

Baguley, Thom; Cahoon, Abbie; Lazareanu, Daniela; Thives Mello, Arthur; Zaneva, Mirela

Working Paper

A comment on "Why Do Children Think Words Are Mutually Exclusive?" (Brody et al., 2024)

I4R Discussion Paper Series, No. 273

Provided in Cooperation with:

The Institute for Replication (I4R)

Suggested Citation: Baguley, Thom; Cahoon, Abbie; Lazareanu, Daniela; Thives Mello, Arthur; Zaneva, Mirela (2025) : A comment on "Why Do Children Think Words Are Mutually Exclusive?" (Brody et al., 2024), I4R Discussion Paper Series, No. 273, Institute for Replication (I4R), s.l.

This Version is available at:

<https://hdl.handle.net/10419/333868>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



No. 273

I4R DISCUSSION PAPER SERIES

A comment on “Why Do Children Think Words Are Mutually Exclusive?” (Brody et al., 2024)

Thom Baguley

Abbie Cahoon

Daniela Lazareanu

Arthur Thives Mello

Mirela Zaneva

December 2025

I4R DISCUSSION PAPER SERIES

I4R DP No. 273

A comment on “Why Do Children Think Words Are Mutually Exclusive?” (Brody et al., 2024)

**Thom Baguley¹, Abbie Cahoon², Daniela Lazareanu³,
Arthur Thives Mello⁴, Mirela Zaneva⁵**

¹Nottingham Trent University, Nottingham/Great Britain

²Ulster University, Belfast/Great Britain

³University of Reading/Great Britain

⁴Brunel University of London/Great Britain

⁵University of Oxford/Great Britain

DECEMBER 2025

Any opinions in this paper are those of the author(s) and not those of the Institute for Replication (I4R). Research published in this series may include views on policy, but I4R takes no institutional policy positions.

I4R Discussion Papers are research papers of the Institute for Replication which are widely circulated to promote replications and meta-scientific work in the social sciences. Provided in cooperation with EconStor, a service of the [ZBW – Leibniz Information Centre for Economics](#), and [RWI – Leibniz Institute for Economic Research](#), I4R Discussion Papers are among others listed in RePEc (see IDEAS, EconPapers). Complete list of all I4R DPs - downloadable for free at the I4R website.

I4R Discussion Papers often represent preliminary work and are circulated to encourage discussion. Citation of such a paper should account for its provisional character. A revised version may be available directly from the author.

Editors

Abel Brodeur
University of Ottawa

Anna Dreber
Stockholm School of Economics

Jörg Ankel-Peters
RWI – Leibniz Institute for Economic Research

A comment on “Why Do Children Think Words Are Mutually Exclusive?” (Brody et al., 2024)

Baguley, Thom (Nottingham Trent University; thomas.baguley@ntu.ac.uk)

Cahoon, Abbie (Ulster University; a.cahoon@ulster.ac.uk)

Lazareanu, Daniela (University of Reading, Henley Business School;
d.m.lazareanu@pgr.reading.ac.uk; dani.lazareanu@gmail.com)

Thives Mello, Arthur (Brunel University of London; Arthur.ThivesMello@brunel.ac.uk)

Zaneva, Mirela (Christ Church, University of Oxford; mirela.zaneva@chch.ox.ac.uk)

(Authors ordered in alphabetical order by last name)

The authors report no conflicts of interest. No financial support was received for this work.

Abstract

An early emerging understanding of information structure in spoken language is demonstrated in Brody et al.’s (2024) study, contrary to the general assumption that young children innately assume new and familiar words to be mutually exclusive. Across 3 experiments, learners up to 2 years of age (N = 106) showed mutual exclusivity if the novel words were spoken with focus, an information-structural marker of contrast. We successfully computationally reproduced the reported results, ran additional robustness tests, and considered potential covariates that were recorded but not used in the original paper (i.e. age and gender). We also assess potentially disproportionate effects of influential participants via Cook’s distance, and refit a model where their data are removed. Further, we consider alternative specifications of random effects (e.g., maximal models) to address convergence issues and obtain interval estimates of experimental conditions’ means that account for random effects. Finally, we examine whether different optimizers address model convergence issues. Across our replication tests, we find that the original results are robust and key estimates remain significant. Of note, our influence analysis and Bayesian analysis showed stronger effects than those originally reported, likely attributable to inattentive/noisy responding to some trials. Overall, our replication efforts should provide increased confidence in the original effects.

1. Introduction

Brody, Feiman and Aravind (2024) question the traditional views of mutual exclusivity in early word learning, which take innate inductive bias as the main driver. Instead, the authors examine young children's understanding of information structure (prosody) and the way this is transmitted (prosodic emphasis or focus). The authors focus on the contrastive information naturally contained within spoken language, testing whether this evokes mutual exclusivity or not. They use a variation of the standard mutual exclusivity paradigm involving a novel object choice task, whereby two novel objects are introduced by a central animated character with *focus* versus with *givenness*. Across 3 Experiments, results show that children ($N = 106$ from across the United States, age range: 24-35 months) only exhibited a mutual exclusivity bias when the novel noun was focused. When it was given (not emphasized), children tended to interpret it (point to) as referring to a previously labeled object. Although prior theories tie mutual exclusivity to word novelty, the present account provides evidence for the important role of information-structural markers, rather than lexical novelty itself, in disambiguating meanings during early language acquisition.

