

A Quantitative Synthesis of Early Language Acquisition Using Meta-Analysis

Molly Lewis¹, Mika Braginsky¹, Sho Tsuji², Christina Bergmann², Page Piccinini²,
Alejandrina Cristia², Michael C. Frank¹

¹ Department Psychology, Stanford University

² Laboratoire de Sciences Cognitives et Psycholinguistique, ENS

Author note

Correspondence concerning this article should be addressed to Molly Lewis, Psychology Department, Stanford University. 450 Serra Mall, Stanford, CA 94305. E-mail: mll@stanford.edu.

Abstract

replicability, etc.

Keywords: replicability, reproducibility, meta-analysis, developmental psychology,
language acquisition

Word count: XXXX

A Quantitative Synthesis of Early Language Acquisition Using Meta-Analysis

Introduction

Psychologists hope to build generalizable theories about human behavior—theories that hold true beyond particulars of an individual study. The field has grown concerned as a result in the face of recent high-profile evidence that an effect observed in one study may not be the same in another (“replicability crisis”; Ioannidis, 2005; Nosek, 2012, 2015). Some of this variability is to be expected, however—the question we should instead be asking is, do the data provide support for the theory, even if they are noisy? Furthermore, to build parsimonious theories of human behavior, we should seek to explain not just individual phenomenon, but entire literatures of research. What is needed, then, is a tool for aggregating noisy data across studies within a phenomenon, as well as a common language for comparing effects across phenomena.

Meta-analytic methods provide a powerful tool for doing just this. The basic unit of meta-analysis—the effect size—provides an estimate of the *size* of an effect, as well as a measure of uncertainty around this point estimate. With such a continuous measure of success, we can apply the same reasoning we use to aggregate noisy measurements over participants in a single study: By assuming each *study*, rather than participant, is sampled from a population, we can appeal to the classical statistical framework to combine estimates of the effect size for a given phenomenon.

This quantitative approach provides a rich tool kit for synthesizing across literatures. By describing different phenomena using the same unit of measurement, we are able to compare effects in different domains. Rather than simply concluding that two effects are both “real,” we can ask more fine-grained questions: Is effect *X* bigger than effect *Y*? Does a moderator influence effect *X* in the same way as effect *Y*? This type of continuous analysis supports building quantitative models, and specifying theories that are more precise and constraining.

In addition to these theoretical motivations, there are practical reasons for conducting

a quantitative synthesis. When planning an experiment, an estimate of the size of an effect on the basis of prior literature can inform the sample size needed to achieve a desired level of power. Meta-analytic estimates of effect sizes can also aid in design choices: If a certain paradigm tends to have overall larger effect sizes than another, the strategic researcher might select this paradigm in order to maximize the power of a study.

In practice, however, the feasibility of this meta-analytic approach relies on the field's commitment to practices that facilitate cumulative science. These practices apply to all stages of the research process. At the stage of experimental planning, researchers must pre-specify analytical decision to limit "researcher" degrees of freedom (Simmons, 2011; Simonsohn, 2014a, 2014b, 2014c). At the stage of completion, researchers should share a result regardless of its significance (Rosenthal, 1979; Fanelli 2012). And, at the stage of sharing, researchers must provide enough information about the method for another lab to conduct a close replication. Critically, reports must also contain complete descriptions of both data and analytical decisions so that effect sizes can be calculated for the purposes of meta-analysis,

In the present paper, we use meta-analytic methods to provide a quantitative synthesis of an entire field of psychological research: language acquisition. We think this field is a particularly informative case study. It may be particularly vulnerable to false findings because running children is expensive (Ioannidis, 2005), and thus:

- sample sizes are small
- replications difficult and rare
- Recent attention about practices in developmental research Peterson (2016)

We have two goals:

- Describe the state of the field in terms of its participation in practices that are prerequisites to cumulative science, and ultimately, a theoretical synthesis
- Provide a preliminary theoretical synthesis of the field

Towards this end, we introduce [Metalab](#).

Method

We analyzed 11 different phenomena in language acquisition. We selected these phenomena in order to describe development at many different levels of the language hierarchy, from the acquisition of prosody and phonemic contrasts, to gaze following in linguistic interaction. This wide range of phenomena allowed us to compare the course of development across different domains, as well as explore questions about the interactive nature of language acquisition (Table 1).

