

Filling in the Mechanisms: How do LMs Learn Filler-Gap Dependencies under Developmental Constraints?

Anonymous ACL submission

Abstract

For humans, filler-gap dependencies require a shared representation across different syntactic constructions. Although causal analyses suggest this may also be true for LLMs (Boguraev et al., 2025), it is still unclear if such a representation also exists for language models trained on developmentally feasible quantities of data. We applied Distributed Alignment Search (DAS, Geiger et al. (2024)) to LMs trained on varying amounts of data from the BabyLM challenge (Warstadt et al., 2023), to evaluate whether representations of filler-gap dependencies transfer between wh-questions and topicalization, which greatly vary in terms of their input frequency. Our results suggest shared, yet item-sensitive mechanisms may develop with limited training data. More importantly, LMs still require far more data than humans to learn comparable generalizations, highlighting the need for language-specific biases in models of language acquisition.

1 Introduction

A major question in language acquisition asks how learners can apply their linguistic knowledge to generalize to unseen information. The Poverty of the Stimulus argument (Chomsky, 1965) claims that children require specialized representations to learn language due to sparse input. Cognitive models of language acquisition have proposed and revised many claims about the nature of the representations humans could use to solve this learning problem (Yang, 2004; Perkins et al., 2022; Pearl, 2023). Language models (LMs), however, are trained on domain-general next-word prediction tasks, yet still posit some generalizations over syntactic structure (Futrell et al., 2019; Warstadt et al., 2020; Linzen and Baroni, 2021; Wilcox et al., 2024). These successes have led some researchers to question the need for language-specific representations altogether (Piantadosi, 2023; Futrell and Mahowald, 2025).

One particular case that evaluates this debate involves *filler-gap dependencies*. Filler-gap dependencies are a type of syntactic relation formed when a constituent (or the filler) is displaced from its canonical position (the gap) and interpreted in another position. These gaps can exist across various constructions including, but not limited to, wh-questions (“*What did the student make _?*”), relative clauses (“*the robot that the student made _*”), and topicalization (“*This robot, the student made _*”).

Filler-gap dependencies are a useful test case for evaluating representations necessary for language acquisition, as they not only require knowledge of hierarchical structure but also the recognition of an empty syntactic position. Many linguistic theories claim that despite superficial differences, these constructions may share an abstract underlying mechanism (Chomsky, 1977; Culicover et al., 1977; Gazdar, 1982; Kaplan and Bresnan, 1982; Postal, 1999). However, results from LMs have led others to argue that statistical patterns over language input may be sufficient (Wilcox et al., 2018, 2024).

These results have since been challenged, highlighting that LMs may not represent filler-gap dependencies in a human-like manner (Bhattacharya and van Schijndel, 2020; Lan et al., 2024; Howitt et al., 2024; Chang et al., 2025). However, most previous work relies largely on evaluating language model probabilities, rather than explicitly finding causal patterns in LM representations. Boguraev et al. (2025) do apply causal interpretability methods (Geiger et al., 2024; Mueller et al., 2025), and identify shared underlying structure across various types of filler-gap dependencies. However, even if Boguraev et al. (2025) show a shared representation of filler-gap dependencies may be learnable in principle, their results do not address whether inductive biases are needed for humans, as they evaluated LMs trained on data that matches nei-

084 ther the content nor the quality of data available to
085 human learners (Wilcox et al., 2025).

086 We run causal interventions on LMs (Boguraev
087 et al., 2025) trained on data from the BabyLM
088 challenge (Warstadt et al., 2023), which reflects
089 material that English-speaking children could be
090 exposed to, up to 12 years of age. Our experi-
091 ments show that LMs may posit a shared represen-
092 tation of filler-gap dependencies across high and
093 low frequency constructions, but this representation
094 is far less general than the ones proposed by lin-
095 guists, showing systematic variability across items
096 and constructions. Importantly, the model fails to
097 posit this representation at the point it is expected
098 to emerge in human learners (Perkins and Lidz,
099 2021). This result suggests that human learners
100 need a combination of domain-general statistical
101 learning mechanisms, along with language-specific
102 inductive biases (Yang, 2004; Portelance and Jasbi,
103 2024).

104 2 Background

105 Filler-gap dependencies are shared across syntac-
106 tic constructions that show different surface forms,
107 even if they may serve different semantic and dis-
108 course functions (Schütze et al., 2015). Looking at
109 the following sentences, (1) is a *wh-question*, while
110 (3) is an example of *topicalization*.

- 111 (1) *Who* did the teacher like _ ?
112 (2) Did the teacher like?*
- 113 (3) *The author*, the teacher liked _ .
114 (4) The teacher liked.*

115 In (1) *who* is the object of *like*, and in (3), *The*
116 *author* is the object of *liked*. These constituents,
117 which are **fillers**, are fronted to form a dependency
118 with the **gaps**, which are unpronounced but marked
119 with _ for readability. Learning the generaliza-
120 tion also involves recognizing where the depen-
121 dency may *not* be valid, such as (2) and (4), which
122 show that (1) and (3) are respectively ungrammat-
123 ical without the fillers. Syntactic configurations
124 called *islands* make extracting a filler ungrammati-
125 cal (Chomsky, 1977).¹ This has made recognizing
126 filler-gap licensing a particularly relevant test case
127 when evaluating syntactic structure in language
128 models.

¹Refer to Wilcox et al. (2024), Howitt et al. (2024), and Chang et al. (2025) for further discussion relating to language models and island constraints.

2.1 LM Surprisal and Filler-Gap Dependencies

129 Many studies use LMs to compute the *surprisal*, or
130 negative log probability, of a word given its context.
131 Surprisal quantifies the effect of processing diffi-
132 culty (Levy, 2008), and evaluating LMs’ surprisals
133 at particular points in a sentence effectively identi-
134 fies which parts are expected to be more difficult
135 to process (Futrell et al., 2019).² Work evaluat-
136 ing LMs on syntactic structure has often relied on
137 comparing two minimal pairs of sentences, such as
138 (1) vs. (2) and (3) vs. (4), where the ungrammati-
139 cal version should typically have a higher surprisal
140 than the grammatical version (Marvin and Linzen,
141 2018; Warstadt et al., 2020; Gauthier et al., 2020).
142

143 In English sentences with embedded *wh*-
144 movement, LSTM language models have shown
145 positive results in recognizing the presence and
146 absence of fillers and gaps (Wilcox et al., 2018).
147 These results have been extended to Transformer
148 models and various island constraints (Wilcox et al.,
149 2024).³

150 Ozaki et al. (2022) evaluated LSTMs across a
151 range of other filler-gap constructions, including
152 topicalization, and found that model performance
153 for each construction is highly correlated with its
154 frequency. They do not present evidence whether
155 this generalization is shared *across* constructions.
156

157 To this end, other approaches involve providing
158 LMs with additional training examples. Simulated
159 priming studies on *wh*-movement (Bhattacharya
160 and van Schijndel, 2020; Prasad et al., 2019) have
161 shown some evidence for a shared representation of
162 filler-gap dependency, but not constraints on the de-
163 pendency. More recent studies have retrained LMs
164 by augmenting their training data with positive ex-
165 amples of a dependency. Lan et al. (2024) show
166 that augmentation improves LMs’ performance on
167 complicated filler-gap constructions (parasitic gaps
168 and across-the-board movement). Extending this
169 approach across constructions, Howitt et al. (2024)
170 adopted their methodology and found that augment-
171 ing LSTMs’ training data with instances of one con-
172 struction (clefting and topicalization) failed to im-
173 prove performance on detecting filler-gap licensing

²However, see Huang et al. (2024) for evidence that LM surprisal cannot reflect the quantitative effects of processing difficulty for particular types of syntactically complex sentences.