In the present report, collaboratively prepared as part of the Replication Games, an initiative supported by the Institute for Replication and Psychological Science, we investigate whether Brody and colleagues' results are computationally reproducible and carry out a number of robustness tests:

- 1) considering the addition of age and gender as covariates,
- 2) assessing sample sensitivity to influential participants and refitting a model without such participants,
- 3) using different optimization algorithms,
- 4) re-expressing the categorical predictors under different contrast schemes
- 5) considering different specifications of the random structures (e.g. maximal models), comparing point and interval estimates between those reported in the paper (derived from raw data), those from the models used for

significance tests in the paper (but not used to derive estimates) and Bayesian models with the appropriate maximal random effects.

We successfully reproduce the results and note only a minor deviation from optimal practices, where two different models were stored under the same name (“m3”). We would like to acknowledge that the original study was first reproduced by the Psychological Science STAR team, who confirmed the computational reproducibility of the results.

Turning to our sensitivity analyses (robustness reproduction), we find that including additional covariates (age, gender) in the models does not change the pattern of significance for the main results. The point estimate for *condition* changes slightly (original $\beta = -12.712$, $se = 4.490$, $p = 0.005$ vs $\beta = -11.850$, $se = 4.934$, $p = 0.0163$ for model with age, vs $\beta = -10.296$, $se = 3.356$, $p = .0022$ for model with gender, vs $\beta = -10.188$, $se = 3.285$, $p = 0.0019$ for model with age + gender; see Table 2). Neither age nor gender were significant factors in our analyses. Next, we used Cook’s distance to examine whether there were any influential participants. 6 participants (associated with 33 trials) were influential (according to a common rule-of-thumb threshold of $4/n$), thus we removed them in a follow-up test. The key estimate for condition was robust, remained significant, and became stronger (original model $\beta = -12.712$, $se = 4.490$, $p = 0.005$ vs model with exclusion of influential participants $\beta = -18.864$, $se = 3.406$, $p < 0.0001$).

In a separate set of tests, we refitted the models using multiple optimizers, including BOBYQA (original), Nelder–Mead, NLMINB, and two NLOpt-based variants (NLOPT-BOBYQA and NLOPT-Nelder–Mead). Across all optimizers, estimates remained unchanged among the models that converged. We further examined robustness to contrast specification by applying reverse-treatment and sum coding, confirming the results are robust to contrast specification. Finally, point and interval estimates for the percentage of new objects selected (the key outcome measure) were derived from raw data whereas significance tests were obtained from a mixed effects generalized linear model with a binomial link.

2. Computational Reproducibility

We summarize the computational reproducibility of the study in regards to the provided data and code in Table 1. The authors provided data in the form of 3 csv files (one for each experiment), alongside a Rmd file with code and the generated html output. The Rmd file contained the complete code, including data wrangling, analysis, and figure generation. The Rmd file ran without errors and produced the same output file as shared by the authors. We checked this against the paper, and the results matched. In other words, we successfully reproduced Experiments 1, 2, 3 (all reported tests). We note that the data provided contains only the included participants, though the authors transparently describe the numbers of excluded participants with reasons in text. The authors have also labeled their data files as ‘cleandata’ (e.g. foxy_exp1_cleandata_wide2.csv). This is why we treat this as ‘cleaned’ and not ‘raw’ data. Similarly, although the provided code includes relevant data wrangling (e.g. wide to long format transitions), it does not contain all relevant cleaning code – for instance, it does not check for missing trials or other exclusion criteria.

	Fully	Partial	No	Notes
Raw data provided			x	Cleaned data provided. Sufficient for replication of results.
Cleaning code provided		x		Code contains some relevant data wrangling but not all (e.g. exclusions)
Analysis data provided	x			
Analysis code provided	x			
Reproducible from raw data			x	Raw data not provided
Reproducible from analysis data	x			

Table 1. Summary table for the provided data and code with regard to reproducibility.

As a minor note, two different models were saved as 'm3' in the original code. This was not an issue in Markdown (which prints outputs in sequence, thus model outputs are visible clearly for each model) but could be an issue if someone later wanted to call the earlier model, which had been overwritten with the later instance of 'm3'. We recommend that the second model is saved with a different name, e.g. 'm3_e3' for anyone reusing these analyses.

2.1 Discrepancies Between Pre-analysis Plan and Article

The authors have pre-registered each of the three experiments, including pre-analysis plan and description of the data collection, procedure, statistical models, inference criteria, and participant exclusion criteria.

