times 3 \times 3 \times 3 + 3 \times 3 \times 3 \times 2 + 3 \times 3 \times 3 = 243$.

Crucially, in order to test for generalization, some of the grammatical utterances were withheld from the exposure corpus. Because the focus was on category X, as defined by the surrounding A and B words, a simplified illustration of the design is shown in Table 1. In Experiment 1, one-third of the $3 \times 3 \times 3 = 27$ possible AXB strings were withheld so that 9 novel strings were available for testing after exposure. During testing, adults heard three types of strings: (a) familiar grammatical utterances that had been presented during the exposure phase, (b) novel grammatical utterances that were withheld from the exposure phase but conformed to the underlying grammar, and (c) ungrammatical utterances that violated the ordering constraints of the grammar (e.g., AXA or BXB).

The measure of grammaticality obtained during the test was based on a 5-point rating scale, with 5 = *highly familiar and heard in the exposure phase* to 1 = *never heard and not a part of the rules that generated the corpus*. If learners acquired the underlying grammar—which entails learning the number of categories (despite utterances varying from three to five words in length) and also learning which words are assigned to the appropriate category—then they should rate familiar and novel grammatical strings highly and rate ungrammatical strings much lower. That is precisely what Reeder et al. (2013)

Table 1 List of AXB strings in Experiment 1 from Reeder et al. (2013). In Experiment 2, the number of exposure was reduced to 9, and the number of withheld strings was increased to 18

Exposure Strings	Withheld Strings
A1 X1 B1	A1 X1 B2
A1 X1 B3	A2 X1 B1
A2 X1 B2	A3 X1 B3
A2 X1 B3	A1 X2 B1
A3 X1 B1	A2 X2 B3
A3 X1 B2	A3 X2 B2
A1 X2 B2	A1 X3 B3
A1 X2 B3	A2 X3 B2
A2 X2 B1	A3 X3 B1
A2 X2 B2	
A3 X2 B1	
A3 X2 B3	
A1 X3 B1	
A1 X3 B2	
A2 X3 B1	
A2 X3 B3	
A3 X3 B2	
A3 X3 B3	

obtained from their adult learners. Figure 1 summarizes this organization of the words into A, X, and B categories from the patterns of strings in the exposure corpus. Importantly, there was no significant difference in ratings of familiar grammatical and novel grammatical strings, indicating that participants generalized fully to the withheld utterances. In a second experiment, Reeder et al. increased the sparsity of the input by withholding two-thirds of the possible AXB strings.

An important feature of the design shown in Table 1 is that despite withholding one-third or two-thirds of the possible AXB strings from the exposure corpus in Experiments 1 and 2, all possible AX and XB bigrams were present. That is, the corpus was balanced in terms of *coverage* of the adjacent words in the categories of interest (A, X, and B). In Experiments 3 and 4, however, the coverage of adjacent words in the corpus was made unbalanced by creating systematic gaps in the bigrams. This allowed us to examine how learners would generalize across these gaps and also to determine when they stopped

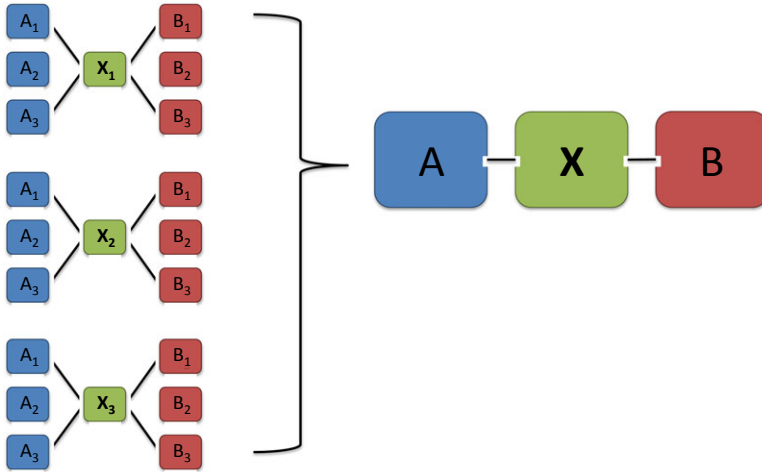


Figure 1 The distribution of nonsense words in Experiment 1 from Reeder et al. (2013). The three words that preceded and followed the three X-words provide robust evidence for categories A, X, and B.

