

# Language Models as Simulations of Early Language Acquisition: analysis of expressive vocabulary

Anonymous EMNLP submission

## Abstract

Large language models (LLMs) have been shown to develop linguistic competence from mere exposure to language content, making them a promising avenue for investigating infants' language learning processes (Lavechin et al., 2023; Chang and Bergen, 2022). Nevertheless, LLMs typically require orders of magnitude more data than children, and language outcomes cannot be directly compared. Here, we introduce *machine-CDI*, a metric based on the learner's output to enable a direct comparison of machines and infants on their expressive vocabulary as a function of input quantity. This metric adapts the Communicative Development Inventories (Fenson et al., 2007; Frank et al., 2017), a normalized inventory of words to quantify child language development, to the evaluation set of language models. We illustrate machine-CDI by comparing the expressive vocabulary in infants and character language models (LSTMs and Transformers) trained on English audiobooks. The results show that language models approximately match the children's learning curves, although Transformers are delayed compared to LSTMs. A further analysis show that the models are more impacted by word frequency than children, with a large delay in acquiring low frequency words for models. This delay is found to be linked to the more general phenomenon of long tail truncation observed in language models, which makes them unable to learn words based on few shot observations. These results shed new light on the principles of language acquisition, and highlights important divergences in how humans and modern algorithms learn to process natural language.

## 1 Introduction

From a cognitive perspective, language models are of theoretical interest to test the distributional hypothesis of language acquisition according to which human children learn from the statistical patterns within language data (Boleda, 2020; Lenci,

2018; Saffran et al., 1996). Previous research has shown that language models can effectively simulate aspects of the language acquisition process, such as phoneme categorization (Lavechin et al., 2023), word acquisition prediction (Chang and Bergen, 2022), and grammatical development (Evanson et al., 2023; Lavechin et al., 2023; Panitto and Herbelot, 2020). However, these studies have predominantly focused on qualitative analyses, often lacking detailed comparisons with real-world human data.

For a more quantitative approach to the simulation of language acquisition, we propose to match learning environment and outcome measures in the following two aspects. First, despite variations in socio-economic factors and cultural settings (Hart et al., 1997; Cristia et al., 2019), current estimates suggest that American English-learning children receive between 300 and 1,000 hours of speech input annually, totaling at most 30 million of words by the age of three. In contrast with modern language models trained on trillions of words (Hart et al., 1997; Cristia et al., 2019), we train models on developmentally plausible input, matched in quantity to the input that children are exposed to. Second, evaluation methods for language models should be consistent with those available from human dataset. Currently, human behavioral data are derived mainly from children's speech production (e.g., CHILDES) or parental reports (Communicative Development Inventories, hereafter *human CDI*) (MacWhinney and Snow, 1985; Fenson et al., 2007). In contrast, language model evaluations often involve zero-shot probing tasks, such as *spot-the-word* (Le Godais et al., 2017) or *grammatical acceptability judgments* (Warstadt et al., 2019), which, although inspired by psycholinguistic methods, are intrinsically different from production-based human data and typically rely on carefully designed probing sets.

To address these issues, we introduce *Machine-*

044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084

085 *CDI*, a lexical benchmark designed for direct compar- 135  
086 ison between infant language acquisition and 136  
087 language modeling. In this study, we (i) train 137  
088 vanilla character LSTMs and Transformers on de- 138  
089 velopmentally plausible data; (ii) introduce a new 139  
090 metric to evaluate how the language models' gener- 140  
091 ations fit into human acquisition patterns from the 141  
092 human CDI (Fenson et al., 2007; Frank et al., 2017), 142  
093 rather than relying on extrinsic evaluations or down- 143  
094 stream tasks; and (iii) provide a comprehensive 144  
095 analysis on generations. Our findings reveal gener- 145  
096 ally comparable curves in expressive vocabulary 146  
097 development between children and LSTMs, while 147  
098 Transformers show learning delays. A detailed ex- 148  
099 amination on missing rates and out-of-vocabulary 149  
100 rate show that this is linked to the more general 150  
101 phenomenon of long tail truncation observed in lan- 151  
102 guage models, which makes them unable to learn 152  
103 words based on few shot observations. These find- 153  
104 ings provide new insights into the principles of 154  
105 language acquisition and highlight important differ- 155  
106 ences in how humans and modern algorithms learn 156  
107 to process natural language. 157