We find no deviations from the pre-analysis plan. We note that there was a non-completion criterion, used for potential exclusion of participants, that was defined as “Does not complete the majority of the experimental items (at least 3/6 usable trials for inclusion).” We did not see code for this included in the shared Rmd file. Nevertheless, data shared included only participants with at least 3/6 trials. In other words, we believe the authors likely checked for this and excluded participants earlier.

We also note in passing that the authors preregistered a significance criterion of “ $p = 0.05$ ”, which we take to mean $p < 0.05$ or perhaps even ≤ 0.05 . This criterion is inconsistently applied in the paper in one instance where the result is described as “*albeit only marginally significant* ($\beta = 4.06, p = .077$)” (Brody et al., 2024, p. 1321).

3. Robustness Reproductions

In this section, we report a series of robustness tests, examining: 1) additional covariates, 2) sample sensitivity, 3) different optimizers, 4) contrast coding, 5) different specifications of the random structures (e.g. maximal models).

3.1 Robustness test: Additional covariates (age and gender)

Here, we focus on the mixed effect logistic regression model for Experiment 1 (with glmer package in R) and consider the addition of covariates. The data contains variables for

participants' age in months (variable: "AgeMonths") and gender (variable: "Gender"). Given the potential rapid and heterogenous developmental changes in early childhood (Lewis et al., 2020; Bergelson & Swingley, 2012), we consider it might be sensible to include age as a covariate. We did not anticipate gender differences in very young children but consider the addition of this variable as a robustness test.

As a note on data cleaning, we found that some values for 'Gender' had an additional space. We ran a simple check for unique values in R: `unique(exp1.wide$Gender)`. We noticed the output showed two different entries for male participants: [1] "F" "M" "M ". We cleaned this by treating values as "M" and "M " as both "M" (i.e., both referring to male participants).

We first examine one of the focal results for Experiment 1, namely for the rate of mutual exclusivity inferences in the focus group: $\beta = -12.712$, $p = 0.005$ (p. 1320 of the original manuscript). In the authors' code notation, this model essentially asks whether two conditions (focus vs givenness) differ from each other. Condition is dummy-coded with the focus condition as the reference level.

In Table 2, we report three further robustness models. These are analogous to the original (base) model but with added covariates: 1) base model + age, 2) base model + gender, 3) base model + age + gender. The estimates for age and gender are not significant in any of these models. The estimate for Condition remains significant in all models, with relatively small changes to the exact estimate value (original $\beta = -12.712$ vs $\beta = -11.850$ for model with age, vs $\beta = -10.296$ for model with gender, vs $\beta = -10.188$ for model with age + gender). In other words, we find this estimate to be robust to the inclusion of covariates.

Model	Intercept (<i>SE</i>) <i>p-value</i>	Condition (<i>SE</i>) <i>p-value</i>	Age (<i>SE</i>) <i>p-value</i>	Gender (<i>SE</i>) <i>p-value</i>
Base (Original)	7.289*** (2.059) <i>p</i> = .0004	-12.712** (4.490) <i>p</i> = .005	-	-

Base + Age	1.418 (9.734) <i>p</i> = .8842	-11.850* (4.934) <i>p</i> = .0163	0.183 (0.300) <i>p</i> = .5425	-
Base + Gender	5.202** (1.967) <i>p</i> = .0082	-10.296** (3.356) <i>p</i> = .0022	-	2.342 (1.713) <i>p</i> = .1717
Base + Age + Gender	0.243* (8.7402) <i>p</i> = .9778	-10.188* (3.285) <i>p</i> = .0019	0.165 (0.289) <i>p</i> = .5686	2.233 (1.703) <i>p</i> = .1897

Table 2. Robustness test considering the addition of covariates. We report the original model (which we refer to as the base model here). We then report three robustness tests, where we add age (in months) as a covariate to the base mode, then gender, then both age and gender. *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$.

3.2 Robustness test: Sample sensitivity with Cook's distance

To assess the impact of potentially influential participants on the mixed-effects model in Experiment 1 (model labelled 'm1' in the authors' code, referred to as 'Base' model in the previous section 3.1 in our report), we looked at Cook's distance (influence.ME package in R). We computed influence measures at the participant level (i.e. used "Participant" as the grouping variable). We looked at two rule-of-thumb thresholds (Bollen & Jackman, 1985; Hair et al., 1998) in order to classify participants as influential, specifically $4/n$ and $4/(n-k-1)$, where n is the number of participants and k is the number of fixed effects in the model. In Figure 1, we show a scatterplot of Cook's distance values. 6 participants were above both of these thresholds (all other participants were below the more liberal threshold). These 6 participants completed a total of 33 trials.