³Results are more mixed in Norwegian (Kobzeva et al., 2023) and Dutch (Suijkerbuijk et al., 2023), which have different filler-gap structures from English.

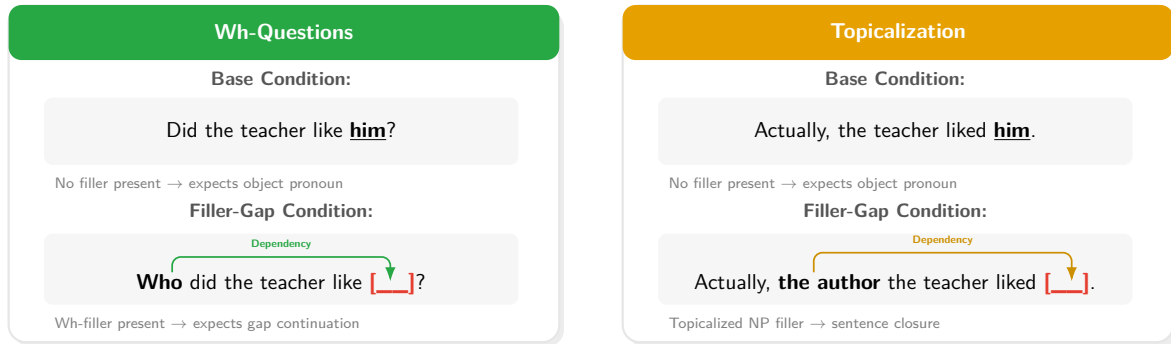


Figure 1: The diagram contrasts Wh-Questions (left, green) and Topicalization structures (right, orange).

few parameters and human-like training data still posit a causal representation of filler-gap dependencies? If so, when does this effect emerge?

- **RQ2:** Are representations localized within examples of the same construction?
- **RQ3:** Do representations transfer from low-frequency to high-frequency constructions?

We create three hypotheses based on these questions. First, the version of the model trained on all 100 million tokens in the BabyLM training corpus *should* learn an abstract filler-gap mechanism that is detectable and transferable by DAS from one construction to another, based on some of the positive results from Chang et al. (2025). However, due to the lack of language-specific inductive biases, we do not predict a strong causal effect before 10 million tokens. This is because the BabyLM checkpoint for 10 million tokens mimics the linguistic input of children from ages 2-5 years (Warstadt et al., 2023), while English-speaking children are able to recognize filler-gap dependencies around 18 months (Perkins and Lidz, 2021) (**H1: Misaligned Emergence Hypothesis**). Second, since LMs likely learn the filler-gap dependency in a piecemeal fashion (Ozaki et al., 2022; Howitt et al., 2024), we expect stronger causal effects for within-construction interventions compared to cross-construction interventions. We also hypothesize that transfer improves when both constructions share the same level of animacy, to replicate (Boguraev et al., 2025)’s findings for lexical boost effects. (**H2: Construction-Specificity Hypothesis**). Third, given the greater prevalence of wh-questions in language corpora, as opposed to topicalization, and existing work showing learning correlates with input frequency (Ozaki et al.,

2022; Boguraev et al., 2025), we expect an asymmetrical and one-way transfer from high-frequency wh-questions to low-frequency Topicalization (**H3: Frequency Modulation Hypothesis**).

4 Methods

4.1 Model

We use the BabyLM-100M model (Warstadt et al., 2023). This model uses the GPT-2-small architecture (Radford et al., 2019) trained on the BabyLM Strict-100M corpus, consisting of a set of approximately 100 million words as training data designed to mimic the total linguistic input received by an English-speaking child until early adolescence (around 12 years of age) (Gilkerson et al., 2017; Warstadt et al., 2023).

The corpus consists of relevant data from the British National Corpus, CHILDES language acquisition database, Switchboard Dialog Act Corpus, subtitles from children’s TV shows, and simplified Wikipedia articles (Charpentier et al., 2025).

Multiple checkpoints for the model across the training process were also released. We evaluate several checkpoints where the model received increasing amounts of input (1M–100M tokens): 10 checkpoints from 0–10M tokens and 9 checkpoints between 10M–100M tokens. Additionally, nine additional checkpoints (100M–1000M) are used in extended analysis in the appendix.⁵

4.1.1 Constructions

We use two filler-gap constructions that have different levels of frequency: matrix wh-questions (high frequency) and topicalization (low frequency), based on Ozaki et al. (2022). This combination

⁵If accepted, code and data for this paper will be publicly released.

allows us to systematically evaluate whether knowledge of high-frequency constructions can be transferred to less frequent ones that may barely be present in children’s input.

Wh-Questions (High Frequency). We use single-clause matrix wh-questions, such as “*What did the doctor do ___?*” and “*What did the student read ___?*” Wh-questions are very common in natural language and exist in child-directed speech (Furrow et al., 1979). Prior work demonstrates that neural LMs readily acquire sensitivity to wh-dependencies in English (Wilcox et al., 2024). This makes Wh-questions a plausible construction for a high-frequency “source” in transfer experiments.

Topicalization (Low Frequency). We employ fronted object topicalization with an optional discourse marker, such as “*The student, the teacher liked ___*” or “*Actually, the book the author read ___*” Topicalization is extremely uncommon in natural language and almost absent from child-directed speech (Roland et al., 2007). Related works find that LMs fail to display the correct behavior for topicalization (Ozaki et al., 2022), even when augmented with examples in the training data (Howitt et al., 2024).⁶ In addition, when analyzing transfer of the filler-gap mechanism across constructions, Boguraev et al. (2025) found topicalization had a low “out-degree” in their transfer network: despite receiving other constructions, topicalization rarely transfers to them. This makes topicalization a strong candidate for a low-frequency “sink” construction for our experiments.

We created minimal pairs following the original template and methodology of Boguraev et al. (2025). Each pair compares a base sentence with no filler-gap dependency with a filler-gap sentence containing a matrix wh-question or a topicalization setup. The model’s expected next-token prediction differs based on the two predicted gaps (marked with `_`):

4.2 Materials

Replicating the experimental structure of Boguraev et al. (2025), we test sentences with both animate and inanimate fillers. Animate fillers use *who* (wh-questions) or NPs with perceived life or agency, such as *the author* (topicalization). In contrast, inanimate fillers use *what* (wh-questions)

⁶Howitt et al. (2024) show that including the discourse marker leads to slight qualitative improvements, but the generalization is only learned in one direction, when the filler is present.

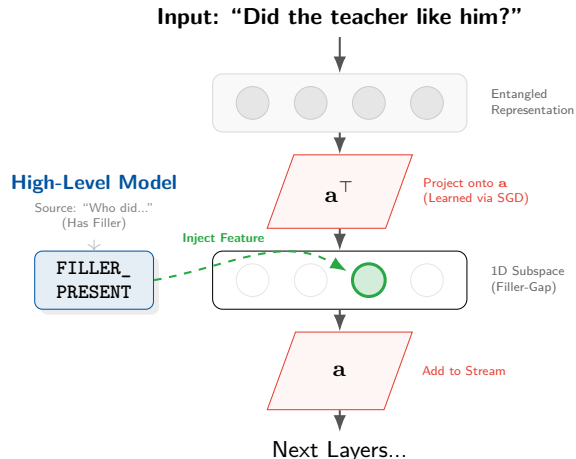


Figure 2: To create a DAS vector, we learn a direction \mathbf{a} to align neural representations with the binary variable `FILLER_PRESENT`. We intervene by projecting the difference between the source and base representations onto \mathbf{a} and injecting it into the base sentence.

or nonliving and nonsentient NPs such as *the book*. This creates four dataset template variants: *wh_animate*, *wh_inanimate*, *topic_animate*, and *topic_inanimate*.

