

Bias Dynamics in BabyLMs: Towards a Compact Sandbox for Democratising Pre-Training Debiasing

Anonymous submission

Abstract

Pre-trained language models (LMs) have, over the last few years, grown substantially in both societal adoption and training costs. This rapid growth in size has constrained progress in understanding and mitigating their biases, especially towards under-represented communities. Since re-training LMs is prohibitively expensive, most debiasing work has focused on post-hoc or masking-based strategies, which often fail to address the underlying causes of bias. In this work, we seek to democratise pre-model debiasing research by using low-cost proxy models, striving to make this research direction accessible to projects outside of the large industry labs. Specifically, we investigate BabyLMs, compact BERT-like models trained on small and mutable corpora that can simulate the bias acquisition and learning dynamics of larger models. We show that BabyLMs display closely aligned patterns of intrinsic bias formation and performance development compared to standard BERT models, despite their drastically reduced size. Furthermore, correlations between BabyLMs and BERT hold across multiple intra-model and post-model debiasing methods. Leveraging these similarities, we conduct pre-model debiasing experiments with BabyLMs, replicating prior findings and presenting new insights regarding the influence of gender imbalance and toxicity on bias formation. Our results demonstrate that BabyLMs can serve as an effective sandbox for large-scale LMs, reducing pre-training costs from over 500 GPU-hours to just over 30 GPU-hours. This provides a way to democratise pre-model debiasing research by enabling faster, more accessible exploration of novel debiasing strategies and the examination of historically under-explored bias topics in service of building fairer LMs.

1 Introduction

The recent dramatic increase in investment in large language models (LLMs) has driven sharp performance gains (Korinek and Vipra 2024) while altering the way research is conducted. Since 2019, parameter counts have grown by roughly three orders of magnitude (Radford et al. 2019; Meta AI 2025), notably inflating experimental costs and straining an academic ecosystem built on small, incremental advances (Cottier et al. 2024; Sathish et al. 2024).

This concentration of compute and data has widened global inequalities in who can participate in LM development. Many research groups, especially in the Global South and marginalised communities, lack the hardware and datasets

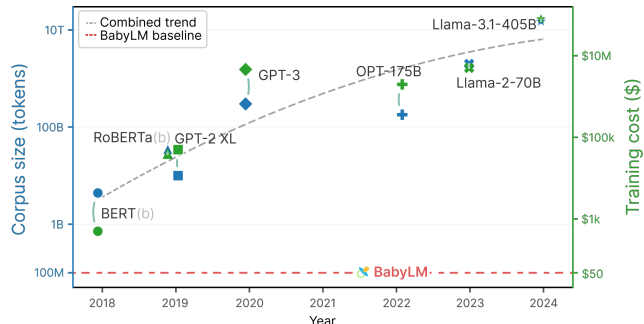


Figure 1: BabyLMs reproduce the bias–performance trends of larger LMs while requiring orders-of-magnitude less data and compute, enabling systematic and democratised study of debiasing.

needed to tackle issues affecting their own communities (Chan et al. 2021).

This problem is particularly pressing in LM bias research, which examines how LLMs treat users differently on the basis of protected attributes – e.g., gender (Zhao et al. 2024), and ethnicity (Field et al. 2021). Despite the ubiquity of LMs, bias-removal research remains relatively linear and, constrained by the high cost of training, focuses mostly on post-hoc debiasing of already trained LMs, altering models only after their internal structures and reasoning circuits have formed (Gallegos et al. 2024). While these methods can help, they often merely mask biases detectable by our simplistic probes rather than fully removing them from the model (Gonen and Goldberg 2019; Gupta et al. 2024).

By contrast, debiasing strategies before or during pre-training remain rare and have become outdated, partly because even a single pre-training run of a relatively simple BERT-style model exceeds 500 GPU-hours on current hardware (Devlin et al. 2019). As a result, only well-resourced institutions can iterate at pre-training time, leaving pre-training debiasing under-explored, slowing methodological progress, and limiting research on the harms faced by marginalised, low-resource groups (Bender et al. 2021; Gelles et al. 2024).

In search of a solution, this work focuses on models developed for the BabyLM Challenge, a research initiative that seeks to create well-performing models trained on small-scale datasets (Warstadt et al. 2023a; Hu et al. 2024a; Charpentier

et al. 2025). Together, these properties show promise in making BabyLMs a valuable tool for democratising debiasing research, combining relevant performance that matches or surpasses BERT, mutable corpora that can be readily experimented with, and much lower pre-training time.

In this paper, we verify this proposition, showing that BabyLMs:

1. *acquire and express biases in a way that is representative of standard LMs*
2. *respond to established debiasing methods similarly to standard LMs*
3. *enable the democratisation of pre-model debiasing research*

2 Background

2.1 BabyLMs

BabyLMs belong to the field of low-resource language models that seeks to democratise NLP with well-performing, affordable LMs (Van Nguyen et al. 2024; Warstadt et al. 2023b).

They are inspired by the fact that human children encounter three to four orders of magnitude less linguistic data than conventional LMs, with Llama-3 using approximately 15 trillion tokens (Grattafiori et al. 2024) while a 13-year-old child may only have encountered about 100 million words (Gilkerson et al. 2017; Linzen 2020). Unlike approaches that reduce parameter counts (Jiao et al. 2019; Hoffmann et al. 2022), the BabyLM task restricts the amount of pre-training data, simulating a more human-like learning environment.

It utilises the aforementioned 100 million words as the training corpus, which consists of transcriptions of child-directed speech, child-oriented texts (e.g., books, subtitles), and other commonly available data (e.g., Wikipedia) (Warstadt et al. 2023b).

2.2 Methods for Evaluating LM Performance

When introducing alternative LM architectures, we need evaluation frameworks to compare performance. **BLiMP** probes core grammatical knowledge via minimal pairs across multiple categories (e.g., syntax, morphology, semantics) (Warstadt et al. 2020). Additionally, the **BabyLM BLiMP supplement** adds five extra probes concerning dialogue and question understanding (Warstadt et al. 2023a). **EWoK** evaluates cognition-inspired world-model knowledge by asking models to match two contexts to two targets, discouraging reliance on surface likelihoods (Ivanova et al. 2024).

Furthermore, **GLUE/SuperGLUE** provide an overview of an LM’s downstream capabilities, and are more compute-intensive to evaluate (Wang et al. 2018, 2019).

There are also more complex evaluation tasks targeted at LLMs concerning properties such as reasoning or word knowledge (Hendrycks et al. 2020; Rein et al. 2024).

2.3 Notable BabyLM Architectures

Throughout the years, numerous LM variants have been submitted to the competition, with the most relevant submissions listed here. The **LTG-BERT** architecture (Samuel et al. 2023)

keeps the BERT backbone but adds NormFormer (Shleifer, Weston, and Ott 2021), gated GELU (Shazeer 2020), disentangled relative positions (He et al. 2020), and span masking (Joshi et al. 2020), outperforming BERT with a much smaller corpus for an identical number of training steps. **GPT-BERT** (Charpentier and Samuel 2024), the current SOTA BabyLM, retains LTG-BERT’s architecture but adds a hybrid objective that combines span masking with causal next-token prediction (BehnamGhader et al. 2024), which further separates it from the classic BERT. Other approaches have also experimented with preprocessing (Cheng et al. 2023) or curriculum learning (Diehl Martinez et al. 2023). However, they underperformed the LTG-based architectures.

2.4 Frameworks for Analysing LM Bias

Blodgett et al. (2020) define bias as systematic patterns in representations or outputs that reinforce social inequalities, resulting in allocational or representational harms. Frameworks for evaluating LM bias range from intrinsic probes to more costly, task-specific extrinsic evaluations, with only a limited correlation observed between the two (Goldfarb-Tarrant et al. 2020; Cao et al. 2022).

StereoSet measures intrinsic bias using Context Association Tests in which the model ranks stereotype, anti-stereotype, and unrelated continuations (Nadeem, Bethke, and Reddy 2020). The bias score is derived from the ratio of examples where the model prefers stereotypical versus anti-stereotypical options, with 50 indicating an unbiased model. Another notable dataset, **CrowS-Pairs** (Nangia et al. 2020), employs entire sentences, utilising stereotype/anti-stereotype sentence pairs, with the score again reflecting the proportion of stereotype sentence choices by the model. Alternative frameworks include SEAT (May et al. 2019), which provides a lightweight intrinsic option, and WinoGender/WinoBias (Rudinger et al. 2018; Zhao et al. 2018a), which target extrinsic gender-based coreference bias.

2.5 Techniques for LM Bias Mitigation

To examine techniques for reducing biases in LMs, we follow Guo et al. (2024) and organise them by application stage. *Post-model approaches*, which leave the model intact and adjust only representations, are cheap and interpretable but are surface-level: **Iterative Nullspace Projection** (INLP) removes linearly decodable protected-attribute signals (Ravfogel et al. 2020) and **Sent-Debias** projects away PCA-estimated bias subspaces (Liang et al. 2020). Nevertheless, non-linear probes frequently reveal residual bias, underscoring that post-hoc fixes are linear and incomplete (Sun et al. 2025; Gonen and Goldberg 2019).

Intra-model methods involve fine-tuning that debiases the full model: increasing **dropout** can disrupt biased representations but risks negative performance impacts (Webster et al. 2020); **Counterfactual Data Substitution** (CDS) rewrites biased instances (e.g., swapping gendered entities), placing them into a new context (Bartl, Nissim, and Gatt 2020; Webster et al. 2018); and **debiasing losses** that motivate the model to debias itself (Park et al. 2023). While bringing some improvement, other studies have demonstrated the brittleness of these methods (Mendelson and Belinkov 2021).

Pre-model methods modify data or pre-training to reduce bias before training even starts. They yield more stable debiasing but are costly to implement and research due to the need to retrain the model (Li et al. 2024; Xie and Lukasiewicz 2023). Key tactics include **Counterfactual Data Augmentation** (CDA) that swaps demographic markers to rebalance either the entire, or fine-tuning, corpus (Lu et al. 2018; Zmigrod et al. 2019; Webster et al. 2020); **toxic-content filtering** that displays a beneficial debiasing effect when applied (BigScience-Workshop et al. 2022; Ranaldi et al. 2024); and **perturbation augmentation**, which stochastically edits sentences along gender, ethnicity, and age axes to produce fairer models (Qian et al. 2022).

3 Experiment Setup

As established in the previous section, pre-model debiasing offers the most stable mitigation, addressing inner representations rather than surface cues. Yet these methods remain under-explored because they require costly re-training of a model from scratch and manipulating large, often improperly specified, corpora. Given the advantages of BabyLM architectures (competitive performance, compact datasets, and inexpensive pre-training), together with established bias- and performance-evaluation frameworks, we argue that BabyLMs can provide a practical sandbox for systematic pre-model debiasing and lower the research barrier.

To validate this proposition, we first need to understand how the debiasing behaviours of standard LMs and BabyLMs are aligned. For this, we need to identify a candidate BabyLM that suitably replicates the bias dynamics of standard LMs like BERT.

3.1 Metrics

We utilise the performance and bias scores from frameworks described in Sections 2.2 and 2.4. For bias, this means using CrowS-Pairs and StereoSet. Both share a mathematically similar approach and a score scale from 0 to 100. Since BabyLMs can be either masked or continuation LMs, lacking next-sentence prediction, we exclude StereoSet’s inter-sentence portion (Ranaldi et al. 2024).

A bias evaluation is not sufficient on its own. Since there is an established relationship between the model’s biases and its performance (Nadeem, Bethke, and Reddy 2020), we estimate the performance through BLiMP, BabyLM BLiMP supplement, and EWoK.

Thus, we obtain three performance metrics and two bias metrics, all of which capture only part of a model’s behaviour. They all share the same scale but probe the model with different sentences and contexts. Therefore, they reveal different parts of its bias and performance profile (Zakizadeh and Pilehvar 2025). To obtain a more comprehensive picture, we average the individual performance scores into a *composite performance* metric and the two bias scores into a *composite bias* metric, with details further discussed in Appendix B.

