

# Sussex Research

## Word learning as the interaction of online referent selection and slow associative learning

Bob McMurray, Jessica Horst, Larissa K Samuelson

### Publication date

01-10-2012

### Licence

This work is made available under the **Copyright not evaluated** licence and should only be used in accordance with that licence. For more information on the specific terms, consult the repository record for this item.

### Citation for this work (American Psychological Association 7th edition)

McMurray, B., Horst, J., & Samuelson, L. K. (2012). *Word learning as the interaction of online referent selection and slow associative learning* (Version 1). University of Sussex.  
<https://hdl.handle.net/10779/uos.23391119.v1>

### Published in

Psychological Review

### Link to external publisher version

<https://doi.org/10.1037/a0029872>

### Copyright and reuse:

This work was downloaded from Sussex Research Open (SRO). This document is made available in line with publisher policy and may differ from the published version. Please cite the published version where possible. Copyright and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners unless otherwise stated. For more information on this work, SRO or to report an issue, you can contact the repository administrators at [sro@sussex.ac.uk](mailto:sro@sussex.ac.uk). Discover more of the University's research at <https://sussex.figshare.com/>

# Word Learning Emerges From the Interaction of Online Referent Selection and Slow Associative Learning

Bob McMurray  
University of Iowa

Jessica S. Horst  
University of Sussex

Larissa K. Samuelson  
University of Iowa

Classic approaches to word learning emphasize referential ambiguity: In naming situations, a novel word could refer to many possible objects, properties, actions, and so forth. To solve this, researchers have posited constraints, and inference strategies, but assume that determining the referent of a novel word is isomorphic to learning. We present an alternative in which referent selection is an online process and independent of long-term learning. We illustrate this theoretical approach with a dynamic associative model in which referent selection emerges from real-time competition between referents and learning is associative (Hebbian). This model accounts for a range of findings including the differences in expressive and receptive vocabulary, cross-situational learning under high degrees of ambiguity, accelerating (vocabulary explosion) and decelerating (power law) learning, fast mapping by mutual exclusivity (and differences in bilinguals), improvements in familiar word recognition with development, and correlations between speed of processing and learning. Together it suggests that (a) association learning buttressed by dynamic competition can account for much of the literature; (b) familiar word recognition is subserved by the same processes that identify the referents of novel words (fast mapping); (c) online competition may allow the children to leverage information available in the task to augment performance despite slow learning; (d) in complex systems, associative learning is highly multifaceted; and (e) learning and referent selection, though logically distinct, can be subtly related. It suggests more sophisticated ways of describing the interaction between situation- and developmental-time processes and points to the need for considering such interactions as a primary determinant of development.

**Keywords:** word learning, fast mapping, referential ambiguity, cross-situational learning, associative learning

**Supplemental materials:** <http://dx.doi.org/10.1037/a0029872.supp>

The word is fundamental to language. Words serve an organizing role in syntactic parsing (MacDonald, Pearlmuter, & Seidenberg, 1994; Tanenhaus & Trueswell, 1995), speech perception (Goldinger, 1998; McClelland, Mirman, & Holt, 2006), and semantic organization (Elman, 2009; Lupyan, Rakison, & McClelland, 2007; Mayor & Plunkett, 2010; Samuelson & Smith, 2000;

Waxman, 2003). Lexical items live at a critical juncture in language processing, linking sound, articulation, syntax, and meaning, and as a result, the acquisition of words has attracted enormous attention (P. Bloom, 2000; Carey, 1978; Fenson et al., 1994; Golinkoff et al., 2000; Mayor & Plunkett, 2010; Xu & Tenenbaum, 2007).

Whether such research examines the growth of the lexicon as a whole (e.g., Fenson et al., 1994; Ganger & Brent, 2004) or conducts microinvestigations of single word learning (e.g., Carey & Bartlett, 1978; Horst & Samuelson, 2008), the fundamental questions concern *word knowledge*: (a) whether children know a word, (b) how they come to know it, and (c) how many words they know. The typical article addressing this starts with the scale: Children acquire about 60,000 words in about 18 years. It then describes why this is so hard. Famously articulated by Quine (1960), in any naming situation there are infinite interpretations for an unknown word. Thus, children face a daunting task of ambiguity resolution that they must solve thousands of times.

Such an article then proposes an explanation for how children solve this problem, but often skips a primary question: What does it mean to know or learn a word? A canonical finding is that toddlers comprehend more words than they produce: Seventy-five

---

Bob McMurray, Department of Psychology and Delta Center, University of Iowa; Jessica S. Horst, School of Psychology, University of Sussex, Brighton, England; Larissa K. Samuelson, Department of Psychology and Delta Center, University of Iowa.

Writing of this article was supported by National Institutes of Health (NIH) Grants DC-008089 to Bob McMurray, RES-000-22-4451 to Jessica S. Horst, and HD-045713 to Larissa K. Samuelson. The content is solely the responsibility of the authors and does not represent the official views of NIH. We would like to thank Keith Apfelbaum and Marcus Galle for helpful discussions on associative learning, Joe Toscano for assistance with the servers and programming, Sarah Kucker and Libo Xhao for new modeling ideas, and Michael Spivey for insight into normalized recurrence.

Correspondence concerning this article should be addressed to Bob McMurray, Department of Psychology, E11 SSH, University of Iowa, Iowa City, IA 52242. E-mail: [bob-mcmurray@uiowa.edu](mailto:bob-mcmurray@uiowa.edu)

percent of 12-month-olds understand *all gone*, but it takes 8 more months before that many say it. When do we consider *all gone* known? Children can often identify the referents of novel words on their first exposure (Mervis & Bertrand, 1994), yet their ability to recognize familiar words develops over some time (Fernald, Pinto, Swingley, Weinberg, & McRoberts, 1998). Again, when do we consider a word known?

Perhaps it is not possible to quantify when, or if, a word is known. If so, the problem of lexical acquisition may be better framed in terms of how children learn to *use* words. After all, we can measure word use directly. Commonly, this idea calls to mind the process of producing words, but we mean something broader. To the extent that a word links sound and meaning, any time that link is used to guide behavior, a word is being used. Thus, word use also includes processes such as comprehending known words, and even determining referents for new words.

If we ignore the uncertainty of knowledge and focus on only children's word use, children must still solve a set of difficult problems. Yet, the concept of using a word does not appear in classic descriptions of word learning. Rather, the focus is on the information needed for learning, the amount that must be gathered, and the difficulties in gathering it. This has led to theoretical views that emphasize knowledge-based processes in accounting for learning but inadvertently deemphasize use (Golinkoff & Hirsh-Pasek, 2006; Golinkoff, Mervis, & Hirsh-Pasek, 1994; Mayor & Plunkett, 2010; Woodward & Markman, 1998; Xu & Tenenbaum, 2007). Our purpose here is to advance an account of the development of word use, both novel and familiar, over multiple timescales. We demonstrate its power with a computational model.

In developing our account, we start by discussing the standard view of word learning and the theoretical tensions surrounding it. We then distill word learning to the minimal computational problem to frame our account. Finally, we present an account based on associative learning and dynamic competition and demonstrate its power to illuminate lexical behavior. These simple mechanisms operate on different timescales: Dynamic competition describes the situation-time process of selecting a referent of a word, whereas associative learning describes the developmental-time process of slowly forming mappings between words and concepts. Critically, the interactions of these timescales can yield emergent power to describe lexical behaviors that does not derive from either one alone.

### The Standard View: Acquiring Lexical Knowledge

Early approaches to word learning used measurement studies to examine the number of words known over development (e.g., L. Bloom, 1973; Dore, Franklin, Miller, & Ramer, 1976; Fenson et al., 1994; Reznick & Goldfield, 1992). A key finding was that word learning appears to accelerate. The source of this is debated (P. Bloom, 2000; Ganger & Brent, 2004; McMurray, 2007; Nazzi & Bertoncini, 2003), but clearly children are efficient learners and may become more efficient over development. This contrasts with the apparent difficulty of word learning. A major obstacle to acquiring words is *referential ambiguity* (Quine, 1960): In any naming event, a novel word can refer to any object present, its properties, the speaker's feelings or intentions for it, an impending action, or something else altogether. Even considering only the

smaller problem of which object or category of objects is being referred to, this is still challenging.

The *constraint approach* offers a metatheory for solving referential ambiguity: Children have (perhaps innate) constraints, principles, or biases that help them infer a word's meaning by providing information not available in the situation (Golinkoff et al., 1994; Woodward & Markman, 1998). The most elementary constraints simply restrict the possible interpretations of a novel word (Markman, 1990), positing that new words refer to whole objects (not parts) or to basic-level categories (rather than super- or subordinate categories). More complex constraints such as social cues may go further, pinpointing the correct referent (Baldwin, 1991; Baldwin, Markman, Bill, Desjardins, & Irwin, 1996; Tomasello, Strosberg, & Akhtar, 1996).

Particularly relevant to the present study is the mutual exclusivity constraint (Markman & Wachtel, 1988) and the similar novel name–nameless category principle (N3C; Mervis & Bertrand, 1994), which describe how children infer the referent of a word on the basis of which other objects they have names for. For example, when presented with a familiar spoon and a novel whisk, children infer that *whisk* refers to the latter, if they know the word *spoon*. The form of this inference has been debated (Grassmann & Tomasello, 2010; Halberda, 2006; Jaswal & Hansen, 2006; Markman & Wachtel, 1988; Mervis & Bertrand, 1994), yet it is clear that children can make inferences that integrate available context with the contents of their lexica.

In the constraint approach, such inferences become the primary route to learning, and some have argued that it is the onset of these constraints or the related social–pragmatic skills that create the sudden acceleration in word learning (Golinkoff & Hirsh-Pasek, 2006; Golinkoff et al., 1994; Markman, 1990; Nazzi & Bertoncini, 2003). However, the constraint approach has been challenged on a number of theoretical and empirical grounds.

### Theoretical Challenges to Constraints

A classic concern with the constraint approach is its ability to scale up. Once children master basic-level terms, how is the taxonomic constraint relaxed to learn superordinates? When can children violate the whole-object constraint to learn properties? And how do children ignore mutual exclusivity to learn synonyms or superordinates? Similarly, the constraint approach does not provide a clear framework for how constraints interact or how conflicts are resolved.

This has led some to cast word learning in terms of general inference processes used for reasoning or social–pragmatic behavior. Mutual exclusivity, for example, applies to concepts (Behrend, Scofield, & Kleinknecht, 2001; P. Bloom & Markson, 1998; Markson & Bloom, 1997; Waxman & Booth, 2000) and can be described by principles such as deductive syllogism (Halberda, 2006). These more general-purpose problem-solving skills may avoid issues related to the relaxation of hard constraints. A powerful formulation of this is to frame the problem in terms of probabilistic (Bayesian) inference (Frank, Goodman, & Tenenbaum, 2009; Xu & Tenenbaum, 2007). In this view, constraints are prior probabilities that interact with context and existing knowledge using the laws of probability to determine the optimal solution. This can handle interacting constraints and allows constraints to be violated when the data permit. Others have cast the problem

in terms of social inference (Akhtar & Martinez-Sussman, 2007; Moore, 2006; Tomasello, 2001). This transforms the problem from determining what a word means to determining what a speaker is referring to. Mutual exclusivity, then, becomes a social inference process in which speakers are expected to follow social conventions and use familiar words (e.g., Diesendruck & Markson, 2001). Social and Bayesian accounts are not mutually exclusive. Both approaches are cast fundamentally in terms of acquiring knowledge about words, but say little about how words are used once (or while) this knowledge is acquired (but see Frank et al., 2009) or how these inferences and knowledge relate to long-term learning.

### Empirical Challenges to Constraints

The foundations of the constraint approach have been shaken by research on four topics: the vocabulary spurt, fast mapping, cross-situational learning, and familiar word recognition.

#### Vocabulary Spurt

The sudden acceleration vocabulary growth has been seen as indirect evidence for constraints by implying their sudden availability or a change in children's approach to word learning. However, Ganger and Brent (2004) argued that if the vocabulary spurt was a singular change, then the velocity profiles of individual children should show a sudden shift in velocity. Yet, for 33 of 38 children, a smoothly accelerating function fit better.

Alternatively, it is possible that principles such as mutual exclusivity are available all the time (Akhtar & Martinez-Sussman, 2007; Markman, Wasow, & Hanson, 2003; Tomasello, 2001), but children simply do not have enough words (or other knowledge) to use them. As the first few words are acquired, these exert greater leverage, allowing the rate of acquisition to increase (cf. Elman et al., 1996; van Geert, 1991). Recently, however, McMurray (2007; Mitchell & McMurray, 2009) demonstrated that acceleration is possible without such mechanisms. As long as learning proceeds in parallel and the distribution of easy and hard words includes few easy words, acceleration is guaranteed. Thus, although a change in underlying learning mechanism or constraints could account for acceleration, acceleration is not evidence for it.

#### Fast Mapping

In her original discussion, Carey (1978) contrasted children's quick mapping of a novel word to a novel referent (which was illustrated with the first demonstrations of mutual exclusivity), with a slower phase of learning the word's full meaning. It is not clear whether this "fast mapping" refers to partial, early stages of learning or purely in-the-moment referent selection (though the word *mapping* implies learning). Nevertheless, if word learning is due to fast mapping, then the act of selecting the referent should result in something being retained.

The most compelling test of this would be to ask first if the child selects the correct referent, and then examine retention when the child is retested in a neutral context. Prior studies failed to do this (see Horst & Samuelson, 2008): Some studies retest with a subsequent trial of the same sort, allowing children to simply solve the problem again (e.g., Mervis & Bertrand, 1994; Wilkinson & Mazzitelli, 2003), and others provide a short review of the name–

object linkages before testing retention (e.g., Goodman, McDonough, & Brown, 1998).

To address this, Horst and Samuelson (2008) presented children two known and one novel object and asked for the referent of a novel name. Children successfully selected the referent via mutual exclusivity. Five minutes later, however, they were unable to map that same to its referent when it was presented with other novel items they had just seen. This suggests that the linkage was not retained. Thus, the use of mutual exclusivity does not necessarily result in long-term learning. It is uncertain whether this generalizes to other task variants or other ages (e.g., Kucker & Samuelson, 2012; Spiegel & Halberda, 2011) or to other constraints. Nevertheless, it suggests that mutual exclusivity may simply bias the child toward the referent in the moment, and is not synonymous with learning. This questions a fundamental assumption of the constraint approach, that resolving referential ambiguity (via constraints) is the same as learning.

#### Cross-Situational Learning

If solving the problem of referential ambiguity is not the same as learning words, how do children do it? One possibility is statistical learning. In any novel naming situation the intended referent may be ambiguous. However, *across* situations there may be only one object consistently paired with a word. For example, while the word *dog* may occur with a dog, a ball, and a leash in one situation, later on it may be heard without ball or leash and with other objects. Over time the referent, dog, is likely to be the most frequently co-occurring object. Thus, at any given time, the child may not need to determine the referent—the child only needs to accumulate co-occurrence statistics to learn the mappings (Horst, McMurray, & Samuelson, 2006; McMurray, Horst, Toscano, & Samuelson, 2009; Siskind, 1996; Smith & Yu, 2008; Yu & Smith, 2007). If true, associative mechanisms (MacWhinney, 1987; Merriam, 1999; Regier, 2005) may suffice for word learning.

This idea had been examined computationally (Horst et al., 2006; McMurray, Horst, et al., 2009; Siskind, 1996), but Yu and Smith offered the first empirical tests. Adults (Yu & Smith, 2007) and infants (Smith & Yu, 2008) were exposed to small artificial lexica that contained such regularities across trials. Both groups successfully learned the word–object linkages from this alone. It is still not known whether such learning can handle categories of objects, and this may change the computational problem. However, this provides an important proof of concept, that statistical or associative learning may proceed without solving referential ambiguity and without constraints (though they may facilitate learning or in the moment language use).

#### Familiar Word Recognition

Finally, by focusing on information used to solve referential ambiguity, the constraint approach has little to say once the child has acquired word–object mappings (familiar words). However, familiar word recognition also changes over development, and it is not clear that this is related to constraints. Fernald et al. (1998) measured the amount of time it took infants to fixate the correct object (in a two-alternative forced-choice looking task) as a measure of the recognition of the word's meaning. This decreased dramatically over development, suggesting a tuning process for



known or recently learned words (Fernald, Perfors, & Marchman, 2006; Fernald et al., 1998). This decrease cannot be accounted for by attentional or oculomotor processes, as infants show no changes in purely visual tasks during this time, and performance in the visual task is not related to their speed of word recognition (Fernald et al., 2006). More importantly, this improvement cannot be characterized as simply refining an existing skill; rather, speed of processing predicts the rate of long-term learning (Fernald et al., 2006) and later linguistic and cognitive outcomes (Marchman & Fernald, 2008), suggesting it is a more fundamental property of word learning and use.

The constraint approach has no way to describe these changes, as it lacks a theory of how words are used. Since such changes are not likely driven by referential ambiguity, under the constraint view they require a separate developmental mechanism. To be clear, few existing accounts make a strong distinction between learning novel and familiar words. Studies such as Merriman, Lipko, and Evey's (2008) have examined how children decide that a word is novel or familiar (implying that such decisions may help engage the right learning or processing strategies). However, this is not an essential component of any major theoretical accounts. Yet at the same time, no major theories address the improvement in familiar word recognition, nor do they seem to have the theoretical tools to do so. Fernald and colleagues' work suggests an important developmental phenomenon that demands an explanation and should be linked to the word learning literature more broadly. If an account can handle both novel and familiar word learning with the same mechanism, this may offer a more parsimonious explanation of word learning.

### A New Direction

Although none of these findings completely rule out constraints, they paint a picture in which the solution to referential ambiguity is subtly independent of long-term learning. These problems are not limited to the constraint approach—any approach focusing exclusively on referential ambiguity and the information used to solve it will struggle to account for these findings. We need an account that emphasizes how novel and familiar words are used and builds from there to understand how this ability develops. Such an account cannot ignore referential ambiguity. However, it must move beyond it to account for development. Our thesis is that we may make more headway by considering behavior at two timescales. Referential ambiguity is a problem that children face in a given situation and must be solved in real time. This differs substantially from the problem of learning and retaining word–object mappings, which may unfold over many situations, and indeed over development.

### Distilling the Word

To develop this account, we first distill word use and word learning to their minimal computational components. We define them in terms of association and activation, processes that are independent of the information that contributes to word recognition and word learning. This distinction is not theoretically novel; it builds on constructs from cognitive development that have been most extensively developed by Munakata, McClelland, and colleagues working in the connectionist paradigm and by Thelen,

Smith, Schöner, and colleagues in dynamic systems theory (Munakata, 1998; Munakata & McClelland, 2003; Munakata, McClelland, Johnson, & Siegler, 1997; Smith, Thelen, Titzer, & McLin, 1999; Thelen, Schöner, Scheier, & Smith, 2001; see also Elman, 1990; Harm & Seidenberg, 1999; McMurray, Horst, et al., 2009; Spencer, Perone, & Johnson, 2009). Our goal here is to translate these concepts to word learning and to use them to develop an account that stands independent of a strongly theoretically connectionist stance, as our ideas are conceptually compatible with other approaches that are distinct from connectionism. This account necessarily oversimplifies many things. We discuss this later. However, it allows us to be precise about mechanisms to frame them computationally.

Word use and learning fundamentally concern the relationship between a phonological pattern and a semantic category. For present purposes, we ignore the complexities of mapping sounds to word forms and assume that the auditory system can identify discrete word forms. This is not trivial, but by the middle of the 2nd year, many of the basic properties of auditory word-form recognition are in place (Fernald, Swingle, & Pinto, 2001; Swingle, 2009; Swingle & Aslin, 2002). Similarly, we assume that infants can analyze a visual scene and categorize referents. This too is not trivial, but again, by the middle of the 2nd year, children appear adept at it (Bauer, Dow, & Hertsgaard, 1995; Behl-Chadha, 1996; Mareschal & Tan, 2007). To be clear, word learning involves mapping words to categories, not merely to individual objects (and there are excellent models that capture aspects of this: Mayor & Plunkett, 2010; Samuelson, 2002). Our goal here is to strip out these important processes to investigate the power of the associations themselves and the real-time processes that operate over them to form the basis of interesting word learning behavior. Thus, we will assume some categorization ability and focus on the mapping between word forms and categories, and for simplicity's sake we will often refer to objects and/or referents when what is meant is a category of objects or referents.

Figure 1 shows our distillation. Circles represent representations, and their shading indicates how strongly each is considered. Figure 1A shows the process of identifying the referent of a word in situation time, in a situation with two visual competitors. Initially, the system starts with every word form under partial consideration or activation (Figure 1A, left side), as nothing has been heard yet. Two objects (*dog* and *tree*) are active, reflecting the visual scene. As the word is heard, the system moves toward considering one word and object (bottom). This shift in activation represents the process of deciding what was heard and what should be attended, as in many interactive activation models (McClelland & Elman, 1986; Spivey, 2007).

Thus, resolving referential ambiguity is a matter of moving from consideration of multiple objects to one. This demands a solution in terms of activation or attention to referents, not learning. Such changes in real-time consideration of the referents could derive from external forces that decrease consideration or activation for incorrect objects or increase consideration of the correct one. In this example, if *dog* is unknown, eye gaze could add consideration for *dog*; mutual exclusivity could reduce consideration of *bug* (if it is known); or a context demanding animacy could rule out the *tree* (see Figure 1B). Such external forces may also include attention processes that facilitate (Fulkerson & Waxman, 2007; Samuelson

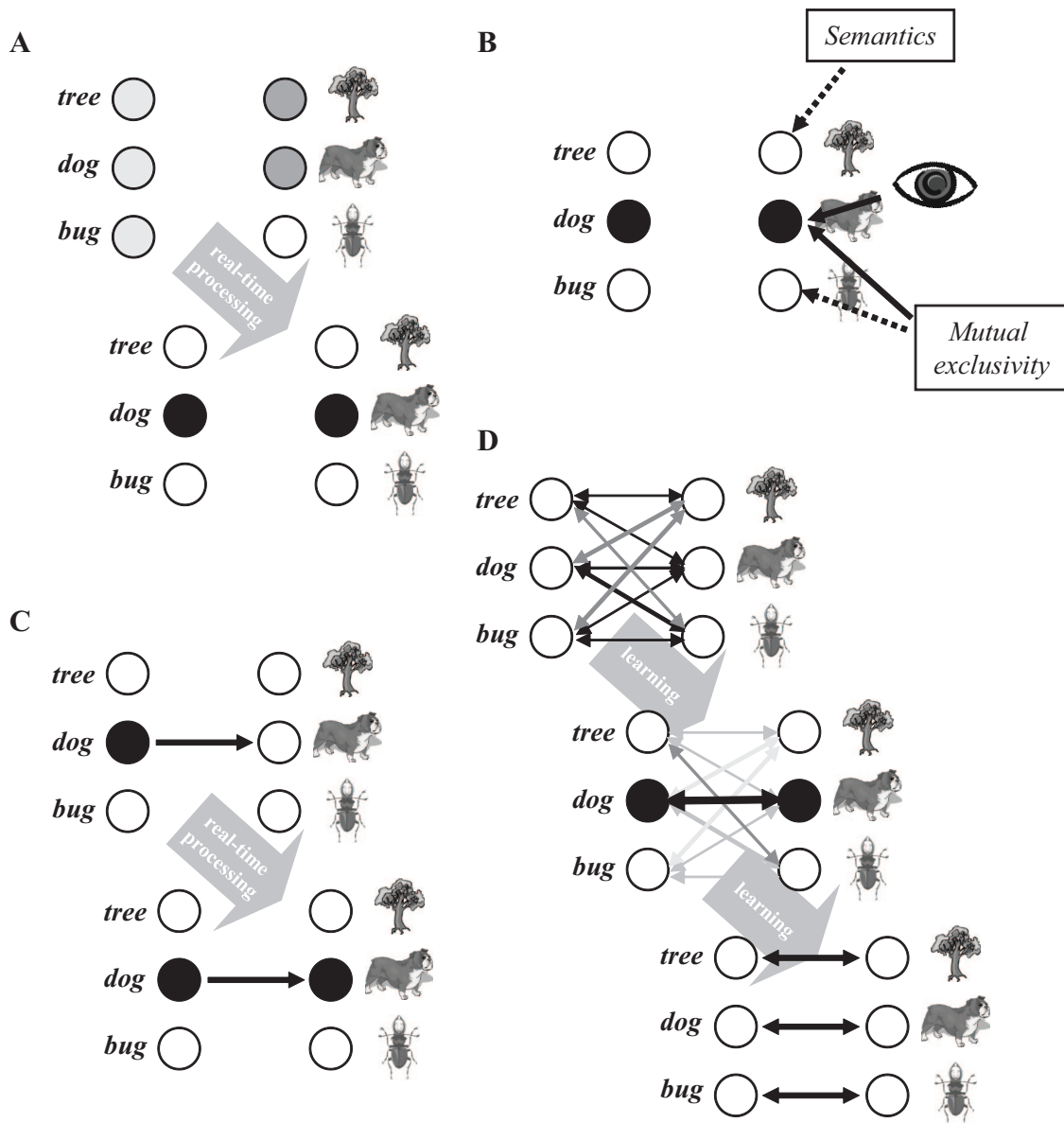


Figure 1. A schematic of word learning. On the left side of each panel are units showing how strongly a particular auditory word form is under consideration; on the right are units showing the strength of a visual object under consideration. (A) Solving the problems of word recognition and referential ambiguity requires a transition from a state in which multiple word forms and object categories are considered to a state in which only one is. (B) Constraints such as pragmatics and mutual exclusivity can act by simply changing the degree of consideration without affecting long-term linkages. (C) Familiar word recognition takes advantage of learned associations to activate object representations from spoken words. (D) Learning is instantiated as long-term linkages between words and objects that are strengthened when both word forms and objects are considered simultaneously.

& Smith, 1998) and/or interfere with word and category learning (e.g., Robinson & Sloutsky, 2004, 2007).

Crucially, none of this has to involve a mapping, nor does it entail learning. As long as the activation of the objects reaches the correct state, there is no need to retain anything—the child has arrived at the right inference. So how does learning occur? Taking typical connectionist assumptions, initially, the system starts with word forms and objects randomly connected—many possible map-

pings are under partial consideration or activation (see Figure 1D, left side), as nothing has been heard or seen yet. Over time, some of these connections will be pruned, and others will be strengthened. For example, the simultaneous consideration of a word (e.g., *dog*) and an object (a dog) could lead to links being strengthened if learning is associative, and the link between *dog* and *tree* (which is not present) being weakened (Figure 1D). Over time, such changes could build a system of links that encompasses many

words and objects. This changing and building of linkages is the result of learning, not external forces such as attention, or pragmatics that guide activation to solve the referential ambiguity problem in the moment. If such links become strong enough, such pathways will be employed when the word is heard again, allowing the word form, *dog*, to activate the appropriate concept without external support (see Figure 1C).

This suggests two distinct processes: (a) the use or recognition of a familiar or novel word (changing activation states) and (b) the changing of the links between the word form and visual referent (learning or the changes of connection weights). These processes can be described on different timescales. The problem of determining the referent of a novel word is a problem of usage. It must be solved very rapidly in situation time, that is, within the context of a single naming event, and it could take advantage of learned mappings or external support. In contrast, the problem of learning is solved over *developmental time*: It is a problem of acquiring lasting linkages between sound patterns and meanings that may take weeks or years. For familiar words, deploying these mappings to understand or produce a word is a situation-time phenomenon, but enhancing the efficiency of this occurs over developmental time. These problems rely on different theoretical mechanisms: Situation-time processes involve changes in activation, or the strength of consideration of particular words and objects, whereas developmental-time learning processes involve changes in knowledge, that is, associations between words and objects (Harm & Seidenberg, 1999; McMurray, Horst, et al., 2009; Munakata & McClelland, 2003; Munakata et al., 1997; Smith et al., 1999; Spencer et al., 2009; Thelen et al., 2001).

Although both problems are fundamentally about matching words and referents, the demands of these tasks are clearly different. The problem of finding a referent in the moment does not necessitate learning. Simply arriving at a state in which one word and one object are under consideration is sufficient, and does not require changing of the strength of the linkage. That is not to say that children can turn off learning, or that learning may not occur in some circumstances. Indeed, it is more parsimonious if learning is “always on.” Rather, we are arguing these are not problems that are solved by learning processes, and it may not matter how much or little is learned in a situation for a child to find a referent. If word–object linkages are ever created in situation time, these linkages do not need to be complete. Thus, children could make use of contextual cues to solve referential ambiguity for their immediate communicative needs, but not necessarily commit to a given mapping from one event.

Conversely, learning does not require the child to solve referential ambiguity. If multiple objects are under consideration, multiple linkages can be laid down. If this is done in small increments over multiple naming events, the more consistent ones could rise to the top, as in cross-situational learning (Smith & Yu, 2008; Yu & Smith, 2007). Such a process would need to be slow. If all available referents are strongly associated with a word in a single event, many erroneous linkages will be considered. If learning is too fast, these linkages could become solidified, permanent—and wrong. This fits with the fact that the average child hears 17,000 words a day (Hart & Risley, 1995), and even at their peak rate of learning, children may acquire only a handful of words in that same period (Sénéchal & Cornell, 1993). Thus, children must learn words slowly.

This framing yields enormous flexibility, as situation-time processes can be optimized for the demands of speaking, comprehending, and inferring, whereas developmental-time processes can be optimized to the demands of learning. Indeed, by moving much of the sophisticated inference of novel word meanings (classically described as constraints) to situation time, it may allow simpler mechanisms of learning to have complex effects, as learning is not entirely independent of such situation-time processes. For example, if the system attends to a referent longer in some circumstances than in others, more learning may result; or if competition between referents resolves faster (in situation time), the system may be able to acquire more unambiguous associations. Similarly, changes in situation-time familiar word recognition could be produced by simply improving the strength of the links between referents and their auditory word forms or by eliminating unnecessary connections.

This approach also addresses the relationship between familiar and novel words. To the extent that any partially formed mappings are available for a novel word, the system may use those partial mappings to increase consideration to the correct object and decrease consideration to erroneous objects—a form of mutual exclusivity. Thus, novel word recognition may take advantage of familiar word processes. Similarly, familiar word recognition may be enhanced by the fact that well-learned words would have stronger associations than newly learned ones—a type of frequency effect. Indeed, the process of tuning these connections to augment familiar words can derive from the same learning mechanisms that establish them. If both novel and familiar words harness the same processes, there is no reason to ignore external support (e.g., pragmatics, semantics) in even familiar word recognition—something of recent interest in adult sentence comprehension (e.g., Chambers, Tanenhaus, & Magnuson, 2004; Hanna & Tanenhaus, 2004; Tanenhaus & Brown-Schmidt, 2008).

Bringing together novel and familiar word learning can unify the literature. It also reframes how we think about classic findings by treating processes that have been described as constraints as descriptors of referent selection, not learning. Finally, it eliminates the classic distinction of fast and slow mapping. Carey (1978) described fast mapping as a special stage of initial learning where links between words and meaning were first forged. In our view, it may reflect purely situation-time referent selection. Indeed, as words are acquired, there are complex changes as connections are built and pruned, but there is no reason to assume the underlying processes are any different the first time. Thus, as the term *fast mapping* conflates learning and situation-time processing, we avoid it. Instead, we use *referent selection* to refer to situation-time behavior in ambiguous naming situations and *learning* to refer to developmental-time changes.