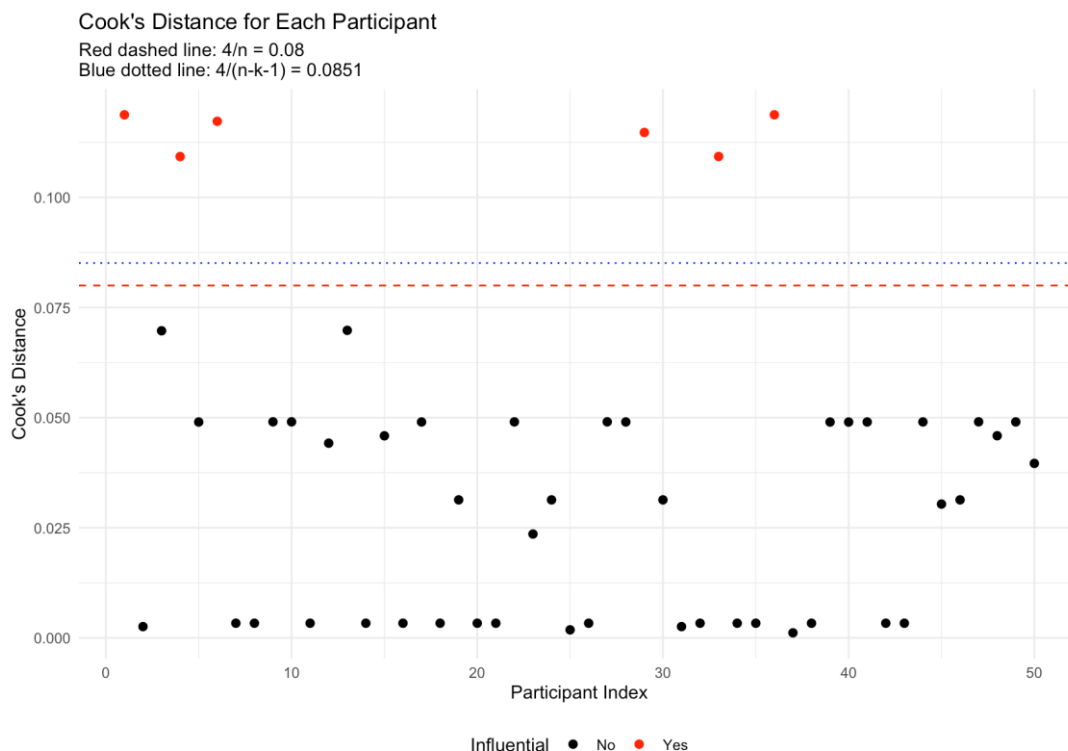


Figure 1. Scatterplot for Cook’s distance. 6 participants (red filled dots) are considered influential, as they are above both considered thresholds. The two thresholds are presented with red and blue dotted lines.

Next, we refit the base model without the influential cases (Table 3). In short, following this exercise, the coefficient for Condition remains significant and becomes more pronounced (original model $\beta = -12.712$ vs. model with exclusion of influential participants $\beta = -18.864$). The result is robust, and in our replication, now stronger. We speculate that the influential participants might have been (somewhat) inattentive responders.

Model	Intercept (<i>SE</i>) <i>p-value</i>	Condition (<i>SE</i>) <i>p-value</i>
Base (Original)	7.289*** (2.059) <i>p = .0004</i>	-12.712** (4.490) <i>p = .005</i>
Base – 6 influential participants removed	9.819*** (2.092) <i>p < .0001</i>	-18.864*** (3.406) <i>p < .0001</i>

Table 3. Robustness test considering sample sensitivity. We report the original model (which we refer to as the base model here). We then report the same model, re-ran with the exclusion of 6 influential participants as determined by Cook's distance. *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$.

3.3 Robustness test: Different optimizers

Fitting each model with different optimization algorithms produced indistinguishable fixed-effect estimates (Table 4), when considering models that converged, indicating that the original authors' results are robust to optimizer choice.