To obtain estimates of effect size, we coded papers reporting experimental data. Within each paper, we calculated a separate effect size estimate for each experiment and age group (“conditions”). In total, our sample includes estimates from 258 papers, 938 different conditions and 11,628 participants. The process for selecting papers from the literature differed by domain, with some individual meta-analyses using more systematic approaches than others. [Simulations here?]

Replicability of the field

A literature is more likely to describe a real effect if studies are randomly sampled from the population of all possible studies that researchers could in principle conduct. This assumption does not mean, however, that there should be *no* variability in effect size across studies: We should expect random variation around the true mean effect size, with smaller studies showing more variability around this mean.

Variability in effect sizes will be biased when this assumption of random study sampling does not hold. Bias may be introduced by the experimenter in a number of ways, including failure to publish null findings (Fanelli, 2010; Rosenthal, 1979; “publication bias”, Rothstein, Sutton, & Borenstein, 2006), analytical flexibility (e.g., “p-hacking”, Simmons, Nelson, & Simonsohn, 2011; Simonsohn, Nelson, & Simmons, 2014), reporting errors, or even fraud. These biases are problematic for theoretical development because they lead to large

Level	Phenomenon	Description	N papers (conditions)
Prosody	IDS preference (Dunst, Gorman, & Hamby, 2012)	Looking times as a function of whether infant-directed vs. adult-directed speech is presented as stimulation.	16 (50)
Sounds	Phonotactic learning (Cristia, in prep.)	Infants' ability to learn phonotactic generalizations from a short exposure.	15 (47)
	Vowel discrimination (native) (Tsuji & Cristia, 2014)	Discrimination of native-language vowels, including results from a variety of methods.	40 (167)
	Vowel discrimination (non-native) (Tsuji & Cristia, 2014)	Discrimination of non-native vowels, including results from a variety of methods.	21 (72)
Phonotactics	Statistical sound learning (Cristia, in prep)	Infants' ability to learn sound categories from their acoustic distribution.	11 (40)
Proto-words	Word segmentation (Bergmann & Cristia, 2015)	Recognition of familiarized words from running, natural speech using behavioral methods.	67 (295)
Words	Mutual exclusivity (Lewis & Frank, in prep.)	Mapping of novel words reflecting children's inference that novel words tend to refer to novel objects.	20 (60)
	Concept-label advantage (Lewis & Long, unpublished)	Infants' categorization judgments in the presence and absence of labels.	16 (100)
	Online word recognition (Frank, Lewis, & MacDonald, 2016)	Online word recognition of familiar words using two-alternative forced choice preferential looking.	12 (32)
Communication	Gaze following (Frank, Lewis, & MacDonald, 2016)	Gaze following using standard multi-alternative forced-choice paradigms.	15 (45)
	Pointing and vocabulary (Colonnaesi et al., 2010)	Longitudinal correlations between declarative pointing and later vocabulary.	25 (30)

Table 1
Overview of meta-analyses in dataset.

but often unknown errors in estimates of the effect size. If bias is present in the literature, estimates of effect size may be poor estimates of the true underlying effect size and thus provide little evidential value. To make theoretical progress, we must therefore distinguish variability in effect sizes due to sample size from variability due to bias.

To assess the replicability of language acquisition phenomena, we conducted several key

analyses: Fail-safe-N, funnel plots, and p-curve. These methods each have limitations, but taken together, they provide converging evidence about the replicability of a literature. We find little evidence of bias in our meta-analyses, suggesting that the language acquisition literature describes real psychological phenomena and can therefore provide the basis for theoretical development.