## 108 2 Related Work 158

### 109 2.1 Language model as distributional learner 159

110 Recently, there has been substantial research ap- 160  
111 plying language models to simulate language ac- 161  
112 quisition. The earliest study by (Rumelhart et al., 162  
113 1986) investigated past tense verb form learning in 163  
114 phoneme-level neural networks, which was later 164  
115 replicated in character-level recurrent neural net- 165  
116 works Kirov and Cotterell, 2018. 166

117 Inspired by "acceptability judgments" in psy- 167  
118 cholinguistic experiments, these models are of- 168  
119 ten evaluated using zero-shot linguistic probes, 169  
120 i.e. comparing the estimated probability of legiti- 170  
121 mate sequences with that of matched implausible 171  
122 ones. Previous studies on infants' language acqui- 172  
123 sition have used this method to probe word-level 173  
124 acquisition ('spot-the-word') and sentence-level ac- 174  
125 ceptability judgments (word: (Chang and Bergen, 175  
126 2022; Vong et al., 2024); syntactic: (Evanson et al., 176  
127 2023)). Notably, linguistic abilities in children and 177  
128 models are tested differently: models are explic- 178  
129 itly tested on next-word prediction using a two- 179  
130 alternative forced-choice metric, while children are 180  
131 implicitly evaluated based on their spontaneous 181  
132 use of linguistic structures during natural speech. 182  
133 This critical difference addresses the necessity for 183  
134 a more comparable metric. 184

135 Additionally, previous studies applying language 136  
137 models to test cognitive hypotheses tend to make 137  
138 an analogy between training dynamics and lan- 138  
139 guage learning process. For instance, Chang and 139  
140 Bergen (2022) has shown the similarity of vocabu- 140  
141 lary learning curves between training steps and age 141  
142 of acquisition(AOA). One follow-up study showed 142  
143 that GPT-2 language models tend to acquire gram- 143  
144 matical knowledge in a sequential order, which 144  
145 corresponds to what has observed from transcripts 145  
146 between children-parents (Evanson et al., 2023). 146  
147 However, most of the studies focus on qualitative 147  
148 analysis by making an analogy between children's 148  
149 developmental stages and training steps. This mis- 149  
150 alignment in time scales makes the trajectories less 150  
151 comparable. Subsequent research addressed this 151  
152 issue by training self-supervised models with vary- 152  
153 ing input sizes to explicitly quantify human's de- 153  
154 velopmental trajectories (Lavechin et al., 2023). 154  
155 Their study demonstrate analogous linear growth 155  
156 patterns in lexical test, initially suggesting the ef- 156  
157 ficacy of language model for vocabulary develop- 157  
158 ment. However, there exists a discrepancy on the 158  
159 evaluation task, with the human reference data rep- 159  
160 resenting the proportion of children knowing the 160  
161 word, whereas the model is measured in probing 161  
162 task accuracy. What's more, the constructed test 162  
163 words do not directly reflect the distributional pat- 163  
164 terns of children's exposure to the words. 164

165 Therefore, the broad motivation of our work is 165  
166 to assess the distributional mechanisms in infants 166  
167 lexical acquisition (Saffran et al., 1996; Romberg 167  
168 and Saffran, 2010) using neural language models as 168  
169 distributional learners. If analogous distributional 169  
170 learning mechanisms were involved in children, 170  
171 then we would expect similar evaluation outcomes 171  
172 from the proposed lexical metrics. 172

### 173 2.2 Word representation in language models 173

174 Language models are typically trained in a way 174  
175 that take as inputs a series of token and output the 175  
176 predicted tokens. In text language models, the pre- 176  
177 dicted tokens can be characters (Xue et al., 2022), 177  
178 entire words (Mikolov et al., 2013; Pennington 178  
179 et al., 2014), or word fragments, for example, byte 179  
180 pair encodings, (Sennrich et al., 2015)). In these 180  
181 two latter cases, the training of the LMs is done 181  
182 in two phases: first a tokenizer is learned, using 182  
183 spaces or punctuation to delimit the word bound- 183  
184 aries, and the training corpus is tokenized; second 184  
185 the LM is trained with the token-prediction objec-

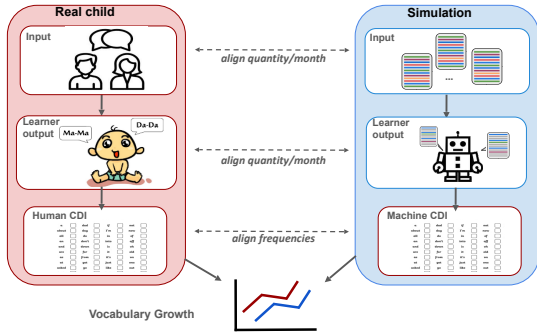


Figure 1: **Overview of Machine-CDI.** Models are fed with linguistic input in matched quantities compared to children and their output is also matched to their outputs. A list of test words is prepared for the machine, matched to the human-CDI in word frequency distribution, yielding comparable vocabulary growth curves.