### Mechanisms of Ambiguity Resolution and Learning

Our *dynamic associative* account makes two claims. First, novel word processing (referent selection) is not distinct from familiar word recognition. Second, and more importantly, word learning is the interaction of situation-time processes that give rise to referent selection and familiar word recognition (word use), and developmental-time processes that give rise to retention and the improvement of these abilities. Here we propose mechanisms. Situation-time processes arise out of dynamic competition between

referents and words, and developmental-time processes arise from associative learning harnessing cross-situational statistics.

### Dynamic Competition

In connectionist models, any time multiple items are considered in parallel, some form of competition is present. Indeed, several models of word learning hypothesize probabilistic representations with this property (MacWhinney, 1987; Merriman, 1999; Regier, 2003, 2005; Xu & Tenenbaum, 2007). These provide a good starting point; however, they do not describe how these probabilities unfold over situation time. Though this is important for modeling behavior such as reaction times or eye movements, one could in principle derive simple linking functions to map these probabilistic representations to reaction times. However, this may not be sufficient. Computational models have shown that when competition unfolds dynamically over time, unexpected effects can occur—gangs of weakly active representations can suddenly inhibit a more active one (Spivey, 2007). Since we want to explore situations such as mutual exclusivity characterized by such ambiguity, and investigate how these unintended consequences can shape learning, it is crucial to implement competition as a dynamic process that unfolds over time.

Dynamic competition has been postulated in a number of domains involving constraint satisfaction and ambiguity resolution including music (Bharucha, 1987), syntactic parsing (MacDonald et al., 1994; McRae, Spivey-Knowlton, & Tanenhaus, 1998), comparison (Goldstone & Medin, 1994), visual scene organization (Vecera & O'Reilly, 1998), visual categorization (Spivey & Dale, 2004), visual search (Spivey, 2007), speech perception (McClelland & Elman, 1986), and language production (Dell, 1986). Many of these problems are clearly relevant for our domain. These computational approaches all incorporate simple neural-like units with graded activation. Activation feeds forward (from perceptual inputs to higher level representations) and backward (from higher level representations) while competing within a level over time. The resulting pattern of activation usually represents the best solution given the constraints imposed by the input, top-down expectations, and structure of the network. Typically, after this process, a single unit is active, and activation for competitors is suppressed, offering a close fit to our distillation of referent selection. There are a number of computational formalisms for this including interactive activation (Dell, 1986; McClelland & Elman, 1986), normalized recurrence (Spivey, 2007), and dynamic field theory (Schutte, Spencer, & Schöner, 2003; Thelen et al., 2001), which share these basic properties.

In our model, referent selection is modeled with dynamic competition through normalized recurrence (McRae et al., 1998; Spivey, 2007). Normalized recurrence has been applied to problems that are related to referent selection including speech perception (McMurray & Spivey, 2000), visual categorization (Spivey & Dale, 2004), and visual search (Spivey, 2007, chapter 8), and has been used to map lexical activation for known words to objects (Spivey-Knowlton & Allopenna, 1997; Spivey, 2007, pp. 187–200), embodying our distillation of word recognition for familiar words. Most importantly, versions of this architecture combine this competition with unsupervised learning (McMurray, Horst, et al., 2009; McMurray & Spivey, 2000).

To model referent selection, words and referents are modeled as localist units. On any trial, one word and multiple objects are active. Words and referents pass activation to a lexical layer (what Spivey, 2007, terms a “decision layer”), and recurrent competition among all three layers forces the network to suppress activation for the objects that do not map to the word. We discuss the motivation for localist representations in the General Discussion. However, at a purely practical level, virtually all of the above referenced competition architectures use localist units, and it is difficult to implement competition in a distributed representation.

### Associative Learning

A number of researchers argue that word learning cannot be associative: The fundamental mechanisms are social (Golinkoff & Hirsh-Pasek, 2006; Nazzi & Bertoncini, 2003), referential or conceptual (Waxman & Gelman, 2009), or constraint based (Woodward & Markman, 1998). Although such accounts describe important sources of information and/or important representational issues, it is not clear what these mean for learning because terms such as social, referential, and conceptual learning do not have clear definitions in learning theory.

Some of these nonassociative accounts still argue that early word learning may be associative (e.g., Golinkoff & Hirsh-Pasek, 2006; Nazzi & Bertoncini, 2003; see also Namy, 2012). This is almost a necessity—there is little lexical knowledge to facilitate mutual exclusivity, and social skills such as the use of eye gaze are still developing (e.g., Moore, 2008). These accounts typically argue that more complex mechanisms such as constraints or social pragmatics take over later. Thus, such accounts posit a discontinuity in the learning process, but even with this discontinuity, they do not offer an explanation for improvements in familiar word recognition.

Such accounts critique a straw man version of associationism in which raw perceptual inputs are directly associated without processing or intervening representations (cf. McMurray, Zhao, Kucker, & Samuelson, *in press*). Indeed, this critique seems to focus on the information that is associated, not in the mechanisms by which the linkages are made (cf. Smith, 2000). In contrast, modern learning theory admits internal representations as a basis of association and allows attention and other factors to shape the strength of these associations (Livesey & McLaren, 2011; Shanks, 2007). This is also central to connectionist learning. Our model does both: associating visual and auditory inputs to a lexical concept and allowing competition to shape their strength. If associative learning uses abstract representations and sophisticated situation-time processing, there is no reason to abandon it after the initial words. Indeed, as we described, by allowing social inference or constraints to shape in-the-moment processes, learning may still be associative at its base while leveraging these richer sources of information.

Under our view, learning is the same whether children are in the so-called association phase or the so-called constraint-based, referential, or social-pragmatic phase. The distinction highlighted by these theories is in terms of the information used during learning and novel word inference, not the learning mechanism. This makes a simple story. Basic learning mechanisms handle the retention of information. Initially co-occurrence may be the only source of information available to them, but later, as the child learns to use



other information in the environment and make more robust decisions about what the referents of words are, these form more precise activation patterns in situation time, which enables richer and faster associative learning—but the associative learning is the same. In this way, the timescale distinction allows us to distinguish and relate these processes. Principles and constraints identify the relevant information in situation time for the purpose of using words, whereas associative learning processes build the correct mappings over developmental time. By buttressing associative learning with dynamic competition to handle the situation-time ambiguity resolution, we may achieve a significantly more powerful word learner, as in any novel naming instance, the learner will have more activation for words and/or referents (enabling stronger associations to be built more quickly) and will have less activation for competing referents, preventing the formation of spurious associations. Crucially, this can be accomplished without having to posit qualitative distinctions between learning mechanisms at different ages.

This framework may also help connect vocabulary learning to classic findings in learning theory that are directly relevant to words. These include phenomena such as the power law of learning (Heathcote, Brown, & Mewhort, 2000; Logan, 1992; Newell & Rosenbloom, 1981; see also Section 2), the role of similarity (Palmeri, 1997; Storkel, Armbrüster, & Hogan, 2006; Swingley & Aslin, 2007; Wifall, McMurray, & Hazeltine, 2012), the role of statistics (Yu & Smith, 2007), and even old phenomena such cue neutralization (Apfelbaum & McMurray, 2011; Bourne & Restle, 1959; Bush & Mosteller, 1951; Rost & McMurray, 2010).

Thus, our goal was to investigate the consequences of associative learning when embedded in this richer framework of multiple timescales and internal representations. Our model uses perhaps the simplest form of associative learning, Hebbian learning. Inputs (words or objects) will be associated with an internal lexical unit if both are active; otherwise, the association decays. As in the case of competition, both implementing and understanding such learning make the most sense with localist units.

Even within this simple approach, there are layers of complexity. First, associations connect word forms and object categories to lexical concepts, not to each other. These lexical concepts function something like lemmas—abstract representations that connect other representations. Their presence means that learning requires at least two connections (word  $\rightarrow$  lexicon; lexicon  $\rightarrow$  object). Second, learning must not just build connections, but also avoid or eliminate unnecessary ones (cf. Regier, 1996). Consider the connections between visual and lexical units (see Figure 2). If the network heard *dog* in the presence of a dog and a bug, the most salient connection is the positive association between the object category *dog* and its lexical unit (thick line). However, this system should also learn that the object category *dog* is not associated with the word *tree* (which was not heard), a negative association (dashed line). It also needs to reduce the negative association between the object category, *tree*, which is not present, and the word, *dog* (dotted line). Thus, for successful learning to occur, it must increase one association and decrease two. In a larger lexicon, there will still be one positive connection, but there will now be hundreds of spurious connections to prune.<sup>1</sup>

What is the source of these spurious associations? As we described, these may exist from the earliest stages of learning. In connectionist models connection weights start from small random

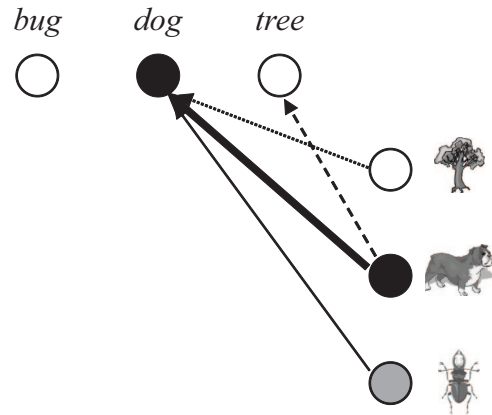


Figure 2. A schematic of the sorts of connections that would need to be acquired or pruned during learning.

values—a necessity in unsupervised learning (e.g., Rumelhart & Zipser, 1986). Some may also be formed during learning when incorrect objects appear with the referent (e.g., *bug*). Either way, these spurious associations will also need to be pruned. This pruning must occur in a way that preserves plasticity for new words. If connections to all words and referents were pruned equally on each naming event, potential positive connections for novel words and categories that have not been heard or seen would be pruned as well, making such words difficult to learn. Rather, we need a form of pruning that preserves potential connections for novel words and referents, but still allows the system to refine its connectivity for familiar words and referents. Thus, this pruning process, which is not often what comes to mind when we think of learning a word, could be an important determinant of development, if only for its massive scale.

## Overview of Architecture

At the broadest level, our dynamic associative model captures short and long timescale dynamics by embedding a model of learning in a model of use. It ignores the complexity of auditory and visual processing, to link word forms to object categories. However, it does not link them directly. Rather, word forms and categories each link to intermediate lexical representations. We implemented this in a hybrid connectionist–dynamic systems approach. Word forms and categories are represented by localist units, which compete in situation time using normalized recurrence to push the network toward a single interpretation. During competition, connections between inputs and the lexical layer are tuned via associative (Hebbian) learning. As we discussed, localist representations are the most transparent way to implement these mechanisms. However, they also offer a theoretical advantage: By stripping out the emergent power of distributed representations, they allow us to isolate these mechanisms and observe their consequence in a more or less pure state. Our goal was to investigate the power of this unique combination of basic mechanisms.

<sup>1</sup> We use the term *prune* here as a vivid metaphor. However we intend a more graded sense in which spurious connections are slowly and gradually reduced, not irrevocably eliminated in one step.

## Relationship to Standard Connectionist Approaches

By situating learning in connection weights and situation-time processing in changing activation, our approach is broadly consistent with classic connectionist thinking on development (Elman, 1990; Harm & Seidenberg, 1999; Munakata, 1998; Munakata & McClelland, 2003; Munakata et al., 1997). In models such as these, such a description can divorce emerging knowledge from the ability to use it in a real task. For example, in Munakata's models, the model may have some latent ability to represent an object under occlusion, but this can be overwhelmed in the moment.

Our model shares these broad properties, though it differs in theoretically important ways. First, as in many models, situation-time processing occurs over recurrent connections between units. However, we argue that competition is the essential element of this situation-time processing (in word learning), something other models (Elman, 1990; Harm & Seidenberg, 1999; Munakata, 1998) have not explored. Second, many models of this sort make the simplifying assumption that each time step corresponds to the presentation of a single input, what Schlesinger and McMurray (in press) term "trial time"<sup>2</sup> (but see Harm & Seidenberg, 1999). In contrast, in this domain it is important to model ongoing processing after the presentation of an input, both to model things such as reaction time and because settling can cause the network to change its interpretation of the input, not just strengthen whatever is already the most active (cf. McClelland & Elman, 1986; Spivey & Dale, 2004). Third, all the prior networks have used more powerful supervised learning, though we argue that unsupervised learning may be fundamental to word learning. Thus, our framework is built on classical connectionist thinking about learning and processing, but makes specific and theoretically motivated decisions about how to use these constructs that have not appeared in prior models of cognitive development.

## Why Another Computational Model of Word Learning?

There are a many computational models of components of word learning, examining topics ranging from the sequencing of phonological material (Gupta & Tisdale, 2009; Sibley, Kello, David, & Elman, 2008) to generalization of category membership (Colunga & Smith, 2005; Li, Farkas, & MacWhinney, 2004; Mayor & Plunkett, 2010; Samuelson, 2002), to embodiment (Roy & Pentland, 2002; Yu, Ballard, & Aslin, 2005). There are several models of the word-referent mapping problem (Frank et al., 2009; Li et al., 2004; MacWhinney, 1987; Mayor & Plunkett, 2010; Merriman, 1999; Regier, 2005; Siskind, 1996; Xu & Tenenbaum, 2007). These models make valuable contributions, highlighting the information that is relevant to the problem (Siskind, 1996), the power of associative mechanisms (Colunga & Smith, 2005; Li et al., 2004; Mayor & Plunkett, 2010; Merriman, 1999; Regier, 2005; Samuelson, 2002), and how constraint-like behavior emerges from simpler systems (Merriman, 1999; Regier, 2005).

A complete analysis of these models is outside the scope of this article (for useful reviews, see Frank et al., 2009; Regier, 2003, 2005), and we are not proposing a competitor to them. Our goal is not to develop a complete model of word learning, but rather to use computational tools to investigate the emergent consequences of theoretical ideas (cf. Schlesinger & McMurray, in press). In that

sense, it is important to address what questions these models have been used to answer, and any limitations that may inhibit their ability to address our questions.

First, by not explicitly capturing both timescales, many models do not succeed in modeling certain phenomena, or are forced to treat problems such as referent selection as developmental-time phenomena. For example, MacWhinney (1987), Merriman (1999), and Regier (2005) incorporated something resembling competition in their probabilistic representations. But they did not incorporate situation-time dynamics, nor distinguish between aspects of the problem that must occur in situation time (e.g., referent selection) and those that occur over developmental time. Consequently, they discussed things such as mutual exclusivity as a limit on learning (e.g., it is difficult to learn a second label for an object), rather than a principle of referent selection. Moreover, without implementing dynamic competition, these models cannot investigate the emergent interactions between competition and learning.

Bayesian models (Xu & Tenenbaum, 2007), in their focus on how interacting constraints lead to accurate inferences, off-load many constraints onto priors (and add new ones). This makes it difficult to understand how these behaviors develop or could arise from simpler processes. More recent Bayesian approaches (Frank et al., 2009) do address independent timescales, acknowledging the demands of both long-term learning and referent selection, and offer an example of how statistical learning can be buttressed by social cues. But lacking situation-time dynamics or a developmentally plausible learning model, such approaches only serve as a metalevel description of the information used for word learning (Jones & Love, 2011).

Regier's (2005) and Mayor and Plunkett's (2010) models are both associative and perhaps closest to ours. In some respects they go further, examining the contributions of auditory and visual similarity. In others they are more limited. At the level of associative learning, both models include rich representations of auditory word forms and basic categories, but do not include abstraction between them, so they may not be able to generalize word knowledge to other processes such as action systems, spatial processes, and orthography. Moreover, competition only arises from probabilistic representations—not true dynamic competition. As a result, they have not investigated the consequences of situation-time competition (though these approaches are likely not opposed to this). Further, these models emphasize issues in word learning that are now understood differently. Regier focused on accelerating vocabulary growth, when McMurray (2007) suggested this is a property of many parallel learning systems. Both models simulate fast mapping by mutual exclusivity as one-shot learning, but Horst and Samuelson (2008) suggested this may not be required—a model could perform well in the moment even if little learning occurred (Horst et al., 2006). Finally, neither deals with referential ambiguity, assuming a form of ostensive naming. This, however, is a limitation of implementation: Such models may be able to cope with many of the phenomena we examine here, given their theoretical overlap.

<sup>2</sup> Though this is a misnomer with respect to the Munakata models, as each trial consists of multiple input presentations (the sequence of visual inputs). The fact that there is not significant cycling of the network after each presentation of the input is the relevant factor here.

Finally, Samuelson (2002) and Colunga and Smith (2005) presented semisupervised associative learning models that do incorporate settling dynamics (though not competition). However, they did not attempt to model the actual learning problem. Rather than examine referential ambiguity, they focused on how associative learning can lead to generalization and to the identification of relevant dimensions for object categorization. Moreover, they did not use the settling dynamics to explain situation-time behaviors such as referent selection.

Thus, our model is consistent with features of many models: It incorporates associative learning (MacWhinney, 1987; Regier, 2003, 2005; Samuelson, 2002) and examines the role of statistical structure (Frank et al., 2009; Siskind, 1996) using graded or probabilistic representations (MacWhinney, 1987; Regier, 2005; Xu & Tenenbaum, 2007). It examines the emergence of constraints (Colunga & Smith, 2005; MacWhinney, 1987; Regier, 2005; Samuelson, 2002), and its architecture combines situation-time and learning processes (Colunga & Smith, 2005; Samuelson, 2002). However, it builds on these notions to examine a central new issue: the emergent power of how mechanisms interact across two timescales. The power of this combination is illustrated by distilling the problem to the point where such mechanisms operate over a minimally informative set of representations (e.g., word forms and concepts stripped of their phonetic, visual, and conceptual processes), and by embedding them within a cross-situational learning framework in which words can only be learned via co-occurrence statistics.

Our goal is not to present a complete model of word learning, but to use our dynamic associative model to learn how these processes interact. This can test the sufficiency of these processes to account for a range of phenomena. More importantly, we can use the model to clarify how these processes are related and develop a theoretical framework on which to base empirical investigations. This allows us to ask whether associative learning can cope with referential ambiguity, and whether children must solve this problem to learn words; how online processes and learning interact; whether processes that underlie familiar word recognition give rise to mutual exclusivity; and whether constraints such as mutual exclusivity emerge without being built in. Answering these will help develop a theoretical approach bigger than any one model.

### Specific Architecture

Our dynamic associative model (see Figure 3) has two layers of localist inputs, for auditory word forms and visual objects. Each auditory unit corresponds to a single word, and each visual unit corresponds to one category of possible referent. During processing, the auditory and visual layers are normalized such that the sum across each layer is 1.0. Thus, if a single node were fully active, its activation would be 1.0; if two were active, each would be .5; and when all nodes are inactive (the resting level), they are set to  $1/N$  where  $N$  is the number of nodes. After normalization, the vector of activations across a layer can be read as the distribution of likelihoods that the auditory (or visual) hypothesis represented by that node is present.

There are no direct connections between auditory and visual units. Interactions occur because both connect to a hidden layer of lexical units. These weights are initially random such that each

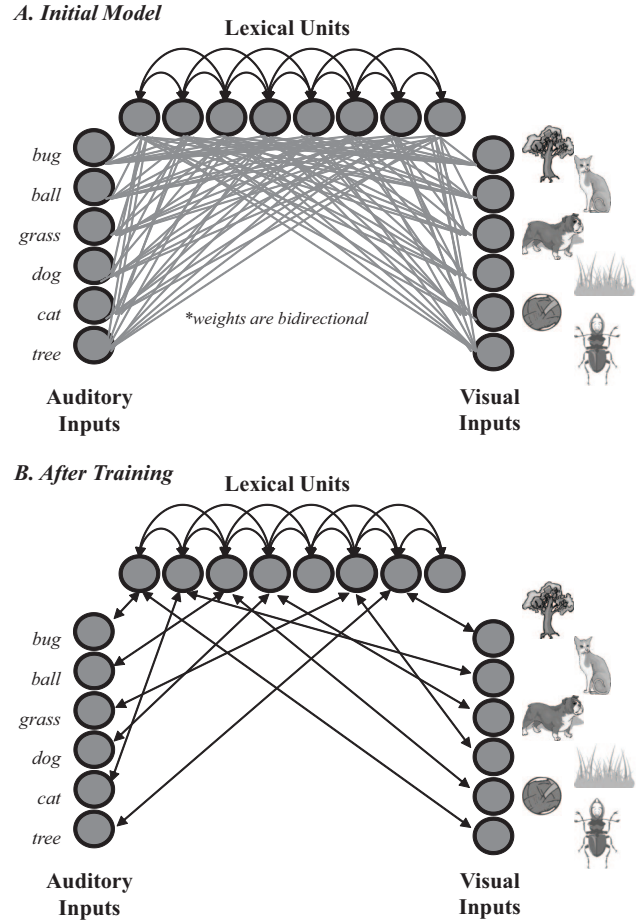


Figure 3. Architecture of the model, both before (A) and after (B) learning.

auditory and visual unit is connected to each lexical unit (with differing strengths; see Figure 1A). However, after learning, these weights generally reflect one-to-one mappings between word forms and lexical units, and between referents and lexical units (see Figure 1B).

The model contains more lexical units than needed to represent all the words. This leads to better learning. Generally, Hebbian learning reinforces existing associations and prunes unnecessary ones. Thus, the first few connections formed during learning are crucial. If the network initially randomly associated two inputs with the same lexical node, this mismatching could be reinforced over subsequent inputs and prevent the network from ever discriminating them. Greater numbers of lexical units make this much less likely (McMurray & Spivey, 2000).

On each trial, a single auditory unit (a name) is activated (set to 1.0). An array of visual units (the objects in the scene) is also activated. Input layers are normalized, and activation from each layer is sent concurrently to the layer of lexical units (Equation 1).

$$\Delta I_x = \left( \sum_{z \in A} w_{xz} a_z + \sum_{z \in V} u_{xz} v_z \right) \quad (1)$$

The change in activation of a lexical unit is based on the net input to that node: the sum of all the auditory units ( $a_z$ , for the  $z$ th



auditory unit) multiplied by their connection weight ( $w_{xz}$ , which connects auditory node  $x$  to lexical node  $z$ ), and the corresponding weighted ( $u_{xz}$ ) activation of the visual input units ( $v_z$ ). This is multiplied by the feed-forward temperature ( $\tau_{ff}$ ), which controls the weighting of the prior activation in setting the new activation (we use the superscripted  $l^{(t)}$  to denote  $l$  at time  $t$ ).

$$l_x^{(t+1)} = l_x^{(t)} + \tau_{ff} \cdot \Delta l_x^{(t)} \quad (2)$$

After updating the activations, lexical units inhibit each other using squared normalization.

$$l_x^{(t+1)} = \frac{(l_x^{(t)})^2}{\sum_{z \in L} (l_z^{(t)})^2} \quad (3)$$

Since activation in the lexicon is always between 0 and 1, the result of this is that highly active units become more active and dominate less active ones.<sup>3</sup> When run repeatedly over several cycles, this approximates winner-take-all competition.

The resulting lexical activation then feeds back to the input layers (Equations 4 and 5). The change in activation of an auditory (or visual) unit,  $y$ , is proportional to the product of its previous activation and the weighted activation ( $w_{yz}$ ) of the lexical units ( $l_z$ ).

$$\Delta a_y = a_y \cdot \sum_{z \in L} l_z w_{yz} \quad (4)$$

$$a_y^{(t+1)} = a_y^{(t)} + \tau_{fb} \cdot \Delta a_y^{(t)} \quad (5)$$

The fact that feedback is multiplied by the current activation means that feedback from the lexical layer only affects *active* input units. This prevents the network from activating perceptual inputs on the basis of top-down evidence alone, and it also introduces nonlinearity into the system. The temperature parameter,  $\tau_{fb}$ , is not required to be the same as for feed-forward activation and was not for the simulations reported here.

After updating, the activations of the auditory and visual layers undergo a small amount of inhibition (Equation 6) and are normalized.

$$a_y^{(t+1)} = \frac{(a_y^{(t)})^\iota}{\sum_{x \in A} (a_x^{(t)})^\iota} \quad (6)$$

Here  $\iota$  represents the degree of inhibition. At  $\iota = 2$ , this would yield strong inhibition as in the lexical layer; at  $2 > \iota > 1$  there is less inhibition; at  $\iota = 1$ , no inhibition; and at  $\iota < 1$  the activation collapses back to the resting state. After this step, activation feeds forward to the lexical layer and the cycle continues. Activation cycles in this way until the lexical layer settles (i.e., the change in lexical activation from time step to time step is close to zero).

Typically, on any trial the network is presented with multiple visual units to simulate a cluttered scene. Throughout cycling, the network partially considers each visual competitor simultaneously, but after many cycles, the competition and feedback generally result in a single object having more activation than the others. This active visual unit is the network's response—the network is allocating more attention to this referent. Thus, recurrent cycling causes the network to settle into an activation pattern across all three layers that reflects the present constraints (the auditory and visual inputs) and partial knowledge (in the weights).

Connections are modified at each cycle with Hebbian learning (Equations 7 and 8). The network increases the strength of the connection between simultaneously active input and lexical units, and decreases the connection in other cases.

$$\begin{aligned} \Delta w_{xy} = & a_x l_y (1 - w_{xy}) \\ & - .5 \cdot (1 - a_x) \cdot l_y w_{xy} \\ & - .5 \cdot a_x (1 - l_y) \cdot w_{xy} \end{aligned} \quad (7)$$

In Equation 7, the first line represents the positive term. If  $a_x$  (auditory unit  $x$ ) and  $l_y$  (activation for lexical unit  $y$ ) are active, the weight is increased proportional to its distance from 1 (its maximum value). The second and third terms are decay terms. A given weight decreases if (a) the input unit is active and the lexical unit is not or (b) the lexical unit is active and the input is not. If neither is active, there is no decay. By restricting weight decay to only connections between units that are actually used at that point in time, the model maintains plasticity in weights connecting input and lexical nodes that are not used. This is crucial for learning new words in the future (for a similar learning rule, see Grossberg, 1976). Weights are updated with Equation 8.

$$w_{xy}^{(t+1)} = w_{xy}^{(t)} + \eta \cdot \Delta w_{xy}^{(t)} \quad (8)$$

Here  $\eta$  is the learning rate and is typically very small—on the order of .0005. This is because learning occurs on each cycle of competition (and with many cycles or input, this will add up). By learning continuously, rather than at the end of processing, we need no homunculus controlling when learning can occur, much as children may not differentiate between training and test trials in the laboratory from other learning opportunities.

An important question is what regularities in the input drive functional learning. We examine one possibility here: co-occurrence between words and referents, or cross-situational statistics (Yu & Smith, 2007). We implemented this style of learning by ensuring that among the set of visual competitors active on any given trial, one was consistently paired with the auditory target while the others were randomly selected.

## General Methods

For most models, a 35-word lexicon was used. Though the network can learn larger lexica (see supplemental materials, Simulation S2), 35 was sufficient to be interesting while allowing the network to run reasonably quickly.<sup>4</sup> Thus, models were initialized with 35 input units and 500 lexical units. Weights were initialized to random values, generally between 0 and .5 (the *wtsize* parameter).

## Training

Models were generally trained for 200,000 epochs, where an epoch is one presentation of a word (though this entails many

<sup>3</sup> This particular form of inhibition instantiates a form of lateral inhibition in which the ability of each unit to inhibit the other units is a function of its proportion of the total activation.

<sup>4</sup> A typical model completed training in 30 min to 2 hr, but for each simulation we typically ran many repetitions of each model in several conditions, requiring several days.



cycles of competition as the word is processed). Although the model can perform quite well after only a handful of exposures to a word, it may take several thousand for the weights to settle on a single strong association. Thus, a long period of training was important for understanding both early and late stages of learning.

On training trials, a single auditory unit was activated, accompanied by the corresponding visual unit and a variable number of competitors. The average number of visual units active across trials is the degree of referential ambiguity present in naming situations for that model. Thus, if on average 14 of 35 units were active, the model faced 40.3% referential ambiguity. On any given trial, the active visual units were determined by first choosing a vector of 35 independent random values between 0 and 1. From this, any unit whose random value was less than the level of referential ambiguity would be active (e.g., for a referential ambiguity of 20%, if the random value for a unit was .2 or lower, this unit would be activated). Consequently, the number of competitors was not constant over trials, though the mean level of ambiguity was.

## Testing

Although many tests of the model were particular to a simulation, most models were tested in the following ways: (a) a simulation of an  $N$ -alternative forced-choice referent selection task, (b) a production task, and (c) an analysis of the weight matrix.

The  $N$ -alternative forced-choice (NAFC) task is similar to what is used with children (e.g., Golinkoff, Hirsh-Pasek, Bailey, & Wenger, 1992; Mervis & Bertrand, 1994). In this procedure, a single auditory unit is activated with  $N$  visual units (typically three). One of these is the target, and the others are randomly selected. Activation settles through the model until the lexical layer stabilizes. The activation in the visual units then determines the model's response. To correctly identify the word, the model must pass activation from the auditory layer to the visual layer, strengthening activation for the correct visual unit and suppressing the competitors. This can be interpreted as the model attending to one object (a look or a reach in a child). This is repeated for each word to estimate the number of words known. This test emphasizes the model's observable behavior in a simulated experiment, rather than analyzes its unobservable intermediate states (the lexical units).

We also created a production task, to simulate laboratory naming tasks, and to simulate the "child says" version of the MacArthur-Bates Communicative Development Inventory (MCDI), which is commonly used to assess which words a child knows. Here we run the model in reverse. This time a single referent is active and the system must select from all the auditory word forms—all the auditory units were partially active ( $1/N$ ). In this task, there are no constraints on which word forms are considered, and there will be typically many more auditory competitors than visual competitors in a comprehension task. Activation settles through the model until lexical layer stabilizes. Here the most active auditory unit serves as the response. This is repeated for each word to estimate the model's productive vocabulary.

In addition to these tasks, we analyzed the weight matrix. This is not comparable to anything that can be done with children, but it assesses learning more directly. For each word, we first examined the auditory  $\rightarrow$  lexical weights to determine which lexical

unit was most strongly connected to that auditory unit. This was repeated for the visual inputs. If this was the same unit, it would imply that the model correctly formed the association between auditory and visual representations for this word, and this word was scored as correct. Occasionally some instantiations of the model selected the same lexical unit for two word-object pairs, implying that both words meant the same thing. Because this was incorrect (in the lexicons used in the first three sections), only one of the two words would be scored as correct. Thus, this analysis asks whether the model has achieved a representation of the word that is close to an idealized state in which each word corresponds to one object (and vice versa).

Each of these tests was often repeated over training, raising the possibility that the model could learn during test. Thus, unless we have noted otherwise, any weight changes that accumulated during testing were discarded before continuing with training.

## Overview of Simulations

Our first goal is to show that this dynamic associative model can account for a variety of empirical results. We do not do this to emphasize the fit of this specific model. Rather, our model is representative of a broad class of interactive and associative mechanisms. Thus, its ability to fit the data really emphasizes the strength of these kinds of approaches, and of approaches operating at two timescales more generally. A second, more important goal is to use the model to understand the processing mechanisms that underlie word learning phenomena such as mutual exclusivity, familiar word recognition, and statistical learning. Thus, our simulations alternate between simulating empirical phenomena and unpacking the models' information processing.

Table 1 shows a summary of all the simulations presented here, with citations to relevant empirical studies. Section 1 examines learning. We ask whether the model can learn under referential ambiguity (Simulation 1.1), and whether it shows differences between comprehension and production (Simulation 1.2). We next ask how much referential ambiguity the model tolerates (Simulation 1.3) and about the shape of vocabulary growth (Simulation 1.4). These demonstrate the basic phenomena of word learning and the role of word use in explaining them.

