

The role of experience in disambiguation during early word learning

Molly Lewis<sup>1, 6</sup>, Veronica Cristiano<sup>2</sup>, Brenden Lake<sup>3</sup>, Tammy Kwan<sup>4</sup>, & Michael C. Frank<sup>5</sup>

<sup>1</sup> University of Chicago

<sup>2</sup> Gallaudet University

<sup>3</sup> New York University

<sup>4</sup> Cognitive Toybox, Inc.

<sup>5</sup> Stanford University

<sup>6</sup> University of Wisconsin-Madison

Author Note

Data from Experiment 2 were previously presented in the Proceedings of the Cognitive Science Society Conference in Lewis & Frank (2013). \*To whom correspondence should be addressed. E-mail: [mollylewis@uchicago.edu](mailto:mollylewis@uchicago.edu)

## Abstract

Young children tend to map novel words to novel objects even in the presence of familiar competitors, a finding that has been dubbed the “disambiguation” effect. This phenomenon is important because it could provide a strong constraint for children in learning new words. But, although the effect is highly robust and widely studied, the cognitive mechanisms underlying it remain unclear. Existing theoretical accounts include a proposal for initial constraints on children’s lexicons (e.g. a principle of mutual exclusivity), situation-specific pragmatic inferences, probabilistic accounts, and overhypothesis account. In the current paper, we have two goals: synthesize the existing body of literature and directly examine the causal role of experience on the effect. We present a synthesis of existing evidence through a meta-analysis of the existing literature, followed by two experiments that examine the relationship between vocabulary development and the disambiguation effect. We conclude by summarizing the empirical landscape, and suggest that multiple mechanisms may underlie the effect.

*Keywords:* mutual exclusivity, disambiguation effect, word learning, meta-analysis

Word count: X

The role of experience in disambiguation during early word learning

## Introduction

A key property of language is that each word maps to a unique concept, and each concept maps to a unique word [Clark (1987); bolinger1977meaning]. Like other regularities in language (e.g., grammatical categories), children cannot directly observe this one-to-one word-concept regularity, yet even very young children behave in a way that is consistent with it. The goal of this paper is develop a theory of the cognitive mechanisms underlying this behavior in children.

Evidence that children obey the one-to-one regularity comes from what is known as the “mutual exclusivity” (ME) effect (we return to the issue of nomenclature below). In a typical demonstration of this effect (Markman & Wachtel, 1988), children are presented with a novel and familiar object (e.g., a whisk and a ball), and are asked to identify the referent of a novel word (“Show me the dax”). Children in this task tend to choose the novel object as the referent, consistent with the one-to-one regularity in language (we refer to this paradigm throughout as the “ME paradigm”). A large body of work has demonstrated that this effect occurs in children across a wide range of ages, experimental paradigms, and populations (Bion, Borovsky, & Fernald, 2013; R.M. Golinkoff, Mervis, Hirsh-Pasek, & others, 1994; J. Halberda, 2003; Markman, Wasow, & Hansen, 2003; Mervis, Golinkoff, & Bertrand, 1994).

The ME effect has received much attention in the word learning literature because the ability to identify the meaning of a word in ambiguous contexts is, in essence, the core problem of word learning. That is, given any referential context, the meaning of a word is underdetermined (Quine, 1960), and the challenge for the world learner is to identify the referent of the word within this ambiguous context. Critically, the ability to infer that a novel word maps to a novel object makes the problem much easier to solve. For example, suppose a child hears the novel word “kumquat” while in the produce aisle of the grocery store. There are an infinite number of possible meanings of this word given this referential context, but the child’s ability to correctly disambiguate would lead her to rule out all

meanings for which she already had a name. With this restricted hypothesis space, the child is more likely to identify the correct referent than if all objects in the context were considered as possible referents.

Additionally, the ME effect has the potential to help the learner acquire words for multiple concepts that can be used to refer to the same object in the world, such as property names and object parts (e.g., “turquoise”, “handle”; Markman and Wachtel (1988))). Consider a child who hears the novel word “turquoise” in the context of a turquoise colored ball. If the child obeys the one-to-one property of language and already knows the word “ball,” the child may assume that “turquoise” refers to a property of the ball, such as color, rather than the ball itself. The one-to-one principle may be particularly useful for learning subordinate (e.g., “dalmation”) and superordinate labels (e.g. “animal”), since each instance of these labels is always consistent with concepts at all levels of the conceptual hierarchy (an observed dalmation is equally consistent with the labels “dalmation,” “dog” and “animal”; e.g., (???)). Unlike for property words, the child will never observe cross-situational evidence that would disambiguate among candidate concepts at different levels of the hierarchy. The one-to-one principle provides one possible route through which children might resolve this inherent ambiguity in word learning.

Despite – or perhaps due to – the attention that the ME effect has received, there is little consensus regarding the cognitive mechanisms underlying it. Does it stem from a basic inductive bias on children’s learning abilities (“constraint and bias accounts,” see below), a learned regularity about the structure of language (“overhypothesis accounts”), reasoning about the goals of communication in context (“pragmatic accounts”), or perhaps some mixture of these? In the current paper, we lay out these possibilities and discuss the state of the evidence. Along the way we present a meta-analysis of the extant empirical literature. We then present two new, relatively large-sample developmental experiments that investigate the dependence of children’s ME inferences on vocabulary (Experiment 1) and experience with particular words (Experiment 2). We end by discussing the emergence of ME inferences

in a range of computational models of word learning. We conclude that:

1. Explanations of ME are not themselves mutually exclusive and likely more than one is at play; \*momen and merriman make this point
2. The balance of responsibility for behavior likely changes developmentally, with basic biases playing a greater role for younger children and learned overhypotheses playing a greater role for older children.
3. All existing accounts put too little emphasis on the role of experience and strength of representation; this lack of explicit theory in many cases precludes definitive tests.
4. ME inferences are distinct from learning.

**A note on terminology.** Markman and Wachtel (1988)’s seminal paper coined the term “mutual exclusivity,” which was meant to label the theoretical proposal that “children constrain word meanings by assuming at first that words are mutually exclusive – that each object will have one and only one label.” (Markman, 1990, p. 66). That initial paper also adopted a task used by a variety of previous authors (including RM Golinkoff, Hirsh-Pasek, Baduini, & Lavalley, 1985; Hutchinson, 1986; Vincent-Smith, Bricker, & Bricker, 1974), in which a novel and a familiar object were presented to children in a pair and the child was asked to “show me the  $x$ ,” where  $x$  was a novel label. Since then, informal discussions have used the same name for the paradigm and effect (selecting the novel object as the referent of the novel word) as well as the theoretical account (an early assumption or bias). This conflation of paradigm/effect with theory is problematic, as other authors who have argued against the theoretical account then are in the awkward position of rejecting the name for the paradigm they have used. Other labels (e.g. “disambiguation” or “referent selection” effect) are not ideal, however, because they are not as specific do not refer as closely to the previous literature. Here we adopt the label “mutual exclusivity” (ME) for the general family of paradigms and associated effects, *without* prejudgment of the theoretical account of these effects.

ME has also been referred to as “fast mapping.” This conflation is confusing at best. In an early study, S. Carey and Bartlett (1978) presented children with an incidental word learning scenario by using a novel color term to refer to an object: “You see those two trays over there. Bring me the *chromium* one. Not the red one, the *chromium* one.” Those data (and subsequent replications, e.g. L. Markson & Bloom, 1997) showed that this exposure was enough to establish some representation of the link between phonological form and meaning that endured over an extended period; a subsequent clarification of this theoretical claim emphasized that these initial meanings are partial (S. Carey, 2010). Importantly, however, demonstrations of retention relied on learning in a case where there was a contrastive presentation of the word with a larger set of contrastive cues (S. Carey & Bartlett, 1978) or pre-exposure to the object (L. Markson & Bloom, 1997).

## Theoretical Views on the ME effect

What are the cognitive mechanisms underlying the ME effect? A number of proposals have been made in the literature, many of which overlap or differ only in subtle ways. Here we briefly describe several influential proposals, highlighting the commonalities and differences across theoretical views.

**Constraint and bias accounts.** Under constraint and bias accounts, children are argued to have a constraint or bias that is innate or emerges after very limited language input. One version of the account, proposed by Markman and colleagues (Markman & Wachtel, 1988; Markman et al., 2003), is that children have a constraint on the types of lexicons considered when learning the meaning of a new word – a “mutual exclusivity constraint.” Under this constraint, children are biased to consider only those lexicons that have a one-to-one mapping between words and objects. Importantly, this constraint is probabilistic and thus can be overcome in cases where it is incorrect (e.g. property names), but it nonetheless serves to restrict the set of lexicons initially entertained when learning the meaning of a novel word. In principle, this constraint could be the result of either

domain-specific or domain-general processes (???). As a domain general property, the ME constraint could be related to other cognitive mechanisms that lead learners to prefer one-to-one mappings (e.g. blocking and overshadowing in classical condition and the discounting principle in motivational research, (???)).

