

1 Estimating age-related change in infants' linguistic and cognitive development using
2 (meta-)meta-analysis

3 Anjie Cao¹ & Michael C. Frank¹

4 ¹ Stanford University

5 ² Konstanz Business School

6 Author Note

7 Add complete departmental affiliations for each author here. Each new line herein
8 must be indented, like this line.

9 Enter author note here.

10 The authors made the following contributions. Anjie Cao: Conceptualization,
11 Writing - Original Draft Preparation, Writing - Review & Editing; Michael C. Frank:
12 Writing - Review & Editing, Supervision.

13 Correspondence concerning this article should be addressed to Anjie Cao, 450 Jane
14 Stanford Way, Stanford, CA 94305. E-mail: anjiecao@stanford.edu

Abstract

Developmental psychology focuses on how psychological phenomena emerge with age. In cognitive development research, however, the specifics of this emergence is often underspecified. Researchers often provisionally assume linear growth by including chronological age as a predictor in regression models. In this work, we aim to evaluate this assumption by examining the functional form of age trajectories across 24 phenomena in early linguistic and cognitive development using (meta-)meta-analysis. Surprisingly, for most meta-analyses, the effect size for the phenomenon was relatively constant throughout development. We investigated four possible hypotheses explaining this pattern: (1) age-related selection bias against younger infants; (2) methodological adaptation for older infants; (3) change in only a subset of conditions; and (4) positive growth only after infancy. None of these explained the lack of age-related growth in most datasets. Our work challenges the assumption of linear growth in early cognitive development and suggests the importance of uniform measurement across children of different ages.

Keywords: keywords

Word count: X

31 Estimating age-related change in infants' linguistic and cognitive development using
32 (meta-)meta-analysis

33 Developmental psychology focuses on how psychological constructs change with age.
34 Throughout the years, many theories have been proposed to characterize and explain how
35 and why developmental changes happen (Bronfenbrenner, 1977; Carey, 2009; Elman, 1996;
36 Flavell, 1994; e.g., Piaget, 1971; Thelen & Smith, 2007). Among these theories, one
37 common assumption is that skills increase with age (positive change assumption): children
38 get better as they get older. Often, researchers treat age as a predictor in linear regression
39 models, and therefore implicitly assume that the constructs of interests follow a linear
40 trajectory (Lindenberger & Pötter, 1998). While both assumptions are widely adopted,
41 especially in early cognitive and language development, their validity is rarely tested.

42 One common approach to evaluating the functional form of age-related changes is
43 through longitudinal studies. Measurements of psychological constructs, when tracked
44 longitudinally, often reveal the age trajectories that violate the linearity assumption. For
45 instance, a longitudinal study that follows the development of executive function (EF) from
46 3 to 5 years-old using a battery of EF tasks show that EF follows a non-linear trajectory
47 over age (Johansson, Marciszko, Brocki, & Bohlin, 2016). Similarly, vocabulary in early
48 childhood, measured by MacArthur-Bates Communicative Development Inventories, also
49 follows the exponential trend rather than the linear trend (Frank, Braginsky, Yurovsky, &
50 Marchman, 2021). In many domains with established measurements, longitudinal research
51 has been used to characterize the functional form of the development (Adolph, Robinson,
52 Young, & Gill-Alvarez, 2008; Cole, Loughheed, Chow, & Ram, 2020; Karlberg, Engström,
53 Karlberg, & Fryer, 1987; McArdle, Grimm, Hamagami, Bowles, & Meredith, 2009; Tilling,
54 Macdonald-Wallis, Lawlor, Hughes, & Howe, 2014). However, longitudinal methods are
55 more rarely applied to experimental studies that identify proposed mechanisms underlying
56 development.

Many important findings in early language and cognitive development are primarily attested in cross-sectional experimental studies. For example, in the language learning domain, many studies have targeted specific mechanisms proposed to underlie how infants acquire specific facets of language. Constructs such as mutual exclusivity (Markman & Wachtel, 1988), statistical learning (Saffran, Aslin, & Newport, 1996), syntactic bootstrapping (Naigles, 1990) and so on, are all attested through decades of experimental evidence acquired through cross-sectional studies. These works are critical to test the causal mechanisms underlying age-related changes, but they are rarely measured in samples with sufficient size and age variation to test the positive change assumption or the assumption of linearity (cf. Frank et al., 2017). In an ideal world, one would run those experiments longitudinally on a large, diverse sample. In practice, this goal is difficult to achieve due to the constraints on both time and financial resources. As a result, the functional forms of age-related changes in critical constructs remain poorly understood.

To address this issue, we turned to meta-analysis. Meta-analysis is a statistical method to aggregate evidence across studies quantitatively. This approach has been widely adopted in many disciplines and subfields, including developmental psychology (Doebel & Zelazo, 2015; e.g. Hyde, 1984; Letourneau, Duffett-Leger, Levac, Watson, & Young-Morris, 2013). Compared to the single study approach, meta-analysis has several advantages. First, it allows us to examine the robustness of the phenomena documented in the literature. By combining results from multiple studies, meta-analysis enhances the statistical power to detect effects that might be too small to identify in individual studies. Second, meta-analysis provides a framework for assessing the consistency of research findings across different contexts (Borenstein, Hedges, Higgins, & Rothstein, 2021; Egger, Smith, & Phillips, 1997). Further, pooling across developmental studies with different cross-sectional samples may yield sufficient variation to explore the functional form of age-related change with greater precision than individual studies.

In this work, we aim to leverage meta-analysis to examine the shape of the

developmental trajectory in key constructs in infant language and cognitive development. Specifically, we use existing meta-analyses from Metalab (<https://langcog.github.io/metalab/>), a platform that hosts community-augmented meta-analyses. Metalab was established to provide dynamic databases publicly available to all researchers (Bergmann et al., 2018). Researchers can deposit their meta-analysis dataset in the platform, and they can also use the dataset for custom analyses (e.g. Cao, Lewis, & Frank, 2023; Lewis et al., 2016). To this date, Metalab contains #FIXME effect sizes from #FIXME different meta-analysis, spanning different areas of developmental psychology. This resource allows us to examine the suitability of meta-analysis as a tool to characterize developmental trajectory – and if suitable, provides insights into how these key constructs develop across the early months of childhood.