We next examine situation-time phenomena. Simulation 2.1 examines the time course of familiar word recognition and Simulations 2.2 and 2.3 examine the use of mutual exclusivity for referent selection. Both are assessed over development and our analysis suggests that even these situation-time phenomena are fundamentally shaped by developmental forces.

Section 3 examines the interaction of situation- and developmental-time processes. Simulation 3.1 examines the apparent independence of these timescales, focusing on referent selection and retention (Horst & Samuelson, 2008). Next, we examine phenomena arguing for more dependent timescales: the interaction of task and mutual exclusivity (Simulation 3.2), and individual differences in familiar word recognition (Simulation 3.3). Finally, Simulation 3.4 manipulates situation-time processes in the model to show that they are necessary for learning.

Section 4 scales the model up in two important ways. First, Simulation 4.1 trains the model on both basic-level and superordinate labels to show that despite the model's use of mutual exclusivity, it can learn multiple labels for a referent. This also

Table 1

*Summary of Findings From the Simulations (With Reference to Relevant Empirical Studies)*

Simulation	Findings
Learning	
1.1: Learning	Model can learn under referential uncertainty. Performance on comprehension exceeds competence due to task constraints. Slow mapping or elaboration effects without semantics. Relevant studies: Capone & McGregor (2005); Smith & Yu (2008); Yu & Smith (2007).
1.2: Comprehension and production	Words can be comprehended earlier than produced. Largely effect of task-competition environment. Relevant studies: Huttenlocher (1974); Reznick & Goldfield (1992).
1.3: Referential ambiguity	Model can learn complete lexicon under high referential ambiguity. Relevant studies: Smith & Yu (2008); Yu & Smith (2007).
1.4: Accelerating learning	Model shows acceleration, as long as learning task is difficult enough and learning is sampled frequently. Relevant studies: Ganger & Brent (2004); McMurray (2007).
Situation-time processes	
2.1: Familiar word recognition	Settling time (reaction time) decreases over development. Effect arises pruning connections between words and incorrect referents. Similar pattern to power law of learning. Acceleration in number of words unrelated to deceleration in reaction time. Relevant studies: Fernald et al. (2006, 1998).
2.2: Fast mapping	Model can fast-map by mutual exclusivity. Function of both online dynamics and the weights set up by the learning and weight decay rule. Relevant studies: Carey & Bartlett (1978); Horst & Samuelson (2008); Markman & Wachtel (1988); Mervis & Bertrand (1994).
2.3: Fast mapping and development	Fast mapping emerges out of changes in weight matrix. Familiarity with visual objects can speed development. Relevant study: Halberda (2003).
Learning-processing interactions	
3.1: Fast mapping and retention	Model fails to retain fast-mapped labels, unless visually familiar. Only a small amount of learning occurs on any fast-mapping event. Relevant studies: Horst & Samuelson (2008); Kucker & Samuelson (2012); Spiegel & Halberda (2011).
3.2: Fast mapping and task	Model succeeds at 3AFC fast-mapping task at earlier points than 5AFC. Both develop over time. Relevant studies: Markman et al. (2003); Mervis & Bertrand (1994).
3.3: Familiar words and individual differences	Model shows stability in reaction time, correlations between reaction time and knowledge. Reaction time at early points in development predicts acceleration in vocabulary growth. Speed of processing is not unitary—emerges out of interactions between processing parameters, learning parameters, and developmental history. Relevant studies: Fernald et al. (2006).
3.4: Processing and learning	Competition is required for cross-situational learning. Feedback slows learning but may be beneficial. Relevant studies: Smith & Yu (2008); Yu & Smith (2007).
Scaling up	
4.1: Superordinate categories	Model can learn superordinate terms for objects in addition to basic level. Basic-level advantage derives from frequency, spreading of associations. Mutual exclusivity does not block learning of second names because it is an online process, not a constraint on learning.
4.2: One-to-one word object mappings	Model can learn when all words have multiple meanings (e.g., polysemy). Model can learn when all objects have multiple labels (e.g., bilingualism). Fast-mapping performance is slightly reduced by polysemy. Fast mapping is significantly degraded when all objects have multiple labels. Relevant studies: Byers-Heinlein & Werker (2009).
Supplement	
S1: Acceleration and word difficulty	Acceleration observed whenever the overall difficulty of the words is high. Varying frequency results in longer period of apparently slow learning.
S2: Temperature and speed of processing	Higher temperatures appear to lead to slower processing for familiar words. Effect derives from learning—higher temperatures offer initially faster settling and, as a result, fewer weights are pruned. Higher temperature also slows learning. Initially faster settling causes system to commit to more erroneous interpretations.
S3: Larger lexica	Model can learn lexica of up to 150 words at high degrees of referential ambiguity (50%, or $M = 75$ competitors per trial).
S4: Slow learning	Manipulated learning rates to see effect on learning. At normal values for typical Hebbian learning ( $>.01$ ), the model fails to learn, but can learn at intermediate and low values. Slow learning prevents model from overcommitting to erroneous mappings.

Note. AFC = alternative forced-choice.

affords the opportunity to examine basic-level advantages in a system that is not hierarchical. Simulation 4.2 generalizes this finding, training the network on multiple referents for a word (e.g., homonyms or polysemy), or multiple labels for a given object (e.g., bilingualism, synonymy). As these can disrupt the one-to-one word–object mappings commonly thought to underlie referent selection by mutual exclusivity, we evaluate both learning and mutual exclusivity.

A number of additional simulations were also conducted and reported in the supplemental materials. These support the simulations presented here and do not introduce new theoretical points. Commented MATLAB code for all simulations is also available in the supplemental materials.

### Section 1: Developmental Time Processes

We first ask whether the model can learn words under referential ambiguity, and whether differences in real-time comprehension and production tasks account for the commonly observed delay in productive vocabulary growth relative to comprehension. Next we assess how much referential uncertainty can be tolerated, and the shape of learning.

#### Simulation 1.1: Learning, Measured by Comprehension

Our first simulation validated that the model can learn under referential uncertainty, using the simulated laboratory comprehension task. We initialized 10 models to learn 35 words under 50% ambiguity (on average 17 competitors were present on every trial—a substantial degree of ambiguity). Every 1000 epochs we tested the model in the 3AFC and 10AFC tasks, and analyzed the weights (Table 2 shows the parameters for all simulations in Section 1).

**Results.** Figure 4 shows the number of words learned as a function of time (training). The perceived rate of learning is largely a function of the task. At 50,000 epochs, the model appears to know most of its lexicon when tested in the 3AFC task, and can do fairly well in the 10AFC task. However, this model has not finished learning. The weight analysis shows that only 1.5 words are known—for most words, the corresponding auditory and visual units do not have strong associations to the same lexical unit. This

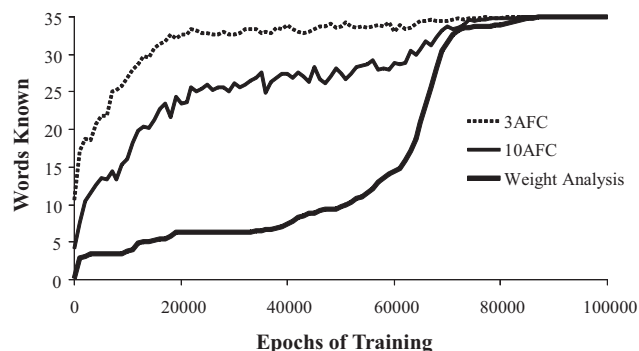


Figure 4. Number of words known over the course of training as measured by 3AFC and 10AFC tasks as well as an analysis of the weights. AFC = alternative forced-choice.

suggests that competition can make the model outperform its stored knowledge, and even after the model can identify the referent, there is still significant learning to do. Some have likened this “slow mapping” process to an elaboration of the meaning, syntax, and phonology (e.g., Capone & McGregor, 2005). This undoubtedly occurs, but this model has no semantics, syntax, or phonology. Thus, the raw word–object mappings may need to undergo a similar process in which competing weights are gradually eliminated, and the correct connections sharpened, even after the model demonstrates understanding of the word (Riches, Tomasello, & Conti-Ramsden, 2005). Crucially, however, we also observe a disconnect between the model’s knowledge (which is poor at early points in training) and its performance (which is simultaneously quite good). Here task constraints, such as the reduced number of competitors, actually allow the model to perform better than its knowledge, in an apparent reversal of the performance competence distinction.

#### Simulation 1.2: Production and Comprehension

Simulation 1.1 suggests a dissociation between what a model “knows” (the associative connections), and what it “does” (performance on comprehension tasks). This raises the possibility that vocabulary growth may appear to follow a different developmental

Table 2  
Parameters for Simulations in Section 1

Parameter	1.1: Comprehension	1.2: Comprehension and production	1.3: Ambiguity	1.4: Acceleration
Input units	35	35	35	35
Lexical units	500	500	500	500
Initial weight size	.5	.5	.5	.5
Learning rate	.0005	.0005	.0005	.0005
Referential ambiguity	.5	.5	.2 – .95	.75
Feed-forward temperature	.01	.01	.01	.01
Feedback temperature	2	2	2	2
Stability point	1e-12	1e-12	1e-12	1e-12
Input inhibition	1.05	1.05	1.05	1.05

*Note.* A number of parameters control the activation flow, rate of learning, and initial conditions of the model. Initial weight refers to the range of values that the connection weights were randomly initialized to (e.g., a random number between 0 and .5). Learning rate affects the amount of weight change for each learning instance.

time course when measured with production and comprehension tasks.

Such differences are observed empirically: Children typically comprehend more words than they can produce (Fenson et al., 1994; Reznick & Goldfield, 1992). Though this is often attributed to memory demands, difficulty planning articulation, or the earlier age at which speech perception develops, in the model, such differences cannot arise from these factors. However, as we have described them, comprehension and production are different tasks, as comprehension requires selecting from a small number of referents, whereas production requires selecting from a vast number of word forms. Thus, this simulation asks whether task differences alone account for part of the delay in production. This is not to say that such differences are artifacts of the tasks used to assess them—they are fundamental to comprehension and production more generally.

Twenty models were trained at 50% referential ambiguity (see Table 1). Every 200 epochs they were tested on a number of comprehension tasks: 3AFC, 5AFC, 10AFC, and 35AFC (as if the entire lexicon were visually present) and a 35AFC production task.

**Results.** Figure 5 shows the number of words learned over time for each task. As before, even when 50% of the lexicon was present on any trial (~17 competitors), the model acquired the full lexicon no matter how it was assessed. Moreover, the 3AFC, 5AFC, and 10AFC tasks suggested that like children, in comprehension tests, the network appears to know more words than in production tests. However, by the end of training the network performs equivalently for both tasks. Interestingly, the production measure matches the estimates based on the weight analysis. Thus, more difficult naming tasks may tap this level of competence.

The number of response alternatives plays a critical role in predicting performance and the apparent rate of learning. To the extent that the response set for comprehension must include fewer objects than the full lexicon, task differences can account for differences between comprehension and production. Similarly, Huttenlocher (1974) described this in information processing terms: Comprehension is a recognition task, while production is a

recall task. Our model instantiates both processes as variants in the same underlying competition dynamics.

This might suggest that the number of response alternatives alone dictates success. However, this is not entirely the case, as a comparison of the 35AFC comprehension and production tasks shows that now comprehension is delayed. This is due to an asymmetry during training. Auditory units are always presented singly, whereas visual units are not. As a result, the network has experience suppressing unnecessary visual units, but has never suppressed unnecessary auditory units.<sup>5</sup> Nonetheless, this small difference illustrates that task differences can arise from differences in both situation-time factors (the number of response alternatives) and developmental-time factors (the history of suppressing competitors).

Critically, however, this distilled account of production versus comprehension suggests that (a) the number of response alternatives during testing can dictate how many words a child appears to know, and (b) this somewhat obvious fact can give rise to differences in production and comprehension vocabulary.

### Simulation 1.3: Learning Under Referential Ambiguity

Thus far, we have held referential ambiguity at 50%. This is substantial, yet we found excellent learning. It is important to determine the robustness of this learning, particularly when buttressed by dynamic competition. Thus, we varied the degree of referential ambiguity from 20% to 95% and trained 10 models at each level. Note that 95% means that on average 33.25 objects were present with the referent in any naming situation, and on 16.5% of the naming instances all 35 words are active. Models received 200,000 training trials, and we assessed performance in the 3AFC and 10AFC tasks as well as a weight matrix analysis every 25,000 epochs.

**Results.** Figure 6A displays the number of words identified in the 3AFC task as a function of referential ambiguity. Chance is 33% (11.7 words). At low levels of noise, the model acquired most words within about 25,000 epochs and learned all of them by 100,000 epochs. At 100,000 epochs an effect of referential ambiguity is seen: The model's performance drops off as with more competitors (though not very far). However, this is overcome with additional training: At 200,000 trials, the model performed at 100% even at 95% ambiguity.

This success is emphasized by our more conservative analysis of the weights (see Figure 6B). Here chance is much lower: The probability of randomly mapping a single auditory and visual unit to the same lexical unit is  $1/500 \times 1/500 = 0.0004\%$ . Moreover, to pass this test, the model cannot rely on competition to arrive at the best guess if the weights are imperfect (as it can in the 3AFC task). By this criterion, it takes much longer to learn a word. At 50,000 epochs, the model's 3AFC performance is good, but its underlying competence (the weight matrix) is far from complete. For example, at 50% referential ambiguity the model has only learned 10 words by this point in time. However, with enough

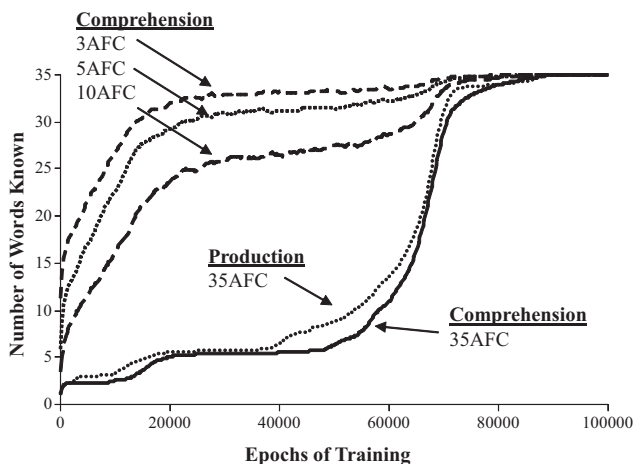


Figure 5. Number of words known by several comprehension and production measures over the course of training. AFC = alternative forced-choice.

<sup>5</sup> Of course, with a more realistic auditory representation, multiple auditory word forms will be active in parallel (e.g., similar sounding words; Allopenna, Magnuson, & Tanenhaus, 1998; Marslen-Wilson, 1987). This may minimize the differences seen when the number of alternatives is equated.



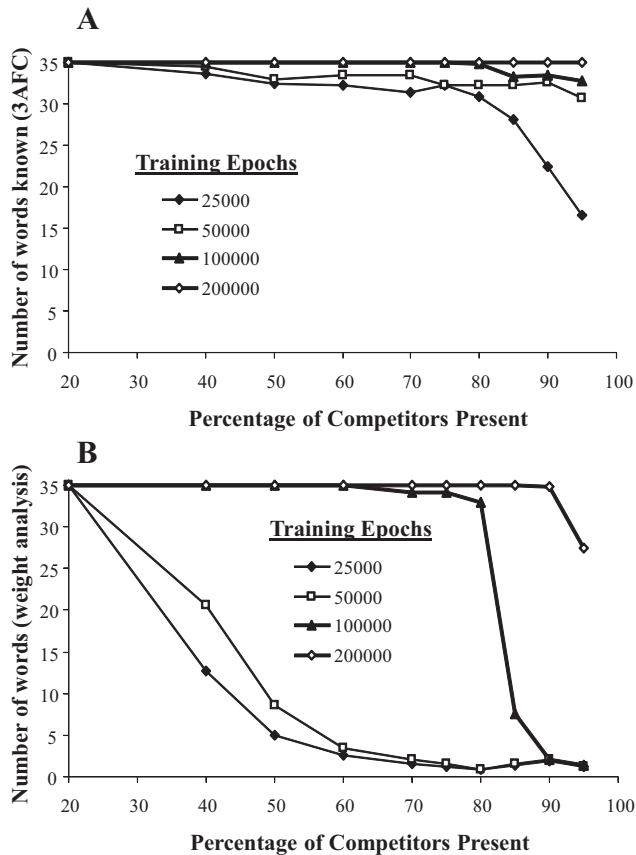


Figure 6. Number of words known as a function of the referential ambiguity and time, as measured in a 3AFC task (A) and via weight analysis (B). AFC = alternative forced-choice.

training, the model performs well on this conservative measure even at the largest degree of ambiguity. With 90% of available referents copresent during learning, all models learned all 35 words; at 95% ambiguity they averaged 27.4 and may not have stabilized yet. This suggests that the simple statistics of co-occurrence (between visual and auditory events) can be extremely powerful. Even when most of the lexicon is available on any given naming situation, the model eventually determines the correct mappings between words and referents. This undercuts claims that associative learning cannot cope with the high degree of referential ambiguity faced by real children. Given sufficient time, such mechanisms may be sufficient when embedded in situation-time competition framework.

#### Simulation 1.4: Acceleration

In the prior simulations, the model appears to start with high rate of acquisition and taper off (e.g., the 3AFC and 10AFC tasks; see Figures 4 and 5). This contradicts the consensus that word learning accelerates (P. Bloom, 2000). Though acceleration was observed in production (Figure 5) and in the weight analysis (Figures 4 and 5), in children comprehension also accelerates (Reznick & Goldfield, 1992). More importantly, the McMurray (2007; Mitchell & McMurray, 2009) analysis suggests that acceleration should be

observed in most parallel learning systems. So, though it is not distinctive of word learning, it was unexpected that it was not consistently observed.

There may be several reasons for this. First, perhaps we are not sampling at a high enough density. The vocabulary explosion typically appears during the 2nd year of life. Given that it takes the model about 100,000 epochs to learn its adult lexicon, the explosion would probably occur in the first 10,000 epochs, though the above simulations only sampled twice during this time window. Second, McMurray (2007) predicted acceleration only when words vary in difficulty such that there are fewer easy words than moderate or difficult words. As a simplifying factor in our simulations, all words were equally difficult—they were equally frequent, and the degree of ambiguity was the same across words.

Third, when measuring children's vocabulary size, we do not subject them to a 3AFC task for each word. Rather, we use a parent questionnaire such as the MCDI (Dale & Fenson, 1996; Fenson et al., 1994), which probably measures something like children's ability to use a word in a variety of contexts, that is, more of an average performance across time for that word. In fact, ongoing work (Mitchell & McMurray, 2008) using a stochastic version of the (McMurray, 2007) model has shown that acceleration is only observable when words require several exposures to learn. When words can be learned in only one exposure, deceleration is guaranteed. Given the heavily constrained 3AFC task, it may only take a small number of repetitions to learn a word by this criterion (particularly at low levels of ambiguity), and we may not see acceleration.

Thus, we ran 10 repetitions of the model with three changes. First, we sampled every 200 epochs. Second, for each word, in addition to the usual 3AFC and 10AFC tests, the model was tested five times and had to get the right answer on at least four (to simulate understanding a word in a variety of contexts). Third, we increased the referential ambiguity to 75% (so all words took longer to learn). We also explored the difficulty distribution by manipulating the frequency of the words such that there were few easy words and many harder ones. This version of the model also showed acceleration. It is discussed in the supplemental materials (Simulation S1).

**Results.** Figure 7 displays the results for the 3AFC and 10AFC tasks when all the words were of equal frequency. The thin lines show when the word was considered known if the model was correct in a single 3AFC or 10AFC task. The thick lines require the model to be correct on four of five trials of these tasks. Requiring

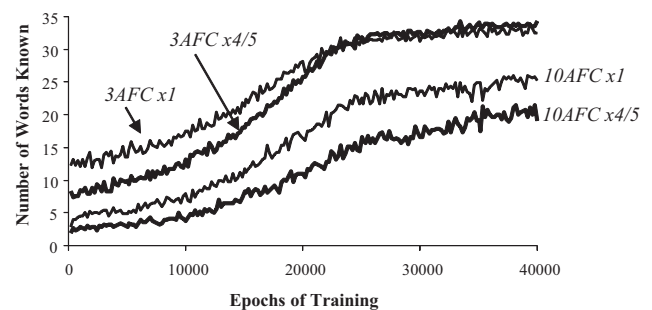


Figure 7. Number of words known as a function of training when all words were of equal frequency. AFC = alternative forced-choice.

a word to pass multiple tests clearly slowed initial performance and led to a steeper learning curve.

There is also more noise in these measurements (from epoch to epoch) than was seen earlier. At this stage in development, the model has not really learned the words, and its performance is affected by fluctuations in the weights of competitors (that are still settling down) and the choice of competitors in the comprehension task (see Adolph, Robinson, Young, & Gill-Alvarez, 2008, for a discussion of sampling issues in development). This suggests that even using standardized measures such as the MCDI the implicit task (the real-time behavior leading up to the measure) may shape the outcomes we measure (see also Sandhofer, Smith, & Luo, 2000).

However, no matter how we measure it, the model undergoes a period of slow learning followed by a period of fast growth (note that the apparent threshold in the 10AFC task is only a momentary plateau—by the end of training, all models learned all 35 words). Thus, the critical factor for acceleration is difficulty—as long as words are fairly difficult to learn, the model shows acceleration. Measurement clearly plays a role, because it acts as sort of a threshold: The number of repetitions required to pass a 3AFC task is less than the number required for a 10AFC task. This is in accord with the (Mitchell & McMurray, 2008) model of the vocabulary explosion.

## Section 1 Discussion

The first simulations suggest that this dynamic associative model can acquire a lexicon under substantial uncertainty. Even with 90% referential ambiguity, the model acquired all the words. However, the model's knowledge is a function of how it is tested. Comprehension tasks with small numbers of competitors show earlier learning than those with larger numbers or production tasks. As a result, we see differences between comprehension and production in a system that does not do either. Similarly, the growth curve is a function of testing: Some tests can show decelerating learning, whereas more realistic assessments show the predicted acceleration.

Virtually all the tasks suggest the model knows more words than are reflected in its connections (knowledge). This offers a paradigm for thinking about the integration of short and long timescale processes. Long-term knowledge (the weights) develops slowly and may be incomplete for substantial portions of development. However, fast competition processes allow the model (or child) to augment these weak representations in the moment and act more intelligently. This throws a novel spin on performance and competence. Typically, children are thought to have better underlying competence than their performance. Indeed, even in connectionist networks that do not make a strong distinction between performance and competence, we still see situations where the networks' apparent knowledge can be overshadowed by in-the-moment task demands (Munakata, 1998; Munakata et al., 1997). Yet here we see the opposite: Situation-time performance compensates for rather lousy competence. This is not a contrived product of the model. In most real situations the environment (including the caregiver) may provide information or support that allows the child to perform significantly above his or her level of competence in many situations (cf. McCabe & Peterson, 1991; Reese & Fivush, 1993). Thus, performance constraints, though usually thought of as impeding

our ability to see the true developmental level of the child, may actually augment it in some circumstances.

These simulations also suggest that referential ambiguity may not be as problematic as typically implied. Even in vastly ambiguous environments, associative learning can be successful over the long term. However, learning must be slow, because the relevant information can only be extracted across naming events—in fact, fast learning may cause the model to overcommit to an incorrect interpretation (as we demonstrate in this model in the supplemental materials, Simulation S4). Here again, however, the interaction of timescales helps: If fast, task-constrained processes buttress poor knowledge, the model can perform well despite imperfect knowledge in its weights.

## Section 2: Situation Time Processes

This section examines children's ability to use their lexica in situation time. The existing literature tends to distinguish familiar and novel word processes. Work on novel words focuses on accuracy: Can children infer, in the moment, the referent of a novel word? For familiar words, the emphasis is on the efficiency or speed of referent selection. Thus, Simulation 2.1 examines developmental changes in the efficiency of children's familiar word recognition, and Simulations 2.2 and 2.3 examine referent selection by mutual exclusivity, a form of referent selection for novel words. Both require children to use available information, in the moment, to identify the referent of a word. Beyond offering a model of these phenomena individually, these simulations also argue that mutual exclusivity emerges from the same processes as changes in familiar word recognition, and they help reveal fundamental properties of learning that make both possible.

### Simulation 2.1: The Development of Word Recognition and the Power Law of Learning

Fernald and colleagues (Fernald et al., 2006; Fernald et al., 1998; Swingley & Aslin, 2000) have examined the time course of children's mapping of familiar words to their referents using fixations. In this paradigm, children see pictures of two objects (e.g., a ball and a car) and are instructed to look at one. The speed at which they fixate the right object is taken as a measure of processing speed, and this tends to decrease over development (Fernald et al., 1998; Hurtado, Marchman, & Fernald, 2007). Fixation time is also correlated with lexicon size (Fernald et al., 2006; Zangl, Klarman, Thal, Fernald, & Bates, 2005) and is stable and predictive among children (Fernald et al., 2006; Marchman & Fernald, 2008). The goal of this simulation was to investigate this computationally, in order to identify the potential loci of these effects.

Ten models were run at three levels of referential ambiguity (25%, 50%, 75%; see Table 3 for parameters). Every 250 epochs, the model was tested on its entire lexicon in a 3AFC task to assess the number of words known and the time it took the model to settle on a referent for each of them. As with Fernald et al. (2006), reaction time (RT) was only saved for trials in which the model selected the correct referent.

**Results: Development of RT.** Figure 8 shows the settling time in cycles as a function of training for each of the three levels of ambiguity. There is a dramatic drop early in training, from 20 to

Table 3  
Parameters for Simulations in Section 2

Parameter	2.1: Word recognition	2.2: Mutual exclusivity	2.3: Development of mutual exclusivity
Input units	35	40	40
Familiar words	35	30	30
Lexical units	500	500	500
Initial weight size	.5	.25	.25
Learning rate	.0005	.0005	.0005
Referential ambiguity	.25, .5, .75	.5	.5
Feed-forward temperature	.01	.01	.01
Feedback temperature	2	2	2
Stability point	1e-12	1e-12	1e-12
Input inhibition	1.05	1.05	1.05

30 cycles at 250 epochs to five to six cycles by the end of training. There are small effects of referential ambiguity. These likely derive from the fact that different levels of referential ambiguity offer the model different amounts of exposure to visual competitors, which may alter how competing associations for a given word can be pruned (as we will discuss, this is a crucial determinant of RT). Optimal learning, then, may require a mix of low-competitor situations (e.g., ostensive naming) to establish the words and rule out strong competitors, and high-competitor situations to improve processing (cf. Horst, Scott, & Pollard, 2010; McMurray et al., in press).

**So what causes the decrease in settling time?** Changes in the dynamics of activation flow cannot account for the decrease in RT because the parameters that control it (temperature and the degree of lateral inhibition) did not change over learning. The only things that did change were the weights. Figure 9 shows a representation of the weights connecting the visual and lexical layers over development. Along the *x*-axis are the 35 visual units. Along the *y*-axis are 35 of the 500 lexical units. The strength of the connection between each unit is represented by the darkness of the patch at their intersection. At the beginning of training (see Figure 9A), these connections are random, and there is no clear structure, but

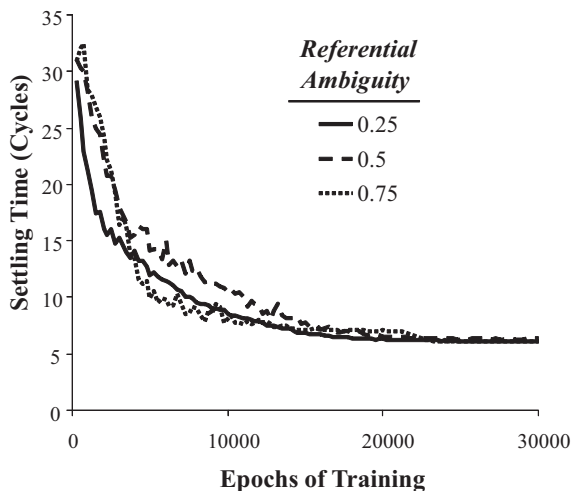


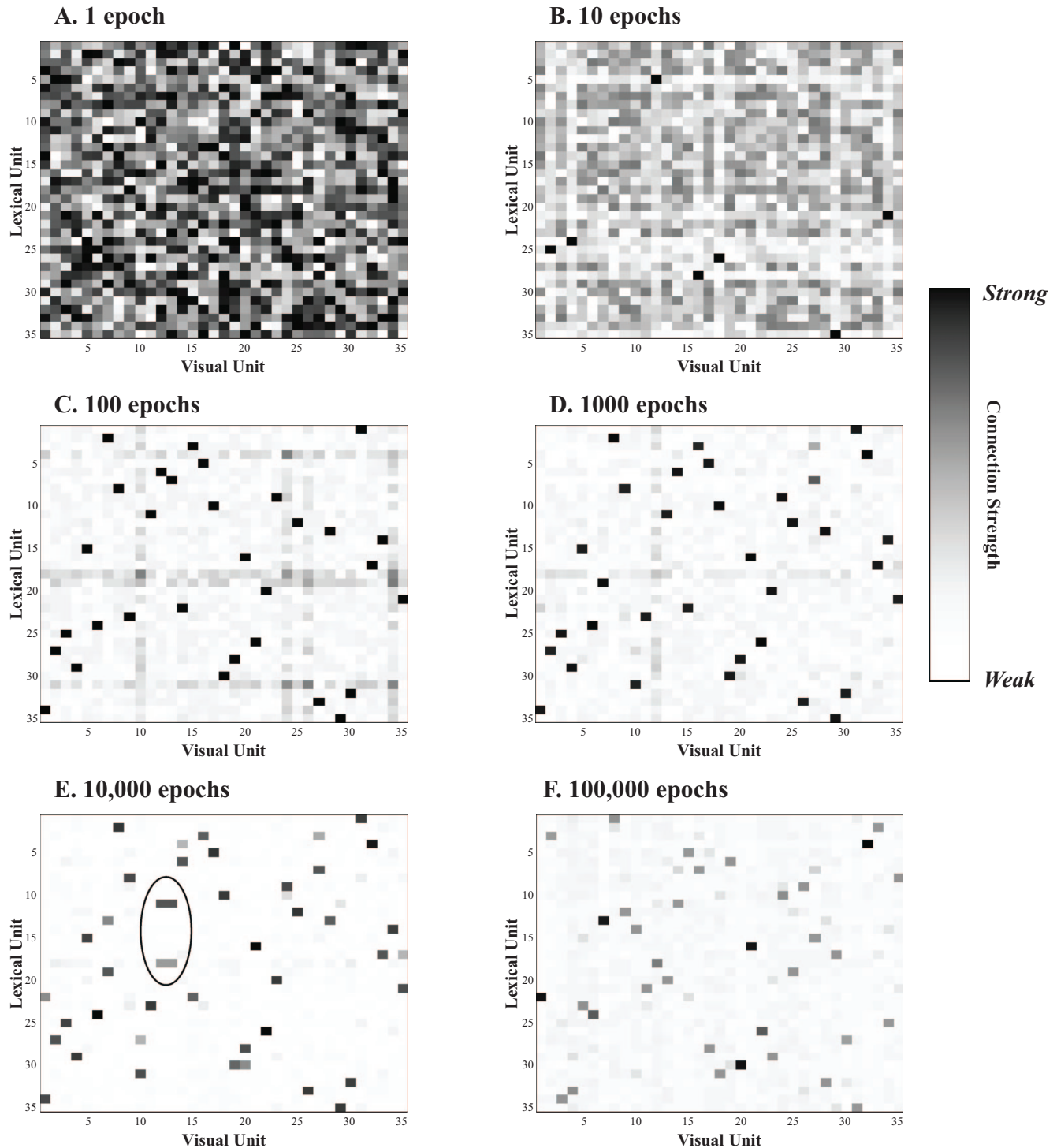
Figure 8. Settling time over the course of training, as a function of referential ambiguity.

even 10 epochs in we can see that the overall strength of the connections has been reduced (though there are a few strong connections for words that have been heard). Over the next 1,000 epochs, unnecessary connections are further pruned and positive ones enhanced. By around 10,000 epochs (see Figure 9D), the model only has a handful of strong connections. However, some of these are ambiguous. For example, Lexical Unit 11 (circled) has strong connections to at least three input units (and those units are also connected to Lexical Unit 17). At 100,000 epochs these competing representations have been eliminated—each visual unit connects to only one lexical unit, and the bulk of the connections are near zero.<sup>6</sup>

We can roughly characterize these changes in terms of the entropy of the weight matrix: A “peaky” distribution of weights characterized by a few strong connections and many weak ones (low entropy) describes a more developed model, whereas a flatter distribution (high entropy) characterizes a less developed model. We tested this by running an additional 10 models with a referential ambiguity level of 50%. Every 500 epochs, models were tested in a 3AFC and 10AFC task, and we evaluated both the auditory and visual weight matrices using three measures. First, we computed entropy, converting weight values to probabilities. By this measure, relatively random weight matrices will have large entropies, and peakier weights will have small entropies. Second, at each point in training, we found the single strongest connection linking each input unit to the lexical layer and recorded its strength as a measure of how strong the positive connections were. Third, we took the average of all the other connections as an indicator of how weak the irrelevant connections were.