3.2 Candidate BabyLM

As noted, our aim is to select a BabyLM that comes the closest behaviour-wise to standard LMs, while still retaining

desirable characteristics, such as low-cost training. Firstly, this requires showing that BabyLMs in general display bias acquisition dynamics, such as verifying that they acquire more biases with increased performance, as observed within larger LMs (Nadeem, Bethke, and Reddy 2020).

We do this by evaluating composite bias and performance metrics for every notable BabyLM and various variants of standard LMs, with all models used listed in Appendix A. This yields Table 1, which shows the correlation between *composite performance* and *composite bias* for both model classes. The strong positive correlation, and its shared strength across BabyLMs and standard LMs, confirms that the overall trend of bias increasing with performance is preserved in BabyLMs. Therefore, BabyLMs can be a valid sandbox for studying debiasing techniques with resources that are feasible for many research groups.

Model Class	N	$r(\text{Composite Performance, Composite Bias})$
BabyLM	9	0.833
Standard	16	0.753

Table 1: Pearson correlation between composite performance and bias for the different model classes

With BabyLM eligibility established, we look at the behavioural distribution of the different models, which is shown in Figure 2 to select the most informative and relevant BabyLMs, which we continue using in our research. The SOTA variant of the LTG-BERT architecture, `ltg-bert-babylm`, is closest to the original BERT. However, this SOTA version has been trained on 1,500 epochs, requiring roughly the same GPU-hours as the original BERT model, making it unsuitable for our democratisation purposes. Thus, we also select its low-resource variant, `ltgbert-100m-2024` (Hu et al. 2024a), trained on just 40 GPU-hours. Although more distant from BERT, it still achieves reasonable performance and exhibits measurable bias, making it a promising low-cost BabyLM. In conclusion, the two candidate models for further study are **LTG-BERT** (`ltg-bert-babylm`) and **LTG-Baseline** (`ltgbert-100m-2024`). LTG-BERT tests whether any BabyLM can sufficiently replicate standard debiasing patterns, while LTG-Baseline tests whether these hold even with such minimal training.

4 Model Viability

Having established that BabyLMs acquire biases in a similar way to standard LMs, the next step is to test whether they also debias comparably. We analyse whether their pre-training corpora are comparable in terms of the biases which can be learned from them. We also investigate how their composite performance and bias change in reaction to a wide selection of intra-model and post-model debiasing techniques, each targeting different parts of the models.

4.1 Corpora

Debiasing strategies, especially pre-model ones, largely entail altering the pre-training corpus; thus, we must demonstrate

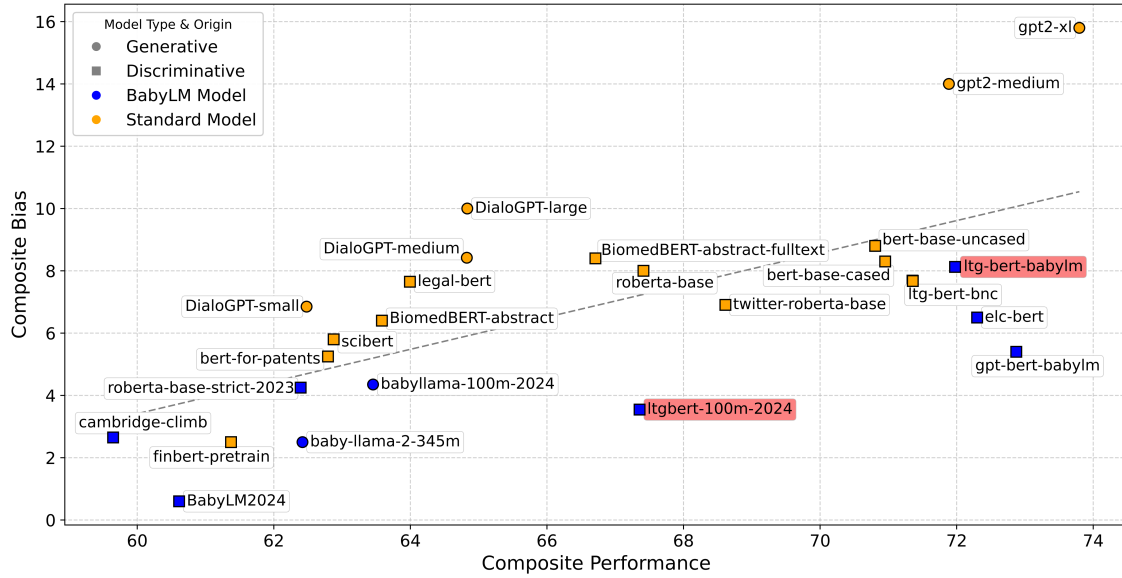


Figure 2: Average performance correlated with average bias for all evaluated LMs, with the average trend being shown and the candidate models being highlighted

that the BabyLM corpus can support the emergence of the same bias dynamics observed in BERT. We therefore examine corpus aspects known to induce LM biases. This section reports the most critical results; the remainder is listed in Appendix C.

In terms of the corpora, we utilise the 2023 BabyLM challenge corpus, used to train LTG-BERT, and closely replicate the unavailable BERT corpus from a Wikipedia dump (Broad 2022) and a BookCorpusOpen re-crawl (Di Liello 2022; Bandy and Vincent 2021), preserving the token ratio between the two corpora and the original $\sim 3\text{B}$ -token size.

With the corpora established, we examine differences in their topical coverage (Chang, Prabhakaran, and Ordonez 2019; Zhao et al. 2019). For each topic category, we compute the percentage of each corpus formed by keywords related to specific subcategories, which were extracted from several bias-oriented studies (Meade, Poole-Dayana, and Reddy 2021; Qian et al. 2022; Sosto and Barrón-Cedeño 2024). Table 2 shows the categories and sub-categories accompanied by several examples.

Gender is the most prominent category. Table 3 shows that the male gender is more represented than the female gender across corpora, with the resulting models also biased in a male-centric direction. This suggests that they should react to gender-focused debiasing in the same way. Across other bias categories, most trends are shared between the corpora. In ethnicity-term frequency, Caucasian terms are consistently most over-represented, Black-related terms second, and Asian strongly last. In religion, Christian terms dominate, while Jewish terms are barely represented in both. LGBTQ+ topics appear equally. The differences lie in the BERT corpus representing Black and Muslim terms much more substantially. Overall, despite BERT containing higher frequencies of biased terms across categories, the BabyLM

corpus largely preserves similar topic ratios, indicating it ought to support debiasing techniques.

Finally, as toxicity and hate-speech are linked to increased biases (BigScience-Workshop et al. 2022), we used established models (Hanu and Unitary AI 2020; Antypas and Camacho-Collados 2023) to label their presence in every sentence in each corpus. Table 5 shows that the BabyLM corpus is more toxic and hateful, even containing blatantly racist and sexualised terms, contrary to its child-aligned disposition.

Overall, compared to the BERT corpus, BabyLM exhibits

Category	Subcategory	Example Terms
Gender	Male	actor, dad, he, ...
	Female	actress, mom, she, ...
Ethnicity	Black	black, african, ...
	Caucasian	caucasian, white, ...
	Asian	asian, china, ...
Religion	Jewish	jewish, torah, ...
	Christian	christian, bible, ...
	Muslim	muslim, quran, ...
LGBTQ+	LGBTQ+	lgbtq, gay, trans, ...

Table 2: Examples of bias-related terms across predefined categories and sub-categories.

Category	BabyLM	BERT
Male	2.07%	1.83%
Female	1.03%	1.30%

Table 3: Gender representation distribution across the corpora









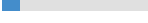
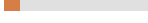




Category	BabyLM	BERT
Black	 0.035%	 0.062%
Caucasian	 0.049%	 0.067%
Asian	 0.014%	 0.034%
Jewish	 0.008%	 0.007%
Christian	 0.034%	 0.065%
Muslim	 0.007%	 0.028%
LGBTQ+	 0.017%	 0.022%

Table 4: Demographic keyword distribution across the corpora





Corpus	Toxic	Hate-speech
BabyLM	 3.12%	 0.55%
BERT	 0.95%	 0.34%

Table 5: Toxicity and hate-speech sentence rates

similar topic frequencies and a higher presence of toxic and hateful sentences, enabling the same bias and debiasing dynamics. Consequently, the BabyLM corpus appears fully capable of supporting both bias acquisition and debiasing similar to those observed with BERT.

4.2 Debiasing behaviour

Building on the finding that the BabyLM corpus aligns with the BERT corpus across all bias-relevant metrics, we test whether the two model classes debias comparably. To analyse behavioural overlaps, we use a broad set of intra- and post-model debiasing techniques, each using different mechanisms and targeting distinct model components. In each test, we compare how debiasing shifts LTG-BERT’s and LTG-Baseline’s composite bias and performance relative to BERT.

Starting with post-model debiasing, we apply two debiasing methods, **Sent-Debias** and **INLP**, both targeting trained models with already frozen encoders. Sent-Debias learns a gender subspace from text and subtracts it from the final hidden representations. INLP trains linear classifiers for protected attributes and iteratively projects out the directions that make those attributes linearly separable. Implementation of both approaches uses a 2.5M-word Wikipedia dump for their debiasing signals (Meade, Poole-Dayana, and Reddy 2021).

Looking at the impact of debiasing in Figure 3, we see that Sent-Debias consistently reduces composite gender bias across all models, although the performance impact varies. INLP remains consistent across architectures regarding its effects on bias. Gender-focused INLP lowers overall bias, while race-focused INLP does not reduce composite bias and also harms accuracy. The gender-focused INLP’s performance effects align with model fit: the over-fitted LTG-BERT benefits from it as a form of regularisation, achieving SOTA performance (Appendix D), the more under-fitted LTG-Baseline loses useful information, and full-data BERT goes without severe penalties. With this behavioural nuance, we conclude that the methods’ impact on bias is consistent across models.

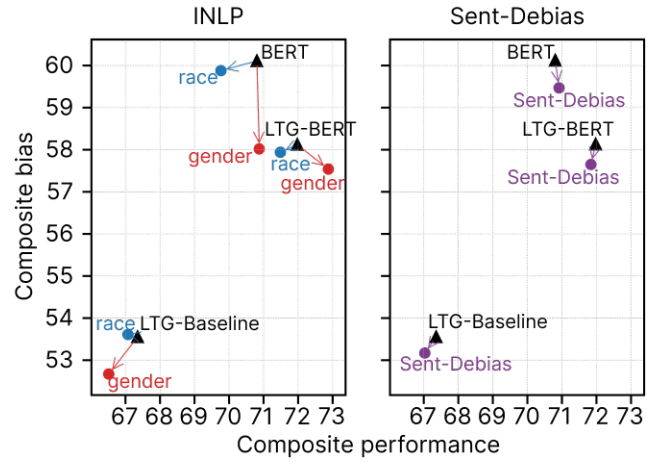


Figure 3: Bias and performance changes caused by post-model debiasing strategies

Next, we probe intra-model debiasing with four strategies. **CDA** balances a 10M-word Wikipedia dump by duplicating sentences containing gendered or racial terms and swapping them (Meade, Poole-Dayana, and Reddy 2021). **CDS** uses the gender-balanced corpus, substituting each gender mention with the opposite to create anti-stereotype contexts (Webster et al. 2018; Bartl, Nissim, and Gatt 2020). A **debiasing-loss** setup trains on the Gender Pronoun Resolution task with Stereotype Neutralisation and Elastic Weight Consolidation (Zhao et al. 2018b; Park et al. 2023). A **dropout** variant increases dropout during continued pre-training on the same Wikipedia dataset.