Fail-safe-N

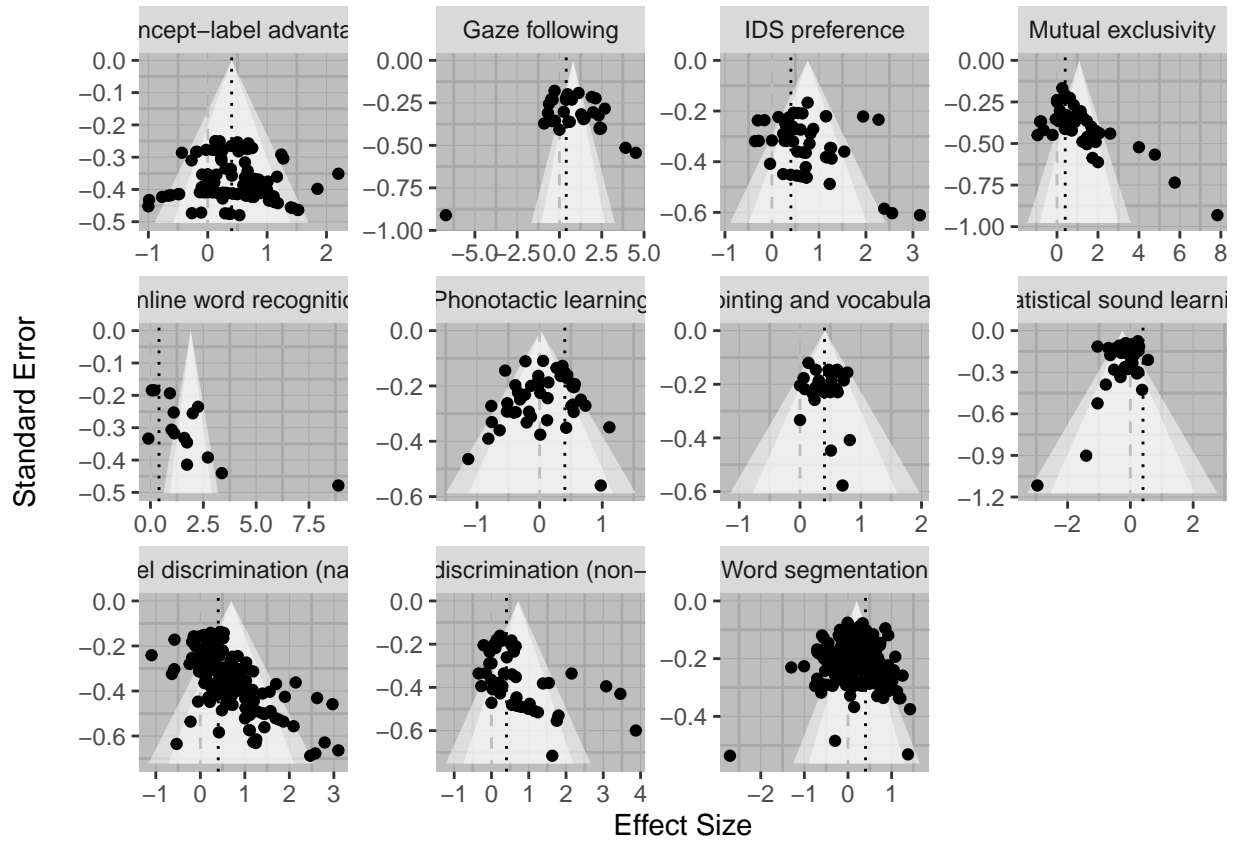
One approach for quantifying the reliability of a literature is to ask, How many missing studies with null effects would have to exist in order for the overall effect size to be zero? This is called the “fail-safe” number of studies. To answer this question, we estimated the overall effect size for each phenomenon (Table 2, column 2), and then used this to estimate the fail-safe-N (Table 2, column 3). This analysis suggests that an unlikely number of studies would have to be “missing” in each literature ($M = 3914$) in order for the overall effect sizes to be 0.