As formulated by Markman and colleagues, the ME constraint operates at the level of extensions (objects), not concepts. For example, the ME constraint says that the labels “policeman” and “cop” – referring to the same entity in the world – are violations of the constraint. Similarly, terms at different levels of the semantic hierarchy that can have the same extensions, such as “animal” and “dog,” are also seen as ME violations. In contrast, these cases are not violations in theories that posit the explanatory construct at the level of concepts (e.g., pragmatic accounts). The distinction between concepts and objects in each theoretical view is important for evaluating whether empirical evidence is consistent with a proposal. Note, however, that in the canonical ME paradigm, where the two referents are both different concepts and objects, the accounts at both levels make identical predictions.

A related proposal to the ME constraint is that children have a bias to map novelty to novelty (???; Novel-Name Nameless-Category principle (N3C); R.M. Golinkoff et al., 1994 (???)). This principle differs from the ME constraint in that the rejection of the familiar object as a potential referent is not part of the inference; instead, children are argued only to map the two novel elements to each other, the novel label and the object (thereby only implicitly rejecting the the familiar object as a referent for the novel label). The N3C principle is argued to be domain-specific to language.

Under a third account, children are motivated to identify objects for which they do not know a label for and fill the “lexical gap” with the novel label (???). Lexical Gap Filling: Merriman & Bowman (1989)

**Probabilistic accounts.** Probabilistic accounts contend that the ME effect does not derive from an explicit representation related to the one-to-one regularity, as proposed by the constraints and bias accounts; rather, under these accounts, the effect is the product of a

word learning system that tracks the frequency of exemplars of words and their referents over time, and then reasons probabilistically about the most likely referent for a novel word within the referential context. (???; Fazly, Alishahi, & Stevenson, 2010; M. C. Frank, Goodman, & Tenenbaum, 2009; McMurray, Horst, & Samuelson, 2012; Regier, 2005).

**Pragmatic accounts.** Under pragmatic accounts, the ME effect derives from reasoning about the intention of the speaker within the referential context (???, ???; Clark, 1987; G. Diesendruck & Markson, 2001). The critical aspect of this account is the claim that children assume that “every two forms contrast in meaning” (Clark, 1988, p. 417), or the “Principle of Contrast.” Clark also argues that speakers hold a second assumption – that speakers within the same speech community use the same words to refer to the same objects (“Principle of Conventionality”). The ME effect then emerges from the interaction of these two principles. That is, the child reasons implicitly: You used a word I’ve never heard before. Since, presumably we both call a ball “ball” and if you’d meant the ball you would have said “ball,” this new word must refer to the new object. Clark (1988, 1990) argues that these two principles are learned, but emerge from a more general understanding that other people have intentions (???, (???)).

**Logical inference accounts.** J. Halberda (2003) argues that the ME effect is the result of domain-general processes used for logical reasoning. Under this proposal, children are argued to be solving a disjunctive syllogism (“A or B, not A, therefore B”) by rejecting labels for known objects. For example, upon hearing the novel label “dax,” the child would implicitly reason that the referent could be either object A or B, and then reject object A because it already has a known label. By deduction, the child would then conclude that “dax” refers to object B. This account shares the same formal reasoning structure as pragmatic accounts, but differs in the underlying source of the key inference: While pragmatic accounts argue that children conclude that object B must be the referent on the basis of reasoning about intention, the logical inference account proposes that this same inference is made on the basis of logical reasoning.



**Over-hypothesis accounts.** Lewis and Frank (2013) suggest that the ME effect could emerge by learning from the statistics of the child’s linguistic. That is, given evidence that words tend to refer to a single concept, the child might develop a learned “overhypothesis” (???) that the lexicon is structured such that each concept is associated with one and only one label. The learning mechanisms are argued to be probabilistic and domain general, while the learned overhypothesis is specific to the structure of the lexicon. The emergent overhypothesis about the structure of the lexicon would be similar to the knowledge a learner is proposed to have under the constraints and biases account.

In order for learning to get off the ground, however, children must notice the one-to-one mapping between a word and a concept in the context of a particular instances of a label’s usage. Lewis and Frank (2013) suggest that this ability could derive from a variety of different mechanisms that make use of the structure of the learning task, such as pragmatic, probabilistic, and logical inference accounts. (???) make a similar proposal, but argue that the overhypothesis is learned primarily from explicit parental corrections (e.g., “that’s an apple, not an orange”).

Under the overhypothesis account, then, the ME effect emerges from multiple mechanisms at two different timescales – one as a function of information about the pragmatic or inferential structure of the communicative context and one as a function of learned higher-order knowledge about how the lexicon is structured. Both mechanisms would then contribute to the inference with different weights across development and across children.

## Theory-Constraining Findings

The literature on the ME effect explores predictions of a variety of theoretical proposals. Here, we highlight a few of the key findings that provide important constraints for a theory of the ME effect.

**Developmental change.** A number of studies provide evidence that the strength of the ME effect increases across development (e.g. ???, (???), Bion et al. (2013)). For example, (Justin Halberda, 2003) tests 14- 16- and 17- mo in the ME paradigm, and finds a pattern of developmental change: 14 mo children are biased to select the familiar object, 16-mo were at chance, and 17-mo were biased to select the novel object, consistent with the one-to-one principle. This evidence suggests that the strength of the ME effect changes across development, though the underlying cause of this developmental change is an open question (an issue we return to below).

**Multilingualism.** Children who are learning multiple languages have been tested in the ME paradigm in order to examine the role of linguistic input in the ME effect. Multilingual children are an interesting test population because the one-to-one mapping between words and concepts is arguably violated in their linguistic input (e.g. a child learning Spanish might know both the words “ball” and “pelota” for the concept ball). Thus, if the ME effect is independent of lexical input, then multilingual children should perform on the ME task siimilar to monolingual children whose input does not violate the one-to-one assumption. In contrast to this prediction, Byers-Heinlein (???) and others find that multilingual children select the novel object at lower rates than monolingual children, suggesting that lexical input plays a role in the strength of the ME effect.

**Speaker-change studies.** Some evidence for pragmatic accounts comes from experiments in which children must reason about the intent of the speaker directly. In one set of experiments (Gil Diesendruck, 2005), children were taught a novel label for a creature that is either a common noun or a proper noun. A speaker, who was not present during the teaching phase, then requests a create by either a novel label. If children are reasoning about the knowledge of the speaker, the pragmatic account predicts that the speaker should know the common name for the known creature (as a competent speaker of the language), but not the proper name. Children show a pattern consistent with this prediction by selecting a novel creature in a 2-AFC task when taught the common noun label, but not a proper noun label.

**Autism.** Children with autism are known to have impairments in reasoning about the intentions of others (Baron-Cohen, Leslie, & Frith, 1986). As such, this population has been tested in the ME paradigm to examine the extent to which reasoning about the intentions of other speakers, as required by the pragmatic view, is a necessary component of the ME effect. Evidence suggests (e.g. (de Marchena, Eigsti, Worek, Ono, & Snedeker, 2011; Preissler & Carey, 2005)) that children with autism select the novel object in the ME task at similar rates to typically developing children, suggesting that pragmatic reasoning is unlikely to be a necessary component for the ME effect.

**Fast mapping + no retention.**

**NF vs. NN.** The canonical ME paradigm involves an object with a known label and an object with an unknown label. In this paradigm, evidence that children are biased to select the novel referent when presented with a novel word is consistent both with accounts that argues that children *reject* the familiar object (e.g. Constraint and bias accounts) and with accounts that children are biased to map the novel word to the novel object (N3C principle). To distinguish between these two types of accounts, researchers have compared the canonical ME paradigm that uses a novel and a familiar object (NF design) to a paradigm that uses two novel objects (NN design). In the NN variant, the child is presented with two novel objects but taught a novel label for one of the objects unambiguously (“This is a zot”). Then, the child is asked to identify the referent of a second novel label (“Can you find the fep?”). If the child relies on novelty alone to identify the correct referent, the ME effect should be absent in the NN design since both objects are novel the child. Instead, there is evidence to suggest that children show the ME effect in both NN and NF designs (???; e.g., G. Diesendruck & Markson, 2001), suggesting that a novelty bias is not sufficient to account for the ME effect.

## Synthesis

To summarize, the empirical findings that a successful theory of ME must account for are: - why the effect is present in young children, but gets larger with development (developmental change); - how language experience supports the effect (multilingualism evidence); - why pragmatic reasoning can support the effect (speaker change evidence), but why it is not necessary (autism evidence).