We acknowledge at the outset that meta-analysis has significant limitations. The quality of a meta-analysis is necessarily constrained by the quality of the existing studies (Simonsohn, Simmons, & Nelson, 2022). If the studies being aggregated are flawed, the conclusions drawn from the meta-analysis will also be questionable. Moreover, one significant issue in interpreting meta-analysis is the heterogeneity among studies. Heterogeneity refers to the variability in study participants, interventions, outcomes, and methodologies. This diversity can make it challenging to aggregate results meaningfully, because differences between studies may reflect true variation in effects rather than a singular underlying effect size (Fletcher, 2007; Higgins & Thompson, 2002; Huedo-Medina, Sánchez-Meca, Marín-Martínez, & Botella, 2006; Thompson & Sharp, 1999). Critically, understanding the source of heterogeneity often requires detailed coding of the potential moderators; this process is frequently hampered by the inadequate reporting standards prevalent in psychological literature, which often leaves essential information for coding these moderators absent (Nicholson, Deboeck, & Howard, 2017; Publications, Journal, & Standards, 2008). In other words, whether meta-analysis can provide insights into the nature of age-related change is dependent upon the quality of the existing literature.

This paper is organized as follows. In the first section, we provide an overview on the estimated general shape of age-related change across the datasets in Metalab. To preview our findings, we found that most datasets showed relatively constant effect size across age. This finding challenges the commonly held linearity assumption and the positive increase assumption. In the second section, we test four hypotheses on why the current meta-analyses failed to reveal age-related changes: (1) age-related selection bias against younger infants; (2) methodological adaptation for older infants; (3) change in only a subset of conditions; and (4) positive growth only after infancy. We found that none of the four explanations provided a satisfying explanation for the lack of age-related change in most meta-analyses.

Datasets

Datasets were retrieved from Metalab. As of February 2024, Metalab hosted 32 datasets in total, with research areas ranging from language learning to cognitive development. All datasets included effect size estimates converted to standardized mean difference (SMD; also known as Cohen's d) as well as estimates of effect size variance and a variety of other moderators (e.g., average age of participants) provided by the contributors. There were 2 desiderata for the datasets to be included in the final analysis:

1. The dataset must describe an experimental (non-correlational) effect that uses behavioral measures, and
2. For a dataset that has already been published, the meta-analytic effect reported in the published form must not be null (i.e., must be significantly different than zero).

Five datasets did not meet the first desideratum (*Pointing and vocabulary (concurrent)*; *Pointing and vocabulary (longitudinal)*; *Video deficit*; *Symbolic play*; *Word segmentation (neuro)*), and one dataset did not meet the second desideratum (*Phonotactic learning*). These datasets were not included in the analysis.

For the remaining 26 datasets, we made the following modifications. Following the organization in the original meta-analysis (Gasparini, Langus, Tsuji, & Boll-Avetisyan, 2021), we separated the Language discrimination and preference dataset into two datasets, one for discrimination and one for preference. We also combined two pairs of datasets because they were testing the same experimental effects: *Gaze following (live)* and *Gaze following (video)* was combined into *Gaze following (combined)*; *Function word segmentation* and *Word segmentation (behavioral)* was combined into *Word segmentation (combined)*. We also replaced the *Infant directed speech preference* dataset with a more up-to-date version reported in Zettersten et al. (2023).

To make the comparison more equivalent to each other, we would run models with the same random effect structure specifications across all datasets. To achieve this goal, we recoded the relevant grouping variables in the datasets with missing grouping variables.

Since we were mostly interested in the age trajectory of these constructs in early childhood, we further trimmed the datasets to include only effect sizes from participants under 36 months of age. This decision did not qualitatively affect our findings as most datasets did not include data above age 36 months. The final analysis included 25 datasets in total. Table 1 presented the names of all the datasets, along with the number of effect sizes and participants included for each dataset.

Methods

All of the statistical analyses were conducted in R. Meta-analytic models were fit using the metafor package (Viechtbauer, 2010). This was an exploratory study in which no hypotheses were pre-registered.

For each dataset, we considered four functional forms as possible candidates for the shape of the developmental trajectory: linear, logarithmic, quadratic, and constant. A linear form is the most common assumption in the literature, whereas logarithmic and

quadratic were chosen to represent sublinear growth and superlinear growth, respectively. The constant form served as a baseline null hypothesis for the other alternative growth patterns. Although other, more complex growth patterns are of course possible, we opted to compare these forms as a first pass. Note that the constant model includes one parameter (an intercept), linear and logarithmic models include two parameters (an intercept and a slope), and the quadratic model includes three parameters (intercept, slope, and quadratic growth term).

For all analyses, we fit multilevel random-effects meta-regression models using nested random intercepts to account for both the testing of individual samples in multiple conditions (e.g., in a between-participants design) and multiple studies within a single paper. Meta-regression models predicted effect sizes (standardized mean difference / Cohen's d) with mean age in months in different functional forms. We fit four meta-regression models in total for each dataset.

Results

Model comparison. Our initial goal was to compare the fit of models with different functional forms for each meta-analysis. Because models differed in their complexity (number of parameters), we extracted the corrected AIC (AICc) for each model. The model with the lowest AICc was considered the baseline model, and all the remaining models were compared against the baseline. The remaining model each received a Δ_{AIC} , which was the difference between the AIC of the model and the AIC of the baseline model. Following standard convention, we treated $\Delta_{AIC} > 4$ as the statistical significance threshold (Burnham & Anderson, 2004). A baseline model was significantly better than an alternative model if and only if the alternative model had $\Delta_{AIC} > 4$.

Surprisingly, the four functional forms could not be meaningfully distinguished in 19 out of 25 datasets.. (This situation typically arises because the data are constant and

hence more complex models with zero parameters fit the data equally well¹). The remaining 6 datasets yielded meaningful contrasts between different functional forms, but the linear form was not the best-fitting form for any dataset. Table 2 shows the model comparison results for each dataset. Figure 1 shows the prediction of each functional form.