Different combinations of the following lexical items are used in respective templates: **Subject NPs** (50 animate nouns: *teacher, doctor, manager, etc.*), **Verbs** (30 transitive verbs: *like, admire, follow, etc.*), **Auxiliaries** for wh-questions (7 verbs: *did, will, could, etc.*), and **Licensing adverbs** for topicalization (25 adverbs: *Actually, Frankly, Surprisingly, etc.*). We verified all materials consist of one token for the model and occur within the BabyLM corpus.⁷

This yielded approximately 21,000 unique sentence pairs for wh-questions and 1,875,000 for topicalization per animacy condition. Topicalization randomly selects both the sentence-initial adverb and the topicalized filler phrase, increasing the amount of unique sentences. However, both pools remain substantially larger than 2000 pairs sampled for DAS training, maintaining sentence diversity for generalization.

4.3 Distributed Alignment Search (DAS)

Distributed Alignment Search (DAS) is a causal intervention method to test if a high-level concept aligns with the internal weights of a language model (Geiger et al., 2024). We can define a mini-

⁷See Nair and Resnik (2023), Giulianelli et al. (2024), and Oh and Schuler (2025) for discussion of tokenization and psycholinguistic applications.

416 mum causal model for a filler-gap dependency using a binary variable: $\text{FILLER_PRESENT} \in \{0, 1\}$.
 417 This variable causally influences gap expectations
 418 and the model’s next-token predictions.
 419

420 Given a *source* sentence with a filler (1) and a
 421 *base* sentence without a filler (1), we can intervene
 422 based on the internal representations of the base
 423 sentence by implanting the learned filler-gap DAS
 424 feature from the source, as seen in figure 2. A
 425 successful implementation should shift the predic-
 426 tion of the model from the base label (*him*) toward
 427 the source label (?). This would suggest that the
 428 filler-gap dependency is encoded at the intervention
 429 site (Wu et al., 2024).

430 4.3.1 Training

431 Following prior work, we use a 1-dimensional vari-
 432 ant of DAS vector (Geiger et al., 2024; Arora et al.,
 433 2024; Boguraev et al., 2025). Given an embed-
 434 ding space with dimensionality d , for each layer ℓ
 435 and token position p , we learn a *direction vector*
 436 $\mathbf{a}_{\ell,p} \in \mathbb{R}^d$ that defines a one-dimensional subspace
 437 in which the filler-gap feature is encoded, between
 438 the model’s representations of the base construc-
 439 tion at ℓ and p ($\mathbf{h}_{\text{base},\ell,p} \in \mathbb{R}^d$) and the source
 440 construction ($\mathbf{h}_{\text{source},\ell,p} \in \mathbb{R}^d$). Once the vector \mathbf{a}
 441 is learned, the intervention projects the difference
 442 between the source and base representations onto
 443 \mathbf{a} and adds it to the base representation:

$$444 \quad \tilde{\mathbf{h}} = \mathbf{h}_{\text{base}} + \mathbf{a}\mathbf{a}^\top(\mathbf{h}_{\text{source}} - \mathbf{h}_{\text{base}}) \quad (1)$$

445 This intervention preserves the orthogonal di-
 446 mensions of the base representation and only mod-
 447 ifies the value along the learned feature direction
 448 \mathbf{a} . The direction is optimized to minimize cross-
 449 entropy loss between the counterfactual predictions
 450 of the model after intervention and the source sen-
 451 tence labels.

452 We use a batch size of 25, 80 training steps per
 453 layer-position combination, and a learning rate of
 454 5×10^{-3} to train each DAS vector. We provide fur-
 455 ther information on hyperparameter selection in
 456 Appendix A. Due to the comparatively low number
 457 of layers in BabyLM’s architecture, we test all 12
 458 layers across 6 token positions aligned to template
 459 slots: prefix (position 0), filler (position 1), auxil-
 460 iary/complementizer (position 2), article (position
 461 3), subject NP (position 4), and verb (position 5).

462 4.3.2 Evaluation Metrics

463 We primarily use the **ODDS** metric to quantify the
 464 magnitude of the causal effect through measuring

465 how much the intervention shifts log-probabilities
 466 toward the counterfactual outcome. This is given
 467 by the formula:

$$468 \quad \text{ODDS} = \log \frac{P(\text{base} \mid \text{clean})}{P(\text{source} \mid \text{clean})} \quad (2)$$

$$469 \quad + \log \frac{P(\text{source} \mid \text{int})}{P(\text{base} \mid \text{int})}$$

469 In Equation (2), ‘clean’ refers to a standard for-
 470 ward pass without an intervention, while ‘int’ refers
 471 to a forward pass where the DAS intervention is
 472 applied.

473 A positive **ODDS** value suggests the intervention
 474 successfully shifts predictions toward the source; in
 475 addition, the higher the **ODDS** values, the stronger
 476 the causal effects.

477 Utilizing empirical findings from Arora et al.
 478 (2024) on the Pythia model family, we establish
 479 the following qualitative thresholds: values near
 480 0 show little to no causal effect (comparable to
 481 random baselines), values in the 3–6 range demon-
 482 strate emerging to moderate causal structure (as
 483 seen in smaller models such as Pythia-14M), and
 484 values greater than 8 indicate strong causal mech-
 485 anisms (as seen in larger models such as Pythia-
 486 6.9B).

487 We primarily report **MAX ODDS**, or the maxi-
 488 mum **ODDS** value across all layers at a given po-
 489 sition. The goal of DAS is to localize the feature
 490 to specific layers, so the maximum represents the
 491 layer-position combination with the most effective
 492 causal effect.

493 4.3.3 Experiments

494 We run two types of experiments: *localization* and
 495 *transfer*, to evaluate generalizations within and
 496 across construction types.