In developing a successful theory, it is important to note that the theoretical accounts of mechanisms underlying the ME effect that have proposed in the literature are not mutually exclusive with each other (???). As pointed out by Markman (???), testing different mechanisms in isolation is the result of an experimental approach to theory building, rather than a reflection of an assumption that there exists one and only one mechanism underlying the effect. That is, in order to identify whether a mechanism is *sufficient* to give rise to the ME effect, logical researchers design experiments in which the ME effect can be observed only if a particular cognitive mechanism is sufficient for the effect. If the effect is observed under these conditions, it provides evidence only that the mechanism is sufficient for the effect, but not that it is necessary and not that other mechanisms are not also sufficient. Indeed, there is reason to think that redundancy in mechanisms for the same behavior is a desirable property of a cognitive system (???).

Instead, in light of the full body of evidence, we argue that multiple mechanisms likely support the ME effect probabilistically. Each child may be making use of multiple mechanisms with varying weights across development and situations, and the relative weights of these different mechanisms may vary across children. For example, learners may be making use of both general knowledge about how the lexicon is structured as well as information about the pragmatic or inferential structure of the task, and both of these sources of information support the ME inference.

## The Current Study

A theory of the ME effect that appeals to multiple cognitive mechanisms is a difficult theory to build. This is because, in order to build such a theory, we must gather empirical evidence that not only describes *that* a mechanism underlies a behavior, but also the degree to which it does and how the contribution of different mechanisms varies across tasks, developmental ages and populations. The goal of the current study is to contribute to building such a theory in two ways. First, we first provide a quantitative synthesis of the current literature related to the ME effect in the form of a meta-analysis. The meta-analysis allows us to gain a clearer picture of the empirical landscape in terms of the magnitude of the effect as well as the role of moderating variables. Second, we present two experiments that examine the causal role of an understudied moderator in the literature – linguistic experience. In Experiment 1, we examine the relationship between vocabulary size and the strength of the ME effect on a large sample of children. We find evidence that children with larger vocabularies tend to show a stronger ME effect, consistent with the notion that language experience influences the ME effect. In Experiment 2, we more directly test the hypothesis that language experience plays a *causal* role in the ME effect, by directly manipulating children’s amount of experience with a word. We find greater experience with the familiar word in the ME paradigm leads to a stronger ME effect. We conclude by re-evaluating a theory of the ME effect in light of our new evidence.

## Meta-analysis

To assess the strength of the disambiguation bias as well a moderating factors, we conducted a meta-analysis on the existing body of literature that investigates the disambiguation effect.

## Methods

**Search strategy.** We conducted a forward search based on citations of Markman and Wachtel (1988) in Google Scholar, and by using the keyword combination “mutual exclusivity” in Google Scholar (September 2013; November 2017).<sup>1</sup> Additional papers were identified through citations and by consulting experts in the field. We then narrowed our sample to the subset of studies that used one of two different paradigms: (a) an experimenter says a novel word in the context of a familiar object and a novel object and the child guesses the intended referent (the canonical paradigm; “Familiar-Novel”), or (b) experimenter first provides the child with an unambiguous mapping of a novel label to a novel object, and then introduces a second novel object and asks the child to identify the referent of a second novel label (“Novel-Novel”). For Familiar-Novel conditions, we included conditions that included more than one familiar object (e.g. Familiar-Familiar-Novel). From these conditions, we restricted our sample to only those that satisfied the following criteria: (a) participants were children (less than 12 years of age)<sup>2</sup>, (b) referents were objects or pictures (not facts or object parts), and (c) no incongruent cues (e.g. eye gaze at familiar object). All papers used either forced-choice pointing or eye-tracking methodology. All papers were peer-reviewed with the exception of two dissertations (Williams, 2009; Frank, I., 1999), but all main results reported below remain the same when these papers are excluded. In total, we identified 43 papers that satisfied our selection criteria and had sufficient information to calculate an effect size.

**Coding.** For each paper, we coded separately each relevant condition with each age group entered as a separate condition. For each condition, we coded the paper metadata (citation) as well as several potential moderator variables: mean age of infants, method (pointing or eyetracking), participant population type, estimates of mean vocabulary size of the sample population from the Words and Gestures form of the MacArthur-Bates Communicative Development Inventory when available (Fenson et al., 2007, MCDI; 1994),

---

<sup>1</sup>Data and analysis code for this and subsequent studies are available in an online repository at: [https://github.com/langcog/me\\_vocab](https://github.com/langcog/me_vocab)

<sup>2</sup>This cutoff was arbitrary but allowed us to include conditions from older children from non-typically-developing populations.

referent type (object or picture), and number of alternatives in the forced choice task. We used production vocabulary as our estimate of vocabulary size since it was available for more studies in our sample. We coded participant population as one of three subpopulations that have studied in the literature: (a) typically-developing monolingual children, (b) multilingual children (including both bilingual and trilingual children), and (c) non-typically developing children. Non-typically developing conditions included children with selective language impairment, language delays, hearing impairment, autism spectrum disorder, and down-syndrome.

In order to estimate effect size for each conditions, we also coded sample size, proportion novel-object selections, baseline (e.g., .5 in a 2-AFC paradigm), and standard deviations for novel object selections,  $t$ -statistic, and Cohen's  $d$ . For several conditions, there was insufficient data reported in the main text to calculate an effect size (no means and standard deviations,  $t$ -statistics, or Cohen's  $d$ s), but we were able to estimate the means and standard deviations though measurement of plots ( $N = 13$ ), imputation from other data within the paper ( $N = 4$ ; see SI for details), or through contacting authors ( $N = 26$ ). Our final sample included 157 effect sizes ( $N_{\text{typical-developing}} = 135$ ;  $N_{\text{multilingual}} = 12$ ;  $N_{\text{non-typically-developing}} = 10$ ).

**Statistical approach.** We calculated effect sizes (Cohen's  $d$ ) from reported means and standard deviations where available, otherwise we relied on reported test-statistics ( $t$  or  $d$ ). Effect sizes were computed by a script, `compute_es.R`, available in the Github repository. All analyses were conducted with the `metafor` package (Viechtbauer & others, 2010) using mixed-effect models with grouping by paper.<sup>3</sup> In models with moderators, moderators variables were included as additive fixed effects. All estimate ranges are 95% confidence intervals.

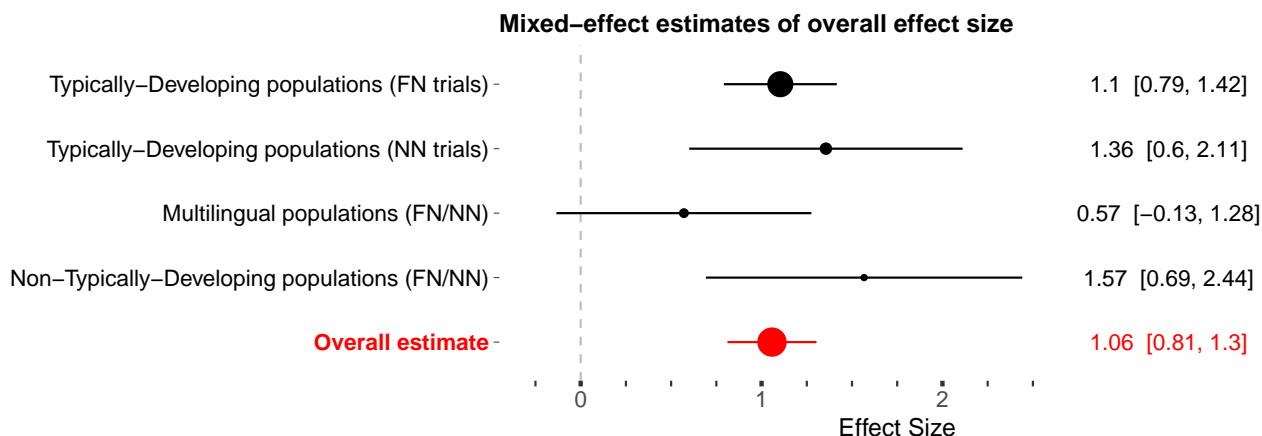
---

<sup>3</sup>The exact model specification was as follows: `metafor::rma.mv(yi = effect_size, V = effect_size_var, random = ~ 1 | paper)`.

## Meta-analytic Analyses

We conducted a separate meta-analysis for four theoretically-relevant conditions: Familiar-Novel trials with typically developing participants, Novel-Novel trials with typically developing participants, conditions with multilingual participants, and conditions with non-typically developing participants.

**Typically-Developing Population: Novel-Familiar Trials.** We first examined effect sizes for the disambiguation effect for typically-developing children in the canonical familiar-novel paradigm. This is the central data point that theories of disambiguation must explain.