Across methods, Figure 4 shows that the models shift in the same direction, differing mainly in magnitude. CDA on gender and race reduces composite bias across models, with stronger gender effects and similar performance-loss trends. CDS again yields near-identical bias reduction across models with modest, method-consistent accuracy costs. Debiasing loss induces the largest accuracy drop and bias decrease, with LTG-Baseline tracking BERT most closely. Dropout yields weaker debiasing while preserving the trends; some noise in the performance–bias-loss ratio is expected since BERT lacks the extra normalisation of LTG models (Shleifer, Weston, and Ott 2021), leading to over-confident logits that drift under perturbation (Gal and Ghahramani 2016; Kong et al. 2020).

These convergent effects show that BabyLMs simulate the behaviour of post- and intra-model debiasing techniques on BERT. We quantify alignment by pairing each model’s performance and bias shifts per method and running canonical correlation analysis (Hotelling 1936). The results in Table 6 show that LTG-Baseline is the most faithful simulator of BERT’s debiasing behaviour, whereas LTG-BERT’s overfitting likely dampens alignment.

5 Pre-Model Debiasing Experiments

In the previous section, we showed that standard LMs and BabyLMs share debiasing dynamics, enabling us to estimate a debiasing method’s impact on a standard LM by applying

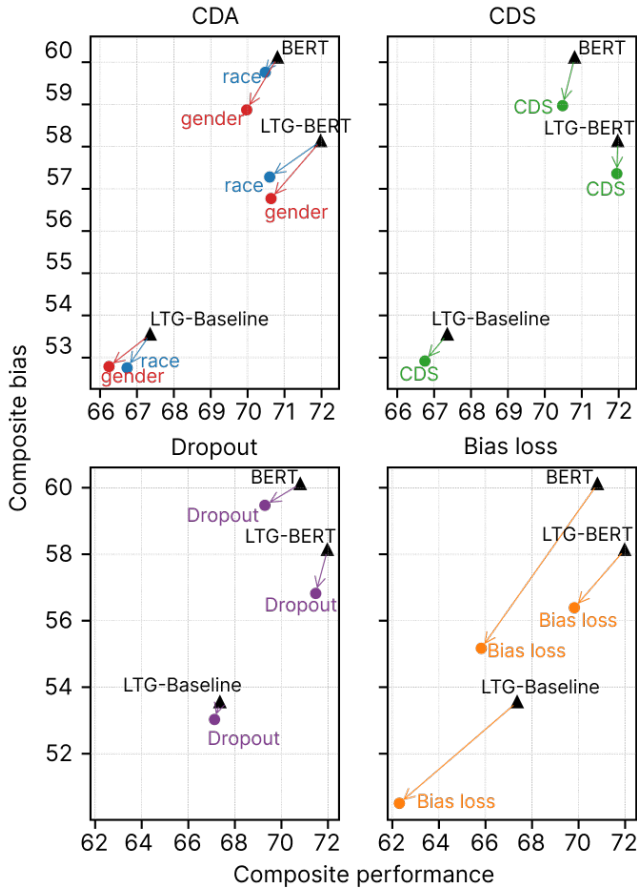


Figure 4: Bias and performance changes caused by intra-model debiasing strategies

Model pair	Correlation(ρ_1)
BERT↔LTG-BERT	0.796
BERT↔LTG-Baseline	0.981
LTG-BERT↔LTG-Baseline	0.772

Table 6: First canonical correlation (ρ_1) between performance and bias shifts across all debiasing methods

it to a BabyLM. We can now utilise this to run pre-model debiasing experiments on LTG-Baseline instead of BERT, reducing the cost from 500 GPU-hours per experiment to just over 30.

Throughout this section, we reinforce that BabyLM mimics BERT’s reported behaviour and show how this benefits future investigation into pre-model debiasing. All experiments utilise the LTG-Baseline architecture, altering only the training corpus. The exact pre-training setup is specified in Appendix E.

Finally, we establish a baseline by training the LTG model on the original BabyLM corpus and tracking performance and bias over time. Within this baseline training, the model picked up bias early, which then stabilised at a steady level, while performance improved more gradually as it learned

richer linguistic structure. The quick uptake of bias indicates that the biases stem directly from the topical imbalances and stereotypes, which most pre-model debiasing techniques target.

5.1 CDA Pre-model Debiasing

As a first experiment, we apply pre-model CDA: for every sentence containing a gendered term, we append a flipped-gender counterpart, increasing corpus size by $\sim 59\%$ and creating a gender-balanced corpus. To isolate whether pre-training changes are due to balancing rather than mere duplication, we run an ablation that duplicates an equal number of random sentences in the BabyLM corpus. Looking at the results in Figure 5, the CDA initially slows the model’s grasp of some linguistic concepts but, with longer training, reaches baseline performance while clearly curbing bias by preventing the steady growth seen in the baseline. On the other hand, the ablation results in the same performance drop while producing only a small bias reduction. Overall, this validates CDA as a sound debiasing strategy and highlights BabyLMs as a cost-effective platform for such controlled experiments.

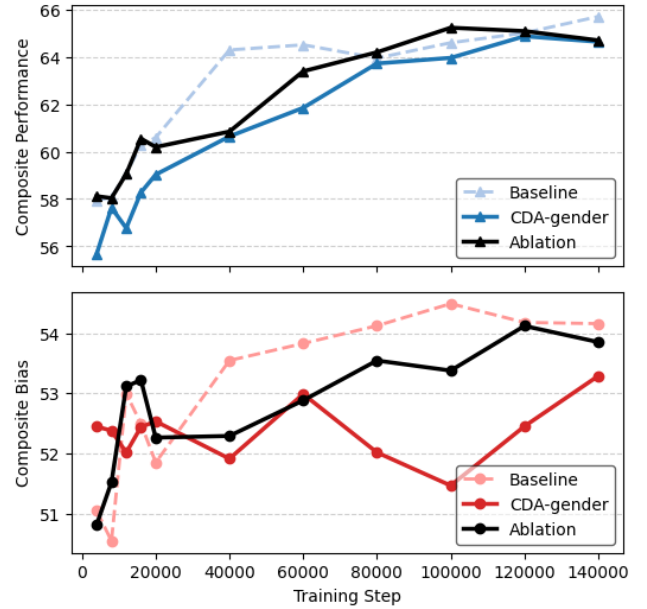


Figure 5: Bias and performance metrics evolution from pre-training on the CDA and CDA-ablation corpora

5.2 Corpus Toxicity Removal

Next, we examine the suggested but unproven claim that corpus toxicity directly drives model bias within LLMs. With 3.39% of BabyLM sentences being toxic or hateful, we seek to push toxicity to 0% via two interventions.

Firstly, we utilise an LLM-in-the-loop (Llama-3.3-70B) to rewrite toxic sentences while preserving text meaning, discarding only 0.42% that could not be detoxified (implementation detailed in Appendix F). This yielded a slight performance gain and a surprisingly small bias drop, likely because it imported the LLM’s style and biases into the corpus.

Second, we dropped all toxic sentences, which significantly reduced bias and slightly harmed performance.

Last, our ablation test removing an equal number of non-toxic sentences matched the performance drop but failed to match the bias decrease, confirming that eliminating toxicity itself, rather than corpus shrinkage, drives the debiasing. Figure 6 summarises these results, establishing a clear link between toxicity and standard bias and highlighting the advantage of our BabyLM-based approach that allows us to identify such behaviour.

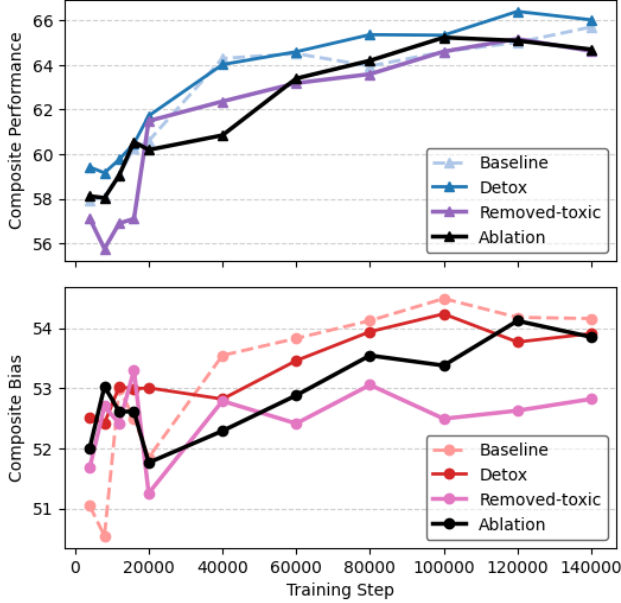


Figure 6: Evolution of bias and performance during pre-training under three corpus strategies: LLM detoxification, removing toxicity, and toxicity-removal ablation

5.3 Perturbation Augmentation

Finally, we evaluate perturbation augmentation (Qian et al. 2022), which uses a perturber LM to rewrite the corpus by randomly swapping demographic references (race and gender), thereby equalising topic–demographic co-occurrence (implementation discussed in Appendix G).

In the original study, pre-training RoBERTa on perturbed data reduced bias substantially with a slight performance gain. Using $\sim 800\times$ fewer GPU-hours, we closely reproduce these trends (Figure 7). The perturbation outperforms CDA in debiasing, likely by introducing greater lexical and syntactic variety and by covering more attributes. Moreover, it boosts performance, with the small gains possibly reflecting improved word order in lower-quality sentences caused by perturbation. Overall, even with this more complex debiasing strategy, BabyLM closely mirrors behaviour reported in far more resource-intensive experiments.

In conclusion, we observed that LTG-Baseline successfully reproduced the expected behaviours in all tested pre-model debiasing tasks. Furthermore, this setup enabled us not only to run the original experiments but also to introduce ablations

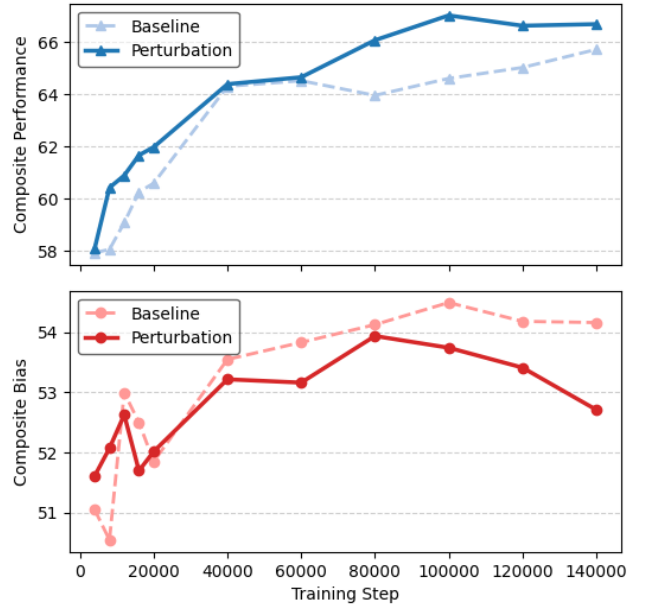


Figure 7: Bias and performance metrics evolution from pre-training on the perturbed corpus

that isolated and verified the impactful components of several debiasing strategies. Consequently, we were able to conduct, under a strict budget, experiments that had not previously been attempted.

Thus, we show that BabyLMs, and possibly other similarly positioned models, can be very versatile and representative tools for ascertaining the effectiveness of debiasing methods, greatly lowering both the cost and the time required to conduct the exploration-stage of debiasing experiments.

6 Findings

With our analysis complete, we notice several crucial trends emerging, which show a direct relevance to our effort to democratise research into pre-model debiasing.

Finding 1: BabyLM corpora and architectures share the same tendencies regarding bias–performance dynamics as standard LMs. Despite far smaller training sets and adapted architectures, BabyLMs acquire linguistic knowledge and biases along the same trajectories as standard BERT-style LMs. Across models, *composite performance* and *composite bias* are strongly and similarly correlated (Table 1). Corpus analysis shows that the BabyLM data, although more grammatically simple, contains the same well-known bias instigators (gender imbalance, topical skew, toxicity, etc.). In particular, male terms are over-represented relative to female terms, Christian terms dominate religious mentions, and toxicity and hate-speech rates are higher than in our reconstructed BERT pre-training corpus. Consequently, the SOTA BabyLMs match BERT on both performance and bias, while lower-resource BabyLMs preserve the same bias–ability correlation. This makes such models viable low-cost sandboxes for bias studies.