- 497 1. **Wh** \rightarrow **Wh (within-construction localiza-)**
 498 **tion):** Train DAS on wh-questions, test on
 499 held-out wh-questions
- 500 2. **Topic** \rightarrow **Topic (within-construction local-**
 501 **ization):** Train DAS on topicalization, test on
 502 held-out topicalization
- 503 3. **Wh** \rightarrow **Topic (forward transfer):** Train DAS
 504 on wh-questions, test on topicalization
- 505 4. **Topic** \rightarrow **Wh (backward transfer):** Train
 506 DAS on topicalization, test on wh-questions

507 The localization experiments (Wh \rightarrow Wh,
 508 Topic \rightarrow Topic) quantify the extent to which

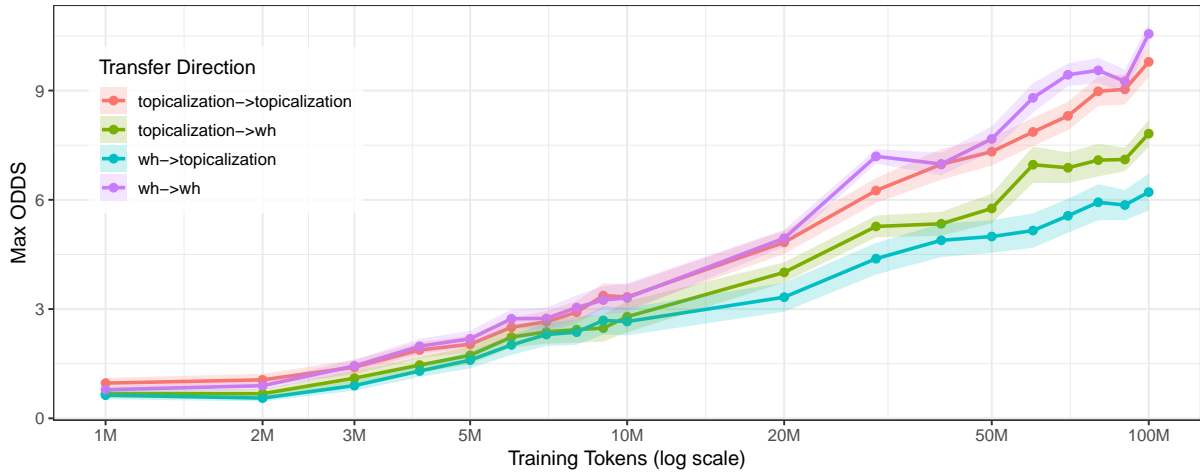


Figure 3: Developmental trajectory of all filler-gap mechanisms across training. Error bands show ± 1 SE across a minimum of 6 seeds. All four conditions show a monotonic increase with training tokens.

the DAS can identify filler-gap representations within each construction type. Likewise, cross-construction transfers (Wh \rightarrow Topic, Topic \rightarrow Wh) test if representations of a given construction are generalizable across the filler-gap constructions. If the transfer is symmetric, ODDS retention should be similar in both directions. If frequency modulates transfer, we predict asymmetrical transfer towards the high-frequency (Wh \rightarrow Topic) direction.

4.4 Statistical Analysis

We fit a linear model predicting MAX ODDS from number of training tokens, transfer direction, and animacy:

$$y = \beta_0 + \beta_t \mathbf{x}_{\text{tokens}} + \beta_d \mathbf{x}_{\text{dir}} + \beta_a \mathbf{x}_{\text{anim}} + \epsilon \quad (3)$$

Where y is MAX ODDS, and \mathbf{x} represents the fixed effects for token count, transfer direction, and animacy. Post-hoc contrasts used estimated marginal means with Holm-Bonferroni correction for 171 pairwise comparisons (Lenth and Piaskowski, 2017). Cohen’s d is used to measure effect sizes using the residual standard deviation.

5 Results

We present the results from 19 BabyLM checkpoints in the developmentally plausible range (1M–100M tokens) across multiple constructions (wh-questions, topicalization) and animacy conditions (animate, inanimate). All experiments were repeated on a minimum of 6 seeds to establish tighter

bounds of confidence. The linear model achieved $R^2 = 0.53$, $F(22, 4537) = 235$, $p < .001$.

5.1 Developmental Trajectory

Figure 3 shows MAX ODDS across training on all four transfer directions. Filler-gap localization significantly increased with training duration ($F(18, 4537) = 264.3$, $p < .001$), from near zero at 1M tokens (MAX ODDS ≈ 0.8) to a robust effect by 100M (MAX ODDS ≈ 10.6 for Wh \rightarrow Wh). Qualitatively, relatively stronger causal effects (MAX ODDS > 8) began emerging after around 50M tokens, while the effects are weaker (MAX ODDS ≈ 3) around 10M tokens.

5.2 Within-Construction Localization

Both constructions displayed successful within-construction localization, increasing with training tokens. At 100M tokens, Wh \rightarrow Wh achieved MAX ODDS = 10.56 (SD = 2.11) and Topic \rightarrow Topic achieved MAX ODDS = 9.79 (SD = 3.24). Compared to the results of the transfer experiments, within-construction localization effects consistently exceeded cross-construction transfer effects (mean difference = 1.33 MAX ODDS, $t = 17.4$, $p < .001$, $d = 0.52$), showing that high performance on one dependency type does not completely generalize to others.

5.3 Cross-Construction Transfer

Unlike the frequency-based prediction in Hypothesis 3, Topic \rightarrow Wh transfer exceeded Wh \rightarrow Topic transfer throughout training (difference = 0.57 MAX ODDS, $t = -5.29$, $p < .001$, $d = -0.22$). At

Contrast	Est.	SE	<i>t</i>	<i>p</i>	<i>d</i>
Within–Across	1.33	0.08	17.4	<.001	0.52
Wh→Topic–Topic→Wh	-0.57	0.11	-5.3	<.001	-0.22
Animate–Inanimate	0.68	0.08	8.6	<.001	0.26

Table 1: Post-hoc comparisons from linear model. Negative asymmetry results suggest Topic→Wh outperformed Wh→Topic.

the 100M checkpoint, Topic→Wh achieved **MAX ODDS** = 7.82 versus Wh→Topic’s **MAX ODDS** = 6.21.

This transfer asymmetry could be due to wh-questions developing a more item-specific representation that transfers less effectively across construction types. Conversely, a more general mechanism may be reflected in topicalization, when it is learned, due to little to no presence in the input. We evaluate this behavior on the model when it is trained on up to 1 billion tokens in Appendix A.3.

5.4 Animacy Effects

We find further evidence of a “lexical boost” effect predicted by Boguraev et al. (2025). DAS transfer was significantly stronger when animacy of training and evaluation matched relative to animacy-mismatched conditions (difference = 0.67 **MAX ODDS**, $t = 4.86$, $p < .001$, $d = 0.28$).

6 Discussion and Conclusions

This study applies causal interventions to determine how a language model learns filler-gap dependencies when provided with developmentally realistic amounts of training data from the BabyLM corpus, determining if findings from larger LMs (Boguraev et al., 2025) apply in this setting. Using DAS, we evaluated BabyLM-100M’s generalizations in four experimental conditions: localization within constructions and transfer across constructions, for two types of filler-gap dependencies, wh-questions and topicalization. Our results show that the model learns a shared representation for filler-gap dependencies, but still requires far more data than children would, and is still highly sensitive to variation across items.

Regarding **RQ1**, which asks both whether and when LMs learn a causal representation, we find evidence in favor of our misaligned emergence hypothesis. Although the BabyLM-100M model showed strong causal effects when trained on the full corpus, they failed to emerge when the model received human-like quantities of training data.

The full corpus was comparable to the input available to English-speaking adolescents (up to around 12 years) (Warstadt et al., 2023), while children exhibit robust knowledge of filler-gap dependencies prior to two years of age (Gagliardi et al., 2016; Atkinson et al., 2018; Perkins and Lidz, 2021). In their description of the BabyLM corpus, Warstadt et al. (2023) report that the 10M checkpoint corresponds to children’s linguistic knowledge between the ages of 2 and 5. If the model *was* learning with a human-like mechanism, we would expect strong causal effects prior to this checkpoint, yet only report weak effects, if any.