185 tive. Tokenization assumes that beginning and ends  
 186 of words can be identified prior learning, while in-  
 187 fants typically acquire language from continuous  
 188 speech without without knowing the correct linguis-  
 189 tic labels like word boundaries a prior (Maye et al.,  
 190 2002). Prior work (Sutskever et al., 2011; Graves  
 191 et al., 2014; Hahn and Baroni, 2019; Nguyen et al.,  
 192 2022; Boldsen et al., 2022; Yu et al., 2024), have  
 193 found that character LMs can learn lexical, syntac-  
 194 tic and semantic representations and do not need  
 195 a prior segmentation in words. In such models,  
 196 words are latent representations instead of being  
 197 explicitly represented, making them a promising  
 198 approach to simulate the process of lexical learn-  
 199 ing.

200 Therefore, we use characters as tokens and leave  
 201 the language models to learn words in an unsu-  
 202 pervised fashion in this work. Instead of applying  
 203 models that are trained on speech or phonemes, we  
 204 start with character language models to take the  
 205 simplified invariant representations of word forms  
 206 as input regardless of the acoustic variability. This  
 207 might provide an upper bound of the overall model  
 208 performance when compared with human data.

### 209 3 Method

210 We follow the reverse engineering approach in  
 211 (Dupoux, 2018) where a simulation jointly models  
 212 the input to the learner, the learner and the outcome  
 213 measure in a quantitative fashion (Figure 1).

## 214 3.1 Metrics

### 215 3.1.1 Machine-CDI

216 We detail two key aspects of Machine-CDI that  
 217 enable direct comparisons on the learning speed of  
 218 human and machine: data quantity alignment and  
 219 evaluation metric alignment.

220 **Aligning input data quantity with infants’ lin-**  
 221 **guistic exposure** Initially, we standardized the  
 222 amount of training data to reflect the estimated  
 223 speech exposure of each child, based on prior re-  
 224 search of an average annual exposure of 1,000  
 225 hours of speech per child (Mendoza and Fausey,  
 226 2021). This number corresponds to an upper bound  
 227 rather than an average when taking into consider-  
 228 ation cross-linguistic variation , with the median  
 229 being 500h/year, and the minimum around 60h/y  
 230 or less(Cristia, 2023). We converted speech du-  
 231 ration into corresponding word counts, assuming  
 232 an average speech rate of 3 words per second, or  
 233 10,000 words per hour as is shown in Figure 2c.  
 234 We then trained models with varying input sizes to  
 235 represent infant cohorts at different developmen-  
 236 tal stages to ensure a realistic approximation of  
 237 linguistic exposure.

238 **Aligning evaluation metric** We base our model  
 239 evaluation task on human CDI (Fenson et al., 2007;  
 240 Frank et al., 2017), a checklist of representative  
 241 word samples used to measure word learning. Our  
 242 evaluation tasks reported binary scores for a word  
 243 set to align with parents’ binary reports (Frank  
 244 et al., 2017).

245 To construct the model’s evaluation sets, we se-  
 246 lected 520 words from the training set by match-  
 247 ing both the number and frequency distribution  
 248 of content words in the human CDI set. To esti-  
 249 mate word frequency of human CDI set, we con-  
 250 catenated and cleaned all the transcripts from the  
 251 English CHILDES Language Data Exchange Sys-  
 252 tem (CHILDES) database, resulting in 14.5 mil-  
 253 lion words of adult speech. For the machine CDI  
 254 word set, frequency was derived from the largest  
 255 language model training set, consisting of approxi-  
 256 mately 1 million word types from 30M words. We  
 257 iteratively minimized the loss function between the  
 258 frequency distributions of the human and machine  
 259 CDI sets.

260 To bridge the gap between parental reports  
 261 and observable language model performance, we  
 262 aligned the vocabulary growth curve derived from  
 263 children’s speech in the CHILDES corpus with hu-  
 264 man CDI scores (see Figure 2b). We constructed

the vocabulary growth curve by calculating the cumulative word frequency for each month and applying a constant count threshold to convert these results into binary scores. Due to data sparsity, the word frequency was recalibrated to approximate monthly speech production based on previous research on child vocalization duration. We selected the count threshold to best fit the human CDI growth curve, assuming that an optimal threshold would closely approximate the observed growth speed (see Figure 2c). Figure 2b shows that, on average, a child is expected to produce a correct word form approximately 60 times per month.