Finding 2: The debiasing behaviour of BabyLMs strongly correlates with that of standard LMs. Building on the first finding, we compared post-model and intra-model debiasing across BERT, LTG-BERT, and LTG-Baseline to test whether BabyLMs share the same debiasing behaviour as standard LMs. Across all debiasing methods, BabyLMs exhibited decreases in bias closely resembling BERT. Greater variance, on the other hand, was observed in performance changes, with post-model approaches being especially sensitive to architectural and pre-training differences between models. Nevertheless, canonical correlation analysis of bias and performance shifts across all methods showed near-perfect alignment between BERT and LTG-Baseline and strong alignment with LTG-BERT. Thus, we show that even a BabyLM trained on 30 GPU-hours reliably proxies debiasing dynamics observed in full-scale models.

Finding 3: BabyLMs enable informative and far more cost-effective pre-model debiasing research. Finally, we demonstrated the impact of our proposed method by running seven pre-training interventions on LTG-Baseline, reproducing established results and enabling new experiments and ablations, using them to both validate our approach and illustrate how it can help us understand and improve pre-model debiasing. We replicated that CDA reduces bias but can slow model learning, and that perturbation augmentation yields larger bias reductions and avoids performance loss. Beyond replication, we directly linked corpus toxicity to downstream bias, making this the first study in a pre-trained model setting to not only to imply but also to directly test the effect of removing toxic sentences on the resulting bias. With additional ablation experiments, we then isolated specific bias causes (gender imbalance, toxicity), showing that debiasing helps primarily by addressing these instigators rather than incidental corpus restructuring. These results position BabyLMs as viable and cost-effective sandboxes for systematic research into pre-model debiasing that can be conducted even with comparatively limited computational resources. This is especially crucial for underserved communities as the BabyLM ecosystem becomes increasingly multilingual (Jumelet et al. 2025), opening pathways to democratising debiasing research beyond English.

7 Conclusion

This work raised the issue that most LM debiasing research focuses only on pre-trained models with already-formed biased circuits. We argued that, to develop effective debiasing strategies, we must first understand how bias emerges in LMs and devise approaches that prevent its original formation. Such research is rare and inaccessible due to its prohibitive cost. To fix this, we proposed investigating debiasing dynamics in well-performing, low-resource, data-efficient models, such as BabyLMs.

Our experiments confirmed that BabyLMs use sufficient data and acquire intrinsic biases comparably to standard LMs when matched for performance. Their debiasing behaviour likewise mirrors other LMs, indicating that BabyLMs’ democratised pre-training setup does not disqualify them. With this, we moved from pre-training costs of

500 GPU-hours to 30 GPU-hours using BabyLM, creating a pathway for affordable debiasing research. With this, we reproduced previously reported results and added findings that solidified toxicity and bias imbalance as the root causes of LM bias.

Overall, our hope is that this research encourages the use of low-cost LMs that facilitate the exploration and mitigation of bias formation, enabling researchers to identify promising methods before committing to costly large-scale experiments. In doing so, it can advance the entire field of LM debiasing while allowing a broader range of under-represented perspectives to participate. This approach holds substantial promise for extending the research capabilities of academic labs and for enabling the exploration of new debiasing topics that were previously too specific or marginal to justify significant computational investment.

8 Limitations

One limitation of this study is the set of metrics it was able to use. BabyLMs limit us to more simplistic bias and performance-evaluation frameworks, since they do not offer sufficiently developed language and world understanding to support advanced extrinsic evaluations. Furthermore, due to the large amount of evaluation throughout the entire paper, frameworks like SuperGLUE, which take 2–3 hours to run on our hardware, are infeasible. In addition, all the metrics we used are English-specific. This is important to note because, even though non-English versions of these benchmarks exist (Névéol et al. 2022; Öztürk et al. 2023), they still do not cover low-resource languages, which is a major obstacle for LM bias research affecting these communities.

Secondly, it is well discussed that bias evaluations represent only a subset of the larger issue, and there might be biases where the different models’ behaviours actively differ but cannot be observed. Thus, we recommend using the composite bias metric to identify overall trends, but there may be a need to investigate or propose more specialised tests when tracing specific types of biases.

Related to this, while BabyLMs show promising alignment, there will be tasks that they cannot replicate. As noted in the paper, we propose them as a tool for *exploration*, helping to identify promising debiasing strategies. When these strategies are identified, they still need to be validated on larger models. Nevertheless, BabyLMs still allow us to skip the costly experimental phase.

Finally, it should be noted that BabyLMs were used as a promising and easily available set of architectures that displayed appropriate properties for the task at hand. We still encourage efforts to explore different corpora and architectures. CLMs, especially in the form of LLMs, cannot be replaced by MLMs. As such, future work should propose a model that democratises debiasing research for CLMs. Likewise, there might be corpora even better suited for testing debiasing methods than the BabyLM one. This study simply shows that the promising results obtained with the LTG-Baseline and the BabyLM corpus provide strong evidence that this path towards the democratisation of debiasing research is possible and promising.

References

- Antypas, D.; and Camacho-Collados, J. 2023. Robust Hate Speech Detection in Social Media: A Cross-Dataset Empirical Evaluation. In *Proceedings of the 7th Workshop on Online Abuse and Harms (WOAH)*, 231–242. Toronto, Canada: Association for Computational Linguistics.
- Bandy, J.; and Vincent, N. 2021. Addressing” documentation debt” in machine learning research: A retrospective datasheet for bookcorpus. *arXiv preprint arXiv:2105.05241*.
- Barbieri, F.; Camacho-Collados, J.; Espinosa Anke, L.; and Neves, L. 2020. TweetEval: Unified Benchmark and Comparative Evaluation for Tweet Classification. In Cohn, T.; He, Y.; and Liu, Y., eds., *Findings of the Association for Computational Linguistics: EMNLP 2020*, 1644–1650. Online: Association for Computational Linguistics.
- Bartl, M.; Nissim, M.; and Gatt, A. 2020. Unmasking Contextual Stereotypes: Measuring and Mitigating BERT’s Gender Bias. In Costa-jussà, M. R.; Hardmeier, C.; Webster, K.; and Radford, W., eds., *Proceedings of the Second Workshop on Gender Bias in Natural Language Processing*.
- BehnamGhader, P.; Adlakha, V.; Mosbach, M.; Bahdanau, D.; Chapados, N.; and Reddy, S. 2024. Llm2vec: Large language models are secretly powerful text encoders. *arXiv preprint arXiv:2404.05961*.
- Beltagy, I.; Lo, K.; and Cohan, A. 2019. SciBERT: A Pre-trained Language Model for Scientific Text. In Inui, K.; Jiang, J.; Ng, V.; and Wan, X., eds., *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 3615–3620. Hong Kong, China: Association for Computational Linguistics.
- Bender, E. M.; Gebru, T.; McMillan-Major, A.; and Shmitchell, S. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? . In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’21, 610–623. New York, NY, USA: Association for Computing Machinery. ISBN 9781450383097.
- BigScience-Workshop; Scao, T. L.; Fan, A.; Akiki, C.; Pavlick, E.; Ilić, S.; Hesslow, D.; Castagné, R.; Luccioni, A. S.; Yvon, F.; et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Blodgett, S. L.; Barocas, S.; Daumé III, H.; and Wallach, H. 2020. Language (Technology) is Power: A Critical Survey of “Bias” in NLP. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5454–5476. Online: Association for Computational Linguistics.
- Broad, N. 2022. English Wikipedia Dump: 2022-03-01 (wiki-20220301-en). Kaggle dataset. Accessed 11 Aug 2025.
- Cao, Y. T.; Pruksachatkun, Y.; Chang, K.-W.; Gupta, R.; Kumar, V.; Dhamala, J.; and Galstyan, A. 2022. On the Intrinsic and Extrinsic Fairness Evaluation Metrics for Contextualized Language Representations. In Muresan, S.; Nakov, P.; and Villavicencio, A., eds., *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 561–570. Dublin, Ireland: Association for Computational Linguistics.
- Chalkidis, I.; Fergadiotis, M.; Malakasiotis, P.; Aletras, N.; and Androutsopoulos, I. 2020. LEGAL-BERT: The Muppets straight out of Law School. In Cohn, T.; He, Y.; and Liu, Y., eds., *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2898–2904. Online: Association for Computational Linguistics.
- Chan, A.; Okolo, C. T.; Terner, Z.; and Wang, A. 2021. The limits of global inclusion in AI development. *arXiv preprint arXiv:2102.01265*.
- Chang, K.-W.; Prabhakaran, V.; and Ordonez, V. 2019. Bias and Fairness in Natural Language Processing. In Baldwin, T.; and Carpuat, M., eds., *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): Tutorial Abstracts*. Hong Kong, China: Association for Computational Linguistics.
- Charpentier, L.; Choshen, L.; Cotterell, R.; Gul, M. O.; Hu, M.; Jumelet, J.; Linzen, T.; Liu, J.; Mueller, A.; Ross, C.; et al. 2025. BabyLM Turns 3: Call for papers for the 2025 BabyLM workshop. *arXiv preprint arXiv:2502.10645*.
- Charpentier, L. G. G.; and Samuel, D. 2024. GPT or BERT: why not both? *arXiv preprint arXiv:2410.24159*.
- Cheng, Z.; Aralikkatte, R.; Porada, I.; Spinoso-Di Piano, C.; and Cheung, J. C. 2023. McGill BabyLM Shared Task Submission: The Effects of Data Formatting and Structural Biases. In Warstadt, A.; Mueller, A.; Choshen, L.; Wilcox, E.; Zhuang, C.; Ciro, J.; Mosquera, R.; Paranjabe, B.; Williams, A.; Linzen, T.; and Cotterell, R., eds., *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*, 207–220. Singapore: Association for Computational Linguistics.
- Cottier, B.; Rahman, R.; Fattorini, L.; Maslej, N.; Besiroglu, T.; and Owen, D. 2024. The rising costs of training frontier AI models. *arXiv preprint arXiv:2405.21015*.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 4171–4186.
- Di Liello, L. 2022. BookCorpusOpen. Hugging Face dataset. Revision: main@edb74e6. Accessed 11 Aug 2025.
- Diehl Martinez, R.; Goriely, Z.; McGovern, H.; Davis, C.; Caines, A.; Buttery, P.; and Beinborn, L. 2023. CLIMB – Curriculum Learning for Infant-inspired Model Building. In Warstadt, A.; Mueller, A.; Choshen, L.; Wilcox, E.; Zhuang, C.; Ciro, J.; Mosquera, R.; Paranjabe, B.; Williams, A.; Linzen, T.; and Cotterell, R., eds., *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*, 112–127. Singapore: Association for Computational Linguistics.
- Ekman, P.; et al. 1999. Basic emotions. *Handbook of cognition and emotion*, 98(45-60): 16.
- Field, A.; Blodgett, S. L.; Waseem, Z.; and Tsvetkov, Y. 2021. A Survey of Race, Racism, and Anti-Racism in NLP. In