RQ2 asks whether the learned representations are specific to particular constructions. We found evidence for this construction-specificity hypothesis because the localization experiments consistently performed better than the transfer experiments. That is, generalization within examples of the same construction was far more effective than transferring representations across constructions. We also replicated the lexical boost effects when animacy gets matched during the intervention, suggesting this feature may transfer across items. Regarding the direction of transfer, as discussed in **RQ3**, we found improved performance generalizing from topicalization to wh-questions, which was the opposite of our predictions in the frequency modulation hypothesis. Future work can determine if this happens because LMs may learn more item-specific representations for wh-questions since they are more frequent, and deploy a more generalized mechanism to represent topicalization.

Overall, instead of positing a single, general representation of filler-gap constructions, LMs learn item and construction-specific representations. Future work should extend DAS to evaluate learning the filler-gap dependency in both directions and sensitivity to island constraints, across more diverse construction types.

When modeling human language acquisition, however, our results show LMs trained solely on next-word prediction are not sufficient. Instead, we emphasize the need to model learning with explicit inductive biases over structured hypothesis spaces. LMs may still play a role in this enterprise, through reflecting inductive biases architecturally (Murty et al., 2023) or specifying possible hypothesis spaces (Misra and Kim, 2024; Portelance and Jasbi, 2024). Thus, cognitive models of language acquisition should build on advances from statistical learning, while acknowledging their limits.

662 Limitations

663 This study only focused on English, limiting the
664 generalizability of these results to other languages
665 where filler-gap dependencies behave differently
666 under LMs (Kobzeva et al., 2023; Suijkerbuijk
667 et al., 2023). Additionally, the training corpus is
668 based on text input alone, while children learn from
669 spoken data, multimodal environments, and social
670 interaction (Meylan et al., 2023; Vong et al., 2024).
671 Since this study focused on extending Boguraev
672 et al. (2025)’s results to a BabyLM-scale model, we
673 did not evaluate whether it could recognize the ab-
674 sence of filler-gap dependencies, and for island con-
675 straints, which have been used in surprisal-based
676 studies (Ozaki et al., 2022; Wilcox et al., 2024;
677 Howitt et al., 2024; Chang et al., 2025). More
678 complex materials would also be useful to ensure
679 results are not associated with confounding factors
680 like punctuation, since we extract representations
681 from periods and question marks. Lastly, although
682 DAS was the best performing method from Arora
683 et al. (2024), measures like Boundless DAS operate
684 over subspaces instead of single dimensions (Wu
685 et al., 2023; Geiger et al., 2024), and could have
686 yielded stronger causal effects.

687 Ethical Considerations

688 This work used publicly available data and models,
689 which are described further in the original publica-
690 tions. We do not foresee any risks associated with
691 this work, as we used the data for their intended
692 purpose to study human language acquisition. Gen-
693 erative AI (GenAI) was used in this project. We
694 used Antigravity⁸ to design plots and refactor code,
695 and Claude Opus 4.5 to refine paper writing for
696 brevity. We never use GenAI for writing text from
697 scratch in this paper. We take complete responsi-
698 bility for any GenAI errors. By discussing GenAI
699 usage here, we aim to encourage other researchers
700 to do the same.

701 References

702 Aryaman Arora, Dan Jurafsky, and Christopher Potts.
703 2024. [CausalGym: Benchmarking causal inter-
704 pretable methods on linguistic tasks](#). In *Proceed-
705 ings of the 62nd Annual Meeting of the Association
706 for Computational Linguistics (Volume 1: Long Pa-
707 pers)*, pages 14638–14663, Bangkok, Thailand. As-
708 sociation for Computational Linguistics.

⁸<https://antigravity.google/>

- Emily Atkinson, Matthew W Wagers, Jeffrey Lidz, 709
Colin Phillips, and Akira Omaki. 2018. Developing 710
incrementality in filler-gap dependency processing. 711
Cognition, 179:132–149. 712
- Debasmita Bhattacharya and Marten van Schijndel. 713
2020. [Filler-gaps that neural networks fail to gen-
714 eralize](#). In *Proceedings of the 24th Conference on
715 Computational Natural Language Learning*, pages
716 486–495, Online. Association for Computational Lin-
717 guistics. 718
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory 719
Anthony, Herbie Bradley, Kyle O’Brien, Eric Hal- 720
lahan, Mohammad Aflah Khan, Shivanshu Purohit, 721
USVSN Sai Prashanth, Edward Raff, et al. 2023. 722
Pythia: A suite for analyzing large language mod- 723
els across training and scaling. In *International
724 Conference on Machine Learning*, pages 2397–2430.
725 PMLR. 726
- Sasha Boguraev, Christopher Potts, and Kyle Mahowald. 727
2025. [Causal interventions reveal shared structure
728 across English filler–gap constructions](#). In *Proceed-
729 ings of the 2025 Conference on Empirical Methods in
730 Natural Language Processing*, pages 25021–25042,
731 Suzhou, China. Association for Computational Lin-
732 guistics. 733
- Chi-Yun Chang, Xueyang Huang, Humaira Nasir, Shane 734
Storks, Olawale Akingbade, and Huteng Dai. 2025. 735
[Mind the gap: How babylms learn filler-gap depen-
736 dencies](#). In *Proceedings of the 2025 Conference on
737 Empirical Methods in Natural Language Processing*,
738 pages 15060–15076. 739
- Lucas Charpentier, Leshem Choshen, Ryan Cotterell, 740
Mustafa Omer Gul, Michael Hu, Jaap Jumelet, Tal 741
Linzen, Jing Liu, Aaron Mueller, Candace Ross, 742
Raj Sanjay Shah, Alex Warstadt, Ethan Wilcox, and 743
Adina Williams. 2025. [Babylm turns 3: Call for
744 papers for the 2025 babylm workshop](#). 745
- Noam Chomsky. 1965. Aspects of the theory of syntax. 746
- Noam Chomsky. 1977. [On Wh-Movement](#). *Formal
747 Syntax*, pages 71–132. Publisher: Academic Press. 748
- Peter W Culicover, Thomas Wasow, and Adrian Akma- 749
jian. 1977. *Formal syntax*. Academic Press. 750
- David Furrow, Katherine Nelson, and Helen Benedict. 751
1979. Mothers’ speech to children and syntactic 752
development: Some simple relationships. *Journal of
753 child language*, 6(3):423–442. 754
- Richard Futrell and Kyle Mahowald. 2025. How linguis- 755
tics learned to stop worrying and love the language 756
models. *arXiv preprint arXiv:2501.17047*. 757
- Richard Futrell, Ethan Wilcox, Takashi Morita, Peng 758
Qian, Miguel Ballesteros, and Roger Levy. 2019. 759
Neural language models as psycholinguistic subjects: 760
Representations of syntactic state. In *Proceedings of
761 NAACL-HLT*, pages 32–42. 762