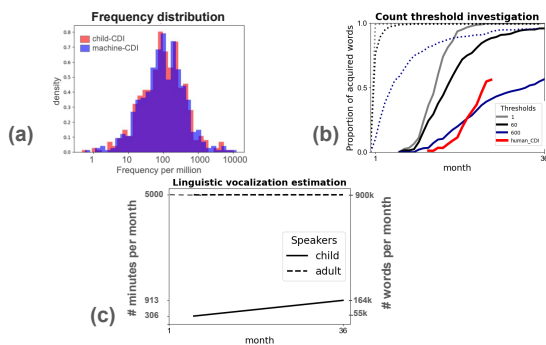


Figure 2: *Calibration method.* a. The distribution of frequency per million of the machine CDI (reference corpus: audiobooks) is matched to that of child-CDI (reference corpus: English-CHILDES). b. To decide whether a word is ‘known’ based on the speaker’s output, we count instances of the word and apply a threshold. In plain blue, resulting vocabulary growth curves for different thresholds. In dotted blue, application of this criterion to the adult’s own input. In red, parental reports of children vocabulary. c. Estimates of monthly parental output (input to the child) and child output.

### 3.1.2 Generation set analysis

To provide a detailed analysis of the generated dataset by language models, we conducted an examination of the *missing rate* and *out-of-vocabulary (OOV) rate*. Specifically, we generated an equal number of words using LSTM and Transformer language models, each trained on a 3.4 million-word corpus. The missing rate quantifies the proportion of words present in the training set but absent in the generated set, indicating the extent to which the generated data replicates the vocabulary of the training set. The OOV rate measures the proportion of words in the generated set that do not appear in the training set, thus assessing the model’s ability to generalize beyond the training data. Additionally, we evaluate the true-

word rate as the proportion of correct word forms in the generated set, to further assess the model’s generalizability. For this evaluation, we use a combination of word lists from CELEX (Van Heuven et al., 2014), the Enchant Library<sup>1</sup>, and Wiktionary<sup>2</sup> as a spelling checker.

To contextualize the characteristics of the model-generated text, we constructed two equal-sized reference sets respectively: the in-domain test set, which are selected from audiobook transcripts that are aligned in genres but not included in the training set, and an out-of-domain set consisting of child-directed speech from the English CHILDES corpus.

### 3.2 Developmentally plausible training set

Following STELA (Lavechin et al., 2022), we built a developmentally plausible training set from the orthographic transcripts of LibriVox English audiobooks (Kearns, 2014), consisting subsets of transcripts of 400h, 800h, 1600h and 3200h respectively. Given our calibration of 1000h/year, this translates into 4.8m, 9.6m, 19.2m, and 38.4m, respectively.

### 3.3 Models

We applied two types of models on developmentally plausible datasets to simulate language acquisition: probabilistic language models (including Long Short-Term Memory models and decoder-only transformers 1), and a non-parametric Bayesian model: Chinese Restaurant Process (Gershman and Blei, 2012).

**Chinese Restaurant Process (CRP)** The non-parametric Bayesian CRP model clusters data by assigning probabilities. A new word joins an existing cluster with a probability proportional to the cluster size and starts a new cluster with a probability proportional to a parameter  $\alpha$  (Gershman and Blei, 2012). We initialized the CRP model using a 3-gram language model derived from the developmentally plausible datasets.

**Neural Network Architectures** We employed two types of neural network architectures: decoder-only LSTMs (Hochreiter and Schmidhuber, 1997) and Transformers (Vaswani et al., 2017). Similar performance from both models would indicate that the learned patterns are robustly present in the data, not artifacts of a specific model architecture.

<sup>1</sup><https://pypi.org/project/pyenchant/>

<sup>2</sup><https://www.wiktionary.org>

For LSTMs, we used a three-layer architecture with an embedding layer of size 200, hidden layers of size 1024, and a feed-forward output layer of size 200, based on prior work (Lavechin et al., 2023). For the Transformer model, we experimented on different attention heads and decoding layers, and ended with 8 attention heads and 6 decoding layers that yielded the optimal perplexity.

### 3.4 Generations

The generation process for the Chinese Restaurant Process (CRP) model uses a concentration parameter,  $\alpha$ , derived from the training data to simulate word occurrences. This approach ensures that each token is generated either as a new word or by incrementing the count of an existing word, thereby creating a corpus that reflects the token distribution dynamics as modeled by the CRP.