Entropy showed strong relationships with both learning and RT. We found a negative correlation between entropy and time during learning ( $R = -.92$ ), and between entropy and the number of words known (10AFC task:  $R = -.77$ ; weight analysis:  $R = -.92$ ). Thus, entropy captures some aspect of overall learning or development. Figure 10A relates entropy to the log of the RT in the 3AFC task. Though the relationship is nonlinear, high entropies (relatively random weights), predict slower RTs (linear:  $R = .43$ ; hyperbolic:  $R = .92$ ).

<sup>6</sup> Note that over development, even the positively associated weights tend to drift downward under this learning rule (they do eventually stabilize; see Figure 10B). This accounts for the fact that the final weight matrix (see Figure 9F) shows lower weights than in some of the earlier panels.



*Figure 9.* A visual representation of the visual  $\rightarrow$  lexical weight matrix over the course of development. On the x-axis are each of the 35 objects; on the y-axis are subsets of the lexical units (including all the ones that are ultimately used). The connection between them is represented by the darkness of the patch.

The strength of the positive connections was not highly predictive of performance on either of our tasks (3AFC:  $R = -.27$ ; 10AFC:  $R = -.41$ ) or settling time ( $R = .31$ ), and analyses of the scatterplots suggested this was not due to a nonlinear relationship

(see Figure 10B). It was moderately related to number of words known ( $R = -.60$ ), but negatively. That is, lower connection strengths tended to indicate more words known. All these results derive from the fact that the positive connections fluctuated over



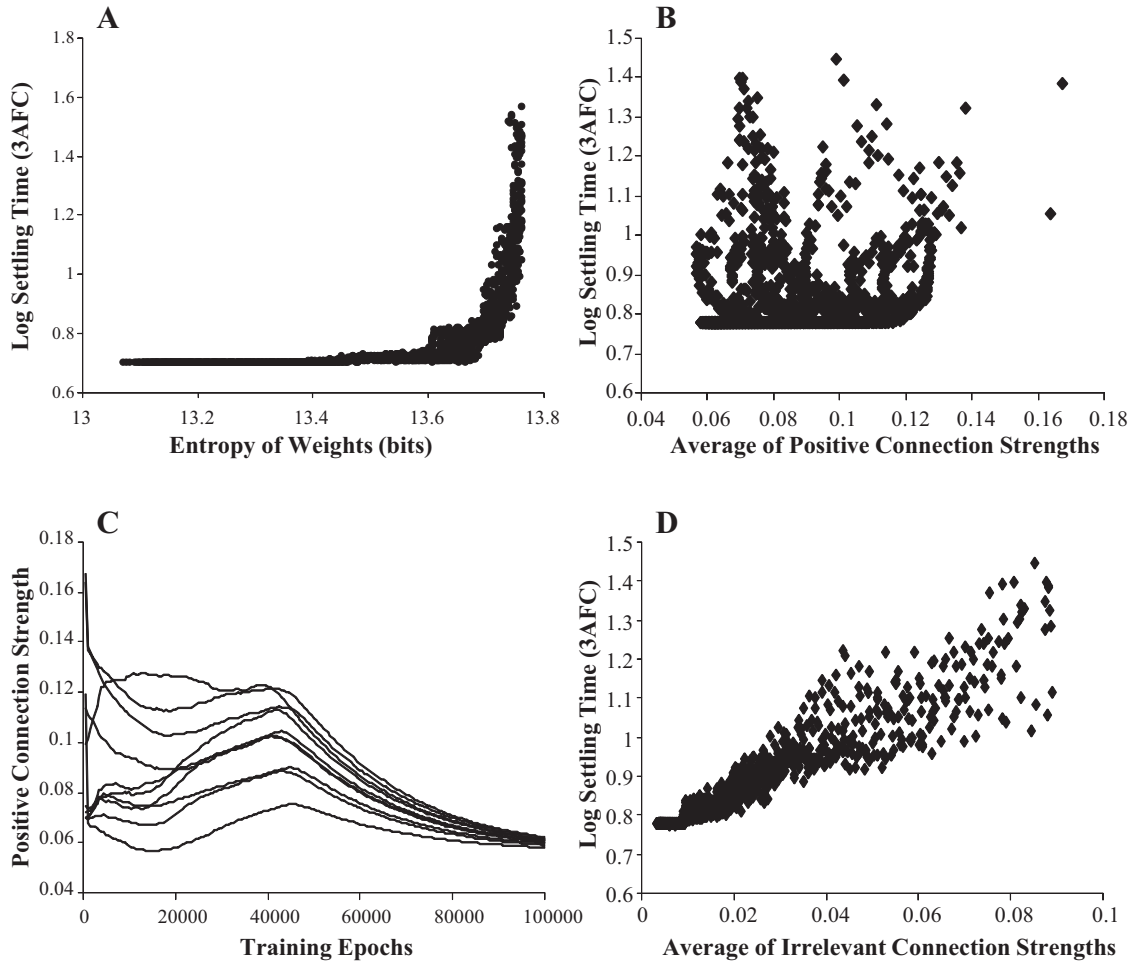


Figure 10. Four measures of the weight matrix related to settling time. (A) Reaction time (RT) as a function of Shannon's entropy. (B) RT as a function of the average of the positive (correct) connections linking each word/object to its correct lexical unit. (C) The strength of the largest positive connection over training. (D) RT as a function of the average of the unused connections. AFC = alternative forced-choice.

training but ultimately decreased slightly (the negative ones decreased a lot more; see Figure 10C). The average of the negative connections was much more predictive. It was highly negatively correlated with performance on the lexical identification tasks (3AFC:  $R = -.94$ ; 10AFC:  $R = -.89$ ) and with the log of settling time ( $R = -.92$ ; see Figure 10D). In both cases, performance increased with smaller irrelevant weights.

As a whole, then, the pruning of unnecessary connections is the driving force behind both acquiring new words and recognizing them faster. Unnecessary connections cause an auditory input to activate multiple lexical units. These lexical units compete, moving in the direction of winner-take-all. Since the network does not settle until this competition is resolved, it is the presence of these momentarily active competitors (driven by unnecessary connections in the weights) that ultimately leads to a longer settling time. This suggests that empirical correlations between RT and vocabulary size (Fernald et al., 2006; Zangl et al., 2005) may be driven by the fact that children who know more words may also have fewer spurious associations at that

point in development. More sophisticated eye movement paradigms may be able to test this by evaluating more precisely degree of competitor activation.

**Acceleration and deceleration.** Developmentally, these simulations suggest a steep decrease in settling times early, followed by a flattening. This pattern is commonly seen in power law or exponential decay function in the literature on general learning principles and has appeared in a variety of motor and cognitive learning tasks (Anderson, 1982; Heathcote et al., 2000; Logan, 1992; Newell & Rosenbloom, 1981; Wifall et al., 2012), suggesting that word learning may operate by similarly general principles. The power law has always been interpreted as demonstrating that learning slows throughout training, which would seem to violate the acceleration commonly observed in word learning. However, our model does both.

To understand this, we looked for a relationship between the greatest change in RT and the change in number of words known. Figure 11A shows the change in RT and words known for the models learning under 25% referential ambiguity. Here both RT

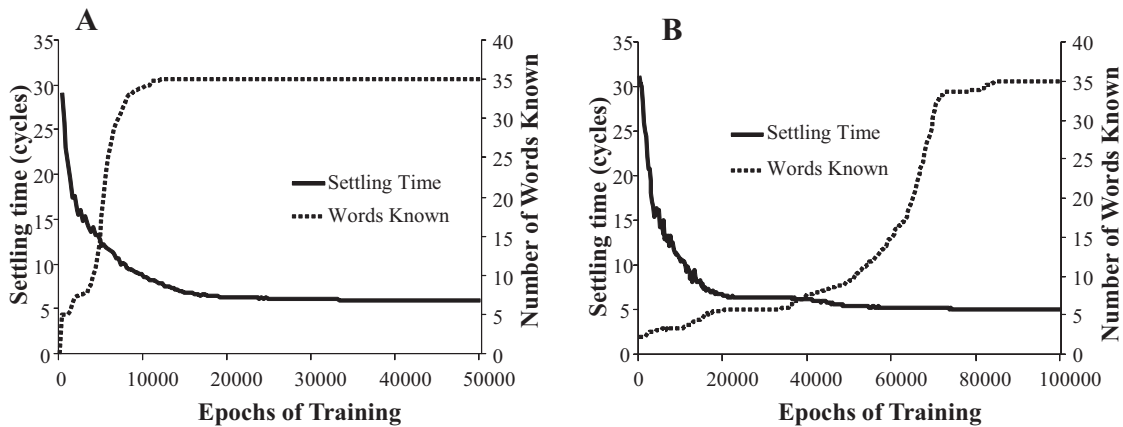


Figure 11. (A) Settling time as a function of training on a log-log scale. (B) Settling time (left axis) and words known (right axis) for model trained under 25% ambiguity.

and words known undergo their greatest change at around 5,000 epochs, implying some fundamental process that affects both at this time. However, Figure 11B shows the same data for the models learning under 50% referential ambiguity. Here the time course of lexical acquisition is pushed much later (peaking between 60,000 and 70,000), but the pattern of RT barely moves. Thus, the timing of changes in RT does not clearly map to changes in the number of words known.

What can explain apparent acceleration and deceleration in learning? As we discussed in the prior section, the best predictor of RT is the magnitude of the spurious connections between auditory and visual competitors and the current lexical unit. These constitute the bulk of the weight matrix. If there are  $35 \times 500 = 17,500$  weights connecting auditory and lexical nodes, only 35 of them are correct—17,465 must be suppressed over learning. However, weight decay is also widespread. Every time a word is heard, thousands of weights are not used and therefore suppressed. Thus, there are quick RT gains to be made for suppressing unnecessary competitor weights, and this can be done in virtually all circumstances. Any word, with any degree of referential ambiguity, will result in some weights being decayed.

The number of words known is more complex. It requires two positive weights (auditory  $\rightarrow$  lexical and visual  $\rightarrow$  lexical) and an absence of competing weights. Neither the irrelevant weights nor the positive weights are singly related to it. Even if the network had a very strong connection between an auditory word form and its lexical unit, if there were other strong competing connections, the word would still not be learned. Thus, in contrast to reducing RT, actually learning a word requires a confluence of events and is much more difficult. This is clear in the correlations observed in the previous section: Words known were best correlated with entropy, a global measure of the weights ( $R = -.91$ ), and less so with either spurious ( $R = -.56$ ) or relevant ( $R = -.60$ ) weights, whereas RT was best correlated with spurious weights ( $R = .91$ ) and less so with the other measures (relevant weights:  $R = .21$ ; entropy:  $R = .42$ ). Thus, despite the fact that the development of both RT and vocabulary size rely on fundamentally the same learning mechanism, the fact that each must be measured through a situation-time measure means that they may tap different aspects of learning, and hence show different learning curves.

**Discussion.** This simulation shows similar results to Fernald et al.'s (1998, 2006) studies: The dynamic associative model's familiar word recognition improves over time. This improvement, which appears as gains in the efficiency of a skill, actually derives from changes in knowledge (connection weights). Crucially, the suppression or pruning of the irrelevant weights is the dominant factor—the bulk of word learning may consist of learning which objects and words do not go together (cf. Regier, 1996). We also show that word learning does not have to differ from general approaches to learning that predict decelerating RTs, even as it shows acceleration in vocabulary size. This underscores the fact that word learning may operate by general learning principles and challenges the utility of drawing strong conclusions based on the shape of vocabulary growth (McMurray, 2007). Critically, each of these measures derives from different changes in the underlying associations, showing the multifaceted nature of association learning in this context. At a broader level, however, the acquisition of word-object linkages, a clearly developmental-time learning phenomenon, is directly implicated in the ongoing development of processing speed, a situation-time measure. In fact, changes in processing derive entirely from a nonobvious component of learning—suppression of irrelevant connections.

## Simulation 2.2: Referent Selection by Mutual Exclusivity

If the development of familiar word recognition derives from the same learning mechanisms as the acquisition of new words, this raises the question of whether the processes that infer the referents of novel words can also arise from these mechanisms. The plethora of proposed constraints and specialized inference processes imply that children deploy additional mechanisms in novel naming situations. Yet, how does the child first determine that the word is novel and then which constraints to apply? Although some have argued for a decision process of some kind (Merriman et al., 2008), an alternative and perhaps more parsimonious account is that these biases emerge out of the same dynamics that give rise to familiar word recognition.

In our dynamic associative model, both novel and familiar words undergo the same competition and associative learning. Given that

early in training all words are novel, the fact that this model is able to learn words at all suggests that it is not a priori necessary to separate novel and familiar word processing to solve the referential ambiguity problem. However, it is not clear whether the model also shows systematic biases in how it interprets novel words.

Simulation 2.2 evaluates mutual exclusivity. Mutual exclusivity is the idea that a novel word cannot refer to a referent that has a previously established word–referent link; that is, words are mutually exclusive. In terms of behavior in referent selection tasks, however, children are said to be following the mutual exclusivity constraint when they exclude objects with known names as the referent of novel words (Mervis & Bertrand, 1994). Note that there are debates over whether the inference children make in such cases is best described as mutual exclusivity (Markman & Wachtel, 1988), a process of matching novel words to novel objects (Mervis & Bertrand, 1994), or another form of inference (Halberda, 2006). Likewise there are debates over whether this constitutes learning (Horst & Samuelson, 2008; Spiegel & Halberda, 2011). Here we are only using mutual exclusivity as a moniker for a behavioral phenomena that happens in the moment when children are confronted with a novel name, a novel object, and several familiar objects. We are not implying any specific inference principle (in our view it arises from competition dynamics), nor are we implying learning—in terms of either an initial, quick link (as the term *fast mapping*, often applied to such situations, does) or a longer, more robust connection. Rather, we will simply refer to the phenomenon of selecting the object that does not have a name when confronted with novel and known objects and a novel name as *referent selection by mutual exclusivity*, or *M.E. reference selection* for short.

Crucially, M.E. reference selection minimally requires a range of available objects, a novel word, and some partially learned weights. All of these are present in this model. Thus, though many of the proposed biases and constraints “live outside” the simple architecture of this model, mutual exclusivity is clearly within the domain of our model, raising the possibility that it could arise in the context of the competition dynamics.

To examine this, 20 networks were initialized with 40 auditory and visual units and 500 output units (see Table 3). Of the 40 input units, 30 were used during training; the remaining 10 novel units were never heard or seen. Thus, by the end of the 100,000 training trials, the model had a lexicon of 30 words, but an additional 10 novel words, whose weights were largely unchanged. After training at 50% referential ambiguity, the models were tested in three ways. For ease of description, we describe these as three-letter strings, with the first letter representing the status of the target. First, we used a 3AFC task with all familiar words (F<sub>1</sub>FF) for comparison with previous models. Second, we tested M.E. referent selection trials, using a novel target, with two familiar objects (N<sub>1</sub>FF). Finally, we used the same configuration of competitors, but with a familiar target (F<sub>1</sub>FN). Here, if the model always selected the novel object, we should see a performance decrement. To construct five novel and familiar trials, the network needed to know at least 10 words, so we tested all 30 words in a production task prior to constructing these test trials, and models without a 10-word vocabulary were not tested.

**Results.** The models ultimately acquired 28.7 of the 30 words (by the production task). Figure 12A shows the models performance on the three primary tests. Models were 99.8% correct on the 3AFC

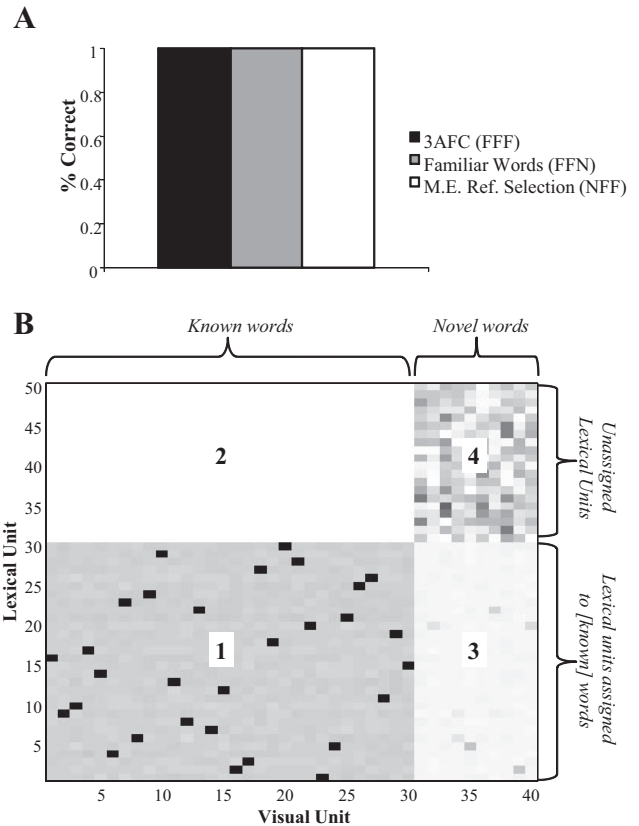


Figure 12. (A) Performance on the 3AFC, familiar word, and mutual exclusivity referent (M.E. ref.) selection trials. (B) Representative weight matrix. Darker patches indicate stronger connections. AFC = alternative forced-choice; F = familiar; N = novel.

task and 100% correct on familiar word trials and M.E. referent selection trials. Thus, fully trained models were effective at identifying familiar words (both with and without the novel object) and at matching the novel word with the unnamed visual object.

Given that the model has no explicit mechanism for mutual exclusivity, how did this emerge? The answer lies in the interaction of the competition and the learning rules. Figure 12B shows the visual-to-lexical weight matrix for a single run of the model after 100,000 training epochs. The x-axis shows the index of each of the 40 visual units (Units 1–30 are familiar words, and Units 31–40 are novel words). The y-axis shows the index of the lexical units. The model had 500 lexical units, but we only show lexical units that were strongly connected to one of the familiar words along with a random sample of 20 output units for the novel words.<sup>7</sup>

For familiar words each object is strongly connected to a single lexical unit. This can be seen in Region 1, which resembles the weight matrices displayed earlier for a trained network. Similarly, Region 2 shows that the connections between these familiar objects and the unused lexical units were eliminated, so it would be difficult to assign a novel word to one of these objects. Region 3

<sup>7</sup> We also changed the order to group the “assigned” lexical units together (a typical run of the model would recruit output units from across the array with no relationship to the order of the inputs).

shows that lexical units that have been assigned to a visual input are also not connected to any other visual input, including the novel objects. Most importantly, however, Region 4 shows that the novel objects all have small and variable connections to the remaining, unused lexical units. As a result, it is highly likely that there will be an associative path that allows activation to spread from a novel word to a novel object, even though these associations are random. Conversely, it is quite unlikely that a novel word could activate known objects.

This particular structure in the weights comes from the weight decay in the learning rule:

$$\begin{aligned}\Delta w_{xy} &= v_x l_y (1 - w_{xy}) \\ &\quad - .5 \cdot (1 - v_x) \cdot l_y \cdot w_{xy} \\ &\quad - .5 \cdot v_x (1 - l_y) \cdot w_{xy}.\end{aligned}\quad (9)$$

Here weights decay in two circumstances. First, they decay whenever a visual unit is on and the corresponding lexical unit is not. Thus, the connections in Region 2 are pruned because the familiar objects have developed connections with a single lexical unit, and lexical inhibition prevents a second one from becoming active. Second, they decay if a lexical unit is active but the visual input is not. Thus, the weights in Region 3 decay because the network has not encountered these objects but has used the corresponding lexical units for other words. The weights connecting novel inputs to unused lexical units never decay because neither class of units is ever active. As a result, connections in Region 4 maintain their original (small, random) values. Then, when a novel word is encountered, these connections permit the network to activate (to various degrees) a large number of lexical units. This activation feeds back to the visual units, but since the familiar objects are not connected to any of these now active lexical units, only the activation of the novel visual objects is amplified. This allows the network to select the correct (novel) object.

Thus, the selection of the novel object is dependent on the learning rule, but not because the network needs to learn something about that object or word. Rather, the weights between the known words and objects and the unused lexical units must decay, and the weights between the novel ones must not in order to create a platform upon which real-time competition dynamics can select the right object. A different type of weight decay (e.g., if all weights decayed on each epoch) would not preserve the right form of the weight matrix. However, learning is not the whole story: This pattern of connectivity could not be harnessed in situation time without the gradual settling process represented by the inhibition and feedback dynamics. Moreover, the model's ability to learn from M.E. referent selection may also depend on this competition–feedback cycle. The model must select a single lexical unit and selectively amplify the novel object in order to eventually turn a word–referent link created during M.E. referent selection into a known word by associating the novel object with the novel word over many instances. Thus, though as a real-time process mutual exclusivity is likely to impact learning, it is really more the product of learning than a mechanism of it. This implies that some types of learning environments may make it more difficult for children to engage in this by eliminating this particular structure of associative weights. This will be examined in Simulation 4.2.

### Simulation 2.3: The Development of M.E. Referent Selection and Visual Familiarity

In some ways, the previous model performs mutual exclusivity too well—children rarely approach 100%. Yet, this was a fully trained adult model, so it was important to examine the model developmentally. There have been few comprehensive developmental investigations of mutual exclusivity. We do know that in a 2AFC task children can succeed at about 18 months, depending on vocabulary size (Markman et al., 2003), but fail at 5AFC novel word tasks until after the vocabulary spurt (Mervis & Bertrand, 1994). Further, Halberda (2003) showed a clear developmental time course with 14- and 16-month-olds failing in a 2AFC looking version of the task but 17-month-olds succeeding. Thus mutual exclusivity is not an innate constraint but develops over time. Given the prior simulation demonstrating the dependence of mutual exclusivity on learning, we investigated this by rerunning the 20 models described in the previous section but measuring performance every 5,000 epochs.

**Results.** Figure 13A shows the results. The lines marked by open diamonds show the model's performance on the 3AFC ( $F_{FF}$ ) and production tasks. These were run regardless of how many words the model knew and show a steady improvement over learning. The models knew enough words to be tested on M.E. referent selection by around 45,000 epochs, and at this point performed at nearly 100%: By the time model knew enough words to be tested on mutual exclusivity, that ability was present. Indeed, runs of this model with fewer novel word trials (hence requiring fewer known words) show even earlier abilities, suggesting that this model may be able to do this task with very little experience.

This was unexpected. Apparently it did not take much learning to set up the right structure in the weight matrix (and given the dramatic changes in irrelevant connections seen in the first 1,000 trials in Figure 9, this may be sufficient). One factor that may moderate this is visual familiarity. Our analysis of the weight matrix suggests that good M.E. referent selection derives from the fact that the novel visual units have never been active to any degree. Yet, most experiments do not use stimuli that are completely unfamiliar. Typical novel objects such as whisks and juicers, though unlikely to be named, have likely been seen before, or may be similar to things that children know. Thus, we ran an additional set of simulations in which the novel visual units were seen (but never named) on some proportion of the trials. The likelihood of seeing a novel object varied from 5% to 50% (since the referential ambiguity rate was 50%, so this last condition was equivalent to the unnamed objects being as familiar as the known objects).

**The effect of familiarity.** Results are shown in Figures 13B–13F. Figure 13B shows the lowest level of familiarity—novel objects appeared 5% of the time—and Figure 13F shows the highest, in which novel objects were as likely to appear as known objects (though never named). The familiarity of the novel objects does not seem to influence responding on the familiar word trials ( $F_{FN}$ )—performance was equally good in all simulations. Importantly, however, even a small amount of visual familiarity impedes M.E. referent selection at early points in development. Figure 13B shows that a marginal amount of familiarity brings initial performance down to 55%, and any more can bring it down to chance. Very quickly after that, performance seems to develop to full



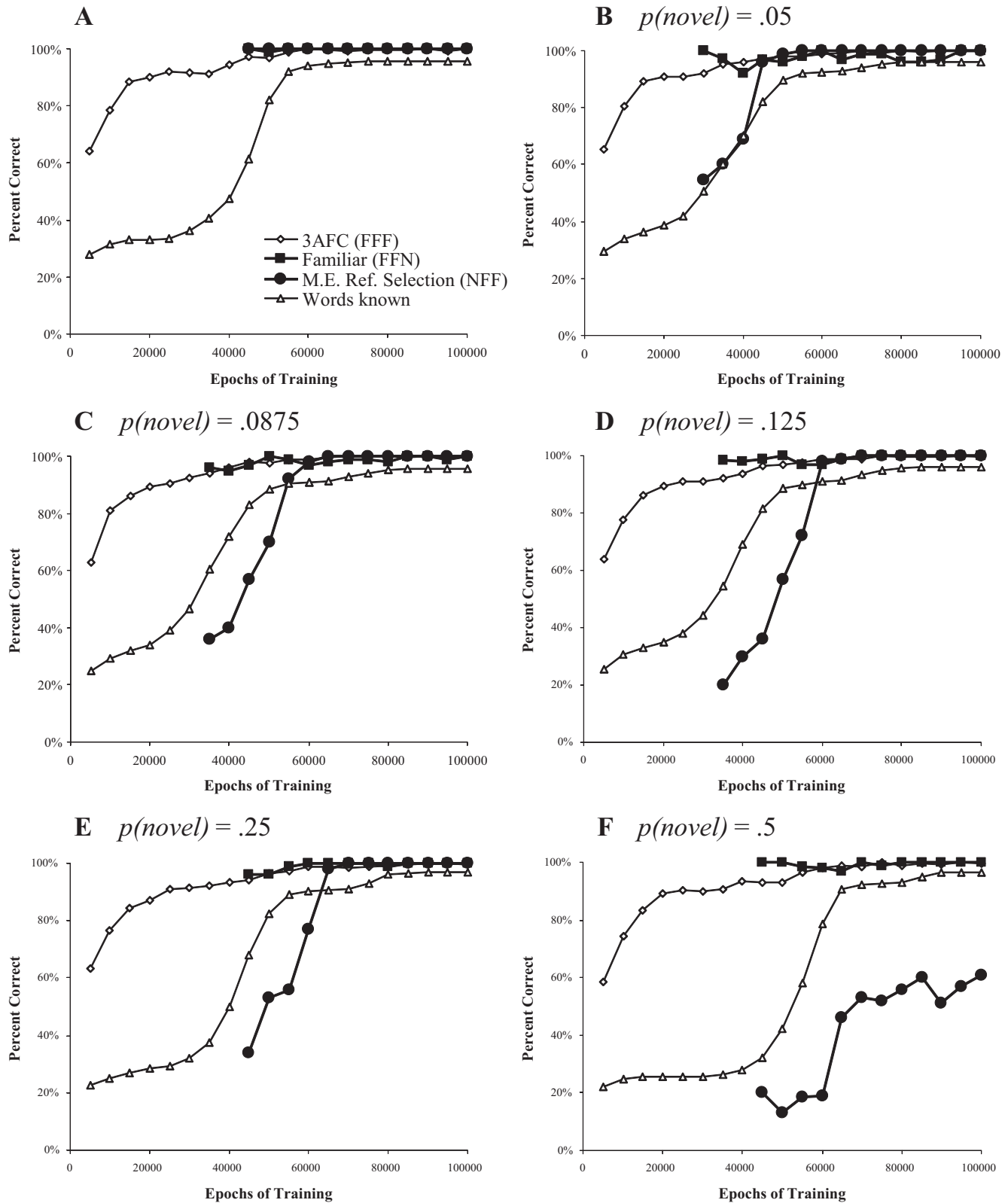


Figure 13. Performance on familiar and mutual exclusivity referent (M.E. ref.) selection trials over development (filled symbols). Also shown are the number of words known (percentage of 35) measured by both the 3AFC and weight analysis (open symbols). Panels represent different likelihoods of the novel words appearing as visual foils. AFC = alternative forced-choice; F = familiar; N = novel.

capacity. However, when the novel foils are highly familiar, M.E. referent selection is never very good. Thus, at least for this model, mutual exclusivity is not solely about lacking a name (as accounts like N3C suggest). The familiarity of the object may play a role as well.

Thus, there may be effects of visual novelty in these tasks that can be seen developmentally—the use of more novel objects could lead to better referent selection by mutual exclusivity. Indeed, Horst, Samuelson, Kucker, and McMurray (2011) showed that in a completely unconstrained referent selection task (three novel objects), children were heavily biased toward objects that were not seen in a brief familiarization. This raises the possibility that at the early stages, mutual exclusivity is really more of a novelty preference, rather than a complex inference. Moreover, beyond visual novelty there may be other ways to slow down the development of M.E. referent selection. For example, we have found in ongoing work with this model that including some quantity of low-ambiguity (ostensive naming) trials along with higher ambiguity trials can give rise to a similar effect (McMurray et al., in press).

## Section 2 Discussion

These simulations capture a number of classic findings in word learning, including the improvement in efficiency of familiar word recognition and referent selection by mutual exclusivity. They show that changes in familiar word recognition, though manifested in online processing, are better characterized by the shape of the learned connections. Moreover, this same learning can give rise to both a deceleration in RT and acceleration in the number of words known. Word learning can be a general learning process. Similarly, referent selection by mutual exclusivity, though it may appear a specialized inference process, can arise out of the same competition dynamics as familiar word recognition, when this process plays out over connections established during learning. This allows us to capture how exposure to objects alters referent selection by mutual exclusivity. More importantly, embedding this within a system that recognizes familiar words and learns word–object linkages allows for a richer explanation.

In both of these simulations, these ostensibly situation-time processes are the product of learning—but not any simple version of learning. With respect to familiar words, the most important predictor of processing speed is how the unnecessary connections decay. Similarly, mutual exclusivity fundamentally relies on a learning rule that describes a particular pattern of weight decay. Thus, suppressing competing associations is essential to multiple aspects of word learning. Similarly to Fernald et al.'s (2006) account of changes in processing speed, Mervis and Bertrand (1994) suggested that the number of words is the critical factor that predicts the onset of M.E. referent selection. However, this does not offer a clear mechanism of change in this context because the competitors are always familiar in mutual exclusivity task. Our model suggests that development has more to do with the pruning of weights (which is likely correlated with the number of words known, and was in the model). Across both simulations, however, the more important message is that to take advantage of the explanatory power inherent in this version of associative learning, we must consider both the positive and negative associations.