Linearity and Positive Increase Assumption. One limitation of the model comparison approach is that it does not quantify growth over time. To further examine the positive increase assumption, we estimated linear meta-regression models and examined the estimates on the age predictor. We found that the slope estimate for age was not significantly different from zero in the majority of the datasets (16/25; Fig 2).

Discussion

We conducted model comparisons to assess the functional forms of age-related change across 25 datasets. Four functional forms—Logarithmic, Linear, Quadratic, and Constant—were largely indistinguishable within most datasets. Notably, in datasets where contrasts were meaningful, linear models received no support, challenging the prevalent linearity assumption for early linguistic and cognitive development. Further, we only detected any positive growth in 8/25 meta-analyses. Past work has successfully revealed age-related changes using meta-analysis (e.g. Best & Charness, 2015; McCartney, Harris, & Bernieri, 1990; Sugden & Marquis, 2017). But in most datasets that we have considered, effect size does not increase with age. Why?

Understanding the lack of developmental change in meta-analytic data

Here we consider four explanations for the lack of age-related change in most of the meta-analyses we examined. First, meta-analyses are susceptible to publication bias

¹ In the situation of a completely constant pattern of effects across age, the maximal difference in model fit would be an AICc of exactly 4 between the constant and quadratic model, reflecting a two-parameter difference.

(Ferguson & Brannick, 2012; Ferguson & Heene, 2012; Francis, 2012; Thornton & Lee, 2000). And the bias could be related to the characteristics of the study, such as the inclusion of younger participants (Kathleen M. Coburn & Vevea, 2015). Consequently, studies with younger participants may have effect sizes that were more inflated, compared to the studies with older participants. The selectivity of publication bias would thus obscure the possible developmental changes in the dataset.

Second, researchers may change methods as infants expand their behavioral repertoire. For example, the high-amplitude sucking paradigm is most likely to be deployed on very young infants, whereas the looking paradigm is most likely to be used on older infants. We did see some evidence for method adaptation in some datasets. For example, in *Language discrimination*, the average age for studies using a sucking paradigm was 0.58 months ($SD = 0.89$), but 5.30 months ($SD = 1.78$) for studies using looking time paradigm. This age-related change in research paradigms could lead to a case of Simpson's paradox: the age-related trend within a single method might be lost when multiple methods are combined (Kievit, Frankenhuys, Waldorp, & Borsboom, 2013; Simpson, 1951).

Third, other methodological factors unrelated to age could also contribute to the lack of developmental effects. 22 of the 25 datasets included in the current analyses has a manuscript associated. Among the manuscripts, 8 identified that the meta-analytic effects were only robust in a subset of the studies. Some of the subsets were identified by certain methodological characters (e.g. in *Syntactic Bootstrapping*, the effect was only present in studies with transitive conditions, Cao & Lewis, 2022), and other subsets were identified by participants characteristics (e.g. in *Familiar word recognition*, the effect was stronger in infants whose primary language exposure was from Romance languages, Carbajal, Peperkamp, & Tsuji, 2021). Perhaps the apparent lack of developmental effects in the current analysis could be attributed to a complex interaction between methodological factors and participant characteristics, rather than a true absence of developmental changes.

Fourth, developmental change in infancy and early childhood might be distinct from one another. Bergelson (2020) has speculated that word comprehension in the looking-while-listening paradigm only shows significant developmental changes after 12 months of age, with infants younger than 12 months showing mostly flat developmental trajectories in this task. This contrast could be attributed to the fact that older infants are not only more experienced compared to younger infants, but also better learners who can more effectively take advantage of the input they receive. Could this pattern generalize to other tasks and domains? There is much evidence suggesting that developmental changes occurring in one domain would have cumulative, cascading effects on changes in other domains (Ahmed, Kuhfeld, Watts, Davis-Kean, & Vandell, 2021; Bornstein, Hahn, Putnick, & Pearson, 2018; Oakes & Rakison, 2019). The outcome of such developmental cascades might not be measurable in the experimental tasks included in the meta-analyses until infants are above 12 months of age.

We investigate each of these explanations in turn, assessing empirical support in our data. We summarise the results of these analyses in Table 3; in brief, no explanation provided traction for more than a small number of datasets.

Age-related publication bias

We first consider whether age-related selection bias can explain the lack of developmental changes in our datasets. If studies with younger infants suffered from publication bias more, then their effect sizes would be more inflated, obscuring possible developmental changes.

Methods. There are many methods to detect publication bias. One of the most common approaches is Egger’s test (Egger, Smith, Schneider, & Minder, 1997), which examines the relationship between the studies’ effect sizes and their precision. A significant result from Egger’s test indicated an asymmetry in the funnel plot, suggesting the presence of publication bias. This method is more sensitive than the rank correlation approach,

another common publication bias detection method (Begg & Mazumdar, 1994). However, Egger’s test can not accommodate predictors other than the study’s precision. As a result, we also turned to the weight-function model developed by Vevea and Hedges (1995). This method detected publication bias by likelihood ratio tests: a bias-corrected model is pitted against the original model to see if the former provides a better fit than the latter. A positive result indicates the presence of publication bias.

To detect age-related publication bias, we split each dataset by the median of the average participant age associated with each effect size (in months). We then run both Egger’s test and the weight-function model on each half of the dataset. We compared the test outcomes from both tests across the two halves of the datasets. For Egger’s test, we used the `regtest` function implemented in `metafor` (Viechtbauer, 2010). For the weight-function model, we used the package `weightr` (Kathleen M. Coburn & Vevea, 2019) and specified random-effect meta-regression models predicting effect sizes with mean age in months.

Results and discussion. Egger’s test was run on all but the 4 datasets in which either half of the datasets contained less than 20 effect sizes. Previous study has shown that Egger’s test has reduced sensitivity in datasets with less than 20 studies (Sterne, Egger, & Smith, 2001). For similar reasons, 7 datasets were excluded in the weight-function analysis.