*Figure 1.* Mixed-effect effect size estimates for all conditions (red) and each of the four theoretically-relevant conditions in our sample. Ranges are 95% confidence intervals. Point size corresponds to sample size. FN = Familiar-Novel trials; NN = Novel-Novel trials.

**Results.** The overall effect size for these conditions was 1.1 [0.79, 1.42], and reliably greater than zero ( $p < .001$ ; Figure 1). The effect sizes contained considerable heterogeneity, however ( $Q = 968.13$ ;  $p < .001$ ).

We next tried to predict this heterogeneity with two moderators corresponding to developmental change: age and vocabulary size. In a model with age as a moderator, age was a reliable predictor of effect size ( $\beta = 0.05$ ,  $z = 11.85$ ,  $p < .001$ ; see Table 1), suggesting that the disambiguation effect becomes larger as children get older. Age of participants was



Table 1

*Meta-analytic model parameters for model including age as a fixed effect. The first model (top) estimates effect sizes for all studies in our sample. The four subsequent models present separate models parameters for four separate conditions. Ranges are 95\% confidence intervals.*

Model	n	term	estimate	Z	p
Overall estimate	157	intercept	-0.18 [-0.47, 0.11]	-1.21	0.23
		age	0.03 [0.03, 0.04]	11.32	<.01
Typically-Developing populations (FN trials)	117	intercept	-0.33 [-0.71, 0.05]	-1.73	0.08
		age	0.05 [0.04, 0.05]	11.85	<.01
Typically-Developing populations (NN trials)	18	intercept	0.06 [-0.8, 0.93]	0.15	0.88
		age	0.03 [0.01, 0.04]	3.55	<.01
Multilingual populations (FN/NN)	12	intercept	0.05 [-0.78, 0.87]	0.11	0.91
		age	0.02 [0, 0.03]	1.77	0.08
Non-Typically-Developing populations (FN/NN)	10	intercept	-0.58 [-2.08, 0.92]	-0.75	0.45
		age	0.04 [0.01, 0.06]	3.15	<.01

*Note.* n = sample size (number of studies); FN = Familiar-Novel; NN = Novel-Novel.

highly correlated with vocabulary size in our sample ( $r = 0.65$ ,  $p < .01$ ), so next we asked whether vocabulary size predicted independent variance in the magnitude of the disambiguation bias on the subset of conditions for which we had estimates of vocabulary size ( $N = 23$ ). To test this, we fit a model with both age and vocabulary size as moderators. Vocabulary size ( $\beta = 0.07$ ,  $z = 2.14$ ,  $p = 0.03$ ), but not age ( $\beta = -0.78$ ,  $z = -1.11$ ,  $p = 0.27$ , was a reliable predictor of disambiguation effect size.

These analyses confirm that the disambiguation phenomenon is robust, and associated with a relatively large effect size ( $d = 1.1$  [0.79, 1.42]). In addition, this set of analyses provides theory-constraining evidence about the mechanisms underlying the effect. In

particular, the finding that vocabulary predicts more variance in effect size, compared to age, suggests that there is an experience related component to the mechanism, independent of maturational development alone.

**Typically-Developing Population: Novel-Novel Trials.** The results from the Familiar-Novel trials point to a role for vocabulary knowledge in the strength of the disambiguation effect. One way in which this vocabulary knowledge could lead to increased performance on the Familiar-Novel disambiguation task is through increased certainty about the label associated with the familiar word: If a child is less certain that a ball is called “ball,” then the child should be less certain that the novel label applies to the novel object. Novel-Novel trials control for potential variability in certainty about the familiar object by teaching participants a new label for a novel object prior to the critical disambiguation trial, where this previously-learned label becomes the “familiar” object in the disambiguation trial. If knowledge of the familiar object is not the only contributor to age-related changes in the disambiguation effect, then there should be developmental change in Novel-Novel trials, as well as Novel-Familiar trials. In addition, if the strength of knowledge of the “familiar” object influences the strength of the disambiguation effect, then the overall effect size should be smaller for Novel-Novel trials, compared to Familiar-Novel trials.

For conditions with the Novel-Novel trial design, the overall effect size was 1.36 [0.6, 2.11] and reliably greater than zero ( $p < .001$ ). We next asked whether age predicted some of the variance in these trials by fitting a model with age as a moderator. Age was a reliable predictor of effect size ( $\beta = 0.03$ ,  $z = 3.55$ ,  $p < .001$ ), suggesting that the strength of the disambiguation bias increases with age. There were no Novel-Novel conditions in our dataset where the mean vocabulary size of the sample was reported, and thus we were not able to examine the moderating role of vocabulary size on these trials.

Finally, we fit a model with both age and trial type (Familiar-Novel or Novel-Novel) as moderators of the disambiguation effect. Both moderators predicted independent variance in disambiguation effect size (age:  $\beta = -0.08$ ,  $z = -0.42$ ,  $p = 0.68$ ; trial-type:  $\beta = 0.04$ ,  $z =$

12.34,  $p < .0001$ ), with Familiar-Novel conditions and conditions with older participants  
tending to have larger effect sizes.

These analyses point to an influence on the disambiguation effect of both development  
(either via maturation or experience-related changes) as well as the strength of the familiar  
word representation. A successful theory of disambiguation will need to account for both of  
these empirical facts.

**Multilingual Population.** We next turn to a different population of participants:  
Children who are simultaneously learning multiple languages. This population is of  
theoretical interest because it allows us to isolate the influence of linguistic knowledge from  
the influence of domain-general capabilities. If the disambiguation phenomenon relies on  
mechanisms that are domain-general and independent of linguistic knowledge, then we  
should expect the magnitude of the effect size to be the same for multilingual children  
compared to monolingual children.

Children learning multiple languages reliably showed the disambiguation effect ( $d =$   
 $1.57 [0.69, 2.44]$ ). We next fit a model with both monolingual (typically-developing) and  
multilingual participants, predicting effect size with language status (monolingual  
vs. multilingual), while controlling for age. Language status was not a reliable predictor of  
effect size ( $\beta = 0.20$ ,  $z = 1.42$ ,  $p = 0.16$ ), but age was ( $\beta = 0.03$ ,  $z = 11.54$ ,  $p < .0001$ ).

These data do not provide strong evidence that language-specific knowledge influences  
effect size, however, the small sample size of studies from this population limit the power of  
this model to detect a difference if one existed.

**Non-Typically-Developing Population.** Finally, we examine a third-population  
of participants: non-typically developing children. This group includes a heterogeneous  
sample of children with diagnoses including Autism-Spectrum Disorder (ASD), Mental  
Retardation, Williams Syndrome, Late-Talker, Selective Language Impairment, and  
deaf/hard-of-hearing. These populations are of theoretical interests because they allow us to  
observe how impairment to a particular aspect of cognition influences the magnitude of the

disambiguation effect. For example, children with ASD are thought to have impaired social reasoning skills (e.g., Phillips, Baron-Cohen, & Rutter, 1998); thus, if children with ASD are able to succeed on disambiguation tasks, this suggests that social reasoning skills are not necessary to making a disambiguation inference.

Overall, non-typically developing children succeeded on disambiguation tasks ( $d = 1.57$  [0.69, 2.44]). In a model with age as a moderator, age was a reliable predictor of the effect, suggesting children became more accurate with age, as with other populations ( $\beta = 0.04$ ,  $z = 3.15$ ,  $p < .001$ ). We were not able to examine the potential moderating role of vocabulary size for this population because there were only 3 conditions where mean vocabulary size was reported.

We also asked whether the effect size for non-typically developing children differed from typically-developing children, controlling for age. We fit a model predicting effect size with both development type (typical vs. non-typical) and age. Development type was a reliable predictor of effect size with non-typically developing children tending to have a smaller bias compared to typically developing children ( $\beta = -0.50$ ,  $z = -2.86$ ,  $p < .0001$ ). Age was also a reliable predictor of effect size in this model ( $\beta = 0.04$ ,  $z = 11.34$ ,  $p < .0001$ ).

This analysis suggests that non-typically developing children succeed in the disambiguation paradigm just as typically developing children do, albeit at lower rates. Theoretical accounts of the disambiguation phenomenon will need to account for how non-typically developing children are able to succeed in the disambiguation task, despite a range of different cognitive impairments.

## Discussion

To summarize our meta-analytic findings, we find a robust disambiguation effect in each of the three populations we examined, as well as evidence that the magnitude of this effect increases across development. We also find that the effect is larger in the canonical Novel-Familiar paradigm compared to the Novel-Novel paradigm, but both designs show

roughly the same developmental trajectory.