## Section 3: The Relationship of Situation- to Developmental-Time Processes

The previous simulations demonstrate that apparently situation-time processes are the product of learning. This section addresses the converse: How do the details of processing impact learning? Simulation 3.1 models data suggesting independence of time-scales: Horst and Samuelson's (2008) work on retention after M.E. referent selection. Simulation 3.2 examines task effects on mutual exclusivity over development. Simulation 3.3 returns to familiar word recognition, and examines longitudinal work showing that recognition time predicts the future rate of acquiring new words. Finally, Simulation 3.4 asks whether learning can occur without processing.

### Simulation 3.1: Referent Selection by Mutual Exclusivity and Retention

Horst and Samuelson (2008) showed that children do not retain words after referent selection by mutual exclusivity. This suggests that this behavior is a situation-time process and not synonymous with learning. Our dynamic associative account is ideal for capturing such effects: Referent selection emerges out of online competition, while learning is slow and may not be able to acquire a word in one exposure. Our model can extend these findings by asking whether anything is retained from referent selection and what circumstances may be necessary to see it.

We simulated Horst and Samuelson (2008) by initializing 20 models with 40 input units, but only training them on 30 words (see Table 4 for parameters). This left 10 novel words and objects that did not receive any training. Five were used to test mutual exclusivity, and the other five were held out. Referential ambiguity for familiar objects (Words 1–30) was set to 50%, with novel objects appearing as competitors 8.75% of the time (but novel names were never heard).

Models were tested on several tasks. First, we assessed which words were known with the production task. Only words that passed this test were used in subsequent testing. Next the model received five novel ( $N_iFF$ ) and five familiar ( $F_iFN$ ) word trials as in Simulations 2.2 and 2.3. However, unlike those simulations, the model learned on these trials, enabling tests of retention for these words. Finally, on retention trials, each of the five novel objects named on the prior M.E. referent selection trials was paired with another novel object and a held-out object ( $N_iNH$ ). Though learning (weight change) occurred throughout the sequence of test trials, the total learning over a single batch of testing was not carried back to training.

**Results.** Figure 14A shows the model's performance on each type of trial over the course of training. As before, familiar and novel word performance was excellent after the emergence of mutual exclusivity at around 45,000 epochs. However, the model was not able to retain the words that were tested in the M.E. referent selection trials, averaging 38% correct retention.

Based on the prior simulations, we were concerned that the visual familiarity of the novel objects may have created this effect. Thus, we replicated these simulations under two conditions: one in which the novel and held-out objects were never seen during training and one in which they occurred frequently ( $p = .375$ ). When the novel objects were completely unfamiliar, results were

Table 4  
Parameters for Simulations in Section 3

Parameter	3.1: Familiar word retention	3.2: Familiar word task	3.3: Longitudinal word recognition	3.4: Learning and competition
Input units	40	40	35	35
Familiar words/objects	30	30	35	35
Novel words/objects	5	10	0	0
Held-out words/objects	5	0	0	0
Lexical units	500	500	500	500
Initial weight size	.25	.25	.5 $\pm$ .025	.5
Learning rate	.0005	.0005	.0005 $\pm$ .00002	.0005
Referential ambiguity	.5	.5	.65 $\pm$ .03	.2, .5, .8
Novel object seen	.0875	.0875, .25		
Feed-forward temperature	.01	.01	.01 $\pm$ .0065	.01
Feedback temperature	2	2	2 $\pm$ .1	2
Stability point	1e-12	1e-12	10 <sup>-12</sup> $\pm$ .25	1e-12
Input inhibition	1.05	1.05	1.05 $\pm$ .01	1.05

similar (see Figure 14B): Referent selection performance was at ceiling (as in Simulation 3.1), but retention was at chance. However, when the novel objects were highly familiar (but unnamed), there was a period during which the model did well in referent selection, but not retention, followed by later points in development when the model could do both (see Figure 14C). This fits with recent work by Kucker and Samuelson (2012) showing that 24-month-old children can retain links created in an M.E. referent selection context if they play with the objects prior to the mutual exclusivity trials. It is also relevant to Spiegel and Halberda's (2011) recent finding that older children (30-month-olds) appear to retain word-object mappings, though in an easier-looking task.

Overall, then, the model fits the pattern of Horst and Samuelson (2008), showing excellent performance in referent selection by mutual exclusivity, but a failure to retain when contextual support is removed. This raises the question of how much, if anything, the model learned from a single mutual exclusivity trial. Though gross performance did not yield evidence of learning, there may be a small amount of learning that is insufficient to drive overt retention.

To assess this, we examined the amount of change in subsets of each weight matrix at various points in training using the root-mean-squared difference between the weights at two points in time (e.g., before and after the mutual exclusivity trials). Weights were divided up (see Figure 14D) into weights connecting lexical units to (a) familiar words, (b) novel words, and (c) held-out words. We then computed the weight change (learning) in each group that occurred during learning and during mutual exclusivity trials. Including the held-out units allows us to determine how much change to expect for completely unused items.

Figures 14E and 14F shows the results. Figure 14E shows the amount of weight change up to the point where the model was tested at 100,000 epochs. There is quite a bit more change in the weights for familiar words (which are being trained) than the novel or held-out words, particularly in the auditory weights. This makes sense: The novel and held-out auditory units are never activated, whereas the novel visual ones occasionally appear as competitors. In contrast, Figure 14F shows the amount of weight change during the mutual exclusivity trials. There was some learning on these trials and generally more learning for the novel words than the others. However, the amount of learning on these trials is far less

than what was learned about those words over the course of training—when they were never heard! It is also far less than what a truly familiar word would have received. Moreover, this small amount of learning is not responsible for the excellent performance in referent selection by mutual exclusivity—when learning was turned off during these trials (Simulation 2.2), the model still performed at 100%.

Nonetheless, this offers a clue to how M.E. referent selection relates to learning. The model learns a little something from each of these trials, and over the course of many such trials, this accumulates to yield complete word learning (see also Horst et al., 2006). But crucially, that tiny amount of learning we observed on that first exposure to a novel word is not different from what would be observed on the second, third, or fourth exposures. Moreover, this learning consists not only of building or maintaining correct associations, but also (and to a much larger extent) of suppressing unnecessary ones. Thus, what happens during this first referent selection is quite different from what earlier views (Carey, 1978) may suggest. Thus, referent selection by mutual exclusivity, though a primarily in the moment process, leaves a small trace in the weights that can accumulate to achieve real knowledge.

In retrospect, the training used in all the simulations thus far likely included many mutual exclusivity trials. Since competitors were randomly chosen on each epoch, there were likely many epochs in which the model knew all the words except the target (or knew more about the competitors than the target). In this way, there is nothing fundamentally different about familiar and novel words.

### Simulation 3.2: Effect of Task

Section 1 suggests that task has a critical effect on the model's performance (e.g., the delay in productive vs. receptive vocabulary). Similarly, mutual exclusivity also has the characteristics of a task effect: The two familiar words constrain the task, permitting the model to perform well despite no knowledge of the novel word. This predicts that task variables such as the number of alternatives may affect mutual exclusivity, particularly early in development.

Mervis and Bertrand (1994) measured referent selection via mutual exclusivity as a function of vocabulary size. They found that only children with relatively large vocabularies (greater than

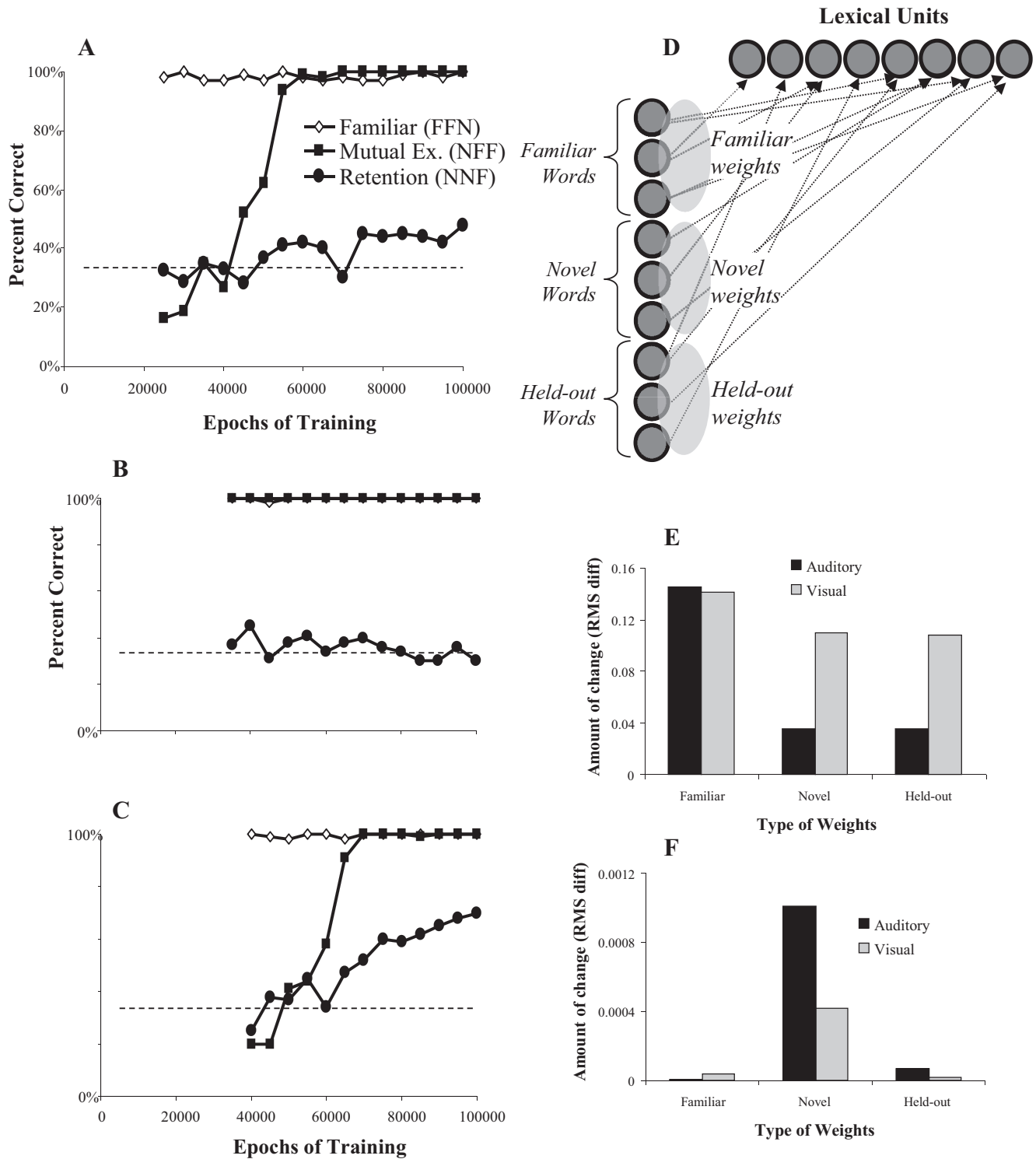


Figure 14. Performance on mutual exclusivity and retention. (A) Performance on familiar word, mutual exclusivity (ex.) referent selection, and retention trials when novel items appeared with a likelihood .0875. (B) Same, but novel items appeared with  $p = 0$ . (C)  $p = .375$ . (D) Diagram of the weight matrix indicating which components were measured. (E) Root-mean-squared (RMS) change in each type of weights during learning. (F) RMS change during the five mutual exclusivity trials. F = familiar; N = novel.



90 words) consistently performed well in the mutual exclusivity task (despite the fact that the familiar objects were known to each child). Moreover, for the children who failed at this task, once their vocabularies reached this level, they too could select the correct referent. Thus, some critical quantity of words may be required before this ability emerges. However, Markman et al. (2003) pointed out that the four- and five-alternative mutual exclusivity tasks used by Mervis and Bertrand may have simply been too difficult. They used a two-alternative task (one familiar and one novel object) and demonstrated that regardless of vocabulary level, all children could succeed at the mutual exclusivity task. Horst et al. (2010) further showed that by 30 months, the number of competitors made little difference in children's performance (though it did affect retention). Thus, task differences may disappear later in development.

To examine this, 20 models learned 30 words at 50% referential ambiguity. As before, 10 novel words were held out for testing. Novel visual units appeared as competitors with a likelihood of 8.75%, to obtain nonceiling performance on mutual exclusivity trials. Every 2,000 epochs, the model was tested in both 3AFC and 5AFC familiar and novel word tasks. For both familiar and novel word tasks, the competitors consisted of two or four familiar objects and one novel object ( $F_1FN$ ,  $F_1FFFN$ ,  $N_1FF$ , or  $N_1FFFF$ ). As before, words were screened with the production task, and models were not tested unless sufficient words were known.

**Results.** Figure 15A shows that, as expected, models performed equivalently and at ceiling on both 3AFC and 5AFC familiar word tasks. In the M.E. referent selection tasks, models performed better in the 3AFC task than the 5AFC task early in development, also as predicted. Figure 15B shows results from a second set of simulations in which the novel objects appeared 25% of the time. It suggests that this task difference was enhanced when the objects were more familiar. In both cases, however, the effect of task was eliminated quickly after M.E. referent selection got off the ground. Thus, as in children, task differences in referent selection by mutual exclusivity are only observed at a narrow window in development.

### Simulation 3.3: Processing Speed and Rate of Acquisition

We have described referent selection by M.E. as primarily a situation-time inference process, but one that is fundamentally based to prior changes in the weights due to learning. We now ask whether such interactions between timescales are also seen in familiar word recognition.

A study by Fernald et al. (2006) offers an intriguing platform for this. They examined 63 infants longitudinally between 12 and 25 months and assessed both the number of words known (using the MCDI), and the infants' speed of processing familiar words using the looking-while-listening task. They found first that speed of processing in this lexical task (RT) is stable across individuals (though decreasing) from month to month. Second, there was a correlation between speed of processing and accuracy in preferential looking. Third, RT was negatively correlated with the number of words known, particularly for the older children (25 months). These first findings could be accounted for by simply assuming that processing speed is a function of learning, much as we showed in Simulation 2.1; learning influences processing.

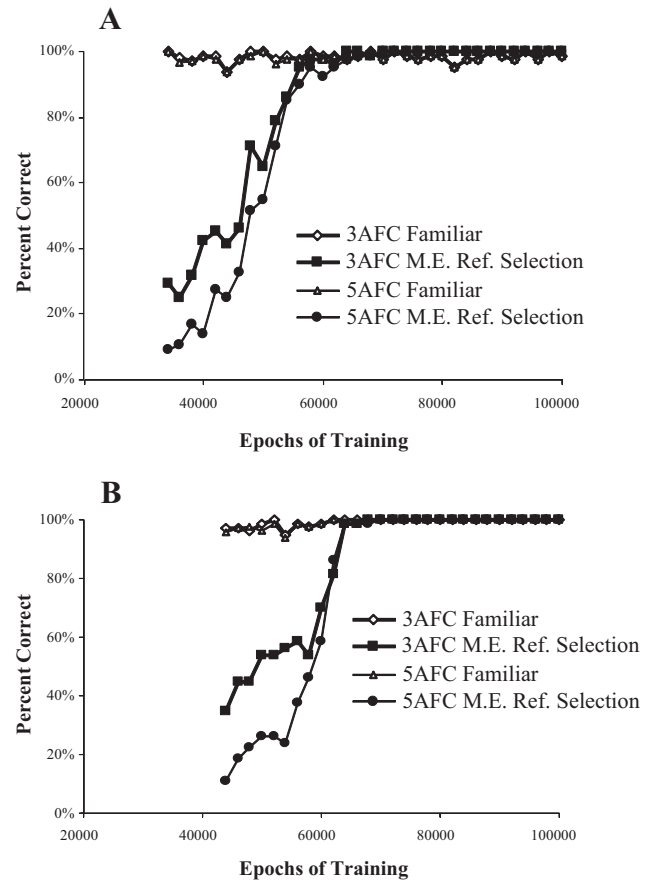


Figure 15. Performance on 3AFC and 5AFC familiar and fast-mapping tasks over training. (A) Novel objects appeared with a likelihood of  $p = .0875$  during training. (B)  $p = .025$ . M.E. ref. = mutual exclusivity referent. AFC = alternative forced-choice.

However, Fernald et al. also found that the children with significantly shorter RTs at 25 months showed more acceleration in the number of words acquired. This suggests the converse, that processing influences learning. Fernald et al. concluded that processing speed may be a property of the child that is fundamentally related to word learning. This motivates a compelling link between online processing and learning. They argued that children who can identify words faster have more resources (or time) to process subsequent words, allowing further opportunities for learning.

This simulation asks whether our model shows these same dependencies. This can validate the model by simulating a complex set of phenomena and extend it to individual differences. It also affords the opportunity to understand the mechanisms that give rise to this particular relationship between learning and processing. In particular, our model does not have the sequential processing demands of real children, so if it still shows such dependencies, it may help illuminate alternative accounts. Moreover, as processing speed can derive from both differences in activation flow and differences in learning, we can start to understand the range of causes that can give rise to this crucial descriptor.

To simulate this, we needed individual differences across models. Though some variation is created by the random initial weights

and the random sequence of training, this was not sufficient. Thus, across simulations we also varied the free parameters of the model as an additional source of variation (cf. McMurray, Samelson, Lee, & Tomblin, 2010). We initialized 100 networks and varied parameters such as the temperatures, the inhibition, and the learning rate by adding Gaussian noise to the means used in prior simulations (see Table 4).

Models were trained on a lexicon of 35 words. They were tested every 1,000 epochs on all 35 words in a 2AFC task, and we recorded both the number of words known (by this measure) and the average settling time. Similar to that in Fernald et al. (2006), settling time was only recorded for trials on which the network answered correctly. Models were also tested on a 20AFC comprehension task. Our analysis follows the findings of Fernald et al. (2006). We assess the stability (within a time slice) of our measures of RT and accuracy, and then the correlation of these measures with the rate of long-term learning.

### Results.

**Stability of RT and accuracy measures.** Table 5 shows the pairwise correlations in RT and accuracy in the 2AFC task for adjacent tests. Both RT and accuracy were highly correlated across time, with an average correlation of .9 (RT) and .57 (Accuracy). These are higher than the correlations found by Fernald et al. (2006), who reported correlations in the range of .2–.4 for RT and .25–.5 for accuracy, but not unexpected for two reasons. First, children have a range of processes outside the model that may introduce variability. Second, our test included all 35 words, whereas Fernald et al. only tested four to eight words—a much smaller sample of a much larger lexicon. As a result, the models' estimates of RT and accuracy were likely to be closer to the true values than behavioral work can derive for children.

Despite these correlations, however, the model is not perfectly stable. Figure 16A shows 2AFC performance over training for six representative runs of the model. Model 09, for example, starts with one of the worst performances, but is quite successful by the end. Model 07 is the worst, but by about 10,000 training trials (Log 4), it tracks quite closely with 11, the best. There is also considerable variation when the models start to perform quite well. Figure 16B shows a similar pattern for RT. Models 07 and 10, for example, start with the worst RT, but end up with the lowest, whereas Models 11 and 13 start low and stay low. Model 09 starts similar to Models 11 and 13, but ends high. Thus, despite these remarkably high pointwise correlations, looking at the whole

time course of development (and sampling at a much higher rate), we see that the underlying instability of the developmental time course is not well captured by the correlations. This supports the kinds of microanalyses advocated by Adolph et al. (2008).

**Relationship of speed to accuracy.** Fernald et al. (2006) also found negative correlations between speed and accuracy at any given month: Children who settled faster got more words correct ( $R = -.3$  to  $-.5$ ). This was also observed in the simulations. Table 6 (Accuracy column) shows the correlations between the number of correct 2AFC trials and the log of the settling time. At 5,000 and 10,000 time steps there was no significant correlation between these two factors. However, this is expected, because the models knew an average of 3.8 words at 5,000 and 4.9 at 10,000 (chance on our 20AFC task would predict 1.75 words). By 15,000 epochs, however, the models knew an average of 7.5 words (by our 20AFC assessment), and correlations between RT and accuracy were significant ( $R = -.30$ ,  $p < .01$ ) and increased throughout development (to  $R = -.85$  by 100,000 epochs). These correlations were negative: Models with lower settling times knew more words.

**Relationship of speed to vocabulary growth.** In the longitudinal study, RT at 25 months predicted the acceleration in word learning across the period studied. To assess this, Fernald et al. (2006) fit quadratic functions to the number of words known at each month. They compared each term of this function to the RT at 25 months and found a significant correlation with the quadratic term (but not the linear term or the intercept), suggesting that settling time is related to acceleration in vocabulary growth.

The model showed the same behavior. For each model, we fit a quadratic function to the number of words known (in the 20AFC task) over the first 25,000 epochs. This roughly corresponds to the period of early learning studied by Fernald et al. (2006)—by 25,000 epochs, the models averaged 11 of the 35 words and were through their first period of acceleration. As Table 6 shows, the quadratic term was correlated with settling time at every time step, whereas the linear and intercept terms were not (except at the first time step). Thus, the model shows the same relationship between speed of processing and vocabulary growth as children.

Fernald et al. (2006) posited that in running speech, children who finish processing a word quickly can move on and learn from subsequent words. This does not seem to be the case here—the model is reset between words. Although this does not rule out this sort of bootstrapping in children, it does suggest that such relationships can arise from other causes. For example, in our model, the parameters controlling settling dynamics (RT) may also alter the networks' ability to resolve referential ambiguity, which could affect learning. Or conversely, as we have shown, settling time is primarily a function of learning, so both effects may derive from the same cause.

**What parameters influence outcome measures?** Fernald et al. (2006) described speed of processing as a fundamental parameter describing variation among children. At the level of description, this is undoubtedly correct, though the underlying mechanisms are not clear. Our model offers a set of candidate parameters, but no single one maps directly to speed of processing. Rather, the speed with which the model processes input is an emergent property that derives from multiple components: Parameters such as the temperature and the degree of inhibition that directly affect the dynamics of settling are clearly important, but

Table 5  
*Correlation of Reaction Time and Accuracy (With Themselves)  
Across Time Slices*

Epochs	Reaction time	Accuracy (2AFC)
5,000–10,000	.77**	.39**
10,000–15,000	.83**	.47**
15,000–20,000	.90**	.54**
20,000–25,000	.94**	.66**
25,000–30,000	.94**	.53**
30,000–50,000	.96**	.54**
50,000–75,000	.92**	.65**
75,000–100,000	.95**	.79**

Note. AFC = alternative forced-choice.

\*\*  $p < .01$ .

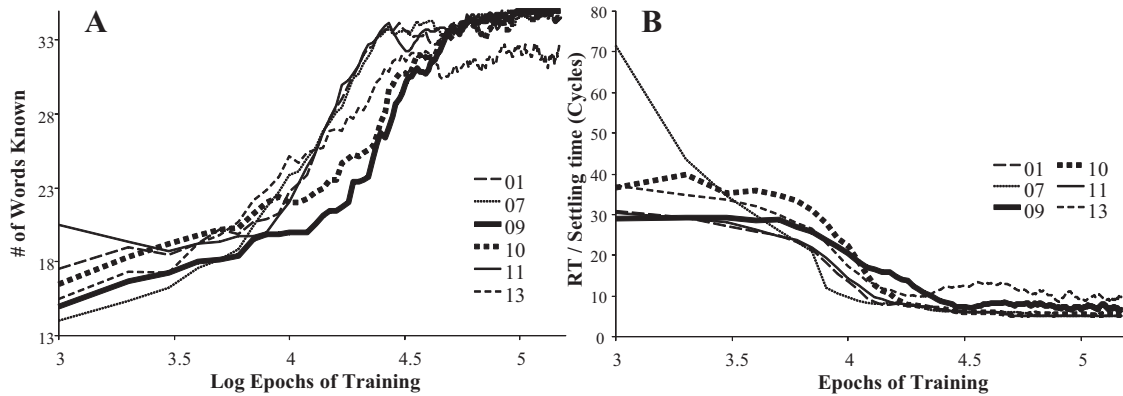


Figure 16. A representative sampling of individual models: (A) number of words known. (B) Settling time. RT = reaction time.

the biggest component may be the weight matrix, a product of learning.

We examined the relationship between our behavioral measures and each of the parameters that were manipulated to create individual differences across development. Correlations are shown in Table 7. A number of factors were significant predictors. The number of visual competitors, for example, was correlated with learning: Models with more competitors learned more slowly. The stability point (e.g., how little change was required to end a trial) was moderately correlated with RT (models with high thresholds tended to settle slower) but inversely related to learning (the models with the lowest thresholds knew the most words). Input inhibition played a role in settling time, but not so much in learning until the end of the simulation at which point models that had stronger inhibition showed more words learned. This suggests that the ability to inhibit competing representations in the visual layer is important both for quickly settling on a target and for ultimately learning the correct mapping. Learning rate was not a strong predictor, but later simulations (see supplemental materials, Simulation S3) suggested that the range of learning rates we tested may not have been sufficient.

Table 6  
*Correlations Between Speed and Accuracy, as Well as the Three Components of the Growth Curve in Overall Words Known (20AFC)*

Time step	Accuracy	Growth function		
		Quadratic	Linear	Intercept
5,000	.10	-.38**	.35**	-.16
10,000	.04	-.34**	.08	.10
15,000	-.30**	-.24*	-.07	.18
20,000	-.35**	-.25*	-.07	.16
25,000	-.48**	-.28**	-.05	.18
30,000	-.52**	-.28**	-.05	.15
50,000	-.64**	-.28**	-.06	.17
70,000	-.77**	-.29**	-.04	.16
100,000	-.85**	-.26**	-.06	.15

Note. AFC = alternative forced-choice.

\*  $p < .05$ . \*\*  $p < .01$ .

The most important predictor was feed-forward temperature, the rate that activation accumulated in the lexical layer from input layers. It was highly correlated with settling time (at all points in development) and with the number of words known at mid- to late points in development. Oddly, however, the correlations with RT were the inverse of what we expected: Models with higher temperatures settled slower (and learned worse). Follow-up simulations (reported in the supplemental materials, Simulation S2) demonstrate that this is due to the fact that models at a high temperature settle slower because they artificially increase activation for the competitors (as well as targets), and thus take longer to suppress them.

Given the complex role of temperature and other parameters contributing to settling time and learning, we suggest that concepts such as “speed of processing” do not reflect a unary dimension of the underlying architecture. Rather, they are emergent on a complex interplay of system dynamics, the performance in the tasks, and the developmental history. Crucially, even the parameters controlling dynamics were correlated with number of words learned (which in turn is a predictor of settling time), and thus many of these effects may be mediated via learning.

**Discussion.** Our dynamic associative model provides a compelling complement to Fernald et al. (2006). The model shows the same stability of RT across development as children, but suggests that there may significant instability when we look closer. It demonstrates the link between RT and accelerating learning but without any simple causal mechanism (e.g., processing capacity). Rather, the relationship derives from the fact that both processing and learning derive from changes in the weight matrix.

More importantly, the model offers a way to interpret RT. RT may not be isomorphic to some elemental individual difference. Rather, it emerges from the interplay of the properties of both learning and the dynamics of competition. These create a fairly stable measure, but one that affects RT and learning at different points in time, as we saw in Simulation 2.1 with respect to acceleration and deceleration. Although processing time is clearly an emergent property of network dynamics and learning, it also reflects individual differences in things such as the learning rate, the temperature, and the like. And because these things also affect learning, it suggests a highly circular and mediating relationship

Table 7  
*Relationship Between Control Parameters and Output Measures in the Network*

Measure	Initial weight size	Learning rate	No. visual competitors	Stability point	Feed-forward temperature	Feedback temperature	Input inhibition
Reaction time							
10,000	-.12	-.15	.09	-.19 <sup>†</sup>	.79**	-.03	-.35**
25,000	-.09	-.09	.14	-.17 <sup>†</sup>	.82**	-.04	-.30**
100,000	-.08	-.10	.09	-.12	.85**	-.07	-.26**
Words known							
10,000	.16	.07	-.31**	-.20*	.04	-.09	-.12
25,000	.22*	.12	-.33**	-.19*	-.31**	.02	.13
100,000	.09	.13	-.09	.27**	-.87**	.09	.29**
2AFC accuracy							
10,000	-.10	.21*	-.32**	-.27**	.01	.02	-.11
25,000	0	.06	-.38**	-.14	-.39**	.20 <sup>†</sup>	.17 <sup>†</sup>
100,000	.09	.19 <sup>†</sup>	-.09	.14	-.77**	.03	.22*
Growth curve							
Quadratic	.02	-.04	.01	.07	-.34**	.06	.24*
Linear	.09	.11	-.15	-.13	.09	.02	-.14
Intercept	-.01	.03	-.09	-.02	.05	-.11	0

Note. AFC = alternative forced-choice.

<sup>†</sup>  $p < .10$ . \*  $p < .05$ . \*\*  $p < .01$ .

between learning and processing. As Simulation 2.1 showed, processing time is determined in large part by the nature of the weight matrix (and the irrelevant connections in particular), so it should not be surprising to find such a relationship. Only if processing time is treated independently of lexical knowledge does this seem surprising. This suggests that explanations of individual differences based on speed of processing (e.g., [Kail, 1994](#)), though perhaps behaviorally stable, may oversimplify the problem, particularly when learning is involved.

### Simulation 3.4: Is Online Processing Required for Learning?

The previous simulations suggest interactions between learning and processing, deriving in part from the nature of the weight matrix and how it influences real-time competition. In the present simulation, we take this to the extreme, asking whether real-time processing is necessary for learning.

Models of unsupervised learning in other domains suggest that unsupervised learning may require some form of competition ([McMurray, Aslin, & Toscano, 2009](#)), and most unsupervised architectures include some form of it (e.g., [Kohonen, 1982](#); [Rumelhart & Zipser, 1986](#)). Perhaps, then, competition is required for learning. If the competition–feedback dynamics allow the model to improve upon ambiguous inputs, it would seem more efficient to use the results of this processing as the basis of association, rather than the more ambiguous inputs to it.

There are three components of competitive real-time processing in this model. First, inhibition between lexical units allows more active units to suppress activation for competitors. Second, feedback between the lexical layer and the inputs helps the network suppress competing inputs (visual competitors) as it makes a decision about the word. Finally, inhibition among input units

helps the network suppress competing inputs. All three contribute to activating the correct lexical and visual units in situation time, but it is not clear whether they are necessary for learning.

Thus, we ran a series of simulations that factorially varied whether feedback, lexical inhibition, and input layer inhibition was used. Each combination was run at three levels of referential ambiguity (20%, 50%, 80%), yielding 24 simulations. This was repeated 10 times for 240 simulations (see [Table 4](#)). Each model was tested every 5,000 epochs. Models without feedback cannot adjust activation in the visual units, making the NAFC tasks useless. Thus, our primary measure was the analysis of the weight matrices.

**Results.** Inhibition at the input layer was a fairly small contributor to learning; thus, we averaged across models with and without it for most of the analyses. [Figure 17A](#) shows the number of words learned over training in each of the four permutations. Lexical inhibition was required for learning. The models without it acquired an average of 1.27 words, whereas all the models that used it acquired all 35 words. Associative learning of this type cannot proceed without the ability to suppress competitors at the lexical level.

The effect of feedback (assuming the presence of lexical inhibition) was more nuanced. The models with feedback (the full model) acquired a few words very quickly, followed by a delay before learning the rest (see [Figure 17B](#)). Models without feedback took longer to get started, but once they did, they quickly outpaced the models with feedback.