Model	Optimizer	Term	Converged	Singular	Estimate	Standard error	Z value	P value	Lower bound 95% CI	Upper bound 95% CI
m1	bobyqa*	ConditionG	TRUE	FALSE	-12.711616193	4.489637809	-2.831323294	0.004635584	-21.511306245	3.911926025
m1	Nelder_Mead	ConditionG	TRUE	FALSE	-12.711624385	4.489451238	-2.831442792	0.004633857	-21.51094877	3.912299919
m1	nlminbwrap	ConditionG	TRUE	FALSE	-12.711625596	4.489750269	-2.831254472	0.004636589	-21.511536091	3.911715031
m1	nloptwrap.NLOPT_LN_NELDERM EAD	ConditionG	TRUE	FALSE	-12.711628091	4.487648289	-2.832581164	0.004617384	-21.50741866	3.915837367
m1	nloptwrap.NLOPT_LN_BOBYQA	ConditionG	TRUE	FALSE	-12.711652071	4.488709631	-2.831916759	0.004626981	-21.50952293	3.913781123
m2	bobyqa*	(Intercept)	TRUE	FALSE	2.9036543231	1.154267466	2.515581887	0.011883608	5.16601854	0.641290094
m2	Nelder_Mead	(Intercept)	TRUE	FALSE	2.971680114	1.306030633	2.275352533	0.022884793	5.531500153	0.411860068
m2	nlminbwrap	(Intercept)	TRUE	FALSE	2.9745513524	1.314118512	2.263533568	0.023602811	5.550223637	0.398879072
m2	nloptwrap.NLOPT_LN_NELDERM EAD	(Intercept)	TRUE	FALSE	2.9744952341	1.314026501	2.263649341	0.023595699	5.549987179	0.399003282
m2	nloptwrap.NLOPT_LN_BOBYQA	(Intercept)	TRUE	FALSE	2.9745154536	1.314125933	2.263493467	0.023605287	5.550202284	0.398828628
m3	bobyqa*	ExperimentEx p2	TRUE	FALSE	0.794416019	1.473467728	0.539147213	0.589785288	-2.093580726	3.682412768
m3	Nelder_Mead	ExperimentEx p2	TRUE	FALSE	0.793612927	1.471884689	0.539181454	0.589761656	-2.091281068	3.678506914
m3	nlminbwrap	ExperimentEx p2	TRUE	FALSE	0.794413792	1.473543682	0.539117912	0.589805498	-2.093731828	3.682559407
m3	nloptwrap.NLOPT_LN_NELDERM EAD	ExperimentEx p2	TRUE	FALSE	0.794451001	1.473538327	0.539145122	0.589786723	-2.093684111	3.682586129

m3	nloptwrap.NLOPT_LN_BOBYQA	ExperimentEx p2	TRUE	FALSE	0.79464938	1.47359031	0.53926072	0.58970696	-	3.68288639
					4	4	5	6	2.09358763	6
m4	bobyqa*	(Intercept)	TRUE	FALSE	4.06016676	2.29452842	1.76949943	0.07681056	-	8.55744247
					4	9	4	1	0.43710895	1
m4	Nelder_Mead	(Intercept)	FALSE	TRUE	3.86323207	1.93496549	1.99653796	0.04587539	0.07069970	7.65576445
					9	7	6	7	4	4
m4	nminbwrap	(Intercept)	TRUE	FALSE	4.06016698	2.29455324	1.76948039	0.07681373	-	8.55749133
					5	2	8	9	0.43715737	9
m4	nloptwrap.NLOPT_LN_NELDERM EAD	(Intercept)	TRUE	FALSE	4.06034040	2.29481436	1.76935462	0.07683471	-	8.55817655
					5	3	4	3	0.43749574	7
m4	nloptwrap.NLOPT_LN_BOBYQA	(Intercept)	FALSE	TRUE	3.86356826	1.93547739	1.99618362	0.04591394	0.07003258	7.65710395
					6	1	1	1	2	2
m5	bobyqa*	(Intercept)	TRUE	FALSE	3.86324035	1.93498834	1.99651866	0.04587749	0.07066319	7.65581750
					2	4	8	6	8	6
m5	Nelder_Mead	(Intercept)	TRUE	FALSE	3.86323238	1.93499339	1.99650934	0.04587851	0.07064534	7.65581943
					6	1	4	2	2	2
m5	nminbwrap	(Intercept)	TRUE	FALSE	3.86324063	1.93502130	1.99648481	0.04588117	0.07059888	7.65588238
					5	1	1	7	6	5
m5	nloptwrap.NLOPT_LN_NELDERM EAD	(Intercept)	TRUE	FALSE	3.86293125	1.93457318	1.99678734	0.04584828	0.07116781	7.65469469
					4	4	7	9	4	4
m5	nloptwrap.NLOPT_LN_BOBYQA	(Intercept)	TRUE	FALSE	3.86320871	1.93497390	1.99651721	0.04587765	0.07065985	7.65575757
					1	7	2	4	3	3
m6	bobyqa*	ExperimentEx p3	TRUE	FALSE	7.80460611	2.16221036	3.60954986	0.00030672	3.56667381	12.0425384
					7	7	9	1	2	2
m6	Nelder_Mead	ExperimentEx p3	TRUE	FALSE	7.80239113	2.15784905	3.61581878	0.0002994	3.57300699	12.0317752
					9	3	4	5	8	8
m6	nminbwrap	ExperimentEx p3	TRUE	FALSE	7.80462622	2.16212569	3.60970051	0.00030655	3.56685987	12.0423925
					9	2	7	1	3	9
m6	nloptwrap.NLOPT_LN_NELDERM EAD	ExperimentEx p3	TRUE	FALSE	7.80535783	2.16248106	3.60944562	0.00030685	3.56689494	12.0438207
					8	9	7	2	3	3
m6	nloptwrap.NLOPT_LN_BOBYQA	ExperimentEx p3	TRUE	FALSE	7.80443715	2.16213011	3.60960568	0.00030666	3.56666212	12.0422121
					6	7	2	3	6	9
m7	bobyqa*	ExperimentEx p3	TRUE	FALSE	-	1.43550048	-	0.70650873	-	2.27304455
					0.54053639	7	0.37654908	1	3.35411734	9
					5	5	9	9	9	9
m7	Nelder_Mead	ExperimentEx p3	FALSE	TRUE	-	0.00121629	-	0	-	-
					0.53696046	6	441.471886	0.53934440	0.53457652	5
					5	6	5	5	5	5
m7	nminbwrap	ExperimentEx p3	FALSE	TRUE	-	1.43129743	-	0.70455512	-	2.26262538
					0.54271758	7	0.37917876	2	3.34806056	8
					9	1	6	6	6	6
m7	nloptwrap.NLOPT_LN_NELDERM EAD	ExperimentEx p3	TRUE	FALSE	-	1.44548528	-	0.71087364	-	2.29733432
					0.53581684	7	0.37068301	4	3.36896800	3
					2	3	3	3	3	3