763	Annie Gagliardi, Tara M Mease, and Jeffrey Lidz. 2016.	Ronald Kaplan and Joan Bresnan. 1982. Lexical-	821
764	Discontinuous development in the acquisition of	Functional Grammar: A Formal System for Gram-	822
765	filler-gap dependencies: Evidence from 15-and 20-	matical Representation. In Joan Bresnan, editor, <i>The</i>	823
766	month-olds. <i>Language Acquisition</i> , 23(3):234–260.	<i>Mental Representation of Grammatical Relations</i> ,	824
		pages 173–281. MIT Press.	825
767	Leo Gao, Stella Biderman, Sid Black, Laurence Gold-	Anastasia Kobzeva, Suhas Arehalli, Tal Linzen, and	826
768	ing, Travis Hoppe, Charles Foster, Jason Phang,	Dave Kush. 2023. Neural Networks Can Learn Pat-	827
769	Horace He, Anish Thite, Noa Nabeshima, Shawn	terns of Island-insensitivity in Norwegian. In <i>Pro-</i>	828
770	Presser, and Connor Leahy. 2020. The Pile: An	<i>ceedings of the Society for Computation in Linguis-</i>	829
771	800gb dataset of diverse text for language modeling.	<i>tics 2023</i> , pages 175–185, Amherst, MA. Association	830
772	<i>arXiv preprint arXiv:2101.00027.</i>	for Computational Linguistics.	831
773	Jon Gauthier, Jennifer Hu, Ethan Wilcox, Peng Qian,	Daria Kryvosheieva, Andrea de Varda, Evelina Fe-	832
774	and Roger Levy. 2020. Syntaxgym: An online plat-	dorenko, and Greta Tuckute. 2025. Different types	833
775	form for targeted evaluation of language models. In	of syntactic agreement recruit the same units within	834
776	<i>Proceedings of the 58th Annual Meeting of the Associ-</i>	large language models.	835
777	<i>ation for Computational Linguistics: System Demon-</i>		
778	<i>strations</i> , pages 70–76.	Nur Lan, Emmanuel Chemla, and Roni Katzir. 2024.	836
779	Gerald Gazdar. 1982. Phrase Structure Grammar. In	Large language models and the argument from the	837
780	Pauline Jacobson and Geoffrey K. Pullum, editors,	poverty of the stimulus. <i>Linguistic Inquiry</i> , pages	838
781	<i>The Nature of Syntactic Representation</i> , Synthese	1–28.	839
782	Language Library, pages 131–186. Springer Nether-		
783	lands, Dordrecht.	Karim Lasri, Tiago Pimentel, Alessandro Lenci, Thierry	840
784	Atticus Geiger, Zhengxuan Wu, Christopher Potts,	Poibeau, and Ryan Cotterell. 2022. Probing for the	841
785	Thomas Icard, and Noah Goodman. 2024. Finding	usage of grammatical number. In <i>Proceedings of the</i>	842
786	alignments between interpretable causal variables	<i>60th Annual Meeting of the Association for Compu-</i>	843
787	and distributed neural representations. In <i>Causal</i>	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	844
788	<i>Learning and Reasoning</i> , pages 160–187. PMLR.	8818–8831.	845
789	Jill Gilkerson, Jeffrey A Richards, Steven F Warren, Ju-	Russell V. Lenth and Julia Piaskowski. 2017. emmeans:	846
790	dith K Montgomery, Charles R Greenwood, D Kim-	Estimated marginal means, aka least-squares means.	847
791	brough Oller, John HL Hansen, and Terrance D Paul.		
792	2017. Mapping the early language environment using	Roger Levy. 2008. Expectation-based syntactic compre-	848
793	all-day recordings and automated analysis. <i>American</i>	hension. <i>Cognition</i> , 106(3):1126–1177.	849
794	<i>journal of speech-language pathology</i> , 26(2):248–	Tal Linzen and Marco Baroni. 2021. Syntactic structure	850
795	265.	from deep learning. <i>Annual Review of Linguistics</i> ,	851
		7:195–212.	852
796	Mario Giulianelli, Luca Malagutti, Juan Luis Gastaldi,	Rebecca Marvin and Tal Linzen. 2018. Targeted syn-	853
797	Brian DuSell, Tim Vieira, and Ryan Cotterell. 2024.	tactic evaluation of language models. In <i>Proceedings</i>	854
798	On the Proper Treatment of Tokenization in Psy-	<i>of the 2018 Conference on Empirical Methods in</i>	855
799	ycholinguistics. In <i>Proceedings of the 2024 Confer-</i>	<i>Natural Language Processing</i> , pages 1192–1202.	856
800	<i>ence on Empirical Methods in Natural Language Pro-</i>		
801	<i>cessing</i> , pages 18556–18572, Miami, Florida, USA.	Stephan C Meylan, Ruthe Foushee, Nicole H Wong,	857
802	Association for Computational Linguistics.	Elika Bergelson, and Roger P Levy. 2023. How	858
803	Sophie Hao and Tal Linzen. 2023. Verb conjugation	adults understand what young children say. <i>Nature</i>	859
804	in transformers is determined by linear encodings	<i>human behaviour</i> , 7(12):2111–2125.	860
805	of subject number. In <i>Findings of the Association</i>	Kanishka Misra and Najoung Kim. 2024. Generat-	861
806	<i>for Computational Linguistics: EMNLP 2023</i> , pages	ing novel experimental hypotheses from language	862
807	4531–4539.	models: A case study on cross-dative generalization.	863
808	Katherine Howitt, Sathvik Nair, Allison Dods, and	<i>arXiv preprint arXiv:2408.05086.</i>	864
809	Robert Melvin Hopkins. 2024. Generalizations	Aaron Mueller, Jannik Brinkmann, Millicent Li, Samuel	865
810	across filler-gap dependencies in neural language	Marks, Koyena Pal, Nikhil Prakash, Can Rager,	866
811	models. In <i>Proceedings of the 28th Conference on</i>	Aruna Sankaranarayanan, Arnab Sen Sharma, Jiud-	867
812	<i>Computational Natural Language Learning</i> , pages	ing Sun, et al. 2025. The quest for the right medi-	868
813	269–279, Miami, FL, USA. Association for Compu-	ator: Surveying mechanistic interpretability for nlp	869
814	tational Linguistics.	through the lens of causal mediation analysis. <i>Com-</i>	870
815	Kuan-Jung Huang, Suhas Arehalli, Mari Kugemoto,	<i>putational Linguistics</i> , pages 1–48.	871
816	Christian Muxica, Grusha Prasad, Brian Dillon, and	Shikhar Murty, Pratyusha Sharma, Jacob Andreas, and	872
817	Tal Linzen. 2024. Large-scale benchmark yields no	Christopher D Manning. 2023. Pushdown layers:	873
818	evidence that language model surprisal explains syn-	Encoding recursive structure in transformer language	874
819	tactic disambiguation difficulty. <i>Journal of Memory</i>		
820	<i>and Language</i> , 137:104510.		