Egger’s test suggested that in 3 out of 25 datasets, there was evidence for publication bias in the younger half but not in the older half (*Audio-Visual Congruence*, *Categorization bias*, *Syntactic bootstrapping*). However, this result was not corroborated by the weight-function analysis. For these three datasets, the weight function analysis did not find evidence for publication bias in either half of the three datasets. This suggests that the significant results found by Egger’s test might be due to factors other than publication bias. The weight-function analysis did find evidence for publication bias in the younger half but not the older half in two datasets: *Mutual exclusivity* (Younger: $\chi^2 = 11.07$, $p < 0.01$; Older: $\chi^2 = 0.02$, $p = 0.89$) and *Vowel discrimination (non-native)* (Younger: $\chi^2 =$

5.18, $p = 0.02$; Older: $\chi^2 = 1.88$, $p = 0.17$). These two datasets yielded significant results for both halves in Egger's test.

Overall, we found little evidence for more severe publication bias among the younger infants. The Egger's test and the function-weight analysis did not yield converging evidence, suggesting that factors other publication bias may be at play in contributing to the results. Interestingly, out of the five datasets that yield significant results for the younger participants, 2 of which were datasets that originally showed significant age-related changes (*Categorization bias*: $\beta = 0, -0.25, 0$, and 0.16 , $SE = 0, 0.34, 0$, and 0.42 , $z = -0.19, -0.72, 0.04$, and 0.38 , $p < 0.01$; *Mutual exclusivity*: $\beta = 0.04, 1.63, 0$, and 1.27 , $SE = 0.01, 0.25, 0$, and 0.15 , $z = 6.01, 6.58, 4.72$, and 8.63 , $p < 0.01$), which was in contrast with the other three datasets in which the age estimates were trending at the negative direction (*Audio-Visual Congruence*: $\beta = -0.02, -0.14, 0$, and 0.33 , $p = 0.39, 0.01, 0.84$, and 0 ; *Syntactic bootstrapping*: $\beta = -0.01, -0.40, 0$, and 0.24 , $p = 0.13, 0.11, 0.16$, and 0.02 ; *Vowel discrimination (non-native)*: $\beta = -0.01, -0.11, 0$, and 0.65 , $p = 0.45, 0.43, 0.55$, and 0). Taken together, we found limited evidence that selective publication bias explains the lack of age-related change across the board.

Methodological adaptation for older infants

In experiments with young children, many design decisions are made to ensure the paradigms are age appropriate (Byers-Heinlein, Bergmann, & Savalei, 2022). For older children, more behavioral measures are available and longer experiments are made possible by increased attention span. As a result, experimenters might test more subtle experimental contrasts. Perhaps the increasing difficulty or subtlety of experimental conditions for older infants mask age-related increase in effect sizes related to a particular construct. For example, imagine that different experimenters wanted to study word learning with 12- and 24-month-olds. The experimenter working with the younger group might choose a paradigm in which only two novel words were taught, while the

experimenter working with the older children might choose to teach four. The resulting effect for older children might be weaker despite overall improvement in the underlying construct.

The accessibility of different methods could also potentially cause an instance of Simpson’s paradox (Kievit et al., 2013). Imagine there were two methods, method A and method B, with the former having lower task demands than the latter. Due to its low task demands, method A would be more likely to be used on younger infants and causes larger effect sizes. In contrast, method B would be more likely to be used on older infants and results in smaller effect sizes. Although the age trend could be positive within each method, when pooling across studies from the two methods, the trend would then be negative, canceling out age-related changes patterns.

Since it is difficult to code for task demands across all studies, we explore whether methodological adaptation influences the developmental trend from the other side: instead of looking at method adaptation with age, we focus on studies using identical methods to test multiple age groups. This subset of datasets should provide the best chance of detecting age-related changes in the absence of methodological variation.

Methods

We first needed to identify the subset of studies in each dataset that satisfy the following two criteria: (1) the same paper tested multiple age groups, and (2) the multiple age groups were all tested using the same experimental design and measure. The first criterion was operationalized as having a paper with multiple age groups with an age difference greater than one month. The second criterion was operationalized based on methodological moderators coded by the original authors and available in MetaLab.

Within the effects selected for each dataset, we calculated Δ_{age} for each effect size. Δ_{age} was the difference between the age associated with a particular effect size and the

minimum age in each subset of the dataset.

19 datasets had subsets of studies fitting our criteria. We focused on the 15 subsets that having 10 and more effect sizes. For each subset, we applied a multilevel meta-regression model using the same nested random intercept as previously described. The model predicts effect sizes based on Δ_{age} . This analysis follows the logic that, if on average there is a greater effect size when the same experiment is conducted with older children relative to younger children, then the relation of effect size to Δ_{age} should be positive.

Results and discussion

We found no significant relationship between Δ_{age} and the effect sizes in any of the dataset (all $p > 0.05$).

This analysis was necessarily constrained by the granularities of the coded moderators. The number of coded methodological moderators ranged from 1 to 9, which means that the experimental design needs to be reduced into at maximum 9 dimensions. However, even at 9 dimensions, it is possible that elements of experiment design influencing task demands were overlooked. For instance, in many domains that use visual stimuli, the particular choice of visual stimuli might significantly vary in complexity (e.g. Cao & Lewis, 2022). Visual complexity has long been proposed as a key factor influencing the task demands (Hunter & Ames, 1988), but stimulus complexity was not coded in any of our meta-analyses. In conclusion, the findings presented here should be interpreted with caution due to potential limitations in the coding of methodological moderators.

Theoretical constraints on effect sizes

Across the 25 datasets, 22 datasets were published through manuscripts in peer-reviewed journals. Among these manuscripts, we found that 8 papers reported that

the meta-analytic effect was significantly stronger in a subset of the data. The subset was often identified by a particular condition in the experimental paradigm (e.g. experiment that shows “giving and taking action” to infants, Margoni & Surian, 2018), or certain characteristics of the participants (e.g. bilingual infants, Tsui, Byers-Heinlein, & Fennell, 2019). In the rest of the data, the meta-analytic effect was either significantly weaker or not present at all. There are many reasons for why the effect would be stronger or only present in a subset of the data. Here, we remain agnostic to the underlying causes for these differences, and leverage these findings to ask: Is it possible that the influence of age was only observable in the subset of the dataset characterized by stronger effect sizes? Perhaps noise in other conditions inadvertently masked age-related changes.