Taken together, these analyses provide several theoretical constraints with respect to the mechanism underlying the disambiguation effect. First, language experience likely accounts for some developmental change. This conclusion derives from the fact that we see a larger effect size in Novel-Familiar trials compared to Novel-Novel trials, and that there is a suggestive correlation between vocabulary size and the strength of the disambiguation effect. Second, independent of familiar word knowledge, the strength of the bias increases across development. This constraint comes from the fact that the bias strengthens across development in the Novel-Novel conditions, and from the fact that there is not a significant impairment to effect in multilingual children (who presumably have less language experience with any particular language). Third, children with a range of different impairments are able to make the inference, suggesting that no single mechanism is both necessary and sufficient for the effect.

These three constraints are consistent with many of individual proposed accounts, as well as a various combinations of them. For example, an effect of language experience on the disambiguation effect via vocabulary knowledge is most consistent with the overhypothesis account, which predicts a stronger learned bias with vocabulary development. However, all four accounts predict developmental change in the NN trials. Under the overhypothesis account, as children are exposed to more language, they develop a stronger learned bias even when the “familiar” word is not previously known; Under the pragmatics account, as children are exposed to more language, they develop more skill in making social inferences, which would led to increased performance on the NN trials; And, under the bias and probabilistic accounts, maturational change could contribute to development in domain-general abilities, leading to a stronger disambiguation inference. Finally, the ability of children to succeed in the disambiguation tasks despite a range of impairments suggests that accounts that rely on a single mechanism, such as pragmatic reasoning or a mutual exclusivity constraint alone, are unlikely to describe the mechanism underlying the disambiguation effect across all children.

In the next section, we gather additional evidence to shed light on the relative contributions of these different mechanisms on the disambiguation effect. In particular, we use experimental methods to more directly examine the relationship between linguistic experience and the disambiguation effect.

### Experiment 1: Disambiguation Effect and Vocabulary Size

Our meta-analysis points to a robust developmental increase in the strength of the disambiguation effect with age. While all four accounts are able to predict this change, only the overhypothesis account predicts that this increase should be directly related to vocabulary knowledge. However, the meta-analytic approach is limited in its ability to measure this relationship since few studies in our sample measure vocabulary size ( $N = 8$ ), and even fewer measure vocabulary size at multiple ages within the same study (Markman et al., 2003;  $N=2$ ; Mather & Plunkett, 2009). In Experiment 1, we therefore aimed to test the prediction that children with larger vocabularies should have a stronger disambiguation bias by measuring vocabulary size in a large sample of children across multiple ages who also completed the disambiguation task. We find that vocabulary size is a strong predictor of the strength of the disambiguation effect across development and that vocabulary size predicts more variance than developmental age.

## Methods

**Participants.** A sample of 226 children were recruited at the Children’s Discovery Museum of San Jose. 72 children were excluded because they did not satisfy our planned inclusion criteria: within the age range of 24-48 months ( $n = 13$ ), completed all trials ( $n = 48$ ), exposed to English greater than 75% of the time ( $n = 37$ ), and correctly answered at least half of the familiar noun control trials ( $n = 55$ ). Our final sample included 154 children ( $N_{\text{females}} = 93$ ).

**Stimuli.** The disambiguation task included color pictures of 14 novel objects (e.g., a funnel) and 24 familiar objects (e.g. a ball; see Appendix). The novel words were the real 1-2

520 syllables labels for the unfamiliar objects (e.g., “funnel”, “tongs”, etc.; See Appendix). Items  
521 in the vocabulary assessment were a fixed set of 20 developmentally appropriate words from  
522 the Pearson Peabody Vocabulary Test (see Appendix; L. M. Dunn, Dunn, Bulheller, &  
523 Häcker, 1965).

524       **Design and Procedure.** Sessions took place individually in a small testing room  
525 away from the museum floor. The experimenter first introduced the child to “Mr. Fox,” a  
526 cartoon character who wanted to play a guessing game (see Fig. X). The experimenter  
527 explained that Mr. Fox would tell them the name of the object they had to find, so they had  
528 to listen carefully. Children then completed a series of 19 trials on an iPad, 3 practice trials  
529 followed by 16 experimental trials. In the practice trials, children were shown two familiar  
530 pictures (FF) on the iPad and asked to select one, given a label. If the participant chose  
531 incorrectly on a practice trial, the audio would correct them and allow the participant to  
532 choose again.

533       The child then completed the test phase. Like the practice trials, each of the test trials  
534 consisted of a word and two pictures, and the child’s task was to identify the referent.  
535 Within participants, we manipulated two features of the task: the target referent (Novel  
536 (Experimental) or Familiar (Control)) and the type of alternatives (Novel-Familiar or  
537 Novel-Novel; NF or NN). On novel referent trials, children were given a novel word and  
538 expected to select the novel object via the disambiguation inference. On familiar referent  
539 trials, children were given a familiar word and expected to select the correct familiar object.  
540 On Novel-Familiar trials, children saw a picture of a novel object and a familiar object (e.g. a  
541 funnel and a bol). On Novel-Novel trials, children saw pictures of two novel objects (e.g. a  
542 pair of tongs and a leak) . The design features were fully crossed such that half of the trials  
543 were of each trial type (Experimental-NF, Experimental-NN, Control-NF, Control-NN).  
544 Trials were presented randomly, and children were only allowed to make one selection.

545       After the disambiguation task, we measured children’s vocabulary in a simple  
546 vocabulary assessment. in which children were presented with four randomly selected images

and prompted to choose a picture given a label. Children completed 2 practice trials followed by 20 test trials.

**Data analysis.** Selections on the disambiguation task were coded as correct if the participant selected the familiar object on Control and the novel object on Experimental trials. We centered both age and vocabulary size for interpretability of coefficients. All models are logistic mixed effect models fit with the lme4 package in R (D. Bates, Mächler, Bolker, & Walker, 2015). Each model was fit with the maximal random effect structure. All ranges are 95% confidence intervals. Effect sizes are Cohen's  $d$  values.

## Results and Discussion

Participants completed the three practice trials (FF) with high accuracy, suggesting that they understood the task ( $M = 0.91$  [0.87, 0.94]).

We next examined performance on the four trial types. Children were above chance (.5) in both types of control conditions where they were asked to identify a familiar referent (Control-NF:  $M = 0.89$ ,  $SD = 0.17$ ,  $d = 2.35$  [2.06, 2.64]; Control-NN:  $M = 0.78$ ,  $SD = 0.25$ ,  $d = 1.14$  [0.9, 1.38]). Critically, children also succeeded on both types of experimental trials where they were required to select the novel object (NF:  $M = 0.84$ ,  $SD = 0.21$ ,  $d = 1.61$  [1.35, 1.87]; NN:  $M = 0.77$ ,  $SD = 0.28$ ,  $d = 0.95$  [0.71, 1.19]).

To compare all four conditions, we fit a model predicting accuracy with target type (F (Control) vs. N (Experimental)) and trial type (NF vs. NN) as fixed effects. We included both target type and trial type as main effects as well as a term for their interaction. There was a main effect of trial type, suggesting that participants were less accurate in NN trials compared to NF trials ( $B = -0.87$ ,  $SE = 0.25$ ,  $Z = -3.51$ ,  $p < .001$ ). The main effect of target type was not significant ( $B = -0.49$ ,  $SE = 0.29$ ,  $Z = -1.69$ ,  $p = 0.09$ ). The interaction between the two factors was marginal ( $B = 0.57$ ,  $SE = 0.36$ ,  $Z = 1.56$ ,  $p = 0.12$ ), suggesting that Novel target trials (Experimental) were more difficult than Familiar target trials (Control) for NF trials but not NN trials.



Table 2

*Parameters of logistic mixed model predicting accuracy on disambiguation trials as a function of trial type (Novel-Familiar (NF) vs. Novel-Novel (NN)), age (months), and vocabulary size as measured by our vocabulary assessment.*

term	Beta	SE	Z	p
(Intercept)	2.01	0.00	2,240.62	<.0001
Vocabulary	5.93	0.00	6,406.33	<.0001
Trial Type (NN)	-0.51	0.00	-564.56	<.0001
Age	0.02	0.00	21.80	<.0001
Vocabulary x Trial Type (NN)	-2.95	0.00	-3,185.91	<.0001
Vocabulary x Age	-0.01	0.00	-9.88	<.0001
Age x Trial Type (NN)	0.02	0.00	18.24	<.0001
Vocabulary x Age x Trial Type (NN)	0.13	0.00	145.54	<.0001

Our main question was how accuracy on the experimental trials changed over development. We examined two measures of developmental change: Age (months) and vocabulary size, as measured in our vocabulary assessment. We assigned a vocabulary score to each child as the proportion correct selections on the vocabulary assessment out of 20 possible. Age and vocabulary size were positively correlated, with older children tending to have larger vocabularies compared to younger children ( $r = 0.43$  [0.29, 0.55],  $p < .001$ ).