875	models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 3233–3247.	928
876		929
877		930
878	Sathvik Nair and Philip Resnik. 2023. Words, subwords, and morphemes: What really matters in the surprisal-reading time relationship? In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 11251–11260.	931
879		932
880		933
881		934
882		935
883	Byung-Doh Oh and William Schuler. 2025. The impact of token granularity on the predictive power of language model surprisal. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4150–4162.	936
884		937
885		938
886		939
887		940
888		941
889	Satoru Ozaki, Dan Yurovsky, and Lori Levin. 2022. How well do lstm language models learn filler-gap dependencies? In <i>Proceedings of the Society for Computation in Linguistics 2022</i> , pages 76–88.	942
890		943
891		944
892		945
893	Lisa Pearl. 2023. Computational cognitive modeling for syntactic acquisition: Approaches that integrate information from multiple places. <i>Journal of Child Language</i> , 50(6):1353–1373.	946
894		947
895		948
896		949
897	Laurel Perkins, Naomi H Feldman, and Jeffrey Lidz. 2022. The power of ignoring: Filtering input for argument structure acquisition. <i>Cognitive Science</i> , 46(1):e13080.	950
898		951
899		952
900		953
901	Laurel Perkins and Jeffrey Lidz. 2021. Eighteen-month-old infants represent nonlocal syntactic dependencies. <i>Proceedings of the National Academy of Sciences</i> , 118(41):e2026469118.	954
902		955
903		956
904		957
905	Steven T Piantadosi. 2023. Modern language models refute chomsky’s approach to language. <i>From fieldwork to linguistic theory: A tribute to Dan Everett</i> , 15:353–414.	958
906		959
907		960
908		961
909	Eva Portelance and Masoud Jasbi. 2024. The roles of neural networks in language acquisition. <i>Language and Linguistics Compass</i> , 18(6):e70001.	962
910		963
911		964
912	Paul M. Postal. 1999. <i>Three Investigations of Extraction</i> . The MIT Press.	965
913		966
914	Grusha Prasad, Marten van Schijndel, and Tal Linzen. 2019. Using priming to uncover the organization of syntactic representations in neural language models. In <i>Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)</i> , pages 66–76.	967
915		968
916		969
917		970
918		971
919		972
920	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	973
921		974
922		975
923		976
924	Douglas Roland, Frederic Dick, and Jeffrey L Elman. 2007. Frequency of basic english grammatical structures: A corpus analysis. <i>Journal of memory and language</i> , 57(3):348–379.	977
925		978
926		979
927		980
	Carson T. Schütze, Jon Sprouse, and Ivano Caponigro. 2015. Challenges for a theory of islands: A broader perspective on ambridge, pine, and lieven. <i>Language</i> , 91(2):31–39.	981
		982
		983
	Michelle Suijkerbuijk, Peter de Swart, and Stefan L Frank. 2023. The learnability of the wh-island constraint in dutch by a long short-term memory network. In <i>Proceedings of the Society for Computation in Linguistics 2023</i> , pages 321–331.	984
		985
	Matthew J Traxler, Kristen M Tooley, and Martin J Pickering. 2014. Syntactic priming during sentence comprehension: Evidence for the lexical boost. <i>Journal of Experimental Psychology: Learning, Memory, and Cognition</i> , 40(4):905.	986
		987
	Wai Keen Vong, Wentao Wang, A Emin Orhan, and Brenden M Lake. 2024. Grounded language acquisition through the eyes and ears of a single child. <i>Science</i> , 383(6682):504–511.	988
		989
		990
	Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2021. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. In <i>The Eleventh International Conference on Learning Representations</i> .	991
		992
	Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mosquera, Bhargavi Paranjabe, Adina Williams, Tal Linzen, et al. 2023. Findings of the babylm challenge: Sample-efficient pretraining on developmentally plausible corpora. In <i>Proceedings of the BabyLM challenge at the 27th conference on computational natural language learning</i> , pages 1–34.	993
		994
		995
	Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R Bowman. 2020. Blimp: The benchmark of linguistic minimal pairs for english. <i>Transactions of the Association for Computational Linguistics</i> , 8:377–392.	996
		997
	Ethan Wilcox, Roger Levy, Takashi Morita, and Richard Futrell. 2018. What do RNN language models learn about filler–gap dependencies? In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 211–221, Brussels, Belgium. Association for Computational Linguistics.	998
		999
	Ethan Gotlieb Wilcox, Richard Futrell, and Roger Levy. 2024. Using computational models to test syntactic learnability. <i>Linguistic Inquiry</i> , 55(4):805–848.	1000
		1001
	Ethan Gotlieb Wilcox, Michael Y Hu, Aaron Mueller, Alex Warstadt, Leshem Choshen, Chengxu Zhuang, Adina Williams, Ryan Cotterell, and Tal Linzen. 2025. Bigger is not always better: The importance of human-scale language modeling for psycholinguistics. <i>Journal of Memory and Language</i> , 144:104650.	1002
		1003
	Zhengxuan Wu, Atticus Geiger, Aryaman Arora, Jing Huang, Zheng Wang, Noah Goodman, Christopher Manning, and Christopher Potts. 2024. <i>pyvene: A library for understanding and improving PyTorch</i>	1004

- 984 [models via interventions](#). In *Proceedings of the 2024*
985 *Conference of the North American Chapter of the*
986 *Association for Computational Linguistics: Human*
987 *Language Technologies (Volume 3: System Demon-*
988 *strations)*, pages 158–165, Mexico City, Mexico. As-
989 sociation for Computational Linguistics.
- 990 Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christo-
991 pher Potts, and Noah Goodman. 2023. Interpretabil-
992 ity at scale: Identifying causal mechanisms in alpaca.
993 *Advances in neural information processing systems*,
994 36:78205–78226.
- 995 Charles D Yang. 2004. Universal grammar, statistics or
996 both? *Trends in cognitive sciences*, 8(10):451–456.

A Appendix

A.1 Hyperparameter Selection

Early experiments found that the default hyperparameters reported in the [Boguraev et al. \(2025\)](#) codebase (batch size 25×16 steps) resulted in undertrained DAS vectors. This is possible because the Pythia 1.4B model has far more parameters and was trained on far more data compared to BabyLM-100M.

We conducted a hyperparameter sweep to determine optimal DAS training parameters. Figure 4 and 5 show **MAX ODDS** across different batch sizes (8, 16, 25, 32) and training steps (40, 60, 80, 100, 120) for the Wh→Wh within-construction condition at the 100M checkpoint. Based on these results, we selected a batch size of 25 with 80 training steps (2000 total samples) for all experiments.

The learning rate was fixed at the default value of 5×10^{-3} used in [Arora et al. \(2024\)](#).

A.2 Animacy Figures

Supplementing the statistical tests for animacy effects in 5.4, we plot the increase in **MAX ODDS** across training split by lexical matching conditions. Figure 6 compares the causal performance when the intervention source and target base sentences share the same animacy status (Animate → Animate) versus differing animacy status (Animate → Inanimate). Results show a consistent gap between the two conditions, suggesting the learned representation may retain sensitivity to lexical features, such as animacy, through the pretraining process.

To separate animacy effects from construction-specific variance, the reported metrics for both Figure 6 and statistical testing are averaged across wh-question and topicalization constructions.

A.3 Beyond Developmental Constraints

We also present our results for the full 1 billion token training trajectory (100M–1000M tokens)⁹ based on additional checkpoints released for the BabyLM model to better understand how filler-gap mechanisms continue to develop beyond developmentally plausible input levels. We see cross-construction generalizations plateau after 100M tokens, and performance begins to overlap for both sets of results. LMs could require more input to learn a representation specific to topicalization, since it rarely shows up in children’s in-

put. These results confirm our overall claims localization within examples of one construction is stronger than transfer across constructions, while showing inconclusive evidence for our third hypothesis regarding interactions between construction frequency and transfer.

⁹The model received 1000M tokens because it was trained on the 100M token dataset for 10 epochs.