Methods. We screened through 22 papers and identified 8 papers that reported a stronger effect on subsets of the data. All subsets had more than 10 effect sizes. For datasets reporting more than one subset as having strong effect, we consider each respectively. In sum, 7 datasets produced 9 subsets that showed stronger effects.

We first investigated whether we could reproduce the original patterns, i.e. the effect sizes in the better halves were indeed stronger than the other halves. We ran the same multilevel meta-regression without any predictor to estimate the meta-analytic effect sizes in each half. Then we ran a Wald test to compare the two estimates by running a fixed-effects meta-regression model predicting effect sizes with the moderator distinguishing the two halves. A significant estimate on the moderator indicates that the meta-analytic effect sizes in both halves are significantly different from one another. We then estimated the slope of the age predictor in a multilevel meta-regression model for each of the subsets with larger effect sizes.

Results and discussion. We did not fully replicate the original findings reported in the original papers: the “better half” identified by the original meta-analysis did not produce significantly stronger effects than the rest of the data in many datasets. We did observe a significantly stronger effect in the remaining 3 datasets: For *Prosocial Agents*,

there was a stronger effect in experimental paradigms showing infants giving-taking actions compared to the studies showing infants other stimuli (Margoni & Surian, 2018, $z = -2.47$, $p = 0.01$); For *Statistical Sound Category Learning*, stronger effect was observed in studies using habituation paradigm compared to other paradigms (Cristia, 2018, $z = -2.42$, $p = 0.02$), and for *Statistical word segmentation*, stronger effect was observed in studies labeled as the conceptual replication of the original work (Black & Bergmann, 2017, $z = 2.51$, $p = 0.01$).

In addition, we did not find constraining our analyses to the “better half” increased the number of significant slope estimates. The two significant slope estimates came from *Mutual Exclusivity* ($\beta = 0.04$, $SE = 0.01$, $z = 4.63$, $p < 0.01$) and *Statistical sound category learning* ($\beta = 0.11$, $SE = 0.05$, $z = 2.23$, $p = 0.03$), which also showed significant slopes in the analyses with the full datasets. Qualitatively, we did see that the estimates increased in magnitude in *Syntactic bootstrapping* ($\beta = 0.01$, $SE = 0.03$, $z = 0.43$, $p = 0.67$) and *Switch task* ($\beta = 0.01$, $SE = 0.03$, $z = 0.27$, $p = 0.79$), but neither reached the statistical significance threshold.

The discrepancy between our analyses and the previously reported finding suggested that the “better half effect” might not be sufficiently robust. This discrepancy could be attributed to the different statistical models we chose – in the original meta-analysis papers, the models tend to differ in their particular specification of the nested random effect structure and/or in the inclusions of moderators. We chose the simplest model with the maximum random effect structure per recommendation (Barr, Levy, Scheepers, & Tily, 2013). This approach ensured fair comparison across all datasets, but it could diminish the strength of the reported effects.

Interestingly, even in the datasets where the better half effect was reproduced, we failed to see a significant age effect in the same datasets (Prosocial agents and Statistical word segmentation) that did not show age-related changes in the original full dataset.

Altogether, this set of analysis suggested that the theoretical constraints on the effect sizes could not adequately explain the lack of age-related change.

Developmental change emerges later

Last but not least, we consider whether there is evidence for discontinuity between the growth patterns in infancy and beyond. Bergelson (2020)’s hypothesis on the development of word comprehension suggests a notable shift post the 12-month mark in infancy. This raises the question of whether such distinctions extend across various tasks. This section aims to delve into these dynamics by only looking at the subset of the dataset with infants older than 12-month-olds.

Methods. Similar to previous analyses, we filtered each dataset to include only studies that reported more than 10 effect sizes that tested infants older than 12 months. 15 datasets met the criteria. We ran the same meta-regressions predicting effect size with mean age in months on this subset, and then we compared the estimates on the age predictor with the same models run on the full datasets.

Results and discussion. If the discontinuity account is true, we should expect to see more significant age effects to emerge on models run on the subset of data with older infants. We found support for this hypothesis in two datasets, *Cross-situational Word Learning* ($\beta = 0.01$, $SE < 0.01$, $z = 2.71$, $p = 0.01$) and *Mispronunciation sensitivity* ($\beta = 0.07$, $SE = 0.01$, $z = 4.69$, $p < 0.01$). In both datasets, there were no age effects in the full datasets, but significant age-related change in the subsets with older infants. However, we also found the opposite patterns. In *Categorization bias* and *Sound symbolism*, there was evidence for age-related change across the entire age range, but no evidence for age-related change in the toddler subset (Both $p > 0.05$).

Genral discussion

How do infants' cognitive and linguistic abilities change with age? In this work, we leveraged a dataset of meta-analyses to evaluate the assumption that these abilities increase positively with age, and that the form of this increase is linear. There was no evidence for linear growth in 16 datasets, and interestingly, in all of these datasets, there was no evidence for any age-related growth at all. In the second section, we investigated four potential explanations for this pattern: (1) age-related selection bias against younger infants; (2) methodological adaptation for older infants; (3) change in only a subset of conditions; and (4) positive growth only after infancy.

Our current work has several limitations. First and foremost, we simply lacked sufficient data to investigate the possible explanations for many domains (see Table 3). In FIXME datasets, when we filtered datasets to answer the corresponding questions, we lacked sufficient data to adequately test our hypotheses. Furthermore, as with many meta-analyses, our datasets also had high heterogeneity, meaning that we can only explain relatively small amounts of the variation among effect sizes.

Our work highlights the importance of improving reporting standards in developmental psychology. Testing moderation of heterogeneity requires consistent coding of moderators across datasets. But surveys of reporting standards show that many potential moderators go unreported. For instance, fewer than half of papers report attrition rate (Nicholson et al., 2017; Raad, Bellinger, McCormick, Roberts, & Steele, 2007). Given these observations, there is a clear need for the developmental psychology community to create and embrace more rigorous and transparent reporting standards. The recently developed framework for reporting demographics information across cultures in developmental psychology is one promising direction moving forwards (Singh et al., 2023).