For LSTMs and Transformers, we employed both unprompted and prompted generation using temperature sampling. Temperature sampling adjusts the output logits by dividing them by a temperature parameter before sampling from the distribution. Higher temperature values make the distribution more uniform, increasing randomness. The number of prompts was matched to the number of sentences generated. Each prompt consisted of 3 words from the audiobook dataset that do not exist in the training set as a whole sequence. We also excluded sequences containing words in the machine CDI test set.

## 4 Results

### 4.1 Human-model comparison

#### Protracted development of language models

The results from our lexical benchmark, illustrated in Figure 3, reveal significant dependencies of vocabulary growth curves on model architectures. The CRP models consistently demonstrate higher vocabulary sizes across all months compared to probabilistic language models. In contrast, probabilistic language models, including LSTMs and Transformers, exhibit slower vocabulary growth, a trend that persists regardless of decoding temperatures and prompt types. This difference highlights a fundamental limitation in their ability to mimic human-like vocabulary expansion.

A closer examination reveals that LSTM models align more closely with human vocabulary growth curves than Transformers, particularly in unprompted generations. Manual analysis of gen-

erated utterances indicates that unprompted generations from Transformers frequently suffer from repeated characters. The architectural influence on model fitness is less pronounced but still present in prompted generations, but the effect is reverse in two architectures: with the LSTM fitness with human vocabulary growth interfered with the prompts; the Transformer’s fitness largely increased by the prompts.

Variations in temperature during the generation process yield similar trends across different experiment settings, with closest fitness observed for temperature settings around 1. Lower temperatures produce more deterministic outputs, while higher temperatures result in more random and noisy generations. These results highlight the importance of model architecture and generation settings in simulating human language acquisition, suggesting that incorporating mechanisms to handle memory and context appropriately could enhance the vocabulary learning capabilities of probabilistic language models.

**Frequency effect** So far, our findings indicate that language models acquire lexical knowledge less effectively than humans, regardless of decoding methods and prompts. One potential reason for this discrepancy is the models’ difficulty with infrequent words. To investigate this, we decomposed the CDI words into six frequency bands, each containing equal number of words, and fitted sigmoid functions for each frequency band (Chang and Bergen, 2022). We then calculated the estimated month for each frequency band to 80% of known words.

As shown in Figure 4, child speech is less influenced by input word frequency than language models, whereas the expressive vocabulary growth of language models is significantly affected by lower-frequency words across all experimental settings. This frequency effect varies by model architectures, with Transformer models require considerably more training data to reproduce the same proportion of infrequent words compared to LSTM models.

Additionally, lower temperatures exacerbate this effect, likely due to the altered probability distribution generated by the output layers of the language models. These observations suggest that while humans learn words more uniformly across frequencies, probabilistic language models struggle with lower-frequency words, and their performance

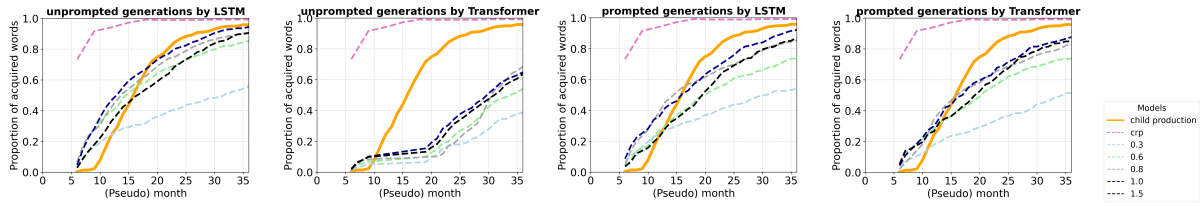


Figure 3: **Vocabulary growth curves for children and models.** Words are considered known if produced more than 60 times. For LSTM and Transformer models, unprompted and prompted generation were sampled at different temperatures. We also plot the curves for an accumulator model and a CRP model.