Figure 3 shows log linear model fits for accuracy as a function of age (left) and vocabulary size (right) for both NF and NN trial types. To examine the relative influence of maturation and vocabulary size on accuracy, we fit a model predicting accuracy with vocabulary size, age, and trial type (Experimental-NN, and Experimental-NF). We included all possible main and interaction terms as fixed effects. Table 2 presents the model

parameters. The only reliable predictor of accuracy was vocabulary size ( $B = 5.93$ ,  $SE = 0$ ,  $Z = 6406.33$ ,  $p < .0001$ ), suggesting that children with larger vocabularies tended to be more accurate in the disambiguation task. Notably, age was not a reliable predictor of accuracy over and above vocabulary size ( $B = 0.02$ ,  $SE = 0$ ,  $Z = 21.8$ ,  $p < .0001$ ).

**Discussion.** Experiment 1 directly examines the relationship between the strength of the disambiguation effect and vocabulary size. We find that the strength of the disambiguation effect is highly predicted by vocabulary size. In addition, we find that the bias is larger for NF trials, compared to NN trials.

The pattern of findings is consistent with meta-analytic estimates of those same effects. Figure 4 presents the data from the experimental conditions in Experiment 1 together with meta-analytic estimates, as a function of age. To compare the experimental data with the meta-analytic data, an effect size was calculated for each participant.<sup>4</sup> As in the meta-analytic models, the effect size is smaller for NN trials compared to NF trials, though the magnitude of this difference is smaller. We also see that the variance is larger for the meta-analytic estimates compared to the experimental data, presumably because there is more heterogeneity across experiments than across participants within the same experiment. The experimental data thus provide converging data with the meta-analysis that there is developmental change in the strength of the bias, and that the effect is weaker for NN trials.

In addition, the data from Experiment 1 provide new evidence relevant to the mechanism underlying the effect: children with larger vocabulary tend to have a stronger disambiguation bias. In principle there are two ways that vocabulary knowledge could support the disambiguation inference. The first is by influencing the strength of the learner's knowledge about the label for the familiar word: If a learner is more certain about the label for the familiar object, they can be more certain about the label for novel object. This account explains the developmental change observed for NF trials. However, this account

---

<sup>4</sup>Because some participants had no variability in their responses (all correct or all incorrect), we used the across-participant mean standard deviation as an estimate of the participant level standard deviation in order to convert accuracy scores into Cohen's  $d$  values.

does not explain the relationship of vocabulary with NN trials, since no prior vocabulary knowledge is directly relevant to this inference. This relationship between vocabulary size and NF size suggests that vocabulary knowledge could also influence the effect by providing evidence for general constraint that there is a one-to-one mapping between words and referents. This empirical fact is consistent with the overhypothesis account.

Importantly, however, data from both the meta-analytic study and the current experiment only provide correlational evidence about the relationship between vocabulary size and the disambiguation inference. In Experiment 2, we experimentally test the hypothesis that the strength of the learner’s knowledge about the familiar object influences the strength of the disambiguation inference, thereby testing one possible route through which vocabulary knowledge may be related to the disambiguation phenomenon.

## Experiment 2: Disambiguation Effect and Familiarity

In Experiment 2, we test a causal relationship between vocabulary size and the disambiguation effect by experimentally manipulating the strength of word knowledge. We do this by teaching participants a label for a novel object and varying the number of times the object is labeled. This manipulation allows us to vary children’s certainty about the label for an object, with objects that have been labeled more frequently associated with high certainty about the label name. The newly, unambiguously labeled object then serves as the “familiar” object in a novel-novel trial. If the strength of vocabulary knowledge about the “familiar” object influences the strength of the disambiguation effect, then we should expect a larger bias when the the familiar object has been labeled more frequently. We find a pattern consistent with this prediction.

## Methods

We report how we determined our sample size, all data exclusions (if any), all manipulations, and all measures in the study.

Table 3

*Demographics of children in Experiment 2.*

Age group	Mean age (months)	Sample size
2	30.99	38
3	40.99	35
4	52.16	37

**Participants.** We planned a total sample of 108 children, 12 per between-subjects labeling condition, and 36 total in each one-year age group (see Table 3). Our final sample was 110 children, ages 25 – 58.50 months, recruited from the floor of the Boston Children’s Museum. Children were randomly assigned to the one-label, two-label, or three label condition, with the total number of children in each age group and condition ranging between 10 and 13.

**Materials.** Materials were the set of novel objects used in de Marchena et al. (2011), consisting of unusual household items (e.g., a yellow plastic drain catcher) or other small, lab-constructed stimuli (e.g., a plastic lid glued to a popsicle stick). Items were distinct in color and shape.

**Procedure.** Each child completed four trials. Each trial consisted of a training and a test phase in a “novel-novel” disambiguation task (de Marchena et al., 2011). In the training phase, the experimenter presented the child with a novel object, and explicitly labeled the object with a novel label 1, 2, or 3 times (“Look at the *dax*”), and contrasted it with a second novel object (“And this one is cool too”) to ensure equal familiarity. In the test phase, the child was asked to point to the object referred to by a second novel label (“Can you show me the *zot*?”). Number of labels used in the training phase was manipulated between subjects. There were eight different novel words and objects. Object presentation side, object, and word were counterbalanced across children.

Table 4

*Parameters of logistic mixed model predicting accuracy on disambiguation trials as a function of age (months) and number of times a label for the familiar object was observed.*

term	B	SE	Z	p
(Intercept)	0.31	0.10	2.94	< .001
Age	0.05	0.01	4.13	< .001
Num. Labels Observed	0.48	0.13	3.75	< .001
Age x Num. Labels Observed	0.02	0.01	1.58	0.11

**Data analysis.** We followed the same analytic approach as we registered in Experiment 1, though data were collected chronologically earlier for Experiment 2. Responses were coded as correct if participants selected the novel object at test. A small number of trials were coded as having parent or sibling interference ( $N = 11$ ), experimenter error ( $N = 2$ ), or a child who recognized the target object ( $N = 4$ ), chose both objects ( $N = 2$ ) or did not make a choice ( $N = 8$ ). These trials were excluded from further analyses; all trials were removed for two children for whom there was parent or sibling interference on every trial. We centered both age and number of labels for interpretability of coefficients. The analysis we report here is consistent with that used in Lewis and Frank (2013), though there are some slight numerical differences due to reclassification of exclusions.

## Results and Discussion

As predicted, children showed a stronger disambiguation effect as the number of training labels increased, and as noise decreased with age (Figure 5).

We analyzed the results using a logistic mixed model to predict correct responses with age, number of labels, and their interaction as fixed effects, and participant as a random

effect. Model results are shown in Table 4. There was a significant effect of age such that older children showed a stronger disambiguation bias and a significant effect of number of labels, such that more training labels led to stronger disambiguation, but the interaction between age and number of labels was not significant.

These data provide causal evidence that the strength of knowledge of the familiar word influences the strength of the disambiguation effect. It thus points to one route through which a child's vocabulary knowledge might influence the disambiguation inference.

## General Discussion

## References

- Baron-Cohen, S., Leslie, A. M., & Frith, U. (1986). Mechanical, behavioural and intentional understanding of picture stories in autistic children. *British Journal of Developmental Psychology*, 4(2), 113–125.
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48. doi:[10.18637/jss.v067.i01](https://doi.org/10.18637/jss.v067.i01)
- Bion, R. A., Borovsky, A., & Fernald, A. (2013). Fast mapping, slow learning: Disambiguation of novel word–object mappings in relation to vocabulary learning at 18, 24, and 30 months. *Cognition*, 126(1), 39–53.
- Carey, S. (2010). Beyond fast mapping. *Language Learning and Development*, 6(3), 184–205.
- Carey, S., & Bartlett, E. (1978). Acquiring a single new word.
- Clark, E. (1987). The principle of contrast: A constraint on language acquisition. *Mechanisms of Language Acquisition*. Hillsdale, NJ: Erlbaum.
- de Marchena, A., Eigsti, I., Worek, A., Ono, K., & Snedeker, J. (2011). Mutual exclusivity in autism spectrum disorders: Testing the pragmatic hypothesis. *Cognition*, 119(1), 96–113.
- Diesendruck, G. (2005). The principles of conventionality and contrast in word learning: An empirical examination. *Developmental Psychology*, 41(3), 451.
- Diesendruck, G., & Markson, L. (2001). Children’s avoidance of lexical overlap: A pragmatic account. *Developmental Psychology*, 37(5), 630.
- Dunn, L. M., Dunn, L. M., Bulheller, S., & Häcker, H. (1965). *Peabody Picture Vocabulary Test*. American Guidance Service Circle Pines, MN.
- Fazly, A., Alishahi, A., & Stevenson, S. (2010). A probabilistic computational model of cross-situational word learning. *Cognitive Science*, 34(6), 1017–1063.
- Fenson, L., Bates, E., Dale, P. S., Marchman, V. A., Reznick, J. S., & Thal, D. J. (2007). *MacArthur-Bates Communicative Development Inventories*. Paul H. Brookes

Publishing Company.