Although competition is clearly required, is there an advantage for feedback? Possibly. It may help the model to acquire a small working vocabulary quickly (see [Figure 17B](#)). It may also benefit online processing. [Figure 17C](#) shows the settling time in a 3AFC task of models with and without feedback. The processing ability



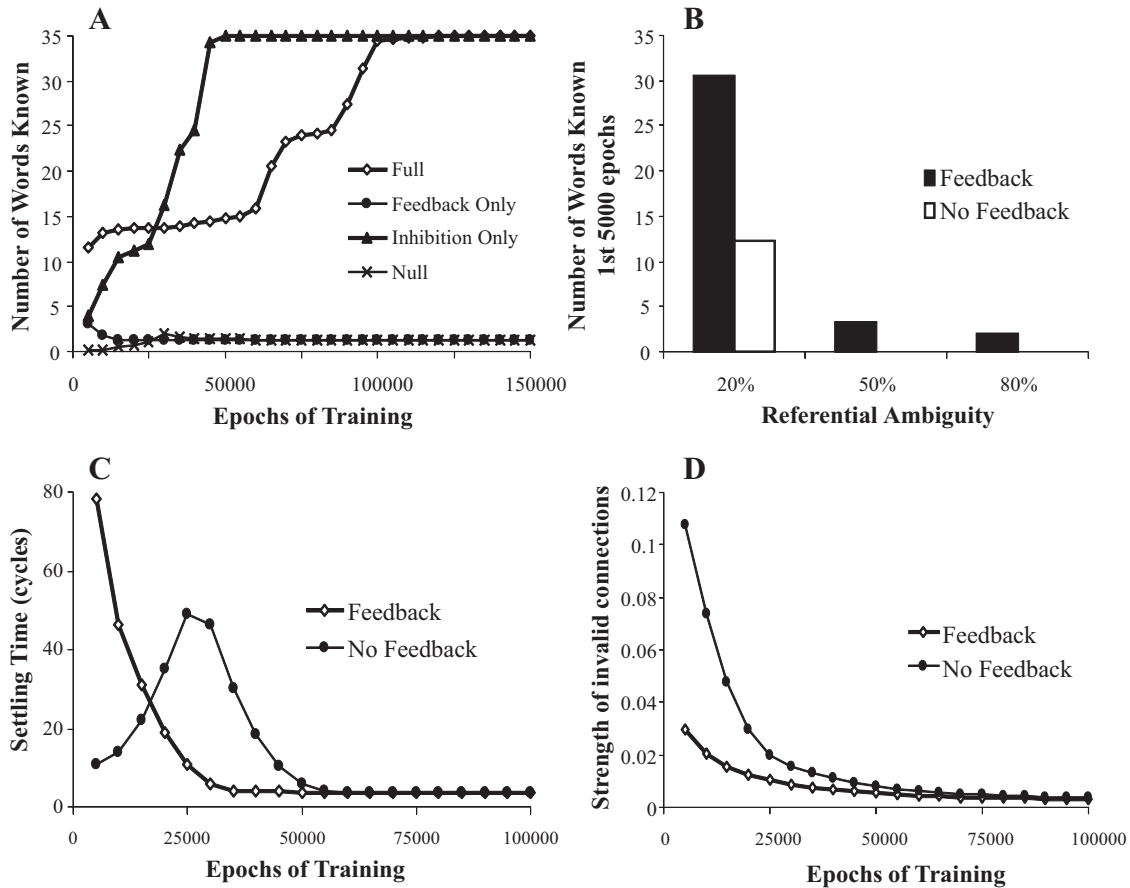


Figure 17. Model performance in Simulation 3.4. (A) Number of words known (via weight analysis) over development for the full model, and those lacking lexical inhibition or feedback. (B) Number of words known in first 5,000 epochs. (C) Average settling time on a single word as a function of lexical feedback and training. (D) Strength of invalid connections over development for models with and without feedback.

of the model with feedback improves much more rapidly and stably than those without. Moreover, though feedback slows the speed of learning, it may improve the quality. Figure 17D shows the average of the weights connecting any given lexical unit to the “incorrect” visual and auditory units. These should be 0 by the end of learning. As can be seen in the figure, situation-time feedback allows the model to suppress these connections faster over developmental time. Finally, feedback enables the model to change activation in the visual layer as a lexical unit is selected, the basis of our “behavioral” tests. Yet, this is not just a computational convenience. For children, word learning must be integrated with behaviors such as selecting referents (focusing attention). It makes more sense to harness the word–object connections established by word learning to guide these behaviors than to rely on a completely feed-forward system that would have to acquire new mappings to do this.

In Simulation 1.1, we described how online processing allows the network to perform better than its partially learned knowledge (weight matrix). However, the current simulations suggest that online processing is much more important than that. Associative learning in this model is simply not possible without some type of lexical inhibition—the model must make a decision about what

word it was hearing. Online processing is not merely shaped by learning, nor does it merely buttress performance. It is essential for, and integrated with, learning.

### Section 3 Discussion

This section simulated three empirical findings regarding the interaction of situation- and developmental-time processes. Across all three, we sought to determine (a) whether situation-time processes are independent of development; (b) whether development, particularly word learning, impacts situation-time processes; and/or (c) whether situation-time processes impact learning.

Our dynamic associative account captured some of the evidence for Objective a: the failure of children to retain recently novel words from mutual exclusivity trials. However, it suggests that though learning may be slow (requiring multiple trials), some learning occurs during referent selection by mutual exclusivity. It also accounted for Objective b: The model’s performance on 3AFC and 5AFC mutual exclusivity tasks was a function of its development. However, unlike prior explanations based on number of words known, the pruning of unnecessary connections was the determining factor. Our examination of Objective c was more

ambiguous. Although we modeled the Fernald et al. (2006) longitudinal work, the model suggested that processing and learning are both emergent from the whole network; however, when we manipulated the component processes of the model in Simulation 3.4, we showed that processing is required for learning.

This last finding suggests the model is not entirely unsupervised—it is self-supervised. The model does not just associate the inputs it sees and hears. Rather, it performs something like an inference process (implemented via competition and feedback) and then uses the output of such a process as the basis of association. The power of this model both to learn under high degrees of ambiguity and to account for a wealth of data speaks to the power of such a scheme.

#### Section 4: One Word/One Object?

Thus far, we have restricted our simulations to “pure” word-object sets in which every object has one label and every label has one referent. However, real language contains many polysemous words (with two meanings). Conversely, most objects can be labeled in multiple ways; most vividly, objects have labels at multiple levels (e.g., basic level, superordinate). In this final section, we begin to probe the limits of such situations in two ways to determine how much of the model’s performance can be attributed to its idealized “language.”

First, in addition to its use in describing real-time referent selection, mutual exclusivity has been described as a hindrance to learning second labels for objects (e.g., Markman et al., 2003; Regier, 2005). As Xu and Tenenbaum (2007) pointed out, this constraint must be relaxed for children to learn a second name for an object (e.g., its superordinate name). Our model shows the first sense of mutual exclusivity (referent selection in ambiguous situations), but it is not clear whether this will also impede learning second labels. This is examined in Simulation 4.1 by training the model on both basic and superordinate labels. Crucially, this allows us to study basic-level advantages in a system that does not represent taxonomy hierarchically.

Second, mapping multiple words to objects disrupts the one-to-one mapping between words and objects. This consistency may be essential both for learning in general and for the development of mutual exclusivity. However, it is unclear whether a purely associative system can generalize a principle across multiple words. Thus, Simulation 4.2 examines two situations that disrupt this mapping: (a) when words can refer to multiple objects (e.g., polysemy) and (b) when the same object can be referred to by multiple words.

#### Simulation 4.1: Multiple Labels

Learning multiple labels for an object may challenge both constraint and associative approaches. Constraints such as the taxonomic constraint or mutual exclusivity must relax to learn properties, synonyms, or other taxonomic categories (e.g., superordinates) for the same object. Similarly, associative learning could commit to a single label for an object and have a hard time linking a second one. Xu and Tenenbaum (2007) argued that to solve this problem, the system must be sensitive to statistical distributions and show graded constraint satisfaction. They argued that Bayesian inference uniquely has these properties. Indeed, the

localist representations and strong inhibition in our dynamic associative account may make it difficult to assign multiple labels to a one word.

Thus, this simulation examined the ability of the model to handle both basic and superordinate names. Models were initialized with 25 auditory word form units and 25 object category units (see Table 8). An additional five auditory units corresponded to five superordinate categories. There were no visual units for these: Each superordinate was associated with five of the 25 objects. On each training trial, we first selected one of the 30 auditory units (25 basic-level and five superordinate units). Thus, the likelihood of hearing a superordinate term was the same as each of the basic-level terms. If the auditory unit was a basic-level name, the corresponding visual unit was active (along with several competitors). If the auditory unit was a superordinate, one of its five corresponding basic-level visual units was activated. Consequently, basic-level names and visual units should be strongly associated with each other, whereas the association between superordinate names and their category members may be smaller (since it will be spread among five objects).

Basic-level performance was assessed by presenting one basic-level name, its referent, and two competitors (for each of the 25 words). Superordinate performance was assessed by selecting a superordinate name along with a target from that category, and competitors from two categories. Each superordinate name was tested five times (once for each member).

**Results.** Figure 18A shows the network’s performance on both tasks. By the end of training, the model mastered both superordinate and basic-level labels performing at 100% on both tasks. Such performance requires that the model acquire multiple names for each object, suggesting that the model displays the necessary flexibility. The model also learned basic-level names before superordinates, an example of the commonly reported advantage for basic-level terms (Rosch, Mervis, Gray, Johnson, & Boyes-Braem, 1976).

**The role of frequency.** One factor that could contribute to the basic-level advantage is frequency. The network was more likely to hear a basic-level name than a superordinate. To examine this,

Table 8  
*Parameters for Simulations in Section 4*

Parameter	4.1: Multiple labels	4.2: Multiple labels and referents
Visual units	25	30–15
Auditory units	30	30–15
Basic-level words	25	
Superordinates	5	0
Words referring to two objects	N/A	0, 5, 10, 15
Objects with two names	N/A	0, 5, 10, 15
Novel visual and auditory units	0	5
Lexical units	500	500
Initial weight size	.5	.5
Learning rate	.0005	.0005
Referential ambiguity	.5	.5
Feed-forward temperature	.01	.01
Feedback temperature	2	2
Stability point	1e-12	1-e-12
Input inhibition	1.05	1.05

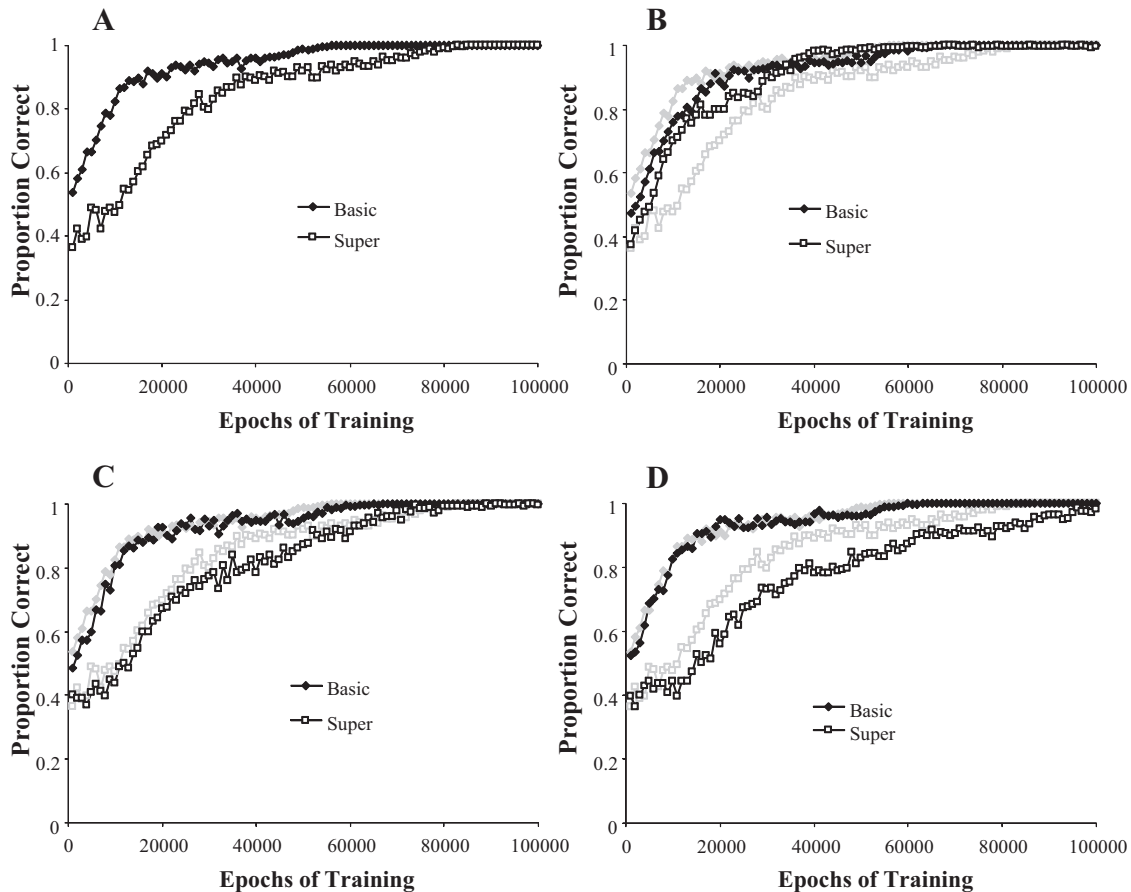


Figure 18. Performance on basic-level and superordinate categorization over development. In all panels, gray curves are the same data as Figure 18A. (A) In model with 25 basic-level terms and five basic-level categories. (B) In the same model in which the probability of a superordinate term was increased to 0.5. (C) Performance in a model with 27 basic-level and three superordinate terms. (D) Performance in that same model when the likelihood of any individual superordinate term was the same as basic-level terms ( $1/30$ ).

we ran an additional simulation in which superordinate and basic-level labels were equally likely by boosting the frequency of individual superordinates (superordinate:  $p = .1$ ; basic:  $p = .02$ ). Figure 18B shows the results (results of the prior simulation are shown in gray). As expected, the network continued to perform well on basic-level categories. This time, superordinate performance was closer to basic level, yet there was still a basic-level advantage.

**Spreading associations.** An additional contributor to the basic-level advantage is the fact that superordinate names must share associations with five objects. For every five exposures to a basic-level name the corresponding object will be seen five times; in contrast, for five exposures to a superordinate name, any given basic-level object will be seen once. Thus, superordinate names may have weaker associations with each of their members. To test this, we conducted a third simulation in which there were 27 basic-level names and only three superordinate terms. Each superordinate category now had nine members. This should enhance the spread of associations and lead to even worse performance. To ensure that frequency was not a factor, we ran two versions of the model. In the first, the overall frequency of superordinate names

was matched to the first simulation in this section where superordinates were as frequent as basic-level terms. That is, the probability of any of the three superordinates was  $5/30$ , or  $.167$ ; thus, the probability of any individual one was  $.167/3 = .055$ . In the second version, the frequency of individual superordinate names was the same as the frequency for individual basic-level names (given the smaller number of superordinates). That is, probability of any of the three superordinates was  $3/30 = .1$ ; thus, the probability of any individual one was  $.1/3 = .033$ . Figures 18C and 18D show the results. Regardless of the superordinate frequency, networks learned the superordinate categories, even with nine members; however, the basic-level advantage in both was larger than in Figure 18A, confirming that the spreading of association contributes to this effect.

**Discussion.** These simulations demonstrate that this dynamic associative model is capable of learning both basic-level and superordinate categories on the basis of co-occurrence statistics alone. Thus, the model is considerably more flexible than the one-word-to-one-object mappings we have largely focused on, and the ability to use mutual exclusivity in the moment does not necessarily constrain learning (nor does it rely on such a constraint).

Given the intermediate lexical representations, our model has two routes to achieve a mapping between one object and two words. It could associate a superordinate word form with multiple lexical units, each of which is associated with a single object; or it could associate a superordinate word with a single lexical unit that is associated with all the category members. It is unclear what the consequences of one or the other are, but a preliminary analysis of several networks' suggested that the latter was the dominant pattern—each superordinate word was associated with a new lexical unit, which in turn was associated with multiple visual units.

More importantly, however, this model illustrates that behaviors such as mutual exclusivity and the use of multiple labels can coexist comfortably in an associative architecture. The key innovation is that mutual exclusivity is not a constraint on learning (as in Regier, 2005; Xu & Tenenbaum, 2007). Rather, it is a constraint on online behavior (referent selection) that has consequences for an unconstrained learning system. In the moment, the network is nudged away from assigning a new name (the superordinate) to a known category (as predicted by mutual exclusivity). However, it must be nudged toward another object in the scene. Across trials, however, the available visual objects are not consistent, so mutual exclusivity never nudges it to consistently select one object for the superordinate name. As a result, much stronger cross-situational statistics take over and establish the correct associations.

#### Simulation 4.2: Violating One Word/One Object

The previous simulation shows that mutual exclusivity need not be a constraint on learning; the model can learn multiple labels for a given word despite the use of mutual exclusivity as an in-the-moment referent selection strategy. The final simulation examines three related issues. First, Simulation 4.1 did not actually test the model's ability to use mutual exclusivity. It is possible that even with the disruption in the one-word/one-object nature of the word-object mapping, the model could still learn, but that its ability to use mutual exclusivity is hampered. Many approaches to M.E. referent selection assume that such inferences are built on a realization by the child that each word refers to one object, and therefore a novel object must have a novel name (Halberda, 2006; Markman & Wachtel, 1988). Although this generalization must be acquired from the statistics of word-object relationships (which our model is sensitive to), our model has no way to store such a principle or strategy. If violating the one-word/one-object assumption impairs the model's ability to use mutual exclusivity, this would be a powerful demonstration that even associative systems can show principled generalization across words.

Second, Simulation 4.1 examined a special case in which objects have multiple names with a clear hierarchy. However, there are also cases in which objects have two equally probable names. An extreme example is bilingual children who learn two words for most objects. An equally important property of real languages is the converse, in which a name can refer to two objects or categories. This property of polysemy is common: Most words have multiple meanings, but it may have different consequences for both learning and M.E. referent selection from the many-names/one-category situation.

To examine this, we ran a series of simulations (see Table 8 for parameters). Two versions were run and varied parametrically. In the *multiple-meanings* models, there were 30 objects, and some number of auditory units referred to two of them, whereas the remainder referred to one. This was varied in increments of five from 30 unique words (30 words each referring to a single object—the equivalent of the prior models) to 0 unique words (15 words, each referring to two of 30 objects). In the *multiple-labels* models, there were 30 words, and some number of visual units had two names, varying from 0 objects with two labels (the equivalent of the prior model) to all 15 objects having two labels.

It was not clear how to test this network with our analysis of the weight matrix, so we conducted a 10AFC task. During testing, foils were restricted such that a word would only have one of its referents present on a test trial (this restriction was not present during training). In addition to this, five auditory and visual units were not trained and used to test M.E. referent selection, as in Simulation 2.2. Parameters (see Table 8) were similar to those of prior simulations, and novel objects appeared (during training) at 17.5% of the 50% referential ambiguity rate (which yielded more realistic M.E. referent selection in the prior simulations).

**Results.** Figures 19A and 19B show performance on the trained words over the last 25,000 epochs of training for both models. This is shown as a function of the number of words with one-to-one mappings, and separately for words that had a unique mapping, and those that did not. In the multiple-meanings model (Figure 19A), all the models learned both types of words well. Though words with one referent were learned slightly better than words mapping to two objects, both types showed accuracy above 95% across all the simulations. Similarly, in the multiple-labels model (Figure 19B), we also see a benefit for learning objects that only have one word, and learning is somewhat lower when no objects have a single label. However, again, performance is excellent, with accuracy in the worst condition at 94.6%.

Figures 19C and 19D show mutual exclusivity performance. Chance (33%) is indicated by the dashed line; the black curve shows performance at the end of training; and the gray curve shows performance early in training. At the end of training, the multiple-meanings model shows no problem in mutual exclusivity—even the model that completely violated the one-word/one-object mapping (with no words with one referent), performed at 97.9%. This implies that “understanding” this systematicity is not a prerequisite for mutual exclusivity in this model. This is underscored by our analysis of the model's mutual exclusivity performance early in training (the gray curve). Here the model with no words referring to one object actually performed better than models with more unique mappings. This was somewhat surprising, and awaits empirical testing, as there are few analyses of the number of polysemous words for which children are exposed to both meanings. However, it powerfully underscores the fact that in this system mutual exclusivity behavior need not rely on a systematic one-word/one-object bias in the input. Rather, as we described in Simulation 2.2., what is required is that the learning rule preserves some pathway through the weights to get from the novel visual to the novel auditory units. Having more than one referent for each word does not disrupt this (since weight decay relies on exposure to the objects and word individually), thus preserving the ability to use mutual exclusivity.



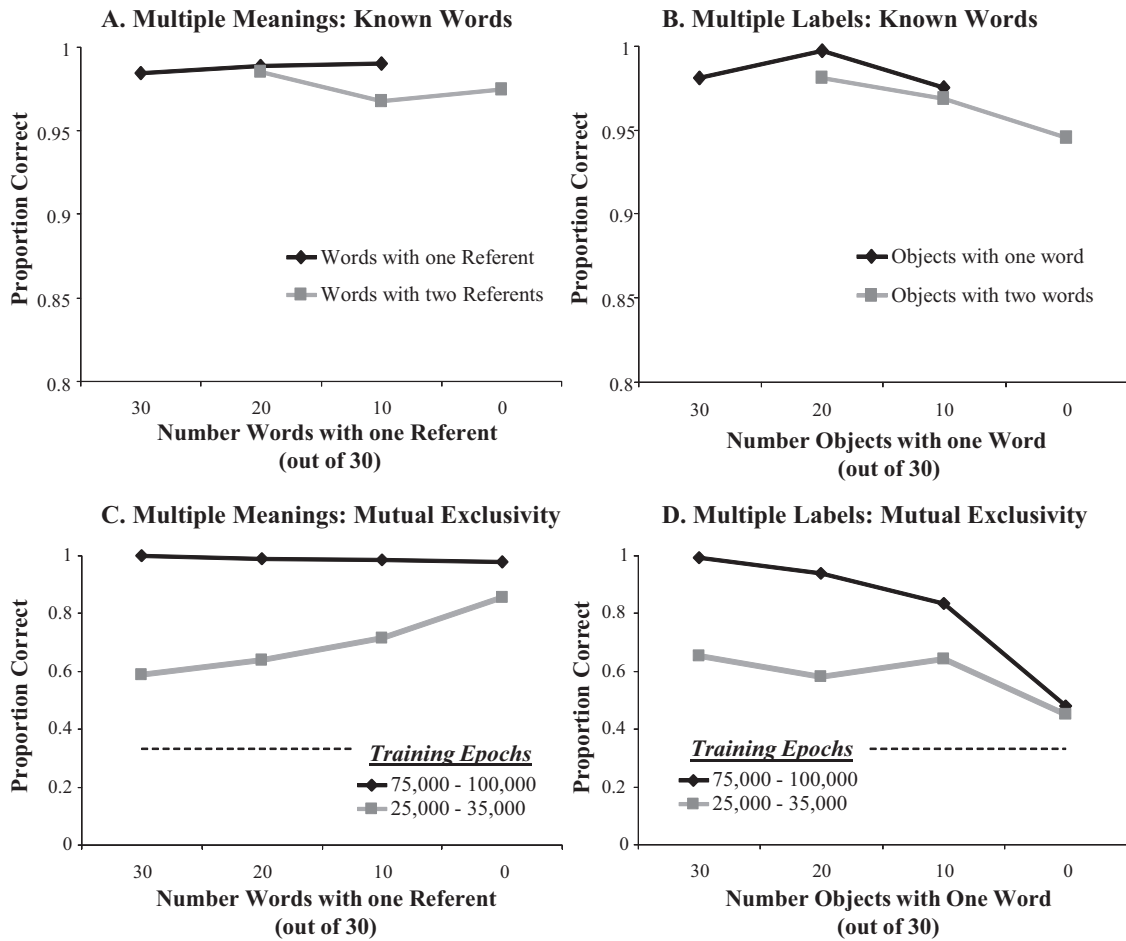


Figure 19. (A) Learning performance of the multiple-meanings model on a 10AFC task during the last 25,000 training trials as a function of the number of words with a single referent, and the type of word. (B) Final learning performance in the multiple-labels model as a function of the number of objects with one word, and the type of object. (C) Performance on a 3AFC mutual exclusivity task as a function of number of one-referent words in the multiple-meanings model. Chance is represented by the dashed line. (D) Mutual exclusivity referent selection performance for the multiple-labels models. AFC = alternative forced-choice.

Intriguingly, however, the multiple-labels model showed a very different pattern of performance. Here we saw that even at the end of training, mutual exclusivity performance dropped, as there were fewer objects with only a single label. Indeed, in the extreme (with no words with one label) the model barely performed above chance.

**Discussion.** These simulations show that first and foremost, disrupting the one-word/one-object mapping does not necessarily disrupt overall learning performance. Though words that had two referents (the multiple-meanings model) and objects that had two words (the multiple-labels model) were learned slightly less well than those with a one-to-one mapping, this performance decrement was negligible. Thus, there is no reason to assume that these principles cannot be scaled up to these more realistic situations, or to situations such as bilingual word learning in which most objects will have multiple names.

Second, mutual exclusivity, as a situation-time process, was largely spared in these models. When words could have two

meanings, mutual exclusivity was fine at every level tested, and appeared to develop faster with more polysemous words. There has been little work looking at the number of words with multiple meanings that children are exposed to or know, and it is unclear where this benefit derives from in the model. However, this counterintuitive prediction may be a useful hallmark of this model. Mutual exclusivity was also largely preserved when models had multiple words for each object (the multiple-labels model). Though performance was degraded with more of these words, it was still above chance (until there were no uniquely mapped words). Thus, situations such as the basic and superordinate situation described above should not pose a problem for referent selection by mutual exclusivity.

Third, in spite of this overall performance, the decrement seen in the multiple-labels model is important. Byers-Heinlein and Werker (2009) suggested that young bilinguals and trilinguals show decrements in using mutual exclusivity. As multilingualism represents the sort of extreme version of a multiple-labels system, our model

offers an explanation for these findings. However, the contrast with the multiple-meanings model, which did not show this decrement, is revealing. It suggests that in an associative system, it is not a strict one-word/one-meaning structure that is necessary. Words may have any number of referents, as long as many objects are largely named by one label. This suggests an important new dimension to principles such as mutual exclusivity and the novel name–nameless category that have been used to describe this behavior.

More importantly, however, these simulations make the broader point that in a way, this ostensibly associative system has derived a principle from across its training experiences, a principle that it can harness in the moment to make decisions about novel words. However, it does so without any real capacity to represent principles such as this. Rather, it seems to have set up its associative weights in such a way that this capacity emerges, in the moment, during real-time competition between words and objects.

### General Discussion

There are two levels on which to evaluate this (or any) model. First, we can consider the range of empirical findings it captures. Second, we can evaluate the theoretical advance made by the model. Does the broader theory tell us something new about word learning?

With respect to the first issue, Table 1 presents a summary of our findings. To briefly summarize, we found that this model could learn words under conditions of very high referential ambiguity. It shows that differences between production and comprehension arise in part due to the fundamental differences in these tasks. It can model the improvement in familiar word recognition over time, including the structure of individual differences, without any need for semantic elaboration, or for bootstrapping type processes. It shows hallmarks of the power law of learning while simultaneously showing accelerating vocabulary growth. It can model referent selection via mutual exclusivity, its development, the lack of retention observed by Horst and Samuelson (2008), and differences in multilinguals (Byers-Heinlein & Werker, 2009). It can also learn multiple names for categories and multiple meanings for words, and it shows a basic-level category advantage.

This model can clearly capture substantial data from very diverse domains of word learning. However, such findings only matter to the extent that they shape our theories of word learning and our intuitions about the nature of the problem. This discussion will focus on these issues. We start by describing some critical limitations of our approach, and ask whether there is an (even) simpler alternative. We then boil down the most important theoretical contributions of this range of work and end by discussing predictions and new directions.

### Limitations

Despite the range of phenomena this dynamic associative model can account for, it is not intended as a complete model of word learning and has a number of limitations. Here we distinguish limitations of the theoretical approach from those of the simple model we developed to illustrate it. There is no reason that similar mechanisms could not be implemented in more complex models to capture an even broader range of behaviors.

First, our model uses localist representations for words and objects, treating each word or category as maximally different from the others. This was deliberately chosen for a number of reasons. Localist representations are easy to interpret and theoretically fairly transparent. More importantly, our twin theoretical claims are best instantiated in this framework. Inhibition between representations is straightforward and easy to implement as a dynamic process in localist representations, and there are no clear ways for doing this with overlapping distributed representations. Similarly, associative learning has long been modeled as linking discrete, localist units, and thus this architecture best captured our approach to learning as well.

Moving to a distributed representation for inputs (e.g., feature vectors for objects) could make it difficult to solve the problem of referential ambiguity. In the current localist scheme, to indicate that multiple objects are present, we simply activate each of their units. However, if objects were represented by a distributed representation across multiple feature vectors, it would be difficult to code more than a few objects—as the feature vectors add up, gradually most, if not all, of the feature units would be active. As a result, the distributed representation for multiple objects simply becomes a concatenation of all features present in a scene. This makes it very difficult to cope with high levels of referential ambiguity. The problem is that in this simple architecture there is no way to bind the features, to know that activity in some features goes together (apples are red and round, blocks are brown and square). This binding problem is a classic issue with distributed representations (Hinton, McClelland, & Rumelhart, 1986), and work must be done to determine whether any of the solutions proposed for this can be integrated with this framework. Perhaps some form of localist (or near localist representations, as seen in topographic maps) is necessary at some level of the system? Additional information in the input (beyond the presence or absence of a feature) may also be helpful, and approaches using spatial location of features to bind them together seem a promising direction (Johnson, Spencer, & Schoner, 2009; Samuelson, Smith, Perry, & Spencer, 2011). In the meantime, however, our goal was to examine the power of competition and associative learning, and localism presented a clear platform in which to do so without solving these historically difficult problems.

Second, our choice of representations completely ignores a fundamental issue—that word learning must link *categories* of words and/or objects. We have attempted to bypass this by assuming that the auditory and visual inputs really represent the output of some other categorization process, but we did not model these processes. However, Hebbian normalized recurrence has already been applied to problems such as speech categorization (McMurray, Horst, et al., 2009; McMurray & Spivey, 2000), and non-learning versions have been used to model visual categorization (Spivey, 2007; Spivey & Dale, 2004). Thus, it may be possible simply to chain together such models. Indeed, work in progress (McMurray et al., in press) has already developed a network that first categorizes visual feature vectors and then uses these categories as the inputs to the network described here. Such an extension may allow the model to capture more of the interesting social, attentional, and conceptual processes that guide children to the right referent in real time, and our preliminary work on this (McMurray et al., in press) suggests it can model complex phenomena by which words sometimes impede visual categorization

(Robinson & Sloutsky, 2004) and sometimes facilitate it (Lupyan et al., 2007).