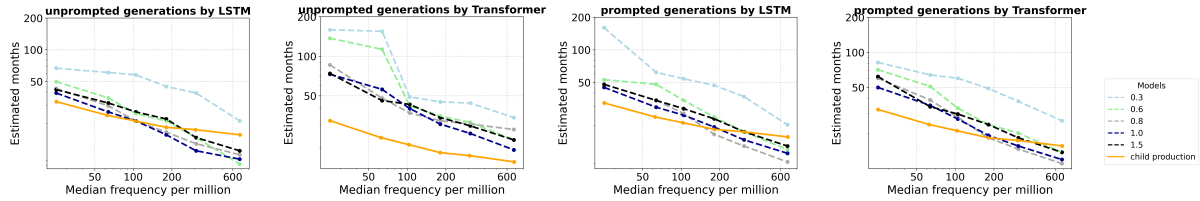


Figure 4: **Estimated age of acquisition as a function of word frequency.** Growth curves for 6 frequency bins of child- and machine-CDI word lists were fitted with a sigmoid and number of months to reach 80% of known words was computed.

441 is further influenced by architectural choices and  
 442 decoding methods.

443 **4.2 Generation set analysis**

444 Notably, the frequency effects observed in the Ma-  
 445 chine CDI are evaluated on a selected subset though  
 446 matched carefully with human-CDI set. It is un-  
 447 clear whether this effect is influenced by the ran-  
 448 domness of sampling. To address this, we expand  
 449 our analysis on the generations using the same size  
 450 of training set, focusing on models trained on 3.4M  
 451 words as a case study. Specifically, we investi-  
 452 gated the potential reasons for the human-model  
 453 discrepancy on vocabulary growth curve. Our anal-  
 454 ysis examined whether the discrepancy arises from  
 455 the omission of infrequent words in the training  
 456 data (missing rate) and whether it can be miti-  
 457 gated by the models generating novel sequences (OOV  
 458 rate).

459 Figure 5 illustrates that the overall missing  
 460 proportion of word types generated by language  
 461 models is higher than that of the reference in-  
 462 domain test set across various settings. This sug-  
 463 gests that language models tend to omit a sig-  
 464 nificant proportion of word types in their gener-  
 465 ated sets. Additionally, lower temperatures result  
 466 in higher missing rates across different experimen-  
 467 tal settings.

468 The comparison with OOV rates reveals a sub-  
 469 stantial gap between the proportion of missing

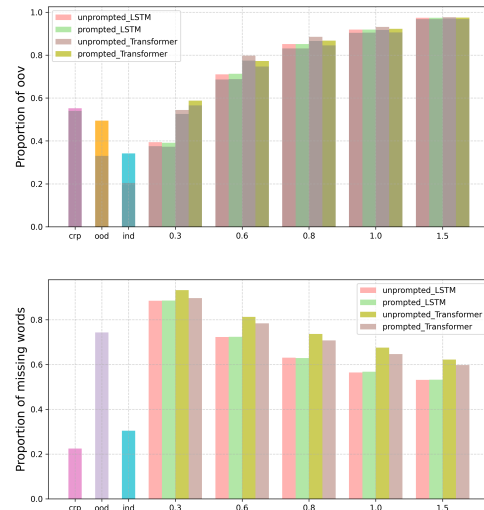


Figure 5: **OOV and word missing rates** Top: the proportion of out-of-vocabulary token types among the generation types; the shaded part shows the proportion of non-words. Below: the proportion of missing token types out of train token types

470 word types and the amount of novel sequence types  
 471 generated. Moreover, Figure 5 highlights a very  
 472 high non-word rate across different generation sets,  
 473 indicating that current models struggle to general-  
 474 ize through compositional rules.

475 We further examined whether the missing words  
 476 are influenced by their frequency in training set.  
 477 Figure 6 shows that most missing words are in  
 478 lower-frequency bands, which indicates LMs' de-

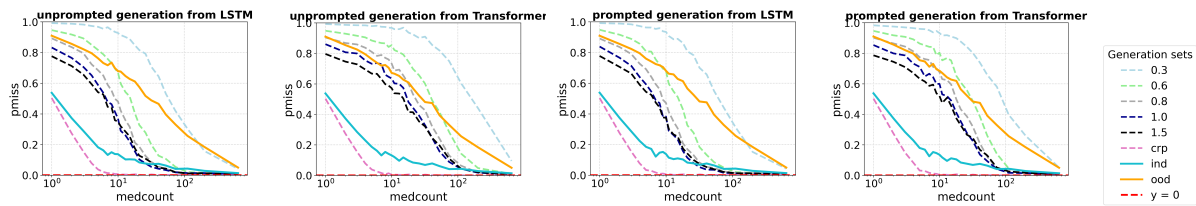


Figure 6: **Missing rates as a function of word count in input corpus.** The subfigures show the missing rates for unprompted and prompted generations from LSTMs and Transformers, respectively. Perfect memorization corresponds to no missing word types from the generation sets, with  $y=0$  as the baseline.