Table 4. Fixed effect estimates for each model using five different optimization algorithms and their associated convergence status. *Note.* *Indicates the optimizer used by the original study authors (bobyqa in all models). M1: Tests whether the two conditions differ. Condition is dummy-coded with the Focus condition as reference. M2: Tests whether the rate of selecting the new referent differs from chance. M3: Tests whether the Given conditions differ between Experiment 1 and Experiment 2. Experiment is dummy-coded with Experiment 1 (Given) as reference. M4: Tests whether the rate of selecting the new referent differs from chance in Experiment 3. M5: Exploratory model testing whether the

rate of selecting the new referent differs from chance without an item intercept. M6: Tests whether the rate of selecting the new object differs between Experiment 2 and Experiment 3. M7: Tests whether the F-conditions differ between Experiment 1 (Focus) and Experiment 3. Experiment is dummy-coded with Experiment 1 (Focus) as reference.

3.4 Robustness test: Contrast coding

Re-expressing the categorical predictors under different contrast schemes (treatment and sum coding) did not change the estimated pairwise log-odds contrasts (Table 5), confirming that the results are robust to different coding schemes.

Model	Coding	Contrast	Estimate (link)	Standard error (link)	Lower bound 95% CI	Upper bound 95% CI
m1	treatment_first*	F - G	12.7116161	4.48963781	3.91208772	21.5111445
m1	treatment_reversed	G - F	-12.711613	4.48521487	-21.502473	-3.9207537
m1	sum_coding	F - G	12.7116194	4.48919876	3.91295146	21.5102872
m3	treatment_first*	Exp1 - Exp2	-0.794416	1.47346773	-3.6823597	2.09352766
m3	treatment_reversed	Exp2 - Exp1	0.79441399	1.47350846	-2.0936095	3.68243749
m3	sum_coding	Exp1 - Exp2	-0.7944144	1.473537	-3.6824939	2.09366504
m6	treatment_first*	Exp2 - Exp3	-7.8046061	2.16221036	-12.042461	-3.5667517
m6	treatment_reversed	Exp3 - Exp2	7.80460708	2.16222682	3.56672039	12.0424938
m6	sum_coding	Exp2 - Exp3	-7.8046079	2.16222106	-12.042483	-3.5667324
m7	treatment_first*	Exp1 - Exp3	0.54053639	1.43550049	-2.2729929	3.35406565
m7	treatment_reversed	Exp3 - Exp1	-0.5405346	1.43547154	-3.3540071	2.27293792
m7	sum_coding	Exp1 - Exp3	0.54053488	1.43548538	-2.2729648	3.35403453

Table 5. Pairwise log-odds contrasts between conditions for each model under three contrast-coding schemes (treatment, reversed treatment, and sum coding). *Note.*

*Indicates the coding used by the original study authors (treatment first in all models).

M1: Tests whether the two conditions differ. Condition is dummy-coded with the Focus condition as reference. M3: Tests whether the Given conditions differ between Experiment 1 and Experiment 2. Experiment is dummy-coded with Experiment 1

(Given) as reference. M6: Tests whether the rate of selecting the new object differs between Experiment 2 and Experiment 3. M7: Tests whether the F-conditions differ between Experiment 1 (Focus) and Experiment 3. Experiment is dummy-coded with Experiment 1 (Focus) as reference.

3.5 Robustness test: Robustness to different random effects specifications

The original authors first fitted models with the structure $y \sim 1 + (1|item) + (1|participant)$ or if a categorical predictor was present $y \sim 1 + Factor + (1|item) + (1|participant)$. These were fitted with the R package lme4 (Bates et al., 2015). However, some of these fits produced warnings in relation to the item effect and models with simpler random effects structures (e.g., excluding item effects) were reported in some cases. This is unsatisfactory for several reasons. First, deriving estimates and tests from different analyses can lead to inconsistencies. Second, the estimates from raw data handle missing trials in a suboptimal way (assuming trials are MCAR; missing completely at random) relative to a mixed effect model (which treats missing trials as MAR; missing at random). Third, the raw data approach ignores the random effects.