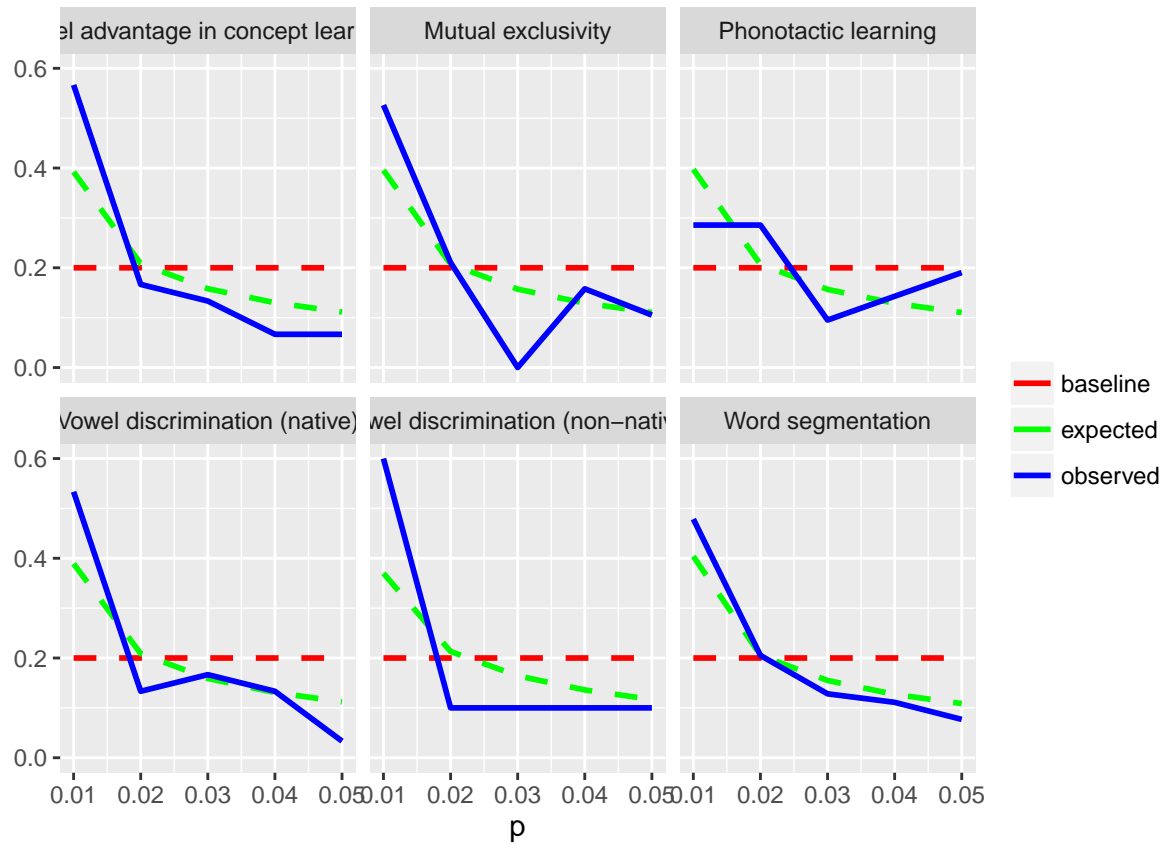
One limitation of this analysis, however, is that it assumes that all reported effect sizes are obtained without p-hacking/analytical flexibility: If experimenters (is this true?)

Funnel Plots

Funnel plots provide a visual method for evaluating whether variability in effect sizes is due only to differences in sample size. Figure XX plots effect sizes versus a metric of sample size, standard error. If there is no bias in a literature, we should expect studies to be randomly sampled around the mean, with more variability for less precise studies.

Figure 1 presents funnel plots for each of our 11 meta-analyses. These plots show evidence of asymmetry (bias) for several of our phenomena (Table 2, column 1). An important limitation of this method, however, is that it is difficult to determine the source of this bias. One possibility is that this bias is due not to experimenter malfeasance, but to true heterogeneity in phenomena (e.g. different ages). P-curves provide a way to address this issue, which we turn to next.



P-curves

(U. Simonsohn, Nelson, & Simmons, 2014; Simonsohn et al., 2014; Simonsohn, Simmons, & Nelson, 2015) bias introduced by meta-analysis in selection (second-order selection bias)

Theoretical Synthesis**OUTLINE****Statistical Approach****METAMETAPLOT****Discussion**

Author Contributions.

Acknowledgments.