Moreover, learning from other fields could provide valuable insights into how to enhance these standards. In biomedical research, numerous reporting standards have been

published and widely adopted Moher et al. (2015). Following these structured guidelines in reporting could significantly increase both the quality and the quantity of information extractable from the original papers, providing more traction for tackling heterogeneity in meta-analysis.

Our work also underscores the importance of multi-laboratory large scale replication projects. The relationship between meta-analysis and multi-laboratory is complicated (Kvarven, Strömmland, & Johannesson, 2020; Lewis, Mathur, VanderWeele, & Frank, 2022). Although the latter approach is much more time- and resource- intensive than the former, it is also much more effective in controlling unwanted heterogeneity and detecting subtle patterns in the data. One prominent example is the comparison between the meta-analysis of Infant directed speech preference (Dunst, Gorman, & Hamby, 2012) and the ManyBabies 1 project on the same topic (Consortium, 2020). Zettersten et al. (2023) found that, after an update to the meta-analysis dataset, both datasets yielded comparable estimated effect sizes ($d = 0.35$), but the age effect was only detected in the MB1 project, not the meta-analysis. That study speculated that our second explanation (methodological variation covarying with age) might account for their studies. In our analysis, we did investigate the methodological adaptation hypothesis in the IDS preference dataset. However, the methodological moderators available for us were limited and we could not incorporate the varying nature of the stimuli into our analysis. This example shows the potential limitations of meta-analyses that rely on aggregated data from studies with varied methodologies. In contrast, multi-laboratory collaboration projects like Manybabies (Visser et al., 2022) can rely on standardized data collection procedure and stimuli, therefore providing a more controlled dataset to answer a specific research question with high power.

It is also worth considering whether the strengths of certain developmental phenomena truly stay constant throughout the first years of life. This counterintuitive possibility casts doubts on the construct validity of the existing measures. Many researchers strive to build on existing experimental procedures and measurements when

they are testing older participants. This then leads to a potentially problematic situation: an experimental paradigm could have high construct validity with participants of a certain age, but low construct validity with participants of different age. As a result, this leads to an interesting conundrum: methodological adaptation could be a source of significant heterogeneity, diminishing the measurable developmental change. But at the same time, paradoxically, it could also be the prerequisite for properly measuring developmental change. An alternative explanation for the true absence of developmental change is the limited sensitivity to change in experimental studies relying on group average for comparison. The group average may stay constant, but there could still be growth in an individual's performance across development (Bornstein, Putnick, & Esposito, 2017). The nuanced nature of developmental change might be best captured by dense, longitudinal data of individual child (e.g. Bergelson et al., 2023; Sullivan, Mei, Perfors, Wojcik, & Frank, 2021).

In sum, our current work presents a surprising finding concerning age-related change in the cognitive and language development literatures in early childhood. Despite decades of research built upon the positive increase and linearity assumptions, we failed to find evidence supporting either in most meta-analyses that we had access to. Our work is not intended to overturn the longstanding developmental theories. Like other researchers, we believe that infants get better across different cognitive and linguistic domains as they get older. Instead, our work aims to highlight the needs for more robust reporting standards and more large-scale multi-laboratory projects that measure children consistently across age groups and over time. Our findings invite the cognitive development community to strengthen our understanding of foundational assumptions via collaborative efforts.

References

- Adolph, K. E., Robinson, S. R., Young, J. W., & Gill-Alvarez, F. (2008). What is the shape of developmental change? *Psychological Review*, 115(3), 527.
- Ahmed, S. F., Kuhfeld, M., Watts, T. W., Davis-Kean, P. E., & Vandell, D. L. (2021). Preschool executive function and adult outcomes: A developmental cascade model. *Developmental Psychology*, 57(12), 2234.
- Altman, D. G., Simera, I., Hoey, J., Moher, D., & Schulz, K. (2008). EQUATOR: Reporting guidelines for health research. *The Lancet*, 371(9619), 1149–1150.
- Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, 68(3), 255–278.
- Begg, C. B., & Mazumdar, M. (1994). Operating characteristics of a rank correlation test for publication bias. *Biometrics*, 1088–1101.
- Bergelson, E. (2020). The comprehension boost in early word learning: Older infants are better learners. *Child Development Perspectives*, 14(3), 142–149.
- Bergelson, E., Soderstrom, M., Schwarz, I.-C., Rowland, C. F., Ramirez-Esparza, N., Hamrick, L. R., et al.others. (2023). *Everyday language input and production in 1001 children from 6 continents*.
- Bergmann, C., Tsuji, S., Piccinini, P. E., Lewis, M. L., Braginsky, M., Frank, M. C., & Cristia, A. (2018). Promoting replicability in developmental research through meta-analyses: Insights from language acquisition research. *Child Development*, 89(6), 1996–2009.
- Best, R., & Charness, N. (2015). Age differences in the effect of framing on risky choice: A meta-analysis. *Psychology and Aging*, 30(3), 688.
- Black, A., & Bergmann, C. (2017). Quantifying infants’ statistical word segmentation: A meta-analysis. *39th Annual Meeting of the Cognitive Science Society*, 124–129. Cognitive Science Society.