It is not clear why estimates were derived from raw data rather than the model used for inference. This could be due to warnings in relation to convergence for some models; alternatively, due to interval estimates from these models relying on asymptotics - whether profiling likelihood or a Wald style (z) approximation is used. It is well known that accurate interval estimates and tests for mixed effects models are not straightforward to obtain. Bolker et al. (2025) suggest that methods for obtaining inferences from single parameters are (from worst to best) Wald z tests, Wald t tests, Likelihood ratio tests and Markov Chain Monte Carlo (MCMC) or parametric bootstrap interval estimates.

For this robustness check we refitted all the models as maximal models using Bayesian mixed effects logistic regression models with the Bernoulli family in the brms package in R (Bürkner, 2017). This uses MCMC estimation to fit the models and is considered one of the gold standard approaches to obtaining accurate inferences. Under this approach missing trials are assumed to be MAR. We used default (weakly informative priors in

brms) to obtain interval estimates: specifically Highest Posterior Density (HPD) intervals. With this prior specification it is expected that the interval estimates should closely match those from a frequentist mixed effects model. The maximal model specification broadly follows recommendations from Barr et al. (2013). The aim is to capture the structure of the data with the effects being modeled. The key models either take the form of a fixed intercept with random effects or a fixed intercept plus factor (either Condition with levels focus vs. given or Experiment with two levels representing experiments being compared). We therefore fit the following random effect structure for intercept only models:

$$y \sim 1 + (1|item) + (1|participant)$$

In these models the focus is on whether the choice of the new item was above chance (50%) for the focus condition or below chance for the given condition. For the models with a fixed factor we fit the structure:

$$y \sim 1 + (1|item) + (1 + Factor||participant)$$

As the factor (Condition or Experiment) always represents a comparison between two groups of participants the common maximal structure $(1 + Factor|participant)$ is inappropriate because it includes a correlation term between the random intercept and dummy coded Factor. As these are independent groups the expectation is that the effects are uncorrelated, but $(1 + Factor||participant)$ in effect fits different random intercepts with distinct variances for each group. This is potentially important because the variances do differ between conditions and experiments (sometimes considerably).

While we fitted maximal Bayesian versions for each of the reported models the overall pattern of results can be summarized by comparing just the three key models M1, M2 and M3 representing the primary analysis for each of the three experiments. Table 6 compare the point and interval estimates for each of the experiments. In each case the outcome is the percent of new objects selected by the participant. These are: i)

estimates derived from raw data (that were reported in the original paper), ii) estimates derived from the models used to derive the hypothesis tests in the original paper (for which point and interval estimates were not reported), and iii) the estimates from the maximal Bayesian models fitted as a robustness check.

3.6. Robustness test comparing point and interval estimates derived from raw data, the reported mixed effect model and a maximal Bayesian model

From Table 6 several clear patterns emerge. The point estimates from the raw data differ from those for the lme4 models used to derive hypothesis tests for each experiment, while the models used to derive hypothesis tests and the maximal Bayesian model estimates are more similar, but not identical. The raw data interval estimates are symmetrical and narrower than the Wald estimates or the Bayesian HPD intervals. Overall, the Bayesian estimates provide stronger support for the hypotheses than either the raw data or models used for hypothesis testing by Brody and colleagues.

	<i>Experiment</i>	<i>E1</i>	<i>E2</i>	<i>E3</i>
	<i>Model Label</i>	<i>M1</i>	<i>M2</i>	<i>M3</i>
		<i>Focus</i>		<i>Focus</i>
Raw data	Percentage new	89.4%	—	86.7%
	95% interval estimate	84.3, 94.3	—	79.8, 93.3
Reported model	Percentage new	99.3%	—	98.3%
	95% interval estimate	96.3, 100.0	—	39.2, 100.0
Bayesian model	Percentage new	99.2%	—	96.0%
	95% interval estimate	95.1, 100.0	—	83.8, 100.0
		<i>Given</i>	<i>Given</i>	
Raw data	Percentage new	24.3%	29.4%	—
	95% interval estimate	17.3, 31.2	23.6, 35.4	—
Reported model	Percentage new	0.4%	5.2%	—

	95% interval estimate	< 0.5, 47.8	0.6, 34.5	—
Bayesian model	Percentage new	3.5%	9.9%	—
	95% interval estimate	< 0.5, 18.9	< 0.5, 33.0	—

Table 6. Robustness test comparing point and interval estimates derived from raw data, the reported mixed effect model and a maximal Bayesian model. *Note.* Raw data estimates are from Brody *et al.* (2024). Raw data interval estimates are from bootstrapping raw data. Reported model estimates are from the labeled model in the Brody *et al.* R script fitted with lme4. Corresponding interval estimates are Wald style normal approximation. Bayesian model estimates are from a fitted model in brms with default priors in R with maximal random effects as appropriate to the model (see full explanation in text). Interval estimates are HPD intervals from brms.