Fenson, L., Dale, P. S., Reznick, J. S., Bates, E., Thal, D. J., Pethick, S. J., . . . Stiles, J.

(1994). Variability in early communicative development. *Monographs of the Society for Research in Child Development*, i–185.

Frank, M. C., Goodman, N. D., & Tenenbaum, J. B. (2009). Using speakers' referential intentions to model early cross-situational word learning. *Psychological Science*, 20(5), 578–585.

Golinkoff, R., Hirsh-Pasek, K., Baduini, C., & Lavalley, A. (1985). What's in a word? The young child's predisposition to use lexical contrast. In *Boston University Conference on Child Language, Boston*.

Golinkoff, R., Mervis, C., Hirsh-Pasek, K., & others. (1994). Early object labels: The case for a developmental lexical principles framework. *Journal of Child Language*, 21, 125–125.

Halberda, J. (2003). The development of a word-learning strategy. *Cognition*, 87(1), B23–B34.

Halberda, J. (2003). The development of a word-learning strategy. *Cognition*, 87(1), B23–B34.

Hutchinson, J. (1986). Children's sensitivity to the contrastive use of object category terms.

Lewis, M., & Frank, M. C. (2013). Modeling disambiguation in word learning via multiple probabilistic constraints. In *Proceedings of the 35th Annual Meeting of the Cognitive Science Society*.

Markman, E. (1990). Constraints children place on word meanings. *Cognitive Science*, 14(1), 57–77.

Markman, E., & Wachtel, G. (1988). Children's use of mutual exclusivity to constrain the meanings of words. *Cognitive Psychology*, 20(2), 121–157.

Markman, E., Wasow, J., & Hansen, M. (2003). Use of the mutual exclusivity assumption by



728 young word learners. *Cognitive Psychology*, 47(3), 241–275.

729 Markson, L., & Bloom, P. (1997). Evidence against a dedicated system for word learning in  
730 children. *Nature*, 385(6619), 813–815.

731 Mather, E., & Plunkett, K. (2009). Learning words over time: The role of stimulus  
732 repetition in mutual exclusivity. *Infancy*, 14(1), 60–76.

733 McMurray, B., Horst, J. S., & Samuelson, L. K. (2012). Word learning emerges from the  
734 interaction of online referent selection and slow associative learning. *Psychological*  
735 *Review*, 119(4), 831.

736 Mervis, C., Golinkoff, R., & Bertrand, J. (1994). Two-year-olds readily learn multiple labels  
737 for the same basic-level category. *Child Development*, 65(4), 1163–1177.

738 Phillips, W., Baron-Cohen, S., & Rutter, M. (1998). Understanding intention in normal  
739 development and in autism. *British Journal of Developmental Psychology*, 16(3),  
740 337–348.

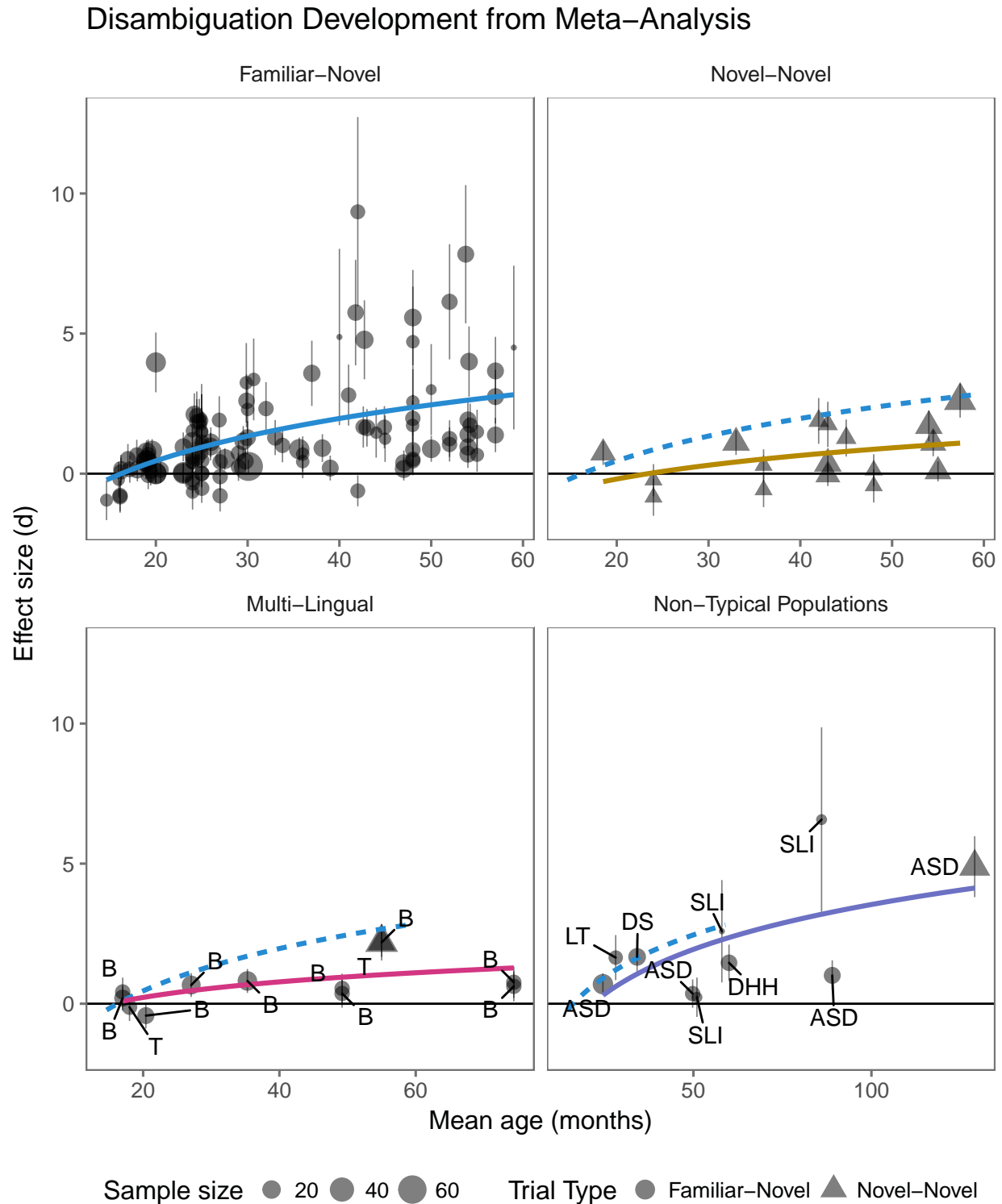
741 Preissler, M., & Carey, S. (2005). The role of inferences about referential intent in word  
742 learning: Evidence from autism. *Cognition*, 97(1), B13–B23.

743 Quine, W. (1960). *Word and object* (Vol. 4). The MIT Press.

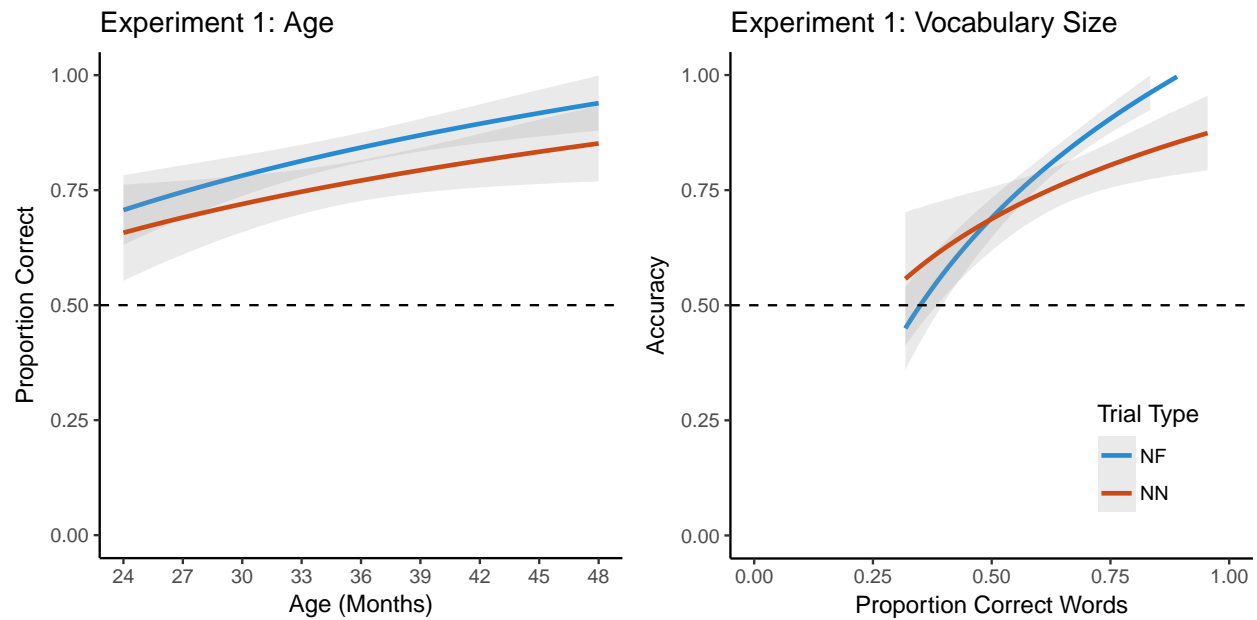
744 Regier, T. (2005). The emergence of words: Attentional learning in form and meaning.  
745 *Cognitive Science*, 29(6), 819–865.