- Borenstein, M., Hedges, L. V., Higgins, J. P., & Rothstein, H. R. (2021). *Introduction to meta-analysis*. John Wiley & Sons.
- Bornstein, M. H., Hahn, C.-S., Putnick, D. L., & Pearson, R. M. (2018). Stability of core language skill from infancy to adolescence in typical and atypical development. *Science Advances*, 4(11), eaat7422.
- Bornstein, M. H., Putnick, D. L., & Esposito, G. (2017). Continuity and stability in development. *Child Development Perspectives*, 11(2), 113–119.
- Bronfenbrenner, U. (1977). Toward an experimental ecology of human development. *American Psychologist*, 32(7), 513.
- Burnham, K. P., & Anderson, D. R. (2004). Multimodel inference: Understanding AIC and BIC in model selection. *Sociological Methods & Research*, 33(2), 261–304.
- Byers-Heinlein, K., Bergmann, C., & Savalei, V. (2022). Six solutions for more reliable infant research. *Infant and Child Development*, 31(5), e2296.
- Cao, A., & Lewis, M. (2022). Quantifying the syntactic bootstrapping effect in verb learning: A meta-analytic synthesis. *Developmental Science*, 25(2), e13176.
- Cao, A., Lewis, M., & Frank, M. C. (2023). A synthesis of early cognitive and language development using (meta-) meta-analysis. *Proceedings of the Annual Meeting of the Cognitive Science Society*, 45.
- Carbajal, M. J., Peperkamp, S., & Tsuji, S. (2021). A meta-analysis of infants' word-form recognition. *Infancy*, 26(3), 369–387.
- Carey, S. (2009). *The origin of concepts*. Oxford University Press.
- Coburn, Kathleen M., & Vevea, J. L. (2015). Publication bias as a function of study characteristics. *Psychological Methods*, 20(3), 310.
- Coburn, Kathleen M., & Vevea, J. L. (2019). *Weightr: Estimating weight-function models for publication bias*.
- Cole, P. M., Loughheed, J. P., Chow, S.-M., & Ram, N. (2020). Development of emotion regulation dynamics across early childhood: A multiple time-scale approach. *Affective*

569 *Science*, 1, 28–41.

570 Consortium, M. (2020). Quantifying sources of variability in infancy research using the
571 infant-directed-speech preference. *Advances in Methods and Practices in Psychological*
572 *Science*, 3(1), 24–52.

573 Cristia, A. (2018). Can infants learn phonology in the lab? A meta-analytic answer.
574 *Cognition*, 170, 312–327.

575 Doebel, S., & Zelazo, P. D. (2015). A meta-analysis of the dimensional change card sort:
576 Implications for developmental theories and the measurement of executive function in
577 children. *Developmental Review*, 38, 241–268.

578 Dunst, C., Gorman, E., & Hamby, D. (2012). Preference for infant-directed speech in
579 preverbal young children. *Center for Early Literacy Learning*, 5(1), 1–13.

580 Egger, M., Smith, G. D., & Phillips, A. N. (1997). Meta-analysis: Principles and
581 procedures. *Bmj*, 315(7121), 1533–1537.

582 Egger, M., Smith, G. D., Schneider, M., & Minder, C. (1997). Bias in meta-analysis
583 detected by a simple, graphical test. *Bmj*, 315(7109), 629–634.

584 Elman, J. L. (1996). *Rethinking innateness: A connectionist perspective on development*
585 (Vol. 10). MIT press.

586 Ferguson, C. J., & Brannick, M. T. (2012). Publication bias in psychological science:
587 Prevalence, methods for identifying and controlling, and implications for the use of
588 meta-analyses. *Psychological Methods*, 17(1), 120.

589 Ferguson, C. J., & Heene, M. (2012). A vast graveyard of undead theories: Publication
590 bias and psychological science's aversion to the null. *Perspectives on Psychological*
591 *Science*, 7(6), 555–561.

592 Flavell, J. H. (1994). *Cognitive development: Past, present, and future*.

593 Fletcher, J. (2007). What is heterogeneity and is it important? *Bmj*, 334(7584), 94–96.

594 Francis, G. (2012). Publication bias and the failure of replication in experimental
595 psychology. *Psychonomic Bulletin & Review*, 19, 975–991.

- Frank, M. C., Bergelson, E., Bergmann, C., Cristia, A., Floccia, C., Gervain, J., et al.others. (2017). A collaborative approach to infant research: Promoting reproducibility, best practices, and theory-building. *Infancy*, 22(4), 421–435.
- Frank, M. C., Braginsky, M., Yurovsky, D., & Marchman, V. A. (2021). *Variability and consistency in early language learning: The wordbank project*. MIT Press.
- Gasparini, L., Langus, A., Tsuji, S., & Boll-Avetisyan, N. (2021). Quantifying the role of rhythm in infants' language discrimination abilities: A meta-analysis. *Cognition*, 213, 104757.
- Higgins, J. P., & Thompson, S. G. (2002). Quantifying heterogeneity in a meta-analysis. *Statistics in Medicine*, 21(11), 1539–1558.
- Huedo-Medina, T. B., Sánchez-Meca, J., Marín-Martínez, F., & Botella, J. (2006). Assessing heterogeneity in meta-analysis: Q statistic or i^2 index? *Psychological Methods*, 11(2), 193.
- Hunter, M. A., & Ames, E. W. (1988). A multifactor model of infant preferences for novel and familiar stimuli. *Advances in Infancy Research*.
- Hyde, J. S. (1984). How large are gender differences in aggression? A developmental meta-analysis. *Developmental Psychology*, 20(4), 722.
- Johansson, M., Marciszko, C., Brocki, K., & Bohlin, G. (2016). Individual differences in early executive functions: A longitudinal study from 12 to 36 months. *Infant and Child Development*, 25(6), 533–549.
- Karlberg, J., Engström, I., Karlberg, P., & Fryer, J. G. (1987). Analysis of linear growth using a mathematical model: I. From birth to three years. *Acta Paediatrica*, 76(3), 478–488.
- Kievit, R., Frankenhuis, W. E., Waldorp, L., & Borsboom, D. (2013). Simpson's paradox in psychological science: A practical guide. *Frontiers in Psychology*, 4, 54928.
- Kvarven, A., Strømmland, E., & Johannesson, M. (2020). Comparing meta-analyses and preregistered multiple-laboratory replication projects. *Nature Human Behaviour*, 4(4),

423–434.

Letourneau, N. L., Duffett-Leger, L., Levac, L., Watson, B., & Young-Morris, C. (2013).

Socioeconomic status and child development: A meta-analysis. *Journal of Emotional and Behavioral Disorders*, 21(3), 211–224.

Lewis, M., Braginsky, M., Tsuji, S., Bergmann, C., Piccinini, P. E., Cristia, A., et al.

(2016). *A quantitative synthesis of early language acquisition using meta-analysis*.