Even if we allow for such chained models, using localist representations of words and objects appears to make the erroneous assumption that similarity among objects or word forms is irrelevant to word learning. This obviously cannot be so. Fortunately, such similarity relationships can be implemented in localist scheme. Minimally, one would expect that similar categories should be partially coactive due to their overlapping inputs. For a real child, when *bug* is heard, similar sounding words such as *bus* will be partially active (Swingley & Aslin, 2002). Similarly, when a bug is identified in the visual scene, similar categories such as ant or spider may be also active. All these competing, partially active categories could significantly raise the number of spurious associations that would have to be considered and eliminated. However, it is not clear whether this would ultimately be problematic: The set of similar-sounding words (*deer* → *gear*, *deal*, *tear*) is distinct from the set of similar-looking objects (*deer* → *horse*, *cow*, *antelope*), so cross-situational statistics may quickly rule out these associations. Conversely, competition dynamics at the level of visual categories could eliminate some of these competitors. Thus, the problem of coactive categories or ambiguous inputs may not be hugely problematic, though it remains to be investigated.

Given all these issues with localist inputs, the self-organizing map approach (Li et al., 2004; Mayor & Plunkett, 2010) may be a natural bridge between localist representations that ignore similarity and distributed representations that naturally reflect it. Such networks capture similarity relations well and are based on similar Hebbian learning rules to our network. However, they also have enough topography and competition that input representations are precise. The simple form of one-step competition they use could likely be modified to be more dynamic, but they may need additional cues (like space) to cope with multiple inputs. Thus, using self-organizing maps as the input to our model, or using our settling dynamics in such approaches, may offer a useful hybrid.

A third limitation is the scale of our simulations. Most of our simulations used only 35 words, which is small by comparison to the real problem. Simulations reported in the supplemental materials (Simulation S3) show that the network can learn 150 words with few modifications. We have not yet tested larger lexica than that, but there is no reason in principle why this would fail. Moreover, larger lexica may create more optimal statistics for learning (Sibley et al., 2008). With a thousand words, the chance of any given competitor appearing is miniscule, so more invariant (correct) associations may pop out quicker. Of course, the larger number of erroneous connections to suppress may also slow learning.

Fourth, our focus on concrete nouns is a limitation, but not problematic. The localist visual units could easily be treated as tags for properties of objects, allowing the model to learn adjectives. Moreover, if the child or model can segment events from the scene (Reynolds, Zacks, & Braver, 2007), visual units could serve as tags for events or actions allowing the model to learn verbs. Siskind (1996) has shown how cross-situational statistics can be used to acquire word meanings from text, potentially enabling this mechanism to be applied to abstract nouns or verbs as well; and Scott and Fisher (2012) have shown that cross-situational statistics could be involved in verb learning as well. More broadly, the fact that lexical representations are situated between multiple layers of

input could allow other sources of information (e.g., conceptual) to interact with existing auditory and visual inputs to guide learning—something that would be difficult to accomplish if auditory units were associated directly to visual ones.

Finally, our model illustrates the properties of learning necessary to give rise to behaviors such as referent selection via mutual exclusivity (the pattern of weight decay). However, our learning rule is one instantiation of a variety of Hebbian rules, and there may be other versions that are capable of learning more quickly or with fewer lexical units. One could also explore the possibility that supervised (error-driven) learning plays a role. Work in language acquisition more generally suggests that children do receive some feedback from caregivers ranging from quite overt error signals to more subtle cues such as the way in which sentences are repeated back, with or without modification (Bohannon & Stanowicz, 1988; Chouinard & Clark, 2003). More directly, in word learning children are corrected in various ways for naming mistakes or incorrect referent selection (Chouinard & Clark, 2003; Gruendel, 1977; see also unpublished evidence cited in Chapman, Leonard, & Mervis, 1986), and there is evidence that children benefit from such feedback in word learning (Chapman et al., 1986; O'Hanlon & Roberson, 2007). Even beyond this, connectionist models suggest that simple prediction error (e.g., hearing a word, predicting which objects are likely to be present, and learning on the basis of the discrepancy) can be an extremely powerful way to use supervised, error-driven learning in an essentially unsupervised context (Elman, 1990). Error-driven learning (even on a handful of learning events) may buttress some of the slowness of purely associative processes. However, error-driven learning is likely to have very different consequences for the conditions under which irrelevant associations are suppressed, which will have ramifications throughout the system. We also will have to explore when, during processing, the error signal is available, and this could have ramifications for learning. Such effects may ultimately be quite diagnostic, allowing us to identify when and where unsupervised and supervised learning contribute to word learning.

At their core these limitations are largely limitations of the simple model we used to explore our broader dynamic associative account. There are clearly more sophisticated competition algorithms, better input representations, and richer approaches to learning that could be incorporated for a more realistic model. However, what is startling is how much of the word learning literature we could to capture by stripping many of these factors out. The delay in productive vocabulary acquisition can be partially accounted for by the nature of the task, without any recourse to articulation, phonology, or perception; relationships between processing time and learning can be accounted for without resource limitations or bootstrapping; so-called slow-mapping effects can be observed without the need for semantic elaboration; a general principle (mutual exclusivity) can emerge in a purely associative framework; and basic-level categorization advantages can emerge with no hierarchical semantics. Though these explanations are only part of the story for these phenomena, we would have missed them in a more complex model. Thus, the broader theory concerning the linkages of real-time competition and associative learning may have much explanatory power by itself.

## Can We Get Any Simpler?

Given this, we might ask if this model could get any simpler. Could this range of processes derive from even more basic principles? There is some impetus to think about word learning in this way; for example, the [McMurray \(2007\)](#) model of the vocabulary spurt modeled learning as simply accumulating points and discovered that acceleration falls out of parallel learning.

In terms of information processing (e.g., Marr's first level of description), perhaps the core of our model is cross-situational learning. We modeled this via Hebbian associations, but one can think of this in even simpler terms as simply co-occurrence counts between words and objects. [Yu and Smith \(2012\)](#) modeled adult cross-situational learning with exactly such a model (their "bare bones" dumb associative model) and found that versions of it performed quite similarly to a hypothesis-testing model. [Medina, Snedeker, Trueswell, and Gleitman \(2011\)](#) have argued that co-occurrence counts are not consistent with the fact that people reap more benefit from low-ambiguity statistics early in training, since by the end of training the accumulated statistics are the same (though we have modeled such effects using our dynamic associative network; [McMurray et al., in press](#)). But given this discussion, it is worth asking whether our approach offers anything substantive over and above a simpler co-occurrence counter.

To some extent, many of our effects can derive from statistics alone. [Yu and Smith \(2012\)](#) have argued that under some circumstances Hebbian learning can compute a co-occurrence matrix quite directly, and our model's ability to learn cross-situationally derives from this. It is likely that our findings of a basic-level advantage (stronger word-object associations for basic-level terms that have one referent than superordinate terms that have multiple) also derive from this. However, interestingly, Yu and Smith's associative model only reached about 40% correct in their most difficult condition (18 words, four presented per trial), which corresponds to a referential ambiguity rate of 17.6%. Our model is at ceiling under these circumstances (though with more training). Of course, our model also includes real-time processing (which consequently enables slower learning), which may enable better learning.

Other findings do not directly fall out of co-occurrence statistics, but would require real-time processing. Differences in word learning or referent selection based on the task (number of alternatives) could not be accounted for with co-occurrence alone. However, they could if co-occurrence statistics were used as the input to some kind of read-out rule (as in [Yu & Smith, 2012](#)). Our finding of decreasing RTs over training would also require some sort of decision rule that converts co-occurrence information into RTs (maybe something like [Ratcliff & Rouder, 1998](#); or [Usher & McClelland, 2001](#)). However, these decision rules would need to be built or tuned by hand to account for things such as exponential decay in RTs, whereas in our model this is an emergent property of our core theoretical principles, competition and learning. Moreover, at this point, the model would start to look quite similar to ours, and it is not clear what would be gained. Perhaps most importantly, these decision rules would ultimately be just a read-out of learning and would not interact substantively with it—that is, the real-time decisions made by the model during training would have no effect on learning. This could be problematic for modeling of [Fernald et al.'s \(2006\)](#) longitudinal work. We have

shown in multiple simulations here how important this linking between real-time and developmental processes can be, and in our other work ([McMurray et al., in press](#)) suggest it may make all the difference to modeling results such as those by [Medina et al. \(2011\)](#). If we are going to have to add some sort of dynamic decision process, we may as well be modeling the interaction of this with learning.

Our work on mutual exclusivity is perhaps the most difficult to account for on the basis of co-occurrence. A novel name will have no co-occurrences data for either the familiar or novel objects in the scene, and thus no principled way to show a bias. As a result, any pure co-occurrence counter would need some decision rule to decide when to go with the novel object. Again, this would have to be built in, to account for these results, rather than emerge out of the same process that recognizes novel words. But how much evidence counts as no evidence? If the model had seen an object once or twice, how would it make this decision? And what about retention? Once the model has a single piece of evidence for a word-object pairing, the most optimal thing is to return to that on retention trials, and yet children apparently do not. And how would this decision rule develop?

Finally, our results on mutual exclusivity when words have more than one referent (or referents have more than one word) suggest that minimally the co-occurrence statistics are not symmetrical: Having two words for an object can hurt mutual exclusivity, whereas the converse does not. This lends credence to our use of internal (lexical) representations to mediate these co-occurrences. But more importantly, the fact that the multilingual models do not show mutual exclusivity suggests that the development of this behavior is sensitive to the word/object statistics in general, and suggests that our associative system can learn a principle (or will fail to, if the statistics demand). It is not clear how this could emerge out of pure co-occurrence counts.

The bottom line: In order to do mutual exclusivity and capture the range of effects we have using only co-occurrence counts, one would have to build it into the model, and build it in such a way so that it can explicitly capture these effects. This treats behaviors such as mutual exclusivity as fundamentally different from familiar word recognition—they use different decision rules—whereas our model shows how both of them emerge from the same competition scheme. This would lead one to a very different theoretical conclusion than what we have shown here.

Thus, our model is doing something more than just co-occurrence statistics. It is the unique interaction of learning and processing, embedded in an environment with such statistics, that enables complex behavior to emerge from these mechanisms. Indeed, reductionistically simplifying the model further would require us to put substantial content in the situation-time decision rules, and it might not be able to model all these effects anyway. More importantly, it may lead to very different theoretical conclusions.

## Theoretical Insights

Although our model was able to capture numerous empirical findings, its strength lies in its ability to highlight new theoretical conclusions about word learning, conclusions that are substantially broader than the rather narrow model we have presented.



**Learning can (and should) be slow.** Most accounts of word learning stress its effortlessness and speed: Children appear to acquire words very rapidly. This is based on phenomena such as so-called fast mapping and the vocabulary explosion. Such learning is seen as difficult for associative accounts, and thus associative accounts of word learning often stress their rapidity (Mayor & Plunkett, 2010; Regier, 2005).

But is word learning really that fast? Children hear approximately 17,000 words a day (Hart & Risley, 1995). By 1 year, when the first word is produced, an average child will have heard 6 million words. Even at the height of the vocabulary explosion, children may show evidence of having acquired 15–20 words per week. So slow learning may be required, just to account for the actual acquisition curves.

One might argue that within these thousands of words the child is hearing, there may be few tokens of the individual words that a child appears to learn in a given week, so this still necessitates fast learning. However, it is important to point out that under associative accounts like this one, much of the problem is simply suppressing competing associations, something that is less dependent on the specific words being heard. So thousands of exposures to *duck* and *goose* may help children improve their ability to process and acquire *chicken*. Despite apparently quick learning in the laboratory, real learning may be fairly slow, but it also may benefit from much broader experience than we normally count—experience with other words, and with objects alone, is all relevant for learning a particular word.

This contrasts with more inferential or hypothesis-testing approaches that consider hypothetical word–meaning in parallel and wait for the right data to update (Medina et al., 2011). In a sense, due to random initial connections, associative systems start by considering all hypotheses (with some variation in strength). A true hypothesis-testing system could never rule them all out, working in serial. However, by suppressing many connections in parallel at each naming instance, this becomes more feasible. In a sense, at each word, there are global, albeit small, changes in the hypothesis space. This is made even more feasible by smarter situation-time processes that allow the child to behave in the moment on the basis of partial knowledge (Simulations 2.2, 2.3, 3.2, 4.3). Considering familiar words in the same framework as novel words only underscores this: Developmental changes in familiar word recognition take years to unfold (e.g., Fernald et al., 1998, 2006; Zangl et al., 2005), even as the child is generally able to select the correct referent fairly early in life. Thus, slow but global changes on each naming instance may yield a fairly robust system that accounts for multiple developmental phenomena.

And why it should be otherwise? The situations that children find themselves in are inherently ambiguous—there are multiple visual referents, and multiple possible interpretations, for any given word. Even much vaunted social cues are not consistently available and do not consistently disambiguate words (Frank et al., 2009). Thus, cross-situational statistics may constitute a good portion of the information available to link words to objects. If this is the case, then slow learning may be more optimal in that it prevents children from committing too strongly to a single (perhaps erroneous) mapping before they have enough data. Indeed, we investigated this in simulations reported in the supplemental materials (Simulation S4) and found quite poor learning when the learning rate of the model was too high.

Our simulations on mutual exclusivity (Simulations 3.1, 3.2, and 4.2) underscore this. Simulation 3.1 captures Horst and Samuelson (2008), showing that the model can use mutual exclusivity for referent selection, but it retains very little from this. The dependence of mutual exclusivity on learning (Simulation 2.1), task configuration (Simulation 3.2), and lexical statistics (Simulation 4.2) argues that this may be an unstable platform for learning, and the fact that multilingual children and models (Simulation 4.2) still learn words suggests it is not required for learning. Indeed, if referent selection by mutual exclusivity was uniquely powerful for learning, one must ask how often children know the name of every item in the visual scene but one. This seems a fairly unlikely event, underscoring the importance of slower, more gradual mechanisms.

But if learning is slow, this raises the question of how children function while they wait for data to accumulate. This is not just theoretical: Slow learning conflicts with the excellent performance we see in many constrained laboratory tasks, and that parents observe every day with their toddlers. Our dynamic associative account suggests that fast, situation-time processes enable children to take advantage of constraints offered by the environment and children's own incomplete mappings to perform impressively in day-to-day and laboratory tasks, even while learning is slow. This was made clear when we compared the model's performance on constrained tasks to its underlying weights: The model's knowledge was incomplete, but it still performed well on NAFC tasks. It also appeared when we compared comprehension to production: Comprehension was necessarily constrained by the response options and was consistently better than less constrained production.

**Even associative learning is multifaceted.** Our approach to associative learning is more complex than commonly considered. This has theoretical consequences beyond our model. First, since auditory and visual units are independently associated with the lexicon, learning on the auditory and visual side can have different effects. This was most apparent in Simulation 4.2, where having two words per object degraded mutual exclusivity performance but having two objects per word did not. Even in a simple associative framework, the relevant statistics (e.g., the fact that children encounter many visual competitors and fewer auditory ones) shape situations in which both auditory and visual units may perform different roles.

This was also seen in Simulations 2.3 and 3.1, which demonstrated how familiarity with the visual objects alone can improve both referent selection by mutual exclusivity and retention (Kucker & Samuelson, 2012), even in the absence of learning any specific content (the objects' properties). This may be an important consideration in future work on the M.E. referent selection task. Similarly, Horst et al. (2011) recently showed that in a referent selection task with completely novel objects pure visual familiarity with some of the objects can bias children's performance.<sup>8</sup> However, in scaling these ideas to the real world, it is clear that novelty may be a major factor in early lexical behavior. Familiarity is graded, and there are few, if any, words that a child has literally no experience with. Thus, our artificial segregation of objects as familiar or novel (both in models and in typical laboratory tasks)

<sup>8</sup> In our model, this novelty bias can derive from something as basic as the random initial connections (that persist for novel objects), again pointing to the importance of pruning (or not).

may not capture the situation in the real world. Crucially, even an associative learning framework, this shows that not all learning needs to involved both ends of the associative link.

This dynamic associative model also shows the surprising importance of suppressing unnecessary associations. This turned out to be the biggest predictor of settling time for familiar word recognition (Simulation 2.1), and the pattern of weight decay was essential for referent selection by mutual exclusivity (Simulation 2.2). This is because the bulk of learning consists of simply suppressing the vast number of irrelevant connections, and this can be done for virtually any naming event. This raises the possibility that even when children select the wrong referent, or do not select a referent at all, they may still be doing useful learning. This nonobvious source of learning has not been considered in prior theoretical, empirical, and computational work.

Finally, our work on referent selection by mutual exclusivity suggests that associations, even among localist inputs, can derive a principle that applies to even novel items. The ability to use mutual exclusivity is not built into this model and clearly develops from the input. More importantly, it can be blocked when there are two words for many objects. In a companion piece (McMurray et al., *in press*), we have conducted an extensive parameter search of the model's ability to use mutual exclusivity and discovered that only models that have real-time competition dynamics and that use internal representations (rather than directly linking words to referents) can do this. Thus, when associative learning is embedded in a more realistic system with both real-time processing and abstract representations, much richer, emergent behavior can arise.

**Learning and processing are quasi-independent.** Our model was built on the theoretical commitment that using words and learning them are different. Learning is accomplished by changing the connections between words, objects, and the lexicon; processing occurs when real-time competition allows activation to flow over those weights to arrive at a solution. This quasi-independence is fundamental. It allows the model to show dissociations between mutual exclusivity and retention (Simulation 3.1) and between performance and knowledge (Simulations 1.1 and 1.2). Moreover, by considering processing independently, we also showed how both novel and familiar word processing can be handled by the same system (Simulations 2.1–2.3).

The quasi-independence has useful functional consequences. By off-loading mutual exclusivity to online processing, it no longer blocks learning of multiple names for a given object (Simulations 4.1 and 4.2). More broadly, learning can be less than perfect because processing can get the child the rest of the way. This is what compensates for slow learning. The best system will be one that uses processes optimized for learning to handle developmental-time learning and processes optimized for in-the-moment demands to handle real-time behavior.

Nonetheless, though computationally these are distinct processes, they are not completely independent in practice: Situation-time processes are dependent on learning. The changes in RT for familiar words derive from learning, and the ability of the model to use mutual exclusivity mapping derives from a weight matrix created by the specific learning rule. Many people describe fast mapping as the sort of initial stages of a slow learning process (e.g., Capone & McGregor, 2005; Carey & Bartlett, 1978; Golinkoff et al., 1992), a sort of incomplete learning. In contrast, we suggest that mutual exclusivity behavior (referent selection) is a

real-time product of the type of learning that has occurred up to that point, but this behavior in turn leaves an associative trace that can build over repetitions to yield word learning. In this light, the learning on the first exposure of a word (what has been termed fast mapping) is no different from subsequent exposures (slow mapping).

Simulation 3.4 showed the converse, suggesting that learning is impossible without processing. The necessity of such competitive processes is implicit in many unsupervised learning models, but its importance has not been highlighted before. Competitive learning (Rumelhart & Zipser, 1986), for example, requires winner-take-all learning; self-organizing maps (Kohonen, 1982; Mayor & Plunkett, 2010) include a competition–interaction kernel; and even the quite unrelated mixture of Gaussians framework for clustering benefits from competition (McMurray, Aslin, et al., 2009). However, though competition is essential, its outcome can be variable. In simulations not reported here, we have found that useful learning still occurs even when competition gets the wrong referent. Thus, the presence of competition is necessary for learning—the system must make a choice. But on any trial, the specific choice is less important.

Finally, Fernald et al.'s (2006) longitudinal study is perhaps the best evidence for the dependence of learning and processing, as they found that children's RTs predicted acceleration in learning. Our model also showed a similar pattern of results, but suggested no simple construct to explain it. Learning rate and settling time were the product of parameters that control learning and processing; the same parameter (e.g., temperature) acts differently depending on the referential ambiguity; and the biggest predictor of RT was the nature of the learned weights. Thus, to understand a functional property of the child such as speed of processing, we must understand the myriad of components of both processing and learning.

**Word learning need not be specialized.** A number of models suggest that acquiring vocabulary may harness general-purpose learning mechanisms. Regier (2005) and Mayor and Plunkett (2010) showed the unexpected power of association learning. Xu and Tenenbaum (2007) and Frank et al. (2009) used general Bayesian inference mechanisms. McMurray (2007) suggested that acceleration is not the hallmark of a specialized system but should be seen virtually everywhere. There is also empirical work showing that similar principles may span word learning and other types of learning, whether they derive from general reasoning strategies that apply to both facts and words (Behrend et al., 2001; P. Bloom & Markson, 1998; Markson & Bloom, 1997; but see Waxman & Booth, 2000) or lower level attention (Samuelson & Smith, 2000) and novelty biases (Horst et al., 2011). Thus, there is a mounting effort to explain vocabulary acquisition in terms of general cognitive and learning processes.

Our dynamic associative account adds to this. First, it shows that novel word inference and familiar word recognition can arise from the same system. Both use the same set of online processes (which themselves are quite general); both operate over the same mappings (weights) that are shaped by the same learning mechanism. There is no need for any sort of monitoring mechanism to route novel words to a more constrained or specialized learning mechanism. Moreover, referent selection by mutual exclusivity can be modeled with something as general as dynamic competition, further emphasizing the generality of these processes.

Second, the patterns of learning observed in vocabulary are not special. At face value, the decelerating learning predicted by the ubiquitous power law of learning conflicts with acceleration during the vocabulary explosion. Our model highlights this conflict, showing both acceleration and deceleration, and both derive from changes in the associative weights (Simulation 2.1). However, the particular changes in the weights that yield gains in RT are not equivalent to those that allow the model to appear as if it acquired a new word: Whereas RT changes largely derive from suppressing unnecessary connections, changes in the number of words known require both the suppression of unnecessary connections and the establishment of the right positive associations. Thus, depending on the measurement (changes in RT or vocabulary size), and the real-time processing that give rise to the behavior, we may reach different conclusions about the shape of learning (accelerative or decelerative), even as the underlying mechanism is the same. This further cements word learning as a general process, but challenges learning theory by suggesting that changes in RT may not fully describe learning.

Finally, our use of associative learning does not entail a particular source of information: Our framework is consistent with information sources such as attentional or social cues, and conceptual structures that have been widely interpreted as special. These exist outside the core lexical mapping system and constrain the settling dynamics (from the outside), or simply determine the type of representations that are associated. Thus, such higher order factors may be fundamental to learning and/or lexical behavior, without needing to be embedded in the learning system. This permits a soft coupling: As children gain sensitivity to things such as speakers' intentions (Moore, 2008), these sources of information gradually play a larger role in shaping online behavior (and through it, learning) without fundamentally restructuring word learning.

## New Directions

The test of any model is its ability to make predictions and highlight new research questions. Given the simplicity of the model, it is not clear that we are in a position to make precise empirical predictions for new tests and paradigms. Nonetheless, the model and the broader theoretical view suggest a number of important new areas of investigation in word learning.

Indeed, during the course of writing this article, a number of predictions from the model were tested empirically. Our finding that M.E. referent selection relies on the relative randomness of the weights connecting the auditory and visual units, and the role of simple (purely visual) experience suppressing them (Simulation 2.1), led us to predict, and confirm, that familiarity with objects may bias children away from selecting them in referent selection tasks (Horst et al., 2011); and our demonstration in Simulation 3.1 that visual familiarity can simultaneously influence retention motivated Kucker and Samuelson (2012), which showed similar results in children. Taken together, these simulations suggest that mutual exclusivity may have two components: (a) a component driven by novelty that leads to excellent referent selection but more retention (b) and something that resembles a constraint satisfaction, which leads to somewhat worse (but still good) referent selection but much better retention. This trade-off should be explored, particularly as these components wax and wane over development.

Similarly, although the simulations in 4.2 match evidence that bilingual children (who have multiple names for many objects) may perform worse in mutual exclusivity tasks (Byers-Heinlein & Werker, 2009), they also suggest the converse—having multiple meanings for many words—may not be problematic. There has been little work on the statistics of word–object mappings in the child's environment and how they relate to behaviors such as mutual exclusivity (analogous to the Perry & Samuelson, 2011, and Samuelson & Smith, 1999, studies of how such statistics predict the shape and material biases), and it is not entirely clear why our model shows this asymmetry, but this is a clear avenue for future work.

More broadly, our work points to a host of issues that the statistics of word–object mappings (co-occurrence statistics) may be involved in. Simulation 4.1 demonstrated how this can give rise to an advantage of basic-level over superordinate category labels. An important part of this is that superordinate terms have their associations spread across multiple objects, whereas basic-level terms are only associated with one. Conversely, we saw in our simulations of familiar word recognition how suppressing irrelevant connections was crucial to improving performance, and that this could only occur when competitors were present (but variable from trial to trial). In this case, the spreading of association across competitors prevents any of them from becoming strongly linked with the target word. This spreading of associations has been invoked in a number of other domains including early speech perception, where spreading associations block children from associating irrelevant talker cues with words (Apfelbaum & McMurray, 2011), and in semantic memory, where “context variability” prevents aspects of the context from serving as a retrieval cue (Steyvers & Malmberg, 2003). More broadly, however, these simulations force us to think more creatively about the co-occurrence statistics of words and objects—our distillation of the problem suggests that they may play a role in numerous domains.

A third avenue of future study is the role of suppressing irrelevant connections. This is pivotal in predicting changes in processing speed and giving rise to referent selection by mutual exclusivity. However, this nonobvious result of learning has not received extensive study. More complex eye movement paradigms may allow us to index the strength of competing associations, to look for correlations with these behaviors, and artificial language paradigms may allow us to manipulate it by temporarily creating strong spurious associations. Indeed, these may ultimately be better predictors of behaviors such as mutual exclusivity and word recognition time. Fitneva and Christiansen (2011) recently demonstrated that in a cross-situational word learning paradigm, participants who looked longer to the incorrect objects during learning showed better performance. This clearly is consistent with the idea that suppressing competing associations is critical for learning and points to a paradigm in which to investigate these issues.

Finally, and most importantly, our approach suggests that even in an associative account, what the child does and the exact sequence of events will matter. For example, the timing of the events in a learning trial could influence whether or not there is sufficient time for competition among referents to resolve, and this would alter learning dramatically. Similarly, the configuration of items on a learning trial (e.g., the number of competitors) and the behavior of the child can both affect learning, by influencing how positive associations are formed (which require a fairly specific



confluence of events) and how negative associations are suppressed (which may be more general and not require a correct response).

## Conclusions

Lexical behavior must fundamentally be considered on two timescales—children learn words over development, but they must also use them here and now. Word learning is not about acquisition of words as a type of knowledge; rather, we must study how children acquire the abilities to recognize and produce words, and infer the meanings of novel words. By embedding learning within a structure of word use, our model offers a unified account for a range of findings in word learning, word recognition, and novel word inference. In this framework, word learning is the simple product of ongoing interactions between developmental-time processes such as associative learning and situation-time processes such as dynamic competition.

## References

- Adolph, K. E., Robinson, S. R., Young, J. W., & Gill-Alvarez, F. (2008). What is the shape of developmental change? *Psychological Review*, 115, 527–543. doi:10.1037/0033-295X.115.3.527
- Akhtar, N., & Martinez-Sussman, C. (2007). Intentional communication. In C. Brownell & C. Kopp (Eds.), *Socioemotional development in the toddler years: Transitions and transformations* (pp. 201–220). New York, NY: Guilford Press.
- Alloppenna, P., Magnuson, J. S., & Tanenhaus, M. K. (1998). Tracking the time course of spoken word recognition using eye movements: Evidence for continuous mapping models. *Journal of Memory and Language*, 38, 419–439. doi:10.1006/jmla.1997.2558
- Anderson, J. R. (1982). Acquisition of cognitive skill. *Psychological Review*, 89, 369–406. doi:10.1037/0033-295X.89.4.369
- Apfelbaum, K. S., & McMurray, B. (2011). Using variability to guide dimensional weighting: Associative mechanisms in early word learning. *Cognitive Science*, 35, 1105–1138. doi:10.1111/j.1551-6709.2011.01181.x
- Baldwin, D. A. (1991). Infants' contribution to the achievement of joint reference. *Child Development*, 62, 875–890. doi:10.2307/1131140
- Baldwin, D. A., Markman, E. M., Bill, B., Desjardins, R. N., & Irwin, J. M. (1996). Infants' reliance on a social criterion for establishing word-object relations. *Child Development*, 67, 3135–3153. doi:10.2307/1131771
- Bauer, P., Dow, G., & Hertsgaard, L. (1995). Effects of prototypicality on categorization in 1- to 2-year-olds: Getting down to basic. *Cognitive Development*, 10, 43–68. doi:10.1016/0885-2014(95)90018-7
- Behl-Chadha, G. (1996). Basic-level and superordinate-like categorical representations in early infancy. *Cognition*, 60, 105–141. doi:10.1016/0010-0277(96)00706-8
- Behrend, D. A., Scofield, J., & Kleinknecht, E. E. (2001). Beyond fast mapping: Young children's extensions of novel words and novel facts. *Developmental Psychology*, 37, 698–705. doi:10.1037/0012-1649.37.5.698
- Bharucha, J. J. (1987). Music cognition and perceptual facilitation: A connectionist framework. *Music Perception*, 5, 1–30.
- Bloom, L. (1973). *One word at a time: The use of single-word utterances before syntax*. The Hague, the Netherlands: Mouton.
- Bloom, P. (2000). *How children learn the meanings of words*. Cambridge, MA: The MIT Press.
- Bloom, P., & Markson, L. (1998). Capacities underlying word learning. *Trends in Cognitive Sciences*, 2, 67–73. doi:10.1016/S1364-6613(98)01121-8
- Bohannon, J. N., & Stanowicz, L. B. (1988). The issue of negative evidence: Adult responses to children's language errors. *Developmental Psychology*, 24, 684–689. doi:10.1037/0012-1649.24.5.684
- Bourne, L. E., & Restle, F. (1959). Mathematical theory of concept identification. *Psychological Review*, 66, 278–296. doi:10.1037/h0041365
- Bush, R. R., & Mosteller, F. (1951). A model for stimulus generalization and discrimination. *Psychological Review*, 58, 413–423. doi:10.1037/h0054576
- Byers-Heinlein, K., & Werker, J. F. (2009). Monolingual, bilingual, trilingual: Infants' language experience influences the development of a word-learning heuristic. *Developmental Science*, 12, 815–823. doi:10.1111/j.1467-7687.2009.00902.x
- Capone, N. C., & McGregor, K. K. (2005). The effect of semantic representation on toddlers' word retrieval. *Journal of Speech, Language, and Hearing Research*, 48, 1468–1480. doi:10.1044/1092-4388(2005)102
- Carey, S. (1978). The child as word learner. In M. Halle, J. Bresnan, & G. A. Miller (Eds.), *Linguistic theory and psychological reality* (pp. 264–293). Cambridge, MA: MIT Press.
- Carey, S., & Bartlett, E. (1978). Acquiring a single new word. In *Papers and reports on child language development* (Vol. 15, pp. 17–29). Stanford, CA: Stanford University, Department of Linguistics.
- Chambers, C. G., Tanenhaus, M. K., & Magnuson, J. S. (2004). Action-based affordances and syntactic ambiguity resolution. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 30, 687–696. doi:10.1037/0278-7393.30.3.687
- Chapman, K. L., Leonard, L. B., & Mervis, C. B. (1986). The effect of feedback on young children's inappropriate word usage. *Journal of Child Language*, 13, 101–117. doi:10.1017/S0305000900000325
- Chouinard, M. M., & Clark, E. V. (2003). Adult reformulations of child errors as negative evidence. *Journal of Child Language*, 30, 637–669. doi:10.1017/S0305000903005701
- Colunga, E., & Smith, L. B. (2005). From the lexicon to expectations about kinds: A role for associative learning. *Psychological Review*, 112, 347–382. doi:10.1037/0033-295X.112.2.347
- Dale, P. S., & Fenson, L. (1996). Lexical development norms for young children. *Behavior Research Methods, Instruments, & Computers*, 28, 125–127. doi:10.3758/BF03203646
- Dell, G. S. (1986). A spreading activation theory of retrieval in sentence production. *Psychological Review*, 93, 283–321. doi:10.1037/0033-295X.93.3.283
- Diesendruck, G., & Markson, L. (2001). Children's avoidance of lexical overlap: A pragmatic account. *Developmental Psychology*, 37, 630–641. doi:10.1037/0012-1649.37.5.630
- Dore, J., Franklin, M. B., Miller, R. T., & Ramer, A. L. (1976). Transitional phenomena in early language acquisition. *Journal of Child Language*, 3, 13–28. doi:10.1017/S0305000900001288
- Elman, J. L. (1990). Finding structure in time. *Cognitive Science*, 14, 179–211. doi:10.1207/s15516709cog1402\_1
- Elman, J. L. (2009). On the meaning of words and dinosaur bones: Lexical knowledge without a lexicon. *Cognitive Science*, 33, 547–582. doi:10.1111/j.1551-6709.2009.01023.x
- Elman, J. L., Bates, E., Johnson, M., Karmiloff-Smith, A., Parisi, D., & Plunkett, K. (1996). *Rethinking innateness: A connectionist perspective on development*. Cambridge, MA: MIT Press.
- Fenson, L., Dale, P. S., Reznick, J. S., Bates, E., Thal, D. J., & Pethick, S. J. (1994). Variability in early communicative development. *Monographs of the Society for Research in Child Development*, 59(5, Serial No. 242). doi:10.2307/1166093
- Fernald, A., Perfors, A., & Marchman, V. A. (2006). Picking up speed in understanding: Speech processing efficiency and vocabulary growth across the 2nd year. *Developmental Psychology*, 42, 98–116. doi:10.1037/0012-1649.42.1.98