479 efficiency in reproducing words in tail distributions.  
 480 Further comparison across different temperatures  
 481 reveals a negative relationship between tempera-  
 482 tures and the proportion of missing words. Simi-  
 483 larly, as shown in Figure 7, OOV words predomi-  
 484 nantly appear in low-frequency bands across differ-  
 485 ent experimental settings. This suggests that while  
 486 language models exhibit some degree of general-  
 487 ization, this effect is minimal and limited to a few  
 488 instances.

## 489 5 Discussion

490 In this study, we assess the distributional mech-  
 491 anisms in infants language acquisition (Saffran  
 492 et al., 1996; Romberg and Saffran, 2010) using  
 493 neural language models as distributional learners.  
 494 Our results demonstrate that a purely distributional  
 495 learner trained on text only approximately repro-  
 496 duce human’s expressive vocabulary growth.

497 We found two main differences. First, the mod-  
 498 els are much more influenced by word frequency  
 499 compared to children. This yield a delay in word  
 500 acquisition for low frequency words. Further anal-  
 501 yses show that very low frequency items (seen less  
 502 than 10 times in the training corpus) tend to be  
 503 overwhelmingly missed by the language models.  
 504 Our findings suggest that while current language  
 505 models approximate the statistical properties of  
 506 their training data, this does not necessarily im-  
 507 ply generating the desired underlying data distri-  
 508 bution across various model architectures and de-  
 509 coding methods. This echoes prior research on  
 510 language model’s memorization, in which a log-  
 511 linear trend between the number of duplicates in  
 512 the training data and the extent of verbatim memo-  
 513 rization (Carlini et al., 2022; Razeghi et al., 2022;  
 514 Kandpal et al., 2022). In contrast, evidence show  
 515 that children can learn new words in a few shot fash-  
 516 ion, suggesting that they may use different learning  
 517 mechanisms (e.g., episodic memory), not available

518 in LMs. Prior study using non-parametric knowl-  
 519 edge to capture long-tail information has shown a  
 520 promising avenue to simulate the episodic memory  
 521 mechanism (Kandpal et al., 2023). Further investi-  
 522 gation needs to be done on cognitive plausibility.  
 523 Second, the models tend to produce a large quantity  
 524 of novel word forms (more than 80% of the word  
 525 forms), the vast majority of which are nonwords.  
 526 This corresponds to the well known tendency of  
 527 LLMs to ‘hallucinate’ (Ji et al., 2023). In only a  
 528 small fraction of the cases, these hallucinations are  
 529 actual words, obtained through the combinatory  
 530 recombination of known words or morphemes. In  
 531 contrast, infants do not produce many nonwords,  
 532 and these nonwords tend to be due to be mispronun-  
 533 ciations of real words.

534 These differences could be interpreted in terms  
 535 both of learning mechanisms and input. Children’s  
 536 linguistic experience is grounded in multi-modal  
 537 experience. Research shows that children as young  
 538 as ten months old learn word-object pairings, map-  
 539 ping novel words onto perceptually salient objects  
 540 (Pruden et al., 2006). By the age of two, they in-  
 541 tegrate social cues such as eye gaze, pointing, and  
 542 joint attention (Çetinçelik et al., 2021). Our find-  
 543 ings suggest that these grounded and interactive  
 544 experiences could impact child word acquisition in  
 545 ways that cannot be fully explained by linguistic  
 546 signals alone. Additionally, the communicative na-  
 547 ture of the language environment provides a more  
 548 dynamic context where infants receive feedback  
 549 from caregivers. Studies on reinforcement learning  
 550 in multi-agent communication tasks highlight the  
 551 importance of these non-distributional properties  
 552 for achieving more human-like natural language  
 553 understanding. For example, research by Chevalier-  
 554 Boisvert et al. (2018), Lazaridou et al. (2016), and  
 555 Zhu et al. (2020) emphasizes the role of interaction  
 556 and feedback in language learning.