- Zong, C.; Xia, F.; Li, W.; and Navigli, R., eds., *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 1905–1925. Online: Association for Computational Linguistics.
- Gal, Y.; and Ghahramani, Z. 2016. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In Balcan, M. F.; and Weinberger, K. Q., eds., *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, 1050–1059. New York, New York, USA: PMLR.
- Gallegos, I. O.; Rossi, R. A.; Barrow, J.; Tanjim, M. M.; Kim, S.; Dernoncourt, F.; Yu, T.; Zhang, R.; and Ahmed, N. K. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, 50(3): 1097–1179.
- Gelles, R.; Kinoshita, V.; Musser, M.; and Dunham, J. 2024. Resource Democratization: Is Compute the Binding Constraint on AI Research? *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(18): 19840–19848.
- Gilkerson, J.; Richards, J. A.; Warren, S. F.; Montgomery, J. K.; Greenwood, C. R.; Kimbrough Oller, D.; Hansen, J. H. L.; and Paul, T. D. 2017. Mapping the Early Language Environment Using All-Day Recordings and Automated Analysis. *American Journal of Speech-Language Pathology*, 26(2): 248–265.
- Goldfarb-Tarrant, S.; Marchant, R.; Sánchez, R. M.; Pandya, M.; and Lopez, A. 2020. Intrinsic bias metrics do not correlate with application bias. *arXiv preprint arXiv:2012.15859*.
- Gonen, H.; and Goldberg, Y. 2019. Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. *arXiv preprint arXiv:1903.03862*.
- Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Vaughan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Guo, Y.; Guo, M.; Su, J.; Yang, Z.; Zhu, M.; Li, H.; Qiu, M.; and Liu, S. S. 2024. Bias in large language models: Origin, evaluation, and mitigation. *arXiv preprint arXiv:2411.10915*.
- Gupta, V.; Narayanan Venkit, P.; Wilson, S.; and Passonneau, R. 2024. Sociodemographic Bias in Language Models: A Survey and Forward Path. In Faleńska, A.; Basta, C.; Costa-jussà, M.; Goldfarb-Tarrant, S.; and Nozza, D., eds., *Proceedings of the 5th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, 295–322. Bangkok, Thailand: Association for Computational Linguistics.
- Hansen, L.; Olsen, L. R.; and Enevoldsen, K. 2023. TextDescriptives: A Python package for calculating a large variety of metrics from text. *Journal of Open Source Software*, 8(84): 5153.
- Hanu, L.; and Unitary AI. 2020. Detoxify.
- He, P.; Liu, X.; Gao, J.; and Chen, W. 2020. DeBERTa: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.
- Hendrycks, D.; Burns, C.; Basart, S.; Zou, A.; Mazeika, M.; Song, D.; and Steinhardt, J. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Hoffmann, J.; Borgeaud, S.; Mensch, A.; Buchatskaya, E.; Cai, T.; Rutherford, E.; Casas, D. d. L.; Hendricks, L. A.; Welbl, J.; Clark, A.; et al. 2022. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*.
- Hotelling, H. 1936. Relations Between Two Sets of Variates. *Biometrika*, 28(3/4): 321–377.
- Hu, M. Y.; Mueller, A.; Ross, C.; Williams, A.; Linzen, T.; Zhuang, C.; Cotterell, R.; Choshen, L.; Warstadt, A.; and Wilcox, E. G. 2024a. Findings of the Second BabyLM Challenge: Sample-Efficient Pretraining on Developmentally Plausible Corpora. In Hu, M. Y.; Mueller, A.; Ross, C.; Williams, A.; Linzen, T.; Zhuang, C.; Choshen, L.; Cotterell, R.; Warstadt, A.; and Wilcox, E. G., eds., *The 2nd BabyLM Challenge at the 28th Conference on Computational Natural Language Learning*, 1–21. Miami, FL, USA: Association for Computational Linguistics.
- Hu, M. Y.; Mueller, A.; Ross, C.; Williams, A.; Linzen, T.; Zhuang, C.; Cotterell, R.; Choshen, L.; Warstadt, A.; and Wilcox, E. G. 2024b. Findings of the second BabyLM challenge: Sample-efficient pretraining on developmentally plausible corpora. *arXiv preprint arXiv:2412.05149*.
- Ivanova, A. A.; Sathe, A.; Lipkin, B.; Kumar, U.; Radkani, S.; Clark, T. H.; Kauf, C.; Hu, J.; Pramod, R.; Grand, G.; et al. 2024. Elements of World Knowledge (EWOK): A cognition-inspired framework for evaluating basic world knowledge in language models. *arXiv preprint arXiv:2405.09605*.
- Jiao, X.; Yin, Y.; Shang, L.; Jiang, X.; Chen, X.; Li, L.; Wang, F.; and Liu, Q. 2019. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*.
- Joshi, M.; Chen, D.; Liu, Y.; Weld, D. S.; Zettlemoyer, L.; and Levy, O. 2020. SpanBERT: Improving Pre-training by Representing and Predicting Spans. *Transactions of the Association for Computational Linguistics*, 8: 64–77.
- Jumelet, J.; Fourtassi, A.; Haga, A.; Bunzeck, B.; Shandilya, B.; Galvan-Sosa, D.; Haznitrana, F. G.; Padovani, F.; Meyer, F.; Hu, H.; et al. 2025. BabyBabelLM: A Multilingual Benchmark of Developmentally Plausible Training Data. *arXiv preprint arXiv:2510.10159*.
- Köksal, A.; Yalcin, O.; Akbiyik, A.; Kilavuz, M.; Korhonen, A.; and Schuetze, H. 2023. Language-Agnostic Bias Detection in Language Models with Bias Probing. In Bouamor, H.; Pino, J.; and Bali, K., eds., *Findings of the Association for Computational Linguistics: EMNLP 2023*, 12735–12747. Singapore: Association for Computational Linguistics.
- Kong, L.; Jiang, H.; Zhuang, Y.; Lyu, J.; Zhao, T.; and Zhang, C. 2020. Calibrated Language Model Fine-Tuning for In- and Out-of-Distribution Data. In Webber, B.; Cohn, T.; He, Y.; and Liu, Y., eds., *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1326–1340. Online: Association for Computational Linguistics.

- Korinek, A.; and Vipra, J. 2024. Concentrating intelligence: scaling and market structure in artificial intelligence*. *Economic Policy*, 40(121): 225–256.
- Li, Y.; Du, M.; Song, R.; Wang, X.; and Wang, Y. 2024. Data-Centric Explainable Debiasing for Improving Fairness in Pre-trained Language Models. In Ku, L.-W.; Martins, A.; and Srikumar, V., eds., *Findings of the Association for Computational Linguistics: ACL 2024*, 3773–3786. Bangkok, Thailand: Association for Computational Linguistics.
- Liang, P. P.; Li, I. M.; Zheng, E.; Lim, Y. C.; Salakhutdinov, R.; and Morency, L.-P. 2020. Towards Debiasing Sentence Representations. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5502–5515. Online: Association for Computational Linguistics.
- Linzen, T. 2020. How Can We Accelerate Progress Towards Human-like Linguistic Generalization? In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5210–5217. Online: Association for Computational Linguistics.
- Lu, K.; Mardziel, P.; Wu, F.; Amancharla, P.; and Datta, A. 2018. Gender Bias in Neural Natural Language Processing. *arXiv preprint arXiv:1807.11714*.
- May, C.; Wang, A.; Bordia, S.; Bowman, S. R.; and Rudinger, R. 2019. On Measuring Social Biases in Sentence Encoders. In Burstein, J.; Doran, C.; and Solorio, T., eds., *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 622–628. Minneapolis, Minnesota: Association for Computational Linguistics.
- Meade, N.; Poole-Dayana, E.; and Reddy, S. 2021. An empirical survey of the effectiveness of debiasing techniques for pre-trained language models. *arXiv preprint arXiv:2110.08527*.
- Mendelson, M.; and Belinkov, Y. 2021. Debiasing Methods in Natural Language Understanding Make Bias More Accessible. In Moens, M.-F.; Huang, X.; Specia, L.; and Yih, S. W.-t., eds., *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 1545–1557. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Meta AI. 2025. The Llama 4 Herd: The Beginning of a New Era of Natively Multimodal AI Innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>. Accessed 29 Jul 2025.
- Mohammad, S.; and Turney, P. 2010. Emotions Evoked by Common Words and Phrases: Using Mechanical Turk to Create an Emotion Lexicon. In Inkpen, D.; and Strapparava, C., eds., *Proceedings of the NAACL HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text*, 26–34. Los Angeles, CA: Association for Computational Linguistics.
- Nadeem, M.; Bethke, A.; and Reddy, S. 2020. StereoSet: Measuring stereotypical bias in pretrained language models. *arXiv preprint arXiv:2004.09456*.
- Nangia, N.; Vania, C.; Bhalariao, R.; and Bowman, S. R. 2020. CrowS-pairs: A challenge dataset for measuring social biases in masked language models. *arXiv preprint arXiv:2010.00133*.
- Névéol, A.; Dupont, Y.; Bezançon, J.; and Fort, K. 2022. French CrowS-Pairs: Extending a challenge dataset for measuring social bias in masked language models to a language other than English. In Muresan, S.; Nakov, P.; and Villavicencio, A., eds., *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 8521–8531. Dublin, Ireland: Association for Computational Linguistics.
- Öztürk, I. T.; Nedelchev, R.; Heumann, C.; Arias, E. G.; Roger, M.; Bischl, B.; and Aßenmacher, M. 2023. How Different is Stereotypical Bias Across Languages? In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 209–229. Springer.
- Park, S.; Choi, K.; Yu, H.; and Ko, Y. 2023. Never Too Late to Learn: Regularizing Gender Bias in Coreference Resolution. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining, WSDM ’23*, 15–23. New York, NY, USA: Association for Computing Machinery. ISBN 9781450394079.
- Qian, R.; Ross, C.; Fernandes, J.; Smith, E. M.; Kiela, D.; and Williams, A. 2022. Perturbation Augmentation for Fairer NLP. In Goldberg, Y.; Kozareva, Z.; and Zhang, Y., eds., *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 9496–9521. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics.
- Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; and Sutskever, I. 2019. Language Models are Unsupervised Multitask Learners.
- Ranaldi, L.; Ruzzetti, E. S.; Venditti, D.; Onorati, D.; and Zanzotto, F. M. 2024. A Trip Towards Fairness: Bias and De-Biasing in Large Language Models. In Bollegala, D.; and Schwartz, V., eds., *Proceedings of the 13th Joint Conference on Lexical and Computational Semantics (*SEM 2024)*, 372–384. Mexico City, Mexico: Association for Computational Linguistics.
- Ravfogel, S.; Elazar, Y.; Gonen, H.; Twiton, M.; and Goldberg, Y. 2020. Null It Out: Guarding Protected Attributes by Iterative Nullspace Projection. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 7237–7256. Online: Association for Computational Linguistics.
- Rein, D.; Hou, B. L.; Stickland, A. C.; Petty, J.; Pang, R. Y.; Dirani, J.; Michael, J.; and Bowman, S. R. 2024. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*.
- Rudinger, R.; Naradowsky, J.; Leonard, B.; and Van Durme, B. 2018. Gender bias in coreference resolution. *arXiv preprint arXiv:1804.09301*.
- Samuel, D.; Kutuzov, A.; Øvrelid, L.; and Velldal, E. 2023. Trained on 100 million words and still in shape: BERT meets British National Corpus. *arXiv preprint arXiv:2303.09859*.