- Fernald, A., Pinto, J. P., Swingle, D., Weinberg, A., & McRoberts, G. W. (1998). Rapid gains in speed of verbal processing by infants in the 2nd year. *Psychological Science*, 9, 228–231. doi:10.1111/1467-9280.00044
- Fernald, A., Swingle, D., & Pinto, J. P. (2001). When half a word is enough: Infants can recognize spoken words using partial phonetic information. *Child Development*, 72, 1003–1015. doi:10.1111/1467-8624.00331
- Fitneva, S. A., & Christiansen, M. H. (2011). Looking in the wrong direction correlates with more accurate word learning. *Cognitive Science*, 35, 367–380. doi:10.1111/j.1551-6709.2010.01156.x
- Frank, M. C., Goodman, N. D., & Tenenbaum, J. (2009). Using speakers' referential intentions to model early cross-situational word learning. *Psychological Science*, 20, 578–585. doi:10.1111/j.1467-9280.2009.02335.x
- Fulkerson, A. L., & Waxman, S. R. (2007). Words (but not tones) facilitate object categorization: Evidence from 6- and 12-month-olds. *Cognition*, 105, 218–228. doi:10.1016/j.cognition.2006.09.005
- Ganger, J., & Brent, M. R. (2004). Reexamining the vocabulary spurt. *Developmental Psychology*, 40, 621–632. doi:10.1037/0012-1649.40.4.621
- Goldinger, S. D. (1998). Echoes of echos? An episodic theory of lexical access. *Psychological Review*, 105, 251–279. doi:10.1037/0033-295X.105.2.251
- Goldstone, R. L., & Medin, D. L. (1994). The time course of comparison. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 20, 29–50. doi:10.1037/0278-7393.20.1.29
- Golinkoff, R. M., & Hirsh-Pasek, K. (2006). Baby wordsmith: From associationist to social sophisticate. *Current Directions in Psychological Science*, 15, 30–33. doi:10.1111/j.0963-7214.2006.00401.x
- Golinkoff, R. M., Hirsh-Pasek, K., Bailey, L. M., & Wenger, N. R. (1992). Young children and adults use lexical principles to learn new nouns. *Developmental Psychology*, 28, 99–108. doi:10.1037/0012-1649.28.1.99
- Golinkoff, R. M., Hirsh-Pasek, K., Bloom, L., Smith, L. B., Woodward, A. L., Akhtar, N., . . . Hollich, G. (2000). *Becoming a word learner: A debate on lexical acquisition*. New York, NY: Oxford University Press. doi:10.1093/acprof:oso/9780195130324.001.0001
- Golinkoff, R. M., Mervis, C. B., & Hirsh-Pasek, K. (1994). Early object labels: The case for a developmental lexical principles framework. *Journal of Child Language*, 21, 125–155. doi:10.1017/S0305000900008692
- Goodman, J. C., McDonough, L., & Brown, N. B. (1998). The role of semantic context and memory in the acquisition of novel nouns. *Child Development*, 69, 1330–1344. doi:10.2307/1132269
- Grassmann, S., & Tomasello, M. (2010). Young children follow pointing over words in interpreting acts of reference. *Developmental Science*, 13, 252–263. doi:10.1111/j.1467-7687.2009.00871.x
- Grossberg, S. (1976). Adaptive pattern classification and universal recoding: I. Parallel development and coding of neural feature detectors. *Biological Cybernetics*, 23, 121–134. doi:10.1007/BF00344744
- Gruendel, J. M. (1977). Referential extension in early language development. *Child Development*, 48, 1567–1576. doi:10.2307/1128520
- Gupta, P., & Tisdale, J. (2009). Does phonological short-term memory causally determine vocabulary learning? Toward a computational resolution of the debate. *Journal of Memory and Language*, 61, 481–502. doi:10.1016/j.jml.2009.08.001
- Halberda, J. (2003). The development of a word-learning strategy. *Cognition*, 87, B23–B34. doi:10.1016/S0010-0277(02)00186-5
- Halberda, J. (2006). Is this a dax which I see before me? Use of the logical argument disjunctive syllogism supports word-learning in children and adults. *Cognitive Psychology*, 53, 310–344. doi:10.1016/j.cogpsych.2006.04.003
- Hanna, J. E., & Tanenhaus, M. K. (2004). Pragmatic effects on reference resolution in a collaborative task: Evidence from eye movements. *Cognitive Science*, 28, 105–115. doi:10.1207/s15516709cog2801\_5
- Harm, M. W., & Seidenberg, M. S. (1999). Phonology, reading acquisition, and dyslexia: Insights from connectionist models. *Psychological Review*, 106, 491–528. doi:10.1037/0033-295X.106.3.491
- Hart, B., & Risley, T. (1995). *Meaningful differences in the everyday experience of young American children*. Baltimore, MD: Brookes.
- Heathcote, A., Brown, S., & Mewhort, D. J. K. (2000). The power law repealed: The case for an exponential law of practice. *Psychonomic Bulletin & Review*, 7, 185–207. doi:10.3758/BF03212979
- Hinton, G. E., McClelland, J. L., & Rumelhart, D. E. (1986). Distributed representations. In D. E. Rumelhart, J. L. McClelland, & PDP Research Group (Eds.), *Parallel distributed processing: Explorations in the microstructure of cognition: Vol. 1. Foundations* (pp. 77–109). Cambridge, MA: MIT Press.
- Horst, J. S., McMurray, B., & Samuelson, L. K. (2006). Online processing is essential for learning: Understanding fast mapping and word learning in a dynamic connectionist architecture. In R. Sun (Ed.), *Proceedings of the 28th Annual Conference of the Cognitive Science Society* (pp. 339–343). Austin, TX: Cognitive Science Society.
- Horst, J. S., & Samuelson, L. (2008). Fast mapping but poor retention in 24-month-old infants. *Infancy*, 13, 128–157. doi:10.1080/15250000701795598
- Horst, J. S., Samuelson, L. K., Kucker, S., & McMurray, B. (2011). What's new? Children prefer novelty in referent selection. *Cognition*, 118, 234–244. doi:10.1016/j.cognition.2010.10.015
- Horst, J. S., Scott, E. J., & Pollard, J. A. (2010). The role of competition in word learning via referent selection. *Developmental Science*, 13, 706–713. doi:10.1111/j.1467-7687.2009.00926.x
- Hurtado, N., Marchman, V. A., & Fernald, A. (2007). Spoken word recognition by Latino children learning Spanish as their first language. *Journal of Child Language*, 34, 227–249. doi:10.1017/S0305000906007896
- Huttenlocher, J. (1974). The origins of language comprehension. In R. Solso (Ed.), *Theories in cognitive psychology: The Loyola Symposium*. Oxford, England: Erlbaum.
- Jaswal, V. K., & Hansen, M. B. (2006). Learning words: Children disregard some pragmatic information that conflicts with mutual exclusivity. *Developmental Science*, 9, 158–165. doi:10.1111/j.1467-7687.2006.00475.x
- Johnson, J. S., Spencer, J. P., & Schoner, G. (2009). A layered neural architecture for the consolidation, maintenance, and updating of representations in visual working memory. *Brain Research*, 1299, 17–32. doi:10.1016/j.brainres.2009.07.008
- Jones, M., & Love, B. C. (2011). Bayesian fundamentalism or enlightenment? On the explanatory status and theoretical contributions of Bayesian models of cognition. *Behavioral and Brain Sciences*, 34, 169–188. doi:10.1017/S0140525X10003134
- Kail, R. (1994). A method for studying the generalized slowing hypothesis in children with specific language impairment. *Journal of Speech Language & Hearing Research*, 37, 418–421.
- Kohonen, T. (1982). Self-organized formation of topologically correct feature maps. *Biological Cybernetics*, 43, 59–69. doi:10.1007/BF00337288
- Kucker, S. C., & Samuelson, L. K. (2012). The first slow step: Differential effects of object and word-form familiarization on retention of fast-mapped words. *Infancy*, 17, 295–323. doi:10.1111/j.1532-7078.2011.00081.x
- Li, P., Farkas, I., & MacWhinney, B. (2004). Early lexical development in a self-organizing neural network. *Neural Networks*, 17, 1345–1362. doi:10.1016/j.neunet.2004.07.004
- Livesey, E., & McLaren, I. (2011). An elemental model of associative

- learning and memory. In E. M. Pothos & A. J. Wills (Eds.), *Formal approaches in categorization* (pp. 153–172). Cambridge, England: Cambridge University Press. doi:10.1017/CBO9780511921322.007
- Logan, G. D. (1992). Shapes of reaction-time distributions and shapes of learning curves: A test of the instance theory of automaticity. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 18, 883–914. doi:10.1037/0278-7393.18.5.883
- Lupyan, G., Rakison, D. H., & McClelland, J. L. (2007). Language is not just for talking: Redundant labels facilitate learning of novel categories. *Psychological Science*, 18, 1077–1083. doi:10.1111/j.1467-9280.2007.02028.x
- MacDonald, M. C., Pearlmutter, N. J., & Seidenberg, M. S. (1994). Lexical nature of syntactic ambiguity resolution. *Psychological Review*, 101, 676–703. doi:10.1037/0033-295X.101.4.676
- MacWhinney, B. (1987). The competition model. In B. MacWhinney (Ed.), *Mechanisms of language acquisition* (pp. 249–308). Hillsdale, NJ: Erlbaum.
- Marchman, V. A., & Fernald, A. (2008). Speed of word recognition and vocabulary knowledge in infancy predict cognitive and language outcomes in later childhood. *Developmental Science*, 11, F9–F16. doi:10.1111/j.1467-7687.2008.00671.x
- Mareschal, D., & Tan, S. H. (2007). Flexible and context-dependent categorization by eighteen-month-olds. *Child Development*, 78, 19–37. doi:10.1111/j.1467-8624.2007.00983.x
- Markman, E. M. (1990). Constraints children place on word meanings. *Cognitive Science*, 14, 57–77. doi:10.1207/s15516709cog1401\_4
- Markman, E. M., & Wachtel, G. (1988). Children's use of mutual exclusivity to constrain the meaning of words. *Cognitive Psychology*, 20, 121–157. doi:10.1016/0010-0285(88)90017-5
- Markman, E. M., Wasow, J. L., & Hanson, M. B. (2003). Use of the mutual exclusivity assumption by young word learners. *Cognitive Psychology*, 47, 241–275. doi:10.1016/S0010-0285(03)00034-3
- Markson, L., & Bloom, P. (1997). Evidence against a dedicated system for word learning in children. *Nature*, 385, 813–815. doi:10.1038/385813a0
- Marslen-Wilson, W. D. (1987). Functional parallelism in spoken word-recognition. *Cognition*, 25, 71–102. doi:10.1016/0010-0277(87)90005-9
- Mayor, J., & Plunkett, K. (2010). A neurocomputational account of taxonomic responding and fast mapping in early word learning. *Psychological Review*, 117, 1–31. doi:10.1037/a0018130
- McCabe, A., & Peterson, C. (Eds.). (1991). *Developing narrative structure*. Hillsdale, NJ: Erlbaum.
- McClelland, J. L., & Elman, J. L. (1986). The TRACE model of speech perception. *Cognitive Psychology*, 18, 1–86. doi:10.1016/0010-0285(86)90015-0
- McClelland, J. L., Mirman, D., & Holt, L. L. (2006). Are there interactive processes in speech perception? *Trends in Cognitive Sciences*, 10, 363–369. doi:10.1016/j.tics.2006.06.007
- McMurray, B. (2007, August 3). Defusing the childhood vocabulary explosion. *Science*, 317–631. doi:10.1126/science.1144073
- McMurray, B., Aslin, R. N., & Toscano, J. C. (2009). Statistical learning of phonetic categories: Insights from a computational approach. *Developmental Science*, 12, 369–378. doi:10.1111/j.1467-7687.2009.00822.x
- McMurray, B., Horst, J. S., Toscano, J. C., & Samuelson, L. K. (2009). Integrating connectionist learning and dynamical systems processing: Case studies in speech and lexical development. In J. P. Spencer, M. S. C. Thomas, & J. L. McClelland (Eds.), *Toward a unified theory of development: Connectionism and dynamic systems theory re-considered* (pp. 218–249). London, England: Oxford University Press. doi:10.1093/acprof:oso/9780195300598.003.0011
- McMurray, B., Samuelson, V. S., Lee, S. H., & Tomblin, J. B. (2010). Eye-movements reveal the time-course of online spoken word recognition language impaired and normal adolescents. *Cognitive Psychology*, 60, 1–39. doi:10.1016/j.cogpsych.2009.06.003
- McMurray, B., & Spivey, M. J. (2000). The categorical perception of consonants: The interaction of learning and processing. *Proceedings of the Chicago Linguistics Society*, 34, 205–220.
- McMurray, B., Zhao, L., Kucker, S., & Samuelson, L. K. (in press). Probing the limits of associative learning: Generalization and the statistics of words and referents. In L. Gogate & G. Hollich (Eds.), *Theoretical and computational models of word learning: Trends in psychology and artificial intelligence*. Hershey, PA: IGI Global.
- McRae, K., Spivey-Knowlton, M. J., & Tanenhaus, M. K. (1998). Modeling the influence of thematic fit (and other constraints) in on-line sentence comprehension. *Journal of Memory and Language*, 38, 283–312. doi:10.1006/jmla.1997.2543
- Medina, T. N., Snedeker, J., Trueswell, J. C., & Gleitman, L. R. (2011). How words can and cannot be learned by observation. *Proceedings of the National Academy of Sciences of the United States of America*, 108, 9014–9019. doi:10.1073/pnas.1105040108
- Merriman, W. (1999). Competition, attention, and young children's lexical processing. In B. MacWhinney (Ed.), *The emergence of language* (pp. 331–358). Mahwah, NJ: Erlbaum.
- Merriman, W. E., Lipko, A., & Evey, J. (2008). How young children judge whether a word is one they know: A dual criterion account. *Journal of Experimental Child Psychology*, 101, 83–98. doi:10.1016/j.jecp.2008.06.001
- Mervis, C. B., & Bertrand, J. (1994). Acquisition of the novel nameless category (N3C) principle. *Child Development*, 65, 1646–1662. doi:10.2307/1131285
- Mitchell, C. C., & McMurray, B. (2008). A stochastic model of the vocabulary explosion. In B. C. Love, K. McRae, & V. M. Sloutsky (Eds.), *Proceedings of the 30th Annual Conference of the Cognitive Science Society* (pp. 1919–1926). Austin, TX: Cognitive Science Society.
- Mitchell, C., & McMurray, B. (2009). On leveraged learning in lexical acquisition and its relationship to acceleration. *Cognitive Science*, 33, 1503–1523. doi:10.1111/j.1551-6709.2009.01071.x
- Moore, C. (2006). Representing intentional relations and acting intentionally in infancy: Current insights and open questions. In G. Knoblich, I. M. Thornton, M. Grosjean, & M. Shiffrar (Eds.), *Human body perception from the inside out: Advances in visual cognition* (pp. 427–442). New York, NY: Oxford University Press.
- Moore, C. (2008). The development of gaze following. *Child Development Perspectives*, 2, 66–70. doi:10.1111/j.1750-8606.2008.00052.x
- Munakata, Y. (1998). Infant perseveration and implications for object permanence theories: A PDP model of the AB task. *Developmental Science*, 1, 161–184. doi:10.1111/1467-7687.00021
- Munakata, Y., & McClelland, J. L. (2003). Connectionist models of development. *Developmental Science*, 6, 413–429. doi:10.1111/1467-7687.00296
- Munakata, Y., McClelland, J. L., Johnson, M., & Siegler, R. (1997). Rethinking infant knowledge: Toward an adaptive process account of successes and failures in object permanence tasks. *Psychological Review*, 104, 686–713. doi:10.1037/0033-295X.104.4.686
- Namy, L. (2012). Getting specific: Early general mechanisms give rise to domain-specific expertise in word learning. *Language Learning and Development*, 8, 47–60. doi:10.1080/15475441.2011.617235
- Nazzi, T., & Bertoncini, J. (2003). Before and after the vocabulary spurt: Two modes of word acquisition? *Developmental Science*, 6, 136–142. doi:10.1111/1467-7687.00263
- Newell, A., & Rosenbloom, P. S. (1981). Mechanisms of skill acquisition and law of practice. In J. R. Anderson (Ed.), *Cognitive skills and their acquisition* (pp. 1–55). Hillsdale, NJ: Erlbaum.
- O'Hanlon, C. G., & Roberson, D. (2007). What constrains children's learning of novel shape terms? *Journal of Experimental Child Psychology*, 97, 138–148. doi:10.1016/j.jecp.2006.12.002

- Palmeri, T. J. (1997). Exemplar similarity and the development of automaticity. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 23, 324–354. doi:10.1037/0278-7393.23.2.324
- Perry, L. K., & Samuelson, L. K. (2011). The shape of the vocabulary predicts the shape of the bias. *Frontiers in Psychology*, 2, 345. doi:10.3389/fpsyg.2011.00345
- Quine, W. V. O. (1960). *Word and object: An inquiry into the linguistic mechanisms of objective reference*. Cambridge, MA: MIT Press.
- Ratcliff, R., & Rouder, J. N. (1998). Modeling response times for two-choice decisions. *Psychological Science*, 9, 347–356. doi:10.1111/1467-9280.00067
- Reese, E., & Fivush, R. (1993). Parental styles of talking about the past. *Developmental Psychology*, 29, 596–606. doi:10.1037/0012-1649.29.3.596
- Regier, T. (1996). *The human semantic potential: Spatial language and constrained connectionism*. Cambridge, MA: MIT Press.
- Regier, T. (2003). Emergent constraints on word-learning: A computational review. *Trends in Cognitive Sciences*, 7, 263–268. doi:10.1016/S1364-6613(03)00108-6
- Regier, T. (2005). The emergence of words: Attentional learning in form and meaning. *Cognitive Science*, 29, 819–865. doi:10.1207/s15516709cog0000\_31
- Reynolds, J. R., Zacks, J. M., & Braver, T. S. (2007). A computational model of event segmentation from perceptual prediction. *Cognitive Science*, 31, 613–643. doi:10.1080/15326900701399913
- Reznick, J. S., & Goldfield, B. A. (1992). Rapid change in lexical development in comprehension and production. *Developmental Psychology*, 28, 406–413. doi:10.1037/0012-1649.28.3.406
- Riches, N. G., Tomasello, M., & Conti-Ramsden, G. (2005). Verb learning in children with SLI: Frequency and spacing effects. *Journal of Speech, Language, and Hearing Research*, 48, 1397–1411. doi:10.1044/1092-4388(2005/097)
- Robinson, C. W., & Sloutsky, V. M. (2004). Auditory dominance and its change in the course of development. *Child Development*, 75, 1387–1401. doi:10.1111/j.1467-8624.2004.00747.x
- Robinson, C. W., & Sloutsky, V. M. (2007). Visual processing speed: Effects of auditory input on visual processing. *Developmental Science*, 10, 734–740. doi:10.1111/j.1467-7687.2007.00627.x
- Rosch, E., Mervis, C. B., Gray, W. D., Johnson, D. M., & Boyes-Braem, P. (1976). Basic objects in natural categories. *Cognitive Psychology*, 8, 382–439. doi:10.1016/0010-0285(76)90013-X
- Rost, G. C., & McMurray, B. (2010). Finding the signal by adding noise: The role of noncontrastive phonetic variability in early word learning. *Infancy*, 15, 608–635. doi:10.1111/j.1532-7078.2010.00033.x
- Roy, D., & Pentland, A. (2002). Learning words from sights and sounds: A computational model. *Cognitive Science*, 26, 113–146. doi:10.1207/s15516709cog2601\_4
- Rumelhart, D., & Zipser, D. (1986). Feature discovery by competitive learning. In D. Rumelhart & J. L. McClelland (Eds.), *Parallel distributed processing: Explorations in the microstructure of cognition* (pp. 151–193). Cambridge, MA: MIT Press.
- Samuelson, L. K. (2002). Statistical regularities in vocabulary guide language acquisition in connectionist models and 15–20-month-olds. *Developmental Psychology*, 38, 1016–1037. doi:10.1037/0012-1649.38.6.1016
- Samuelson, L. K., & Smith, L. B. (1998). Memory and attention make smart word learning: An alternative account of Akhtar, Carpenter, and Tomasello. *Child Development*, 69, 94–104.
- Samuelson, L. K., & Smith, L. B. (1999). Early noun vocabularies: Do ontology, category structure and syntax correspond? *Cognition*, 73, 1–33. doi:10.1016/S0010-0277(99)00034-7
- Samuelson, L. K., & Smith, L. B. (2000). Grounding development in cognitive processes. *Child Development*, 71, 98–106. doi:10.1111/1467-8624.00123
- Samuelson, L. K., Smith, L. B., Perry, L. K., & Spencer, J. P. (2011). Grounding word learning in space. *PLoS ONE*, 6(12), e28095. doi:10.1371/journal.pone.0028095
- Sandhofer, C. M., Smith, L. B., & Luo, J. (2000). Counting nouns and verbs in the input: Differential frequencies, different kinds of learning? *Journal of Child Language*, 27, 561–585. doi:10.1017/S0305000900004256
- Schlesinger, M., & McMurray, B. (in press). Modeling matters: What computational models have taught us about cognitive development. *Journal of Cognition and Development*.
- Schutte, A. R., Spencer, J. P., & Schöner, G. (2003). Testing the dynamic field theory: Working memory for locations becomes more spatially precise over development. *Child Development*, 74, 1393–1417. doi:10.1111/1467-8624.00614
- Scott, R. M., & Fisher, C. (2012). 2.5-year-olds use cross-situational consistency to learn verbs under referential uncertainty. *Cognition*, 122, 163–180. doi:10.1016/j.cognition.2011.10.010
- Sénéchal, M., & Cornell, E. H. (1993). Vocabulary acquisition through shared reading experiences. *Reading Research Quarterly*, 28, 360–374. doi:10.2307/747933
- Shanks, D. R. (2007). Associationism and cognition: Human contingency learning at 25. *Quarterly Journal of Experimental Psychology*, 60, 291–309. doi:10.1080/17470210601000581
- Sibley, D. E., Kello, C. T., David, P. C., & Elman, J. L. (2008). Large-scale modeling of wordform learning and representation. *Cognitive Science*, 32, 741–754. doi:10.1080/03640210802066964
- Siskind, J. M. (1996). A computational study of cross-situational techniques for learning word-to-meaning mappings. *Cognition*, 61, 39–91. doi:10.1016/S0010-0277(96)00728-7
- Smith, L. B. (2000). Avoiding associations when it's behaviorism you really hate. In R. M. Golinkoff, K. Hirsh-Pasek, L. Bloom, L. B. Smith, A. L. Woodward, N. Akhtar, . . . G. Hollich (Eds.), *Becoming a word learner: A debate on lexical acquisition* (pp. 169–174). New York, NY: Oxford University Press.
- Smith, L. B., Thelen, E., Titzer, R., & McLin, D. (1999). Knowing in the context of acting: The task dynamics of the A-not-B error. *Psychological Review*, 106, 235–260. doi:10.1037/0033-295X.106.2.235
- Smith, L. B., & Yu, C. (2008). Infants rapidly learn word-referent mappings via cross-situational statistics. *Cognition*, 106, 1558–1568. doi:10.1016/j.cognition.2007.06.010
- Spencer, J. P., Perone, S., & Johnson, J. S. (2009). Dynamic field theory and embodied cognitive dynamics. In J. P. Spencer, M. S. C. Thomas, & J. L. McClelland (Eds.), *Toward a unified theory of development: Connectionism and dynamic systems theory re-considered* (pp. 86–118). London, England: Oxford University Press. doi:10.1093/acprof:oso/9780195300598.003.0005
- Spiegel, C., & Halberda, J. (2011). Rapid fast-mapping abilities in 2-year-olds. *Journal of Experimental Child Psychology*, 109, 132–140. doi:10.1016/j.jecp.2010.10.013
- Spivey, M. J. (2007). *The continuity of mind*. New York, NY: Oxford University Press.
- Spivey, M. J., & Dale, R. (2004). On the continuity of mind: Toward a dynamical account of cognition. In B. H. Ross (Ed.), *The psychology of learning and motivation: Advances in research and theory* (Vol. 45, pp. 87–142). San Diego, CA: Elsevier. doi:10.1016/S0079-7421(03)45003-2
- Spivey-Knowlton, M. J., & Allopenna, P. (1997, August). *A computational account of the integration of linguistic and visual information during spoken word recognition. Paper presented at the Computational Psycholinguistics Conference*, Berkeley, CA.
- Steyvers, M., & Malmberg, K. J. (2003). The effect of normative context variability on recognition memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 29, 760–766. doi:10.1037/0278-7393.29.5.760



- Storkel, H. L., Armbrüster, J., & Hogan, T. P. (2006). Differentiating phonotactic probability and neighborhood density in adult word learning. *Journal of Speech, Language, and Hearing Research*, 49, 1175–1192. doi:10.1044/1092-4388(2006/085)
- Swingle, D. (2009). Onsets and codas in 1.5-year-olds' word recognition. *Journal of Memory and Language*, 60, 252–269. doi:10.1016/j.jml.2008.11.003
- Swingle, D., & Aslin, R. N. (2000). Spoken word recognition and lexical representations in very young children. *Cognition*, 76, 147–166. doi:10.1016/S0010-0277(00)00081-0
- Swingle, D., & Aslin, R. N. (2002). Lexical neighborhoods and the word-form representations of 14-month-olds. *Psychological Science*, 13, 480–484. doi:10.1111/1467-9280.00485
- Swingle, D., & Aslin, R. N. (2007). Lexical competition in young children's word learning. *Cognitive Psychology*, 54, 99–132. doi:10.1016/j.cogpsych.2006.05.001
- Tanenhaus, M. K., & Brown-Schmidt, S. (2008). Language processing in the natural world. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 363, 1105–1122. doi:10.1098/rstb.2007.2162
- Tanenhaus, M. K., & Trueswell, J. C. (1995). Sentence comprehension. In J. L. Miller & P. D. Eimas (Eds.), *Handbook in Perception and Cognition: Vol. 11. Speech, language, and communication* (pp. 217–262). New York, NY: Academic Press.
- Thelen, E., Schöner, G., Scheier, C., & Smith, L. B. (2001). The dynamics of embodiment: A dynamic field theory of infant perseverative reaching errors. *Behavioral and Brain Sciences*, 24, 1–34. doi:10.1017/S0140525X01003910
- Tomasello, M. (2001). Perceiving intentions and learning words in the second year of life. In M. Tomasello & E. Bates (Eds.), *Language development: The essential readings* (pp. 111–128). Malden, MA: Blackwell.
- Tomasello, M., Strosberg, R., & Akhtar, N. (1996). Eighteen-month-old children learn words in non-ostensive contexts. *Journal of Child Language*, 23, 157–176. doi:10.1017/S0305000900010138
- Usher, M., & McClelland, J. L. (2001). The time course of perceptual choice: The leaky, competing accumulator model. *Psychological Review*, 108, 550–592. doi:10.1037/0033-295X.108.3.550
- van Geert, P. (1991). A dynamic systems model of cognitive and language growth. *Psychological Review*, 98, 3–53. doi:10.1037/0033-295X.98.1.3
- Vecera, S. P., & O'Reilly, R. C. (1998). Figure-ground organization and object recognition processes: An interactive account. *Journal of Experimental Psychology: Human Perception and Performance*, 24, 441–462. doi:10.1037/0096-1523.24.2.441
- Waxman, S. R. (2003). Links between object categorization and naming: Origins and emergence in human infants. In D. H. Rakison & L. M. Oakes (Eds.), *Early category and concept development: Making sense of the blooming, buzzing confusion* (pp. 213–241). New York, NY: Oxford University Press.
- Waxman, S. R., & Booth, A. E. (2000). Principles that are invoked in the acquisition of words, but not facts. *Cognition*, 77, B33–B43. doi:10.1016/S0010-0277(00)00103-7
- Waxman, S. R., & Gelman, S. (2009). Early word-learning entails reference, not merely associations. *Trends in Cognitive Sciences*, 13, 258–263. doi:10.1016/j.tics.2009.03.006
- Wifall, T., McMurray, B., & Hazeltine, R. E. (2012). *Similarity impairs motor learning: Implications for the power law of learning?* Manuscript submitted for publication.
- Wilkinson, K. M., & Mazzitelli, K. (2003). The effect of “missing” information on children's retention of fast-mapped labels. *Journal of Child Language*, 30, 47–73. doi:10.1017/S0305000902005469
- Woodward, A. L., & Markman, E. M. (1998). Early word learning. In W. Damon (Ed.), *Handbook of child psychology: Vol. 2. Cognition, perception, and language* (pp. 371–420). Hoboken, NJ: Wiley.
- Xu, F., & Tenenbaum, J. B. (2007). Word learning as Bayesian inference. *Psychological Review*, 114, 245–272. doi:10.1037/0033-295X.114.2.245
- Yu, C., Ballard, D. H., & Aslin, R. N. (2005). The role of embodied intention in early lexical acquisition. *Cognitive Science*, 29, 961–1005. doi:10.1207/s15516709cog0000\_40
- Yu, C., & Smith, L. B. (2007). Rapid word learning under uncertainty via cross-situational statistics. *Psychological Science*, 18, 414–420. doi:10.1111/j.1467-9280.2007.01915.x
- Yu, C., & Smith, L. B. (2012). Modeling cross-situational word-referent learning: Prior questions. *Psychological Review*, 119, 21–39. doi:10.1037/a0026182
- Zangl, R., Klarman, L., Thal, D., Fernald, A., & Bates, E. (2005). Dynamics of word comprehension in infancy: Developments in timing, accuracy, and resistance to acoustic degradation. *Journal of Cognition and Development*, 6, 179–208. doi:10.1207/s15327647jcd0602\_2

Received April 6, 2011

Revision received July 5, 2012

Accepted July 24, 2012 ■