Phenomenon	d	fail-safe-N	funnel skew	p-curve skew	power
IDS preference	0.71 [0.53, 0.89]	3762	1.88 (0.06)		
Phonotactic learning	0.04 [-0.09, 0.16]	45	-1.08 (0.28)	-1.52 (0.06)	0.14
Vowel discrimination (native)	0.6 [0.5, 0.71]	9536	8.98 (0)	-5.42 (0)	0.67
Vowel discrimination (non-native)	0.66 [0.42, 0.9]	3391	4.13 (0)	-3.24 (0)	0.78
Statistical sound learning	-0.14 [-0.27, -0.02]	Inf	-1.87 (0.06)		
Word segmentation	0.19 [0.14, 0.24]	5372	2.17 (0.03)	-9.67 (0)	0.56
Mutual exclusivity	1 [0.68, 1.33]	6417	6.08 (0)		
Concept-label advantage	0.4 [0.29, 0.51]	3928	0.31 (0.76)	-6.15 (0)	0.69
Online word recognition	1.89 [0.81, 2.96]	2843	2.92 (0)		
Gaze following	0.84 [0.26, 1.42]	2641	-1.69 (0.09)		
Pointing and vocabulary	0.41 [0.32, 0.49]	1202	0.59 (0.55)		

Table 2

Summary of replicability analyses. d = Effect size (Cohen's d) estimated from a random-effect model; fail-safe- N = number of missing studies that would have to exist in order for the overall effect size to be $d = 0$; funnel skew = test of asymmetry in funnel plot using the random-effect Egger's test (Stern & Eggers, 2005); p-curve skew = test of the right skew of the p-curve using the Stouffer method (Simonsohn, Simmons, & Nelson, 2015); power = power to reject the null hypothesis at the 5% significance level based on the p-curve (Simonsohn, Nelson, & Simmons, 2014); Brackets give 95% confidence intervals, and parentheses show p-values.

References. Bergmann, C., & Cristia, A. (2015). Development of infants' segmentation of words from native speech: A meta-analytic approach. *Developmental Science*.

Dunst, C., Gorman, E., & Hamby, D. (2012). Preference for infant-directed speech in preverbal young children. *Center for Early Literacy Learning*, 5(1).

Fanelli, D. (2010). Positive Results Increase Down the Hierarchy of the Sciences. *PLoS ONE*, 5(4), e10068–10.

Frank, M. C., Lewis, M. L., & MacDonald, K. (in press). A performance model for early word learning. In *Proceedings of the 38th Annual Conference of the Cognitive Science Society*. Retrieved from

http://langcog.stanford.edu/papers_new/frank-2016-underrev.pdf

Lewis, M., & Frank, M. C. (in prep). Multiple routes to disambiguation.

Peterson, D. (2016). The Baby Factory: Difficult Research Objects, Disciplinary Standards, and the Production of Statistical Significance. *Socius: Sociological Research for a Dynamic World*, 2(0), 1–10.

Rosenthal, R. (1979). The file drawer problem and tolerance for null results. *Psychological Bulletin*, 86(3), 638.

Rothstein, H. R., Sutton, A. J., & Borenstein, M. (2006). *Publication bias in meta-analysis: Prevention, assessment and adjustments*. John Wiley & Sons.

Simmons, Nelson, L. D., & Simonsohn, U. (2011). False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant. *Psychological Science*, 22(11), 1359–1366.

Simonsohn, U., Nelson, L. D., & Simmons, J. P. (2014). P-curve and effect size correcting for publication bias using only significant results. *Perspectives on Psychological Science*, 9(6), 666–681.

Simonsohn, Nelson, L. D., & Simmons, J. P. (2014). P-curve: A key to the file-drawer. *Journal of Experimental Psychology: General*, 143(2), 534.

Simonsohn, Simmons, J. P., & Nelson, L. D. (2015). Better p-curves. *Simonsohn, Uri, Joseph P. Simmons, and Leif D. Nelson (Forthcoming), "Better P-Curves," Journal of Experimental Psychology: General*.

Sterne, J. A., & Egger, M. (2005). Regression methods to detect publication and other bias in meta-analysis. *Publication Bias in Meta-Analysis: Prevention, Assessment, and Adjustments*, 99–110.

Tsuiji, S., & Cristia, A. (2014). Perceptual attunement in vowels: A meta-analysis.

Developmental Psychobiology, 56(2), 179–191.