- Sathish, V.; Lin, H.; Kamath, A. K.; and Nyayachavadi, A. 2024. Llempower: Understanding disparities in the control and access of large language models. *arXiv preprint arXiv:2404.09356*.
- Shazeer, N. 2020. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*.
- Shleifer, S.; Weston, J.; and Ott, M. 2021. Normformer: Improved transformer pretraining with extra normalization. *arXiv preprint arXiv:2110.09456*.
- Sosto, M.; and Barrón-Cedeño, A. 2024. Queerbench: Quantifying discrimination in language models toward queer identities. *arXiv preprint arXiv:2406.12399*.
- Sun, L.; Mao, C.; Hofmann, V.; and Bai, X. 2025. Aligned but Blind: Alignment Increases Implicit Bias by Reducing Awareness of Race. *arXiv preprint arXiv:2506.00253*.
- Van Nguyen, C.; Shen, X.; Aponte, R.; Xia, Y.; Basu, S.; Hu, Z.; Chen, J.; Parmar, M.; Kunapuli, S.; Barrow, J.; et al. 2024. A survey of small language models. *arXiv preprint arXiv:2410.20011*.
- Wang, A.; Pruksachatkun, Y.; Nangia, N.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32.
- Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. R. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.
- Warstadt, A.; Mueller, A.; Choshen, L.; Wilcox, E.; Zhuang, C.; Ciro, J.; Mosquera, R.; Paranjabe, B.; Williams, A.; Linzen, T.; and Cotterell, R. 2023a. Findings of the BabyLM Challenge: Sample-Efficient Pretraining on Developmentally Plausible Corpora. In Warstadt, A.; Mueller, A.; Choshen, L.; Wilcox, E.; Zhuang, C.; Ciro, J.; Mosquera, R.; Paranjabe, B.; Williams, A.; Linzen, T.; and Cotterell, R., eds., *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*, 1–34. Singapore: Association for Computational Linguistics.
- Warstadt, A.; Mueller, A.; Choshen, L.; Wilcox, E.; Zhuang, C.; Ciro, J.; Mosquera, R.; Paranjabe, B.; Williams, A.; Linzen, T.; and Cotterell, R., eds. 2023b. *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*. Singapore: Association for Computational Linguistics.
- Warstadt, A.; Parrish, A.; Liu, H.; Mohananey, A.; Peng, W.; Wang, S.-F.; and Bowman, S. R. 2020. BLiMP: The Benchmark of Linguistic Minimal Pairs for English. *Transactions of the Association for Computational Linguistics*, 8: 377–392.
- Webster, K.; Recasens, M.; Axelrod, V.; and Baldridge, J. 2018. Mind the GAP: A Balanced Corpus of Gendered Ambiguous Pronouns. *Transactions of the Association for Computational Linguistics*, 6: 605–617.
- Webster, K.; Wang, X.; Tenney, I.; Beutel, A.; Pitler, E.; Pavlick, E.; Chen, J.; Chi, E.; and Petrov, S. 2020. Measuring and reducing gendered correlations in pre-trained models. *arXiv preprint arXiv:2010.06032*.
- Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; Davison, J.; Shleifer, S.; von Platen, P.; Ma, C.; Jernite, Y.; Plu, J.; Xu, C.; Scao, T. L.; Gugger, S.; Drame, M.; Lhoest, Q.; and Rush, A. M. 2020. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 38–45. Online: Association for Computational Linguistics.
- Xie, Z.; and Lukasiewicz, T. 2023. An empirical analysis of parameter-efficient methods for debiasing pre-trained language models. *arXiv preprint arXiv:2306.04067*.
- Yuan, S.; Nie, E.; Kouba, L.; Kanger, A. Y.; Schmid, H.; Schütze, H.; and Färber, M. 2025. LLM in the Loop: Creating the PARADEHATE Dataset for Hate Speech Detoxification. *arXiv preprint arXiv:2506.01484*.
- Zakizadeh, M.; and Pilehvar, M. T. 2025. Blind Men and the Elephant: Diverse Perspectives on Gender Stereotypes in Benchmark Datasets. *arXiv preprint arXiv:2501.01168*.
- Zhao, J.; Ding, Y.; Jia, C.; Wang, Y.; and Qian, Z. 2024. Gender bias in large language models across multiple languages. *arXiv preprint arXiv:2403.00277*.
- Zhao, J.; Wang, T.; Yatskar, M.; Cotterell, R.; Ordonez, V.; and Chang, K.-W. 2019. Gender Bias in Contextualized Word Embeddings. In Burstein, J.; Doran, C.; and Solorio, T., eds., *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 629–634. Minneapolis, Minnesota: Association for Computational Linguistics.
- Zhao, J.; Wang, T.; Yatskar, M.; Ordonez, V.; and Chang, K.-W. 2018a. Gender bias in coreference resolution: Evaluation and debiasing methods. *arXiv preprint arXiv:1804.06876*.
- Zhao, J.; Wang, T.; Yatskar, M.; Ordonez, V.; and Chang, K.-W. 2018b. Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods. In Walker, M.; Ji, H.; and Stent, A., eds., *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, 15–20. New Orleans, Louisiana: Association for Computational Linguistics.
- Zmigrod, R.; Mielke, S. J.; Wallach, H.; and Cotterell, R. 2019. Counterfactual Data Augmentation for Mitigating Gender Stereotypes in Languages with Rich Morphology. In Korhonen, A.; Traum, D.; and Màrquez, L., eds., *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 1651–1661. Florence, Italy: Association for Computational Linguistics.

A Evaluated Models

In order to run our experiment, we must analyse a sufficiently large set of BabyLMs and standard LMs, so that we can observe meaningful trends.

Regarding BabyLMs, we collect every notable BabyLM that is available online with executable code, utilising our overview in Section 2.3 (Samuel et al. 2023). It should be

noted that some models require altered data input pipelines, making them impractical. Likewise, several promising papers never released their models. To ensure comparability, we take only models from the strict track. The final list is shown in Table 7.

For standard LMs, we also need to observe how bias scales with performance. Testing just a few primary models would not establish a trend. To resolve this, we examine the popular architectures used to create BabyLMs together with other variants of these architectures trained on different corpora (e.g. legal documents (Chalkidis et al. 2020), Twitter (Barbieri et al. 2020), scientific papers (Beltagy, Lo, and Cohan 2019)). All these models are listed in Table 8.

B Model Alignment

Running all models through our pipelines, we obtain Pearson correlations between the bias metrics and the performance metrics. Figures 8 and 9 show the correlation results for BabyLMs and standard LMs respectively. In both groups, performance and bias metrics are strongly positively correlated, while the two bias scores have a positive but more limited correlation with each other. BabyLMs show especially noisy behaviour since some of them have very low word knowledge needed to display biases. This leads us to establish the composite metrics, which more completely capture the overall behaviour.

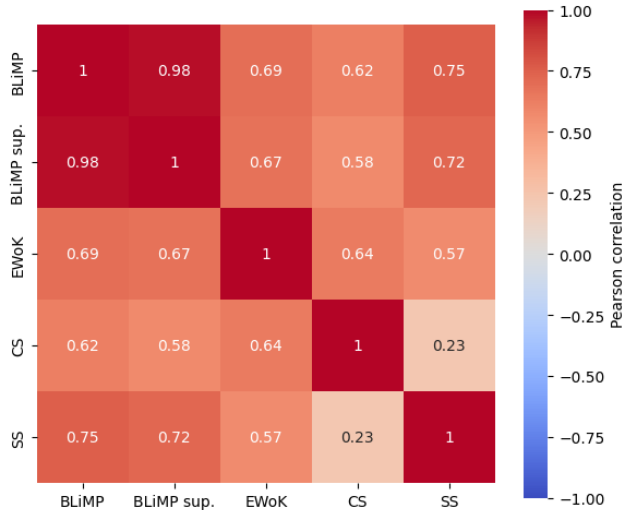


Figure 8: Pearson correlations (BabyLM models, $N = 9$)

C Corpora

This section details all other experiments conducted with the aim of exploring and comparing the corpora of BabyLMs and standard LMs.

C.1 Structural and Syntactic Metrics

We start the analysis by examining the basic linguistic metrics present in the corpora, trying to understand the composition

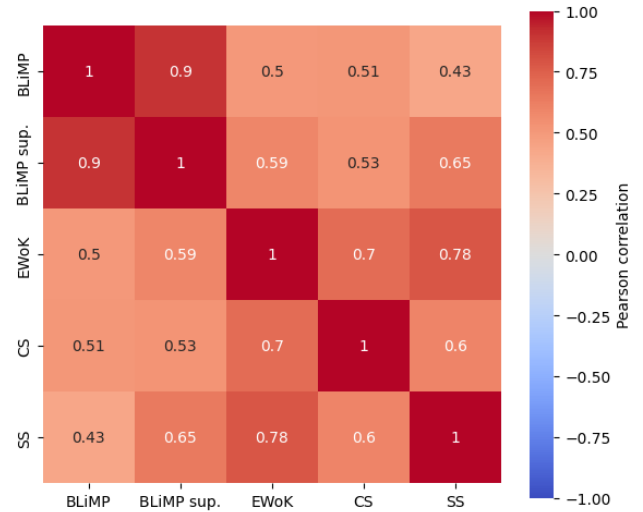


Figure 9: Pearson correlations (standard models, $N = 16$)

and complexity of each dataset. This might help us understand whether the BabyLM corpus is sufficiently coherent to form a full range of biases.

To do this, we leverage the TextDescriptives library by Hansen, Olsen, and Enevoldsen (2023), which allows us to derive more than sixty syntactic and discourse-level statistics. Because all corpora exceed the toolkit’s token limit for text processing, we split the text into chunks of 8192 tokens, retaining large enough samples to assess the more structural metrics while keeping the problem tractable. In the end, the outputs per chunk are averaged into final values representing the entire corpus. This ensures that no corpus is penalised for its size while capturing the structural and syntactic profile needed for analysis.

Part-of-Speech Profiles PoS statistics give us a compact high-level overview, informing us about the average type of the texts present in each corpus. We are especially interested in the frequency of function words (e.g., pronouns, determiners), which organise discourse, versus the frequency of content words (e.g., nouns, verbs), which carry meaning.

Figure 10 shows us that BabyLM displays a lower frequency of nouns and pronouns, signalling a simpler and more narrative-driven corpus.

Figure 11 serves to further highlight this contrast. BabyLM’s proportion of pronouns to nouns is, unlike that of BERT, almost at parity. Thus, in the BabyLM corpus more tokens are used to refer to individuals rather than describing unique objects or concepts, again indicating that it is linguistically simpler than the BERT corpus.

Linguistic Complexity and Readability Since it appears that the main difference is text complexity, we investigate the readability metrics.

We choose to use the Flesch–Kincaid Grade Level (FKGL), which is computed from the average sentence length and the average number of syllables per word. Its final value then corresponds to grades in the US school system, ranging

Model key	Hugging Face ID
BabyLM2024	jdebene/BabyLM2024
elc-bert	lgcharpe/ELC_BERT_baby_100M
cambridge-climb	cambridge-climb/baseline-roberta_pre_lay er_norm-model
gpt-bert-babylm	ltg/gpt-bert-babylm-base
ltg-bert-babylm	ltg/ltg-bert-babylm
ltgbert-100m-2024	babylm/ltgbert-100m-2024
roberta-base-strict-2023	babylm/roberta-base-strict-2023
babylm-100m-2024	babylm/babylm-100m-2024
baby-llama-2-345m	JLTastet/baby-llama-2-345m

Table 7: Examined BabyLMs with their Hugging Face IDs

Model key	Hugging Face ID
BiomedBERT-abstract	microsoft/BiomedNLP-BiomedBERT-base-unc ased-abstract
BiomedBERT-abstract-fulltext	microsoft/BiomedNLP-BiomedBERT-base-unc ased-abstract-fulltext
DialoGPT-large	microsoft/DialoGPT-large
DialoGPT-medium	microsoft/DialoGPT-medium
DialoGPT-small	microsoft/DialoGPT-small
bert-base-cased	bert-base-cased
bert-base-uncased	bert-base-uncased
bert-for-patents	anferico/bert-for-patents
BNC bert	ltg/ltg-bert-bnc
finbert	yyanghkust/finbert-pretrain
gpt2-medium	openai-community/gpt2-medium
gpt2-xl	openai-community/gpt2-xl
legal-bert	nlpaueb/legal-bert-base-uncased
roberta-base	FacebookAI/roberta-base
scibert	allenai/scibert_scivocab_uncased
twitter-roberta-base	cardiffnlp/twitter-roberta-base

Table 8: Examined standard models with their Hugging Face IDs

between 0 and 18. With the results shown in Table 9, we can once again see that BERT has the linguistically more advanced texts and BabyLM has the simpler ones. The FKGL places the BabyLM corpus between 5th and 6th grades, which corresponds to students aged 10-12 and thus aligns with what would be expected for BabyLM.

Unsurprisingly, when investigating the type-token ratio shown in Table 10 to ascertain the richness of the vocabulary, we see that BERT has the more diverse vocabulary. BabyLM is again much simpler. This is likely due to its enforced child-alignment.

Coherence Finally, we use the first-order coherence to examine how tightly each text corpus adheres to a topic. In Table 11, we see that the corpora score highly, with both being similarly consistent.