Lewis, M., Mathur, M. B., VanderWeele, T. J., & Frank, M. C. (2022). The puzzling relationship between multi-laboratory replications and meta-analyses of the published literature. *Royal Society Open Science*, 9(2), 211499.

Lindenberger, U., & Pötter, U. (1998). The complex nature of unique and shared effects in hierarchical linear regression: Implications for developmental psychology. *Psychological Methods*, 3(2), 218.

Margoni, F., & Surian, L. (2018). Infants' evaluation of prosocial and antisocial agents: A meta-analysis. *Developmental Psychology*, 54(8), 1445.

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

McArdle, J. J., Grimm, K. J., Hamagami, F., Bowles, R. P., & Meredith, W. (2009).

Modeling life-span growth curves of cognition using longitudinal data with multiple samples and changing scales of measurement. *Psychological Methods*, 14(2), 126.

McCartney, K., Harris, M. J., & Bernieri, F. (1990). Growing up and growing apart: A developmental meta-analysis of twin studies. *Psychological Bulletin*, 107(2), 226.

Moher, D., Shamseer, L., Clarke, M., Ghersi, D., Liberati, A., Petticrew, M., ... Group, P.-P. (2015). Preferred reporting items for systematic review and meta-analysis protocols (PRISMA-p) 2015 statement. *Systematic Reviews*, 4, 1–9.

Naigles, L. (1990). Children use syntax to learn verb meanings. *Journal of Child Language*, 17(2), 357–374.

Nicholson, J. S., Deboeck, P. R., & Howard, W. (2017). Attrition in developmental

psychology: A review of modern missing data reporting and practices. *International Journal of Behavioral Development*, 41(1), 143–153.

Oakes, L. M., & Rakison, D. H. (2019). *Developmental cascades: Building the infant mind*. Oxford University Press.

Piaget, J. (1971). *The theory of stages in cognitive development*.

Publications, A., Journal, C. B. W. G. on, & Standards, A. R. (2008). Reporting standards for research in psychology: Why do we need them? What might they be? *The American Psychologist*, 63(9), 839.

Raad, J. M., Bellinger, S., McCormick, E., Roberts, M. C., & Steele, R. G. (2007). Brief report: Reporting practices of methodological information in four journals of pediatric and child psychology. *Journal of Pediatric Psychology*, 33(7), 688–693.

Saffran, J. R., Aslin, R. N., & Newport, E. L. (1996). Statistical learning by 8-month-old infants. *Science*, 274(5294), 1926–1928.

Schulz, K. F., Altman, D. G., & Moher, D. (2010). CONSORT 2010 statement: Updated guidelines for reporting parallel group randomised trials. *Journal of Pharmacology and Pharmacotherapeutics*, 1(2), 100–107.

Simonsohn, U., Simmons, J., & Nelson, L. D. (2022). Above averaging in literature reviews. *Nature Reviews Psychology*, 1(10), 551–552.

Simpson, E. H. (1951). The interpretation of interaction in contingency tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 13(2), 238–241.

Singh, L., Barokova, M. D., Baumgartner, H. A., Lopera-Perez, D. C., Omane, P. O., Sheskin, M., et al.others. (2023). A unified approach to demographic data collection for research with young children across diverse cultures. *Developmental Psychology*.

Sterne, J. A., Egger, M., & Smith, G. D. (2001). Investigating and dealing with publication and other biases in meta-analysis. *Bmj*, 323(7304), 101–105.

Sugden, N. A., & Marquis, A. R. (2017). Meta-analytic review of the development of face discrimination in infancy: Face race, face gender, infant age, and methodology

677 moderate face discrimination. *Psychological Bulletin*, 143(11), 1201.

678 Sullivan, J., Mei, M., Perfors, A., Wojcik, E., & Frank, M. C. (2021). SAYCam: A large,
679 longitudinal audiovisual dataset recorded from the infant's perspective. *Open Mind*, 5,
680 20–29.

681 Thelen, E., & Smith, L. B. (2007). Dynamic systems theories. *Handbook of Child*
682 *Psychology*, 1.

683 Thompson, S. G., & Sharp, S. J. (1999). Explaining heterogeneity in meta-analysis: A
684 comparison of methods. *Statistics in Medicine*, 18(20), 2693–2708.

685 Thornton, A., & Lee, P. (2000). Publication bias in meta-analysis: Its causes and
686 consequences. *Journal of Clinical Epidemiology*, 53(2), 207–216.

687 Tilling, K., Macdonald-Wallis, C., Lawlor, D. A., Hughes, R. A., & Howe, L. D. (2014).
688 Modelling childhood growth using fractional polynomials and linear splines. *Annals of*
689 *Nutrition and Metabolism*, 65(2-3), 129–138.

690 Tsui, A. S. M., Byers-Heinlein, K., & Fennell, C. T. (2019). Associative word learning in
691 infancy: A meta-analysis of the switch task. *Developmental Psychology*, 55(5), 934.

692 Vandenbroucke, J. P., Elm, E. von, Altman, D. G., Gøtzsche, P. C., Mulrow, C. D.,
693 Pocock, S. J., ... Initiative, S. (2007). Strengthening the reporting of observational
694 studies in epidemiology (STROBE): Explanation and elaboration. *Annals of Internal*
695 *Medicine*, 147(8), W–163.

696 Vevea, J. L., & Hedges, L. V. (1995). A general linear model for estimating effect size in
697 the presence of publication bias. *Psychometrika*, 60, 419–435.

698 Viechtbauer, W. (2010). Conducting meta-analyses in r with the metafor package. *Journal*
699 *of Statistical Software*, 36(3), 1–48.

700 Visser, I., Bergmann, C., Byers-Heinlein, K., Dal Ben, R., Duch, W., Forbes, S., et
701 al.others. (2022). Improving the generalizability of infant psychological research: The
702 ManyBabies model. *Behavioral and Brain Sciences*, 45.

703 Zettersten, M., Cox, C. M. M., Bergmann, C., Tsui, A., Soderstrom, M., Mayor, J., et

704 al.others. (2023). *Evidence for infant-directed speech preference is consistent across*
705 *large-scale, multi-site replication and meta-analysis.*