746 Viechtbauer, W., & others. (2010). Conducting meta-analyses in r with the metafor package.  
747 *J Stat Softw*, 36(3), 1–48.

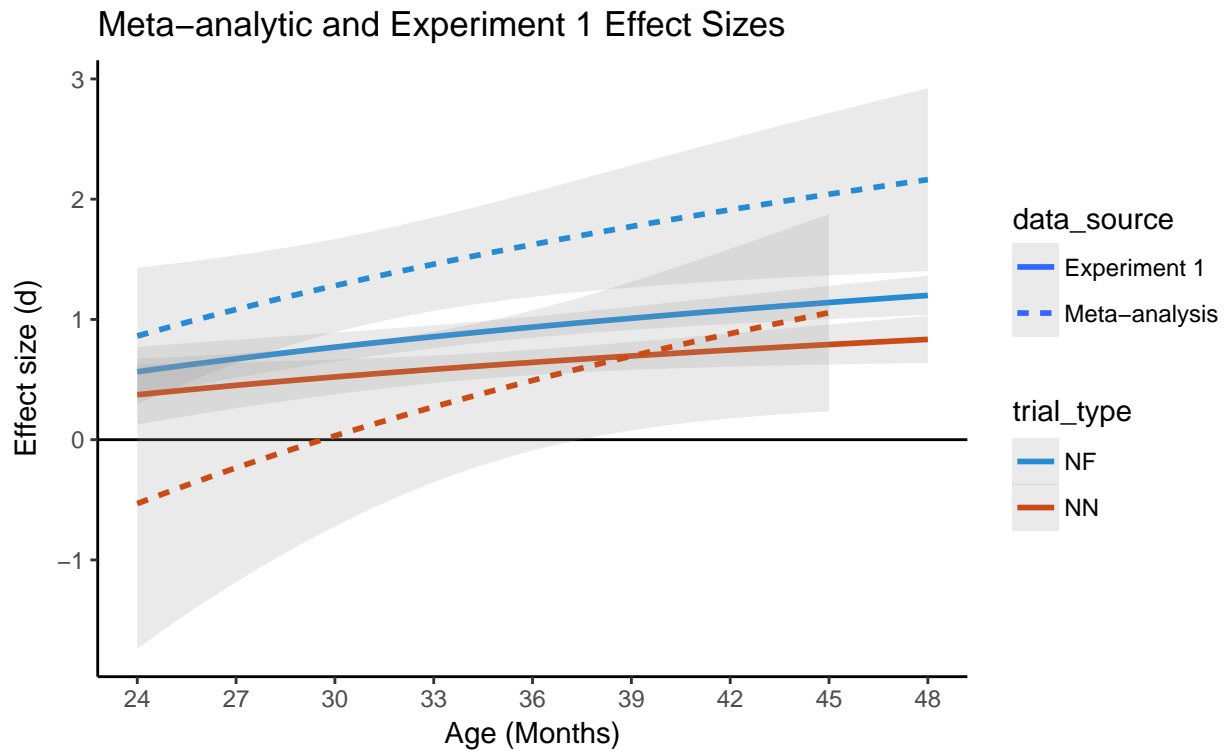
748 Vincent-Smith, L., Bricker, D., & Bricker, W. (1974). Acquisition of receptive vocabulary in  
749 the toddler-age child. *Child Development*, 189–193.



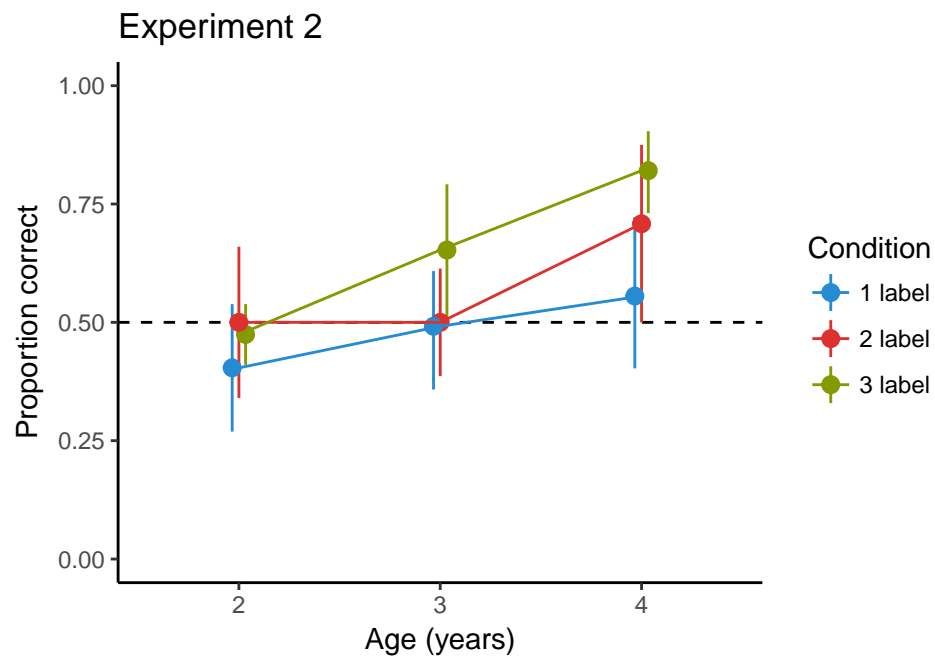
*Figure 2.* Developmental plots for each moderator. Ranges correspond to 95% confidence intervals. Model fits are log-linear. Point size corresponds to sample size, and point shape corresponds to trial type (Familiar–Novel vs. Novel–Novel). Note that the x-axis scale varies by facet. B = bilingual; T = trilingual; LT = late-talker; ASD = autism spectrum disorder; DS = down syndrome; SLI = selective language impairment; DHH = deaf/heard-of-hearing.



*Figure 3.* Experiment 1 results. Accuracy as a function of age (months; left) and vocabulary size (proportion correct on vocabulary assessment; right). Blue corresponds to trials with the canonical novel-familiar disambiguation paradigm, and red corresponds to trials with two novel alternatives, where a novel label for one of the objects is unambiguously introduced on a previous trial. The dashed line corresponds to chance. Ranges are 95% confidence intervals.



*Figure 4.* Meta-analytic data and data from experimental trials in Experiment 1 as a function of age. Effect sizes for Experiment 1 data are calculated for each participant, assuming the across-participant mean standard deviation as an estimate of the participant level standard deviation. Ranges are 95% confidence intervals.



*Figure 5.* Accuracy data for three age groups across three different conditions. Conditions varied by the number of times the child observed an unambiguous novel label applied to the familiar object prior to the critical disambiguation trial. The dashed line corresponds to chance. Ranges are 95% confidence intervals.

## Appendix

**Vocabulary Assessment Items (Exp. 1).**

1. hatchet
2. elephant
3. flamingo
4. duck
5. hug
6. broccoli
7. panda
8. hexagon
9. parallelogram
10. carpenter
11. drum
12. chef
13. bear
14. harp
15. vase
16. globe
17. triangle
18. vegetable
19. beverage
20. goat

**Familiar Words (Exp. 1).**

1. bottle
2. cup

- 774 3. spoon  
775 4. bowl  
776 5. apple  
777 6. cookie  
778 7. banana  
779 8. pretzel  
780 9. ball  
781 10. shoe  
782 11. flower  
783 12. balloon  
784 13. guitar,  
785 14. bucket

786 ### Novel Words (Exp. 1) 1. kettle 2. ladle 3. whisk 4. tongs 5. radish 6.  
787 leek 7. bok choy 8. kumquat 9. rudder 10. beaker 11. funnel 12. disk 13. bung 14.  
788 cam 15. chestnut 16. dulcimer 17. fig 18. ginger 19. gourd 20. longan 21. luffa 22.  
789 okra 23. pipette 24. sieve