The most striking feature is the difference in the point estimates for the raw data and the estimates from either of the mixed effects models. This discrepancy requires some explanation. In a mixed effect model the estimates exhibit shrinkage. In essence this shifts (“shrinks”) estimates towards mean values and reduces the impact of atypical responses (those with large residuals). Shrinkage makes the model slightly worse at describing the observed data but improves out of sample generalization. The degree of discrepancy in this study is however larger than typically seen. This is in part because of the small sample size but largely because of the nature of the data and sample. With children in this age range (2 - 3 years) some are likely inattentive and essentially responding at random on one or more trials. The exclusion criteria removed participants fewer than 3 usable trials out of six. Figure 2 in Brody *et al.* (2024, p. 1320) indicates that typical children selected the new object near ceiling on focus trials and near floor on given trials (with the modal response being 100% or 0% in many cases). Thus the model estimates are shrunk towards the more consistent and (presumably) less distracted or inattentive children. The analysis of raw data also exacerbates any differences because of missing trials. Averaging percentages per participants weights trials from children with missing data more than those with complete cases (who are also likely the less distracted participants). This is also consistent with the influence

analysis reported earlier (where removing the higher influence data points strengthened the effects).

(Shrinkage also occurs for the items, but item variance was negligible, and it is unlikely to have had much impact on the estimates in any of the models. Estimation issues in the original studies were in fact largely singularity warnings when the item variance was estimated as zero or negative.

In terms of inferences, the interval estimates from the maximal Bayesian model are more accurate than the CIs derived from the raw data and or the Wald CIs for reasons outlined above. They provide somewhat stronger support for the main findings of Brody et al. (2024) than the reported interval estimates and tests. This is most apparent for M3 in Table 6. The reported test of this effect in Experiment 3 is a Wald test of the log odds of the intercept vs. zero (i.e., selecting the new object above the 50% chance). This is non-significant with $\beta = 4.06$, $z = 1.769$, $p = .0768$ (and a corresponding interval estimates that includes 50%). A more accurate likelihood ratio test for the same model gives $\chi^2(1) = 15.73$, $p = .0073$. This is also supported by the Bayesian HPD interval of [83.8, 100] which comfortably excludes 50%. We can also derive a Bayesian approximation to a p value for this test. A one-sided frequentist p value can be considered an approximation to the Bayesian posterior probability of an effect in the opposite direction to that observed (Casella & Berger, 1987). For this test, this probability is .001 and thus a Bayesian approximation to the two-sided p value (for comparison) is $p = .002$.

4. Conclusion

Our replication efforts confirm the computational reproducibility and robustness of Brody et al.'s (2024) original results. We successfully reproduced key effects, and those remained significant across our robustness tests, including the addition of covariates, the removal of influential participants, varying optimizers or contrast coding, and fitting maximal Bayesian models. Notably, in two of our tests (influence and Bayesian analyses) we found support for stronger effects than the ones originally reported by the authors. We

speculate this could be because of slightly inattentive responding on some trials. These stronger effects should not be interpreted as a critique against the authors, who had transparently pre-registered sensible exclusion criteria for inattentive responding. Instead, among attentive children, it should be interpreted as further support for the conditional effect of mutual exclusivity. Altogether, our replication tests give us increased confidence in the results reported by Brody and colleagues.

References

- 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. <https://doi.org/10.1016/j.jml.2012.11.001>
- Bates, D., Maechler, M., Bolker, B., Walker, W. (2015). Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software*, 67(1), 1-48.
- Bergelson, E., & Swingley, D. (2012). At 6–9 months, human infants know the meanings of many common nouns. *Proceedings of the National Academy of Sciences* 109(9), 3253–3258.
- Bollen, K. A., & Jackman, R. W. (1985). Regression diagnostics: An expository treatment of outliers and influential cases. *Sociological Methods & Research*, 13(4), 510-542
- Bolker, B. *et al.* (2025, 19 July). GLMM FAQ. <https://bbolker.github.io/mixedmodels-misc/glmmFAQ.html> Retrieved 20 November 2025.
- Brody, G., Feiman, R., & Aravind, A. (2024). Why do children think words are mutually exclusive?. *Psychological Science*, 35(12), 1315-1324.
- Bürkner, P.-C. (2017). brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal of Statistical Software*, 80(1), 1-28.
- Casella, G., & Berger, R. L. (1987). Reconciling Bayesian and Frequentist Evidence in the One-Sided Testing Problem. *Journal of the American Statistical Association*, 82(397), 106–111.

Hair, J. F., Anderson, R. E., Tatham, R. L., & Black, W. C. (1998). *Multivariate data analysis*. Englewood cliff. New jersey, USA, 5(3), 207-2019.

Lewis, M., Cristiano, V., Lake, B. M., Kwan, T., & Frank, M. C. (2020). The role of developmental change and linguistic experience in the mutual exclusivity effect. *Cognition* 198, 104–191.