Overall, the BabyLM corpus is markedly simpler in text complexity yet maintains coherence comparable to BERT. Consequently, if it contains sufficient bias-inducing material, it is likely adequate to impart those biases to the model.

C.2 Toxicity, Sentiment, and Emotions

Furthermore, we need to understand the tone of the corpora and the context in which they present their different topics. While no exact link has been created, there is some evidence of the fact that displaying information in a negative or toxic light pushes the model to develop stronger biases (BigScience-Workshop et al. 2022). As such, we must examine the presence of toxicity, hate-speech, sentiment, and the overall emotional composition throughout the corpora.

Emotion Scores Starting with emotion detection, we are interested in measuring Ekman’s core emotions: joy, sadness, fear, surprise, anger, and disgust (Ekman et al. 1999). To this end, we take advantage of the *NRC Emotion Lexicon*, which links individual words to various emotions (Mohammad and Turney 2010), allowing us to infer the overall emotional tone of the corpus. The normalised *emotion score* of a corpus is calculated by dividing the number of words expressing a specific emotion by the total number of eligible words. This *emotion score* is computed for all six emotions across all our

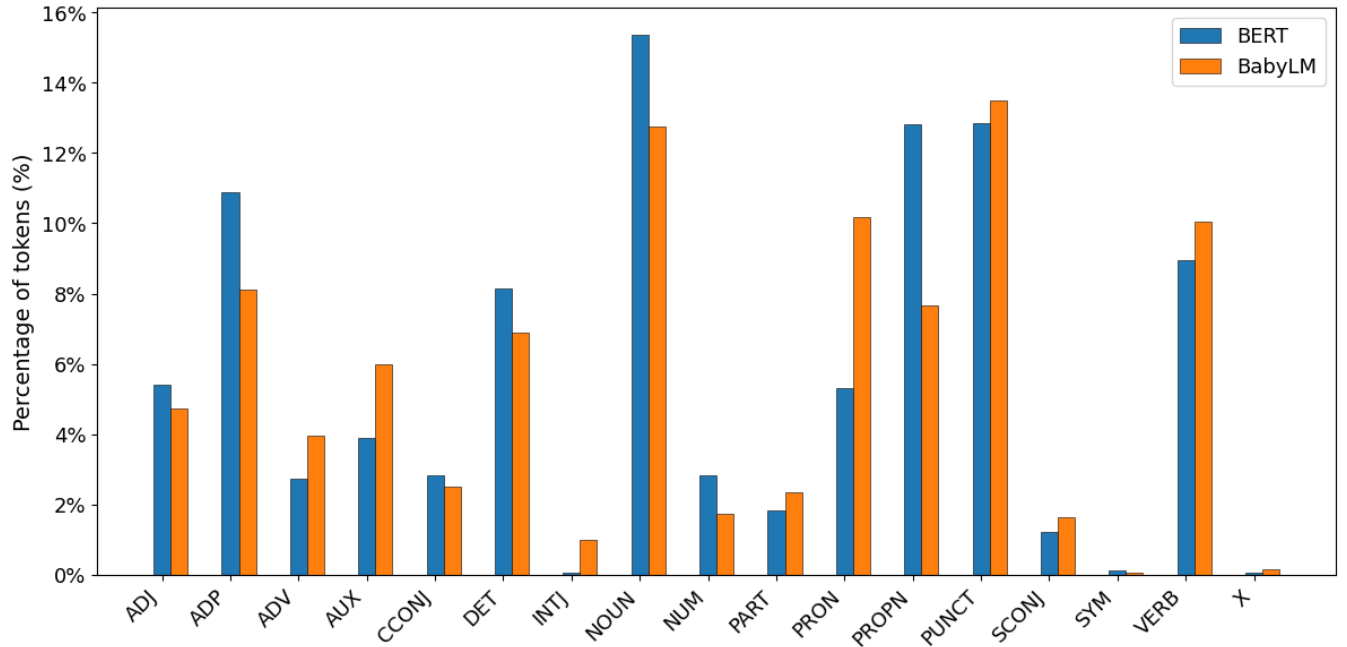


Figure 10: Part-of-Speech distribution across corpora

Corpus	Avg. syllables per word	Avg. word length	Avg. sentence length	FKGL
BabyLM	1.25	4.01	15.91	5.40
BERT	1.38	4.57	19.84	8.43

Table 9: Readability and length statistics for the BERT and BabyLM corpora.

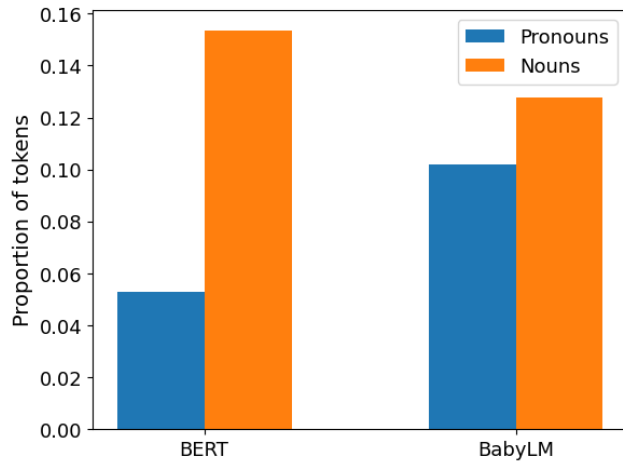


Figure 11: Pronoun vs. noun distribution across corpora

corpora.

The results of the lexical analysis can be seen in Figure 12. While the emotions are mostly balanced across the corpora, BabyLM is consistently more emotional, especially in terms of joy and surprise. However, in both cases, the difference

Corpus	Type-token ratio
BabyLM	0.33
BERT	0.53

Table 10: Type-token ratio for the BERT and BabyLM corpora

Corpus	First-order coherence
BabyLM	0.811
BERT	0.802

Table 11: First-order coherence scores for the BERT and BabyLM corpora

is mild and both are positive emotions, which are not the carriers of bias.

Sentiment Works like the one by (Köksal et al. 2023) have indicated that sentiment is passed from corpora into the biases of models, forging positive or negative connections that might end up strengthening stereotypes. Because of this, we measure the sentiments in each corpus, both the overall ones and those connected to specific topics.

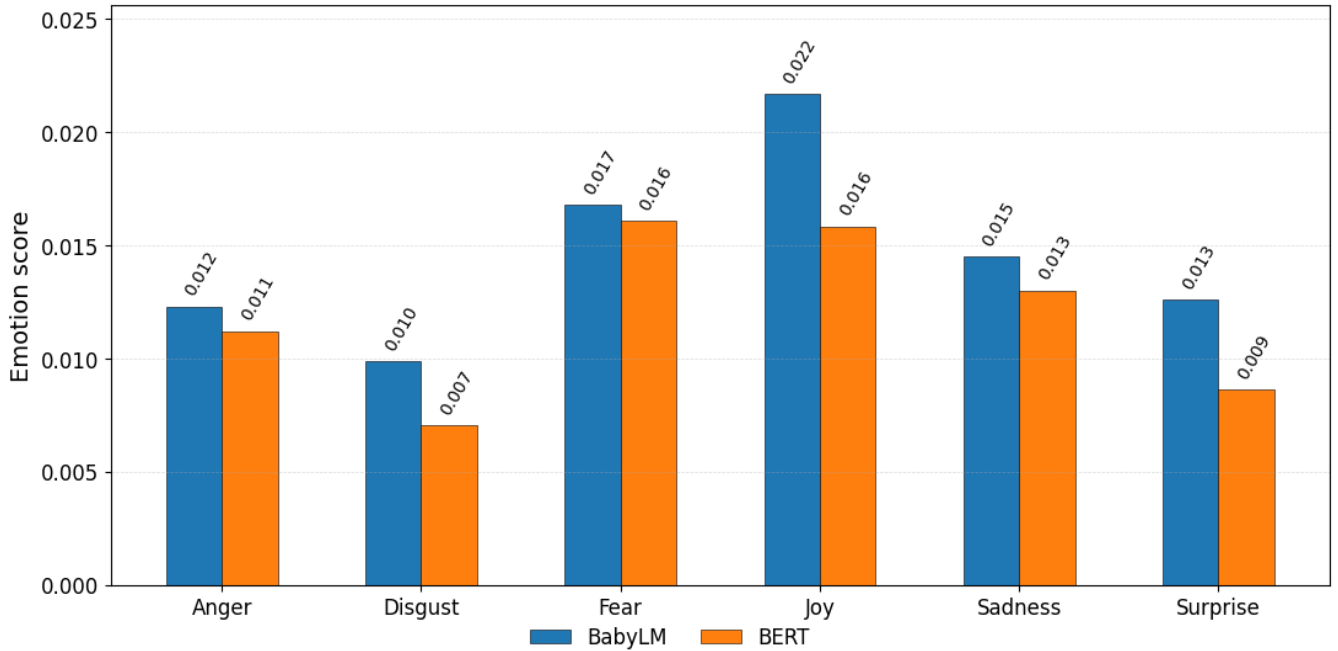


Figure 12: Distribution of emotions across the corpora using the emotion score obtained by lexical analysis

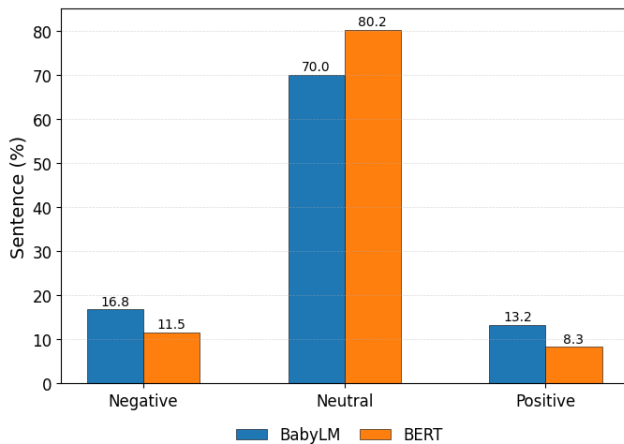


Figure 13: Distribution of sentences per sentiment across corpora

We use the established RoBERTa-based solution by Barbi-eri et al. (2020), which is one of the most popular publicly available sentiment analysis models. With it, we classify each corpus as positive, negative, or neutral. Finally, we compute the *sentiment score* in the same fashion we computed the *emotion score*, only switching from word-based to sentence-based evaluation.

Figure 13 shows the overall sentiment profile of the corpora. The BabyLM corpus is consistently more emotional, reinforcing that it is well-positioned to transfer the biases into the model despite its child-alignment.

Examining the results more closely, Figure 14 displays the

sentiment of sentences containing a bias-inducing word from a specific category for each corpus. To calculate this, we take all eligible sentences per topic and subtract the percentage of negative sentences from the percentage of positive sentences. Here, we can see that the BabyLM corpus keeps being mostly more negative than the BERT corpus, with the exception of the Muslim topic. While it carries more negativity, the ratios between topics are similar, meaning that it retains the same bias orientation. The only true divergence is higher negativity towards LGBTQ+ topics.

In summary, BabyLM proves to be more emotional and negative than BERT, while retaining the same larger trends, meaning that the corpora are largely compatible. There is no evidence that BabyLM should not be able to teach model biases.

C.3 Specifications Regarding Toxicity

While this topic has already been covered in Section 4.1, we want to provide a validation of the result. Given that we have established that the BabyLM corpus mostly contains simple sentences, the toxicity and hate-speech results in Table 5 raise the question of whether the toxicity and hate-speech labels are overly sensitive. To investigate this, we sampled 200 such sentences, identifying that the vast majority are indeed toxic, containing slurs or highly aggressive expressions. For illustration, we list some sampled examples in Table 12 (**Contains explicit offensive statements**). These examples display blatant racist and sexualised terms, showing that, despite its child-alignment, BabyLM contains diverse and toxic texts.