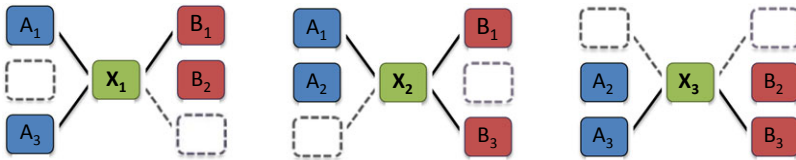


Figure 2 The distribution of nonsense words in Experiment 3 from Reeder et al. (2013). Only two words preceded and followed the three X-words, providing systematic gaps in the evidence for categories A, X, and B.

generalizing. As shown in Figure 2, one-third of the AX or XB bigrams never appeared in the exposure corpus. In contrast to Experiments 1 and 2, adults in the test phase now showed reliable evidence of rating novel grammatical strings as less acceptable than familiar grammatical strings. This indicates that with systematic gaps in the input, learners restrict generalization. Perhaps most revealing, when the same input corpus with systematic gaps was repeated three times, learners judged novel grammatical strings as even less acceptable, indicating that the more reliable the gaps are in the input, the more likely learners are to judge those strings as ungrammatical. This pattern of results shows that learners within the same paradigm can shift from broad generalization

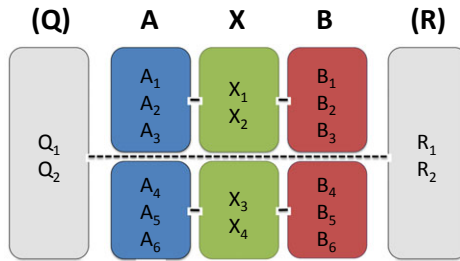


Figure 3 The distribution of nonsense words in the subcategorization experiment from Reeder et al. (2009). Separate sets of three words preceded and followed two subsets of X-words, providing robust evidence for two A, X, B subcategories.

to more restricted and lexically specific learning, depending on the pattern of distributional information they receive.

Two additional experiments lend further support to the robustness and flexibility of the distributional learning mechanism employed by adults in these grammatical category-learning experiments. Schuler, Reeder, Newport, and Aslin (2014) modified the (Q)AXB(R) design by introducing variation in the frequency of words within each category (approximating the Zipfian distribution of word frequencies seen in natural languages). This word-frequency variation presents an additional challenge to the unsupervised learner: does the sparser evidence of certain word combinations signal a significant “gap” in the corpus (perhaps because these combinations are ungrammatical), or does the sparser evidence simply reflect the variation in word frequency? That is, do learners rely on combinatorial probabilities rather than the absolute frequency of word sequences? The results suggest that adults learners are able to do this more complex computation: they rated the withheld novel grammatical strings just as highly as the familiar grammatical strings regardless of word frequencies, despite a 3:1 ratio of word frequencies within each category.

Perhaps the strongest test of grammatical category learning comes from Reeder, Newport, and Aslin (2009, 2014). Their design was similar to the foregoing experiments, except that the grammar was divided into two subsets, analogous to subcategories in natural languages (e.g., transitive and intransitive verbs or feminine vs. masculine nouns). As shown in Figure 3, half of the AXB strings had one set of A, X, and B words and the other half of the AXB strings had a different set of A, X, and B words—though all sentences had the same basic format and could begin and end with the same Q and R words. If learners simply treated this corpus as a slightly more complicated version

of gaps due to sparsity—or if they confused or did not remember the precise combinations of A, X, and B words in the corpus—they would incorrectly generalize to novel strings that “crossed the subcategory divide” (e.g., $A_1X_2B_5$). On the other hand, if learners were overly sensitive to the precise A, X, and B distributional information, it might affect their judgments of gaps present *within* a subcategory, leading them to rate a novel grammatical string that observed the subcategory structure as less acceptable. The results of this experiment were quite clear: adults learned to distinguish the subcategories, and yet they were still able to generalize appropriately within each subcategory. That is, when presented with strings that *cross* a subcategory, learners rated them as less acceptable than grammatical strings. However, when presented with novel grammatical strings from within a subcategory, they rated them as highly as familiar grammatical strings. These results demonstrate, in a single experiment, that adult learners *both* restrict generalization and generalize appropriately from the same corpus; that is, they extract the correct underlying grammatical structures, given distributional information that represents these structures.