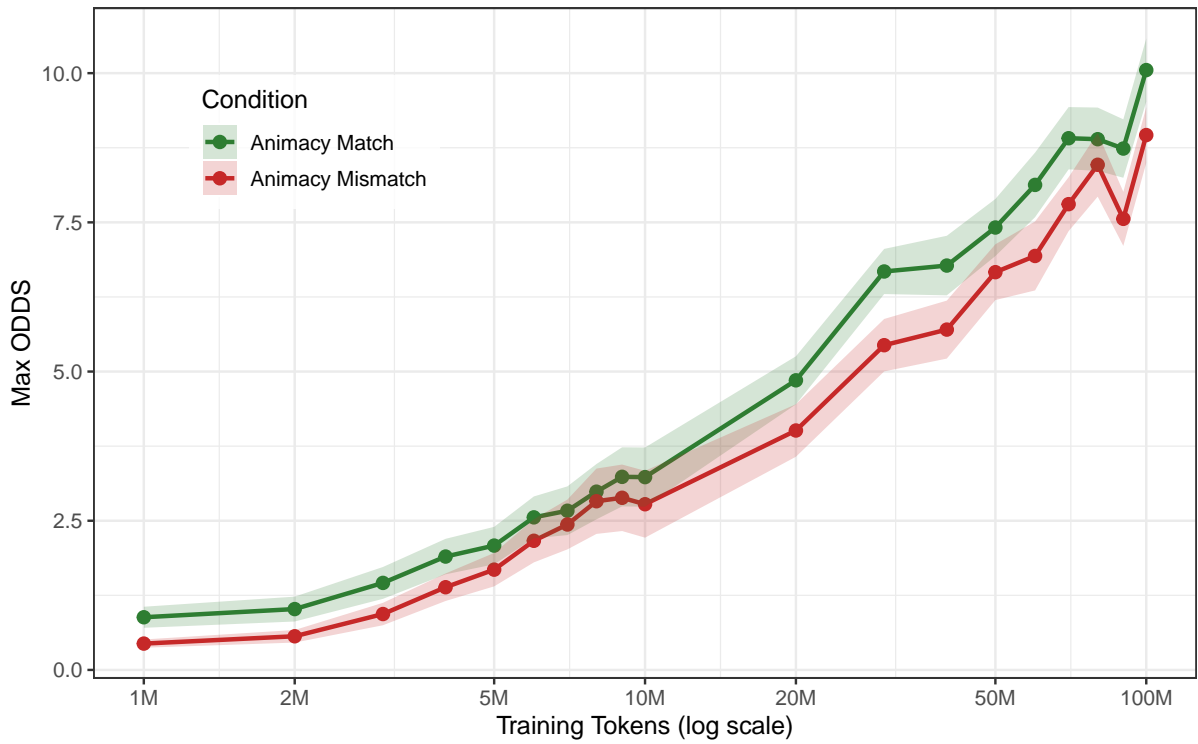


Figure 6: Developmental trajectory of lexical boost across training. Error bands show ± 1 SE across a minimum of 2 seeds.

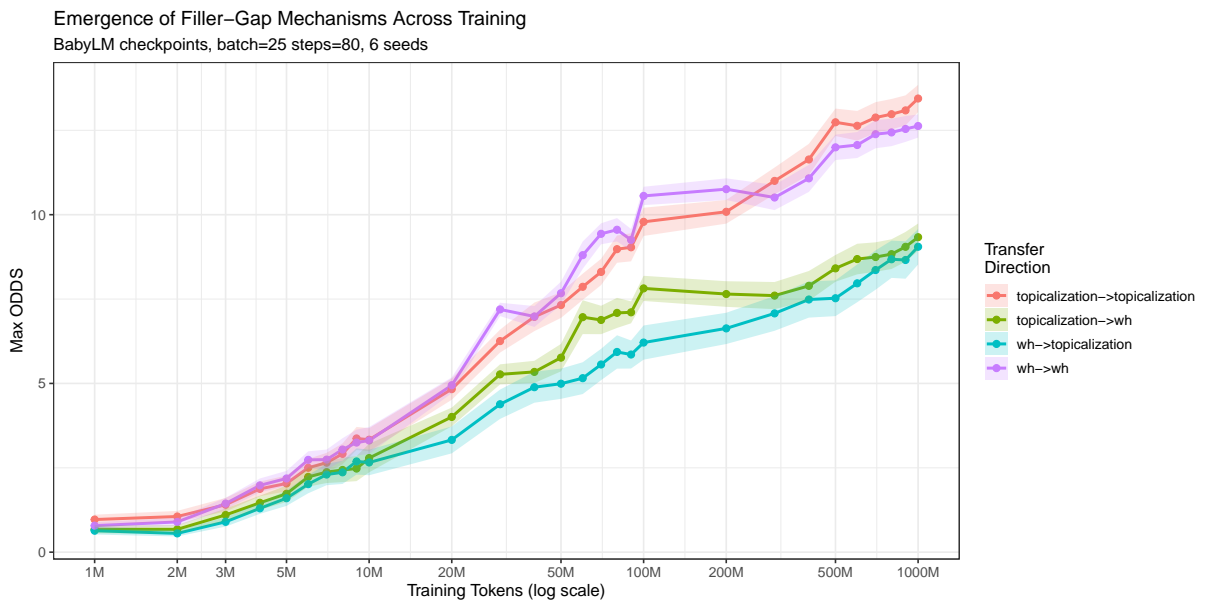


Figure 7: Full developmental trajectory from 1M to 1000M tokens. Filler-gap mechanisms continue to improve but begin to plateau around 500M–700M tokens.

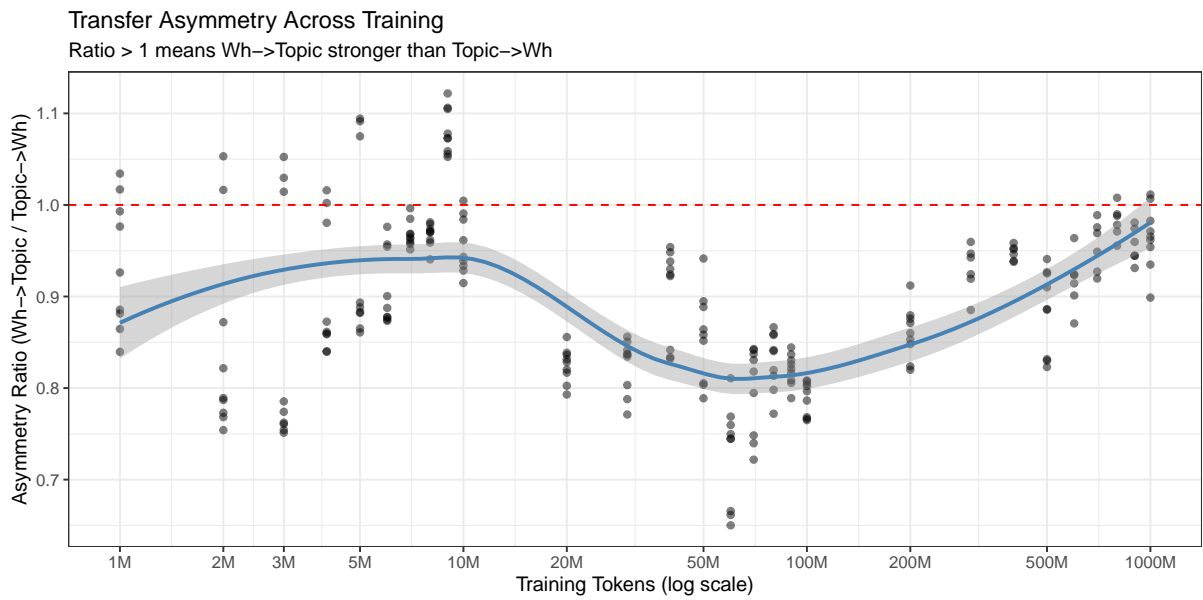


Figure 8: Transfer asymmetry across full training range. The asymmetry ($\text{Topic} \rightarrow \text{Wh} > \text{Wh} \rightarrow \text{Topic}$) persists and slightly increases at later checkpoints.