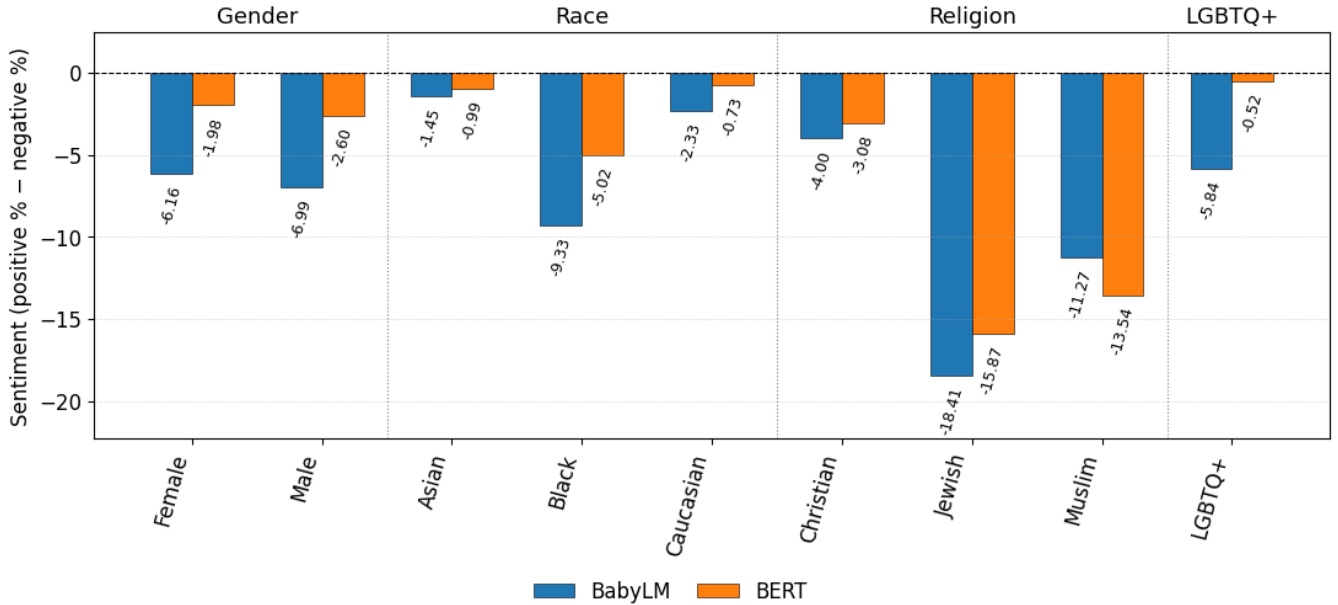


Figure 14: Sentiment stances to the selected topics across corpora

Sentence
Her earhole ain't big enough for fucking!
cause they don't want to look like morons too.
He was a black nigger what?

Table 12: Sampled examples of toxic sentences from the BabyLM corpus

D INLP Reduces Performance Penalty for Over-Fitted BabyLMs

The INLP experiment has also revealed an important piece of behaviour from LTG-BERT that concerns the BabyLM challenge itself. As mentioned, INLP helped improve LTG-BERT’s performance, likely by reducing the issues connected to over-fitting on its training corpus.

This is notable because LTG-BERT is the best-performing BabyLM-class MLM, with the SOTA BabyLM being GPT-BERT, which was trained with the same training data-to-epoch ratio. Thus, if INLP debiasing makes LTG-BERT surpass the SOTA model on our performance tasks, as shown in Figure 15, it could even yield further improvements to other over-fitted BabyLMs, generating further performance improvements.

E Implementing LTG-Baseline Pre-training

We must define the exact specifics of training our own LTG-Baseline model. This requires us to prepare the corpus, define a training script, and specify the exact hyperparameters. Following the specifications from the authors (Hu et al. 2024b), we use the predefined BabyLM corpus established in Section 4.1, and adopt the LTG architecture by Samuel et al. (2023) as the starting point. However, an issue with further

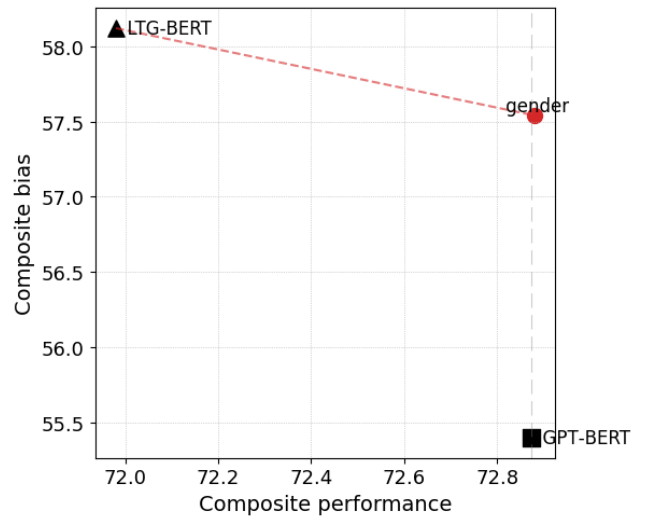


Figure 15: Performance improvement of INLP-debiased LTG-BERT compared with the GPT-BERT BabyLM

implementation is that the authors have not published the LTG-Baseline training script or its hyperparameters.

To resolve this, we reached out directly to the authors, who supplied us with all the necessary details. They used the standard MLM training script by Wolf et al. (2020) with hyperparameters, which are provided in Table 13.

Furthermore, for purposes of model comparison, epochs cease to be a representative unit when we alter the corpus. To remedy this, we note that, in the standard training, 20 epochs translate to roughly 140,000 training steps. Thus, we compare all of our models on the interval from 0 to 140,000

Hyperparameter	Value
max_seq_length	128
per_device_train_batch_size	128
num_train_epochs	20
learning_rate	5e-4
adam_beta1	0.9
adam_beta2	0.999
adam_epsilon	1e-8
max_grad_norm	1
warmup_steps	0

Table 13: LTG-Baseline pre-training hyperparameters

steps, removing any unfair advantage stemming from an expanded corpus.

F LLM-in-the-loop Debiasing Implementation

The LLM-in-the-loop detoxification approach has garnered interest because it allows us to remove toxicity without removing information from a corpus (Yuan et al. 2025). The method uses an LLM to rewrite any toxic sentence in a non-toxic way whilst preserving the sentence’s meaning. Thus, we are able to remove toxicity without disrupting coherence.

For the purposes of implementation, we selected Llama-3.3-70B to detoxify the sentences (Grattafiori et al. 2024). It was picked due to being a well-behaved and high-performing open-source model. In order to identify a well-performing detoxification prompt, we sampled 20 toxic sentences and tested multiple prompt variants, with the selected one, together with an example, shown in Figure 16.

G Perturbation Augmentation Implementation

Perturbation Augmentation debiasing technique by Qian et al. (2022) relies on the idea of using a pre-trained LM (perturber) to rewrite the corpus, randomly swapping every demographic reference for a different one.

In practice, this means that when supplied with a chunk of text, a target word, and a target topic, the perturber changes the gender, age, or race of the subject in the chunk to a new one randomly selected from the set list in Table ?? . If applied to the entire corpus, this ensures uniform distribution of topics. An example of this kind of perturbation is shown in Figure 17.

To apply this perturbation to the BabyLM corpus, we use the perturber trained by the authors, ensuring the quality of the result. Unfortunately, we also need the target words, which the authors have not provided. Nevertheless, we resolved this by extracting a noisier list of target words from the perturber’s training data.

During the perturbation process, we split the corpus into 128-token chunks, extract all target words from the chunk, randomly select one of them together with a random sub-category, and perturb the chunk. We only focus on race and

gender perturbations as age proved to generate more erroneous changes. If the chunk does not change, we repeat this process with another word until we change the chunk or exhaust the target words. Thus, the noisy ineffective words do not disrupt the experiment. With the entire corpus transformed, Table 14 details the distribution of changes.

Category	Subcategory	Perturbed Chunks
Overall Corpus	Any Change	84.9%
	Gender	79.0%
	Race	5.9%
Gender Perturbation	Non-binary	27.9%
	Woman	26.8%
	Man	24.2%
Race Perturbation	Pacific-Islander	1.1%
	Native-American	1.0%
	White	1.0%
	Asian	1.0%
	Black	1.0%
	Hispanic	1.0%

Table 14: Change distribution in the perturbed corpus

H Computational Resources

For all debiasing experiments, covering both fine-tuning and pre-training, we utilised a server with four A100 GPUs. Most evaluation tasks were conducted on a single NVIDIA T4 GPU. In total, including test runs, the experiments consumed approximately 600 GPU-hours. No hyper-parameter search experiments were conducted for any of the experiments.

I Code Repository

The code for all experiments presented in this paper is available here *LINK HIDDEN DUE TO ANONYMOUS SUBMISSION*.

J Reproducibility Checklist

- This paper:
 - Includes a conceptual outline and/or pseudocode description of AI methods introduced (**yes**/partial/no/NA)
 - Clearly delineates statements that are opinions, hypotheses, and speculation from objective facts and results (**yes**/no)
 - Provides well-marked pedagogical references for less-familiar readers to gain background necessary to replicate the paper (**yes**/no)
- Does this paper make theoretical contributions? (**no**/yes)
- Does this paper rely on one or more datasets? (**yes**/no)

If yes, please complete the list below:

 - A motivation is given for why the experiments are conducted on the selected datasets. (**yes**/partial/no/NA)
 - All novel datasets introduced in this paper are included in a data appendix. (**partial**/yes/no/NA)

Prompt:

Each line below is a sentence that has been labeled as toxic or hateful. Rewrite each line so that it is not toxic or hateful, without changing its meaning. Return the cleaned sentences in the same order, one per line, with no numbering or comments:

Input:

...
It's a load of fucking crap to be quite honest as far as I'm concerned!
...

Output:

...
To be honest, I'm really disappointed with this, it's not what I expected at all.
...

Figure 16: Detoxification prompt and its effect on an example

See what **catherine's** gonna wear? Look. Look joseph. See she has a ladybug! See the ladybug? Very nice. Purple ladybug. And now what do we need to get for **catherine**?

Target word

Catherine

Target topic

Gender

Man

See what **Cameron's** gonna wear? Look. Look joseph. See **he** has a ladybug! See the ladybug? Very nice. Purple ladybug. And now what do we need to get for **Cameron**?

Figure 17: Perturbation showcase

- All novel datasets introduced in this paper will be made publicly available upon publication with a license allowing free research use. (yes/partial/no/NA)
 - All datasets drawn from the existing literature are accompanied by appropriate citations. (yes/no/NA)
 - All datasets drawn from the existing literature are publicly available. (yes/partial/no/NA)
 - Datasets that are not publicly available are described in detail, with justification. (yes/partial/no/NA)
4. Does this paper include computational experiments? (yes/no)
- If yes, please complete the list below:
- Number/range of values tried per (hyper-)parameter and selection criteria are reported. (NA/yes/partial/no)
 - Code for data preprocessing is included in the appendix. (partial/yes/no)
 - Source code for conducting and analyzing experiments is included. (partial/yes/no)
 - Code will be released publicly upon publication with a permissive license. (yes/partial/no)
 - Code includes comments with implementation details and paper references. (partial/yes/no)
 - Seed setting methods for stochastic algorithms are described. (partial/yes/no/NA)
 - Computing infrastructure (hardware/software specs) is reported. (yes/partial/no)
 - Evaluation metrics are formally described with motivations. (yes/partial/no)
 - Number of runs per result is specified. (partial/yes/no)
 - Performance analysis includes variation, confidence, or distributions. (partial/yes/no)
 - Significance of performance differences is assessed with statistical tests. (partial/yes/no)
 - Final (hyper-)parameter settings are listed.

(yes/partial/no/NA)

Notes: The camera-ready version of the paper will include an anonymised code repository containing the code needed to reproduce all reported debiasing experiments. The dataset augmentations are described, and the code required to augment the BabyLM dataset is provided in the repository. However, we do not include the datasets themselves.