Of course, adults are not the target population for studies of language acquisition. Thus it is important to extend these studies of artificial grammar learning to children and eventually to infants. We have completed similar experiments with 5- to 7-year-old children using slightly simpler AXB grammars (Newport, Schuler, Reeder, Lukens, & Aslin, 2014). To engage the children, the artificial grammar listening task was embedded in an “alien language” context with interesting but irrelevant video displays to maintain the children’s interest. To ensure that the children were actually listening (and not just attending to the video displays), they had to perform a one-back monitoring task in which sentences were intermittently repeated and the children had to indicate when the alien repeated himself. Despite these design modifications and the challenges of testing children, their performance on the rating tasks after exposure was remarkably similar to the performance of adults across the series of experiments. They judged novel grammatical strings as just as acceptable as familiar grammatical strings when the proportion of withheld strings was one-third or two-thirds but the corpus was balanced. And when the balance of the bigrams contained systematic gaps, like adults they judged novel grammatical strings as less acceptable than familiar grammatical strings. Thus, young children show the same overall pattern of distributional learning as adults—generalizing or withholding generalization depending on the patterning of the input.

Models of Category Learning and Generalization

An important aspect of these studies of grammatical category-learning in adults and children is that they illustrate how a single mechanism of distributional learning can lead to what appears to be “surface statistical learning” under some circumstances and “abstract rule learning” under other circumstances. Yet all of the experiments were, to the naïve learner, identical. It is implausible that participants engage a fundamentally different mechanism in one experiment than in another, and even more implausible that they use different mechanisms for specific test items (as in the subcategorization experiment). Thus, our hypothesis is that there is a *single* distributional learning mechanism that exhibits a *gradient* of generalization, depending on a “rational” interpretation of the patterning of the input. When the bulk of the evidence supports generalization—because it is plausible that gaps are due to chance, given their inconsistent distribution or low frequency—then learners judge novel strings as grammatical. But when the evidence begins to raise doubt (implicitly, of course) about whether the gaps in the input strings are random—that is, when gaps are highly systematic and consistent over a large input corpus—learners begin to restrict generalization.

How might we think about such a gradient mechanism of generalization? Recall that a key feature of the design of the Reeder et al. (2013) and follow-up experiments was variation in the distribution of bigrams (i.e., AX and XB word-pairs). Perhaps such a lexical bigram model is sufficient to account for the observed variations in how strongly adults generalize (or not). In an extensive analysis of the lexical bigram model, Qian, Reeder, Aslin, Tenenbaum, and Newport (2014) have shown that it works quite well even when two-thirds of the strings are withheld during the exposure phase when the coverage of bigrams is balanced. The lexical bigram model even does a reasonable job of fitting the data from the subcategorization experiment because the two sets of bigrams are nonoverlapping. But the lexical bigram model begins to fail when the input becomes supersparse, and particularly when bigrams are not fully shared across words. For example, if a new X-word appears in only one of the bigram contexts that are shared by other X-words, the lexical bigram model treats this new word as a lexical exception, not fully belonging to the X-word category and not generalizing to the other X-word contexts in which it has not been seen. In contrast, if there is a recurring cluster of contexts that most X-words appear in and if a new word appears in one of these, adult learners will generalize the whole set of contexts to the new word, as if it were fully a part of the X-word category. Children show the same pattern. One way to account for this finding is to create a model that has a “latent” level (i.e., an unobserved

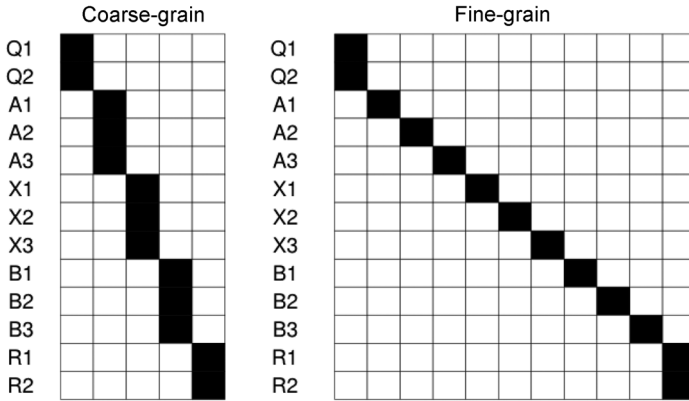


Figure 4 Model predictions from Qian et al. (2014) when a bias parameter is allowed to vary from a coarse to a fine granularity in how gaps in the input are fit to the underlying category structure.

category level), rather than generalizing solely on surface bigrams. In general, of course, adding more parameters to a model is not a principled way to better fit a complex set of data. The key insight of the Qian et al. model, however, is that by adding a latent category level and allowing the model to adjust the “grain” of how broad or narrow the categories are, a *single* model can account for the gradient nature of the generalization observed across all of the Reeder et al. experiments without custom fitting the parameters of the model to each experiment (see Figure 4).

Of course, it remains to be seen whether the Qian et al. (2014) model provides a realistic account of grammatical category learning in much younger learners (i.e., infants and toddlers). Although there is ample evidence that infants can form categories, it is not clear if they have the kind of flexibility seen in adults and older children to utilize a single mechanism of category learning to adjust the process of generalization in a gradient manner. But in the absence of clear evidence for multiple mechanisms of distributional learning, we favor our current working hypothesis that a single mechanism, with a latent category level and parameters that take into account the patterning of the input, provides the most parsimonious and flexible model of statistical language learning—at least for learning word categories as well as segmenting words from the speech stream. In ongoing work we address the ways in which this type of model must be expanded to cover more complex aspects of language, as well as other types of complex serial learning.

References

- Aslin, R. N., & Newport, E. L. (2012). Statistical learning: From acquiring specific items to forming general rules. *Current Directions in Psychological Science, 21*, 170–176.
- Aslin, R. N., Saffran, J. R., & Newport, E. L. (1998). Computation of conditional probability statistics by 8-month-old infants. *Psychological Science, 9*, 321–324.
- Baldwin, D., Andersson, A., Saffran, J., & Meyer, M. (2008). Segmenting dynamic human action via statistical structure. *Cognition, 106*, 1382–1407.
- Bonatti, L. L., Peña, M., Nespor, M., & Mehler, J. (2005). Linguistic constraints on statistical computations: The role of consonants and vowels in continuous speech processing. *Psychological Science, 16*, 451–459.
- Chomsky, N. (1957). *Syntactic structures*. The Hague, Netherlands: Mouton.
- Chomsky, N. (1959). A Review of B. F. Skinner's verbal behavior. *Language, 35*, 26–58.
- Creel, S. C., Newport, E. L., & Aslin, R. N. (2004). Distant melodies: Statistical learning of non-adjacent dependencies in tone sequences. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 30*, 1119–1130.
- de la Mora, D. M., & Toro J. M. (2013). Rule learning over consonants and vowels in a non-human animal. *Cognition, 126*, 307–312.
- Eimas, P., Siqueland, E., Jusczyk, P. W., & Vigorito, P. (1971). Speech perception in the human infant. *Science, 171*, 303–306.
- Endress, A. D., & Bonatti, L. L. (2007). Rapid learning of syllable classes from a perceptually continuous speech stream. *Cognition, 105*, 247–299.
- Fiser, J., & Aslin, R. N. (2002). Statistical learning of new visual feature combinations by infants. *Proceedings of the National Academy of Sciences, 99*, 15822–15826.
- Frank, M. C., Tenenbaum, J. B., & Gibson, E. (2013). Learning and long-term retention of large-scale artificial languages. *PloS one, 8*(1), e52500.
- Gentner, T. Q., Fenn, K. M., Margoliash, D., & Nusbaum, H. C. (2006). Recursive syntactic pattern learning by songbirds. *Nature, 440*, 1204–1207.
- Gerken, L. A. (2006). Decisions, decisions: Infant language learning when multiple generalizations are possible. *Cognition, 98*, B67–B74.
- Gervain, J., Macagno, F., Cogoi, S., Peña, M., & Mehler, J. (2008). The neonate brain detects speech structure. *Proceedings of the National Academy of Sciences, U.S.A., 105*, 14222–14227.
- Graf Estes, K. M., Evans, J. L., Alibali, M. W., & Saffran, J. R. (2007). Can infants map meaning to newly segmented words? Statistical segmentation and word learning. *Psychological Science, 18*, 254–260.
- Kirkham, N. Z., Slemmer, J. A., & Johnson, S. P. (2002). Visual statistical learning in infancy: Evidence for a domain general learning mechanism. *Cognition, 83*, B35–B42.

- Johnson, E. K., & Jusczyk, P. W. (2001). Word segmentation by 8-month-olds: When speech cues count more than statistics. *Journal of Memory and Language*, *44*, 548–567.
- Johnson, S. P., Fernandes, K. J., Frank, M. C., Kirkham, N. Z., Marcus, G. F., Rabagliati, H., et al. (2009). Abstract rule learning for visual sequences in 8- and 11-month-olds. *Infancy*, *14*, 2–18.
- Marcus, G. F., Fernandes, K. J., & Johnson, S. P. (2007). Infant rule learning facilitated by speech. *Psychological Science*, *18*, 387–391.
- Marcus, G. F., Vijayan, S., BandiRao, S., & Vishton, P. M. (1999). Rule learning in 7-month-old infants. *Science*, *283*, 77–80.
- McMurray, B. (2007). Defusing the childhood vocabulary explosion. *Science*, *317*, 631.
- Mintz, T. H. (2002). Category induction from distributional cues in an artificial language. *Memory & Cognition*, *30*, 678–686.
- Mintz, T. H. (2003). Frequent frames as a cue for grammatical categories in child directed speech. *Cognition*, *90*, 91–117.
- Mintz, T. H., Newport, E. L., & Bever, T. G. (2002). The distributional structure of grammatical categories in speech to young children. *Cognitive Science*, *26*, 393–424.
- Murphy, R. A., Mondragon, E., & Murphy, V. A. (2008). Rule learning by rats. *Science*, *319*, 1849–1851.
- Newport, E. L., & Aslin, R. N. (2004). Learning at a distance: I. Statistical learning of non-adjacent dependencies. *Cognitive Psychology*, *48*, 127–162.
- Newport, E. L., Schuler, K. D., Reeder, P. A., Lukens, K., & Aslin, R. N. (2014). *Learning grammatical categories in an artificial language by 5- to 7-year-olds using distributional information*. Manuscript in preparation.
- Qian, T., Reeder, P. A., Aslin, R. N., Tenenbaum, J., & Newport, E. L. (2014). *Rational generalization via rational categorization*. Manuscript in preparation.
- Reeder, P. A., Newport, E. L., & Aslin, R. N. (2009). The role of distributional information in linguistic category formation. In N. Taatgen & H. van Rijn (Eds.), *Proceedings of the 31st annual meeting of the Cognitive Science Society* (pp. 2564–2569). Austin, TX: Cognitive Science Society.
- Reeder, P. A., Newport, E. L., & Aslin, R. N. (2013). From shared contexts to syntactic categories: The role of distributional information in learning linguistic form-classes. *Cognitive Psychology*, *66*, 30–54.
- Saffran, J. R., Aslin, R. N., & Newport, E. L. (1996). Statistical learning by 8-month-old infants. *Science*, *274*, 1926–1928.
- Saffran, J. R., Johnson, E. K., Aslin, R. N., & Newport, E. L. (1999). Statistical learning of tone sequences by human infants and adults. *Cognition*, *70*, 27–52.
- Saffran, J. R., Pollak, S. D., Seibel, R. L., & Shkolnik, A. (2007). Dog is a dog is a dog: Infant rule learning is not specific to language. *Cognition*, *105*, 669–680.

- Schuler, K. D., Reeder, P. A., Newport, E. L., & Aslin, R. N. (2014). *The effects of uneven frequency information on artificial linguistic category formation in adults*. Manuscript in preparation.
- Shukla, M., Nespors, M., & Mehler, J. (2007). An interaction between prosody and statistics in the segmentation of fluent speech. *Cognitive Psychology*, *54*, 1–32.
- Shukla, M., White, K. S., & Aslin, R. N. (2011). Prosody guides the rapid mapping of auditory word forms onto visual objects in 6-month-old infants. *Proceedings of the National Academy of Sciences*, *108*, 6038–6043.
- Skinner, B. F. (1957). *Verbal behavior*. Acton, MA: Copley.
- Swingle, D. (2005). Statistical clustering and the contents of the infant vocabulary. *Cognitive psychology*, *50*, 86–132.
- Thiessen, E. D., & Saffran, J. R. (2003). When cues collide: Use of stress and statistical cues to word boundaries by 7- to 9-month-old infants. *Developmental psychology*, *39*, 706–716.
- Toro, J. M., Sinnett, S., & Soto-Faraco, S. (2005). Speech segmentation by statistical learning depends on attention. *Cognition*, *97*, B25–B34.
- Toro, J. M., & Trobalon, J. (2005). Statistical computations over a speech stream in a rodent. *Perception & Psychophysics*, *67*, 867–875.
- Turk-Browne, N. B., Jungé, J. A., & Scholl, B. J. (2005). The automaticity of visual statistical learning. *Journal of Experimental Psychology: General*, *134*, 552–564.
- Wu, R., Gopnik, A., Richardson, D. C., & Kirkham, N. Z. (2011). Infants learn about objects from statistics and people. *Developmental Psychology*, *47*, 1220–1229.
- Yu, C., Ballard, D. H., & Aslin, R. N. (2005). The role of embodied intention in early lexical acquisition. *Cognitive Science*, *29*, 961–1005.