557 In this paper, we have described how lexical eval-

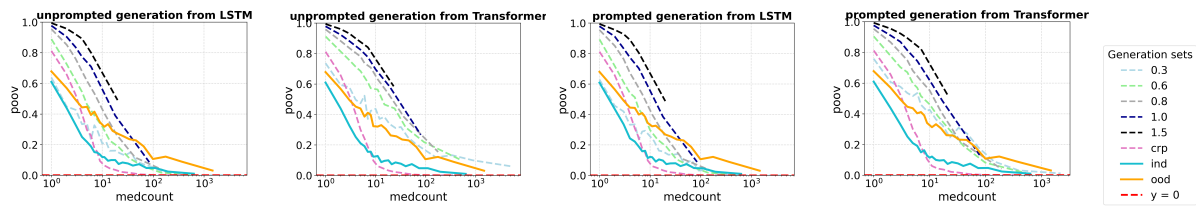


Figure 7: **OOV rates as a function of word count in generated corpus.** The subfigures show the OOV rates from unprompted and prompted generations from LSTMs and Transformers, respectively. This score reflects the novelty of the generation sets.

558 uation metrics have been carefully designed to evaluate  
 559 language models trained on developmentally plausible  
 560 text corpus. Notably, we only focused on the word  
 561 form inspection, which might inflate the model  
 562 performance. Even with the upper bound of the  
 563 model performance, the models are delayed for  
 564 expressive vocabulary. And we found a stable  
 565 frequency effect across different language model  
 566 architectures and decoding settings. We show that  
 567 this is linked to the more general phenomenon of  
 568 long tail truncation observed in language models,  
 569 which makes them unable to learn words based  
 570 on few shot observations. These results shed new  
 571 light on the principles of language acquisition, and  
 572 highlights important divergences in how humans  
 573 and modern algorithms learn to process natural  
 574 language.

## 575 6 Limitations

576 One limitation lies in discrepancies between infant  
 577 behavioral measures derived from parental reports  
 578 and those evaluated on our model. This difference  
 579 potentially accounts for differences observed  
 580 between estimates from CHILDES transcripts and  
 581 parental reports. On the one hand, the model’s  
 582 evaluation test, focusing on word form segmenta-  
 583 tion from input data. In contrast, parental criteria  
 584 may involve the proper usage of the test words,  
 585 which typically involve a broader scope of linguis-  
 586 tic knowledge on semantic and syntactic levels.  
 587 On the other one hand, the CHILDES transcripts,  
 588 though pre-processed carefully to remove all the  
 589 human annotations, it is possible that incomplete  
 590 word forms are completed and normalized by hu-  
 591 man annotators, which may cause the inflated lan-  
 592 guage performance in the transcript. Also, some  
 593 subsets interview procedure’s richness may boost  
 594 children’s expressivity beyond everyday speech,  
 595 potentially leading parents to underestimate vocabu-  
 596 lary in CDI inventories. What’s more, CHILDES

597 transcripts aggregate data from multiple children,  
 598 while parental reports are averaged on single child.  
 599 Also, we calibrated the word counts based on the  
 600 estimated vocalization length. This might result in  
 601 duplicated counts on children’s production. Nev-  
 602 ertheless, all these differences might inflate the  
 603 lexical scores obtained from transcripts. Notably,  
 604 we apply the exactly the same post-process on the  
 605 model’s generation and also compare model’s gen-  
 606 erations with CHILDES references. Therefore, this  
 607 might exert a trivial influence on CHILDES-model  
 608 difference.

609 Another limitation lies in the usage of character-  
 610 level input rather than speech input. Characters  
 611 preserve the invariant form of words, and space or  
 612 punctuation indicate word boundaries. Hence the  
 613 models we tested correspond to an upper bound  
 614 of what could be found with some realistic mod-  
 615 els based on speech inputs (Lavechin et al., 2023,  
 616 2024), where word forms are variable and not de-  
 617 limited with clear boundaries. Further studies are  
 618 needed to evaluate speech-LMs (Lakhotia et al.,  
 619 2021; Nguyen et al., 2024) and address the techni-  
 620 cal difficulty of transcribing the speech output of  
 621 such models in a format that can be applied to our  
 622 machine-CDI benchmark.

## 623 Ethics Statement

624 Use of human data: While we did not collect any  
 625 new human data ourselves, many of our analy-  
 626 ses involved the use of prior datasets within the  
 627 CHILDES database. All of these datasets were  
 628 collected in accordance with IRB policies at the  
 629 institutions of the data collectors, and all followed  
 630 standard practices in obtaining informed consent  
 631 and deidentifying data.



632	<b>References</b>		683
633	Maria Boldsen, Per Anker Jensen, Lars Kai Hansen, and Ole Winther Andersen. 2022. Perceptual representations vs character embeddings in cross-lingual analysis. <i>arXiv preprint arXiv:2203.12345</i> .		684
634			685
635			686
636			687
637	Gemma Boleda. 2020. Distributional semantics and linguistic theory. <i>Annual Review of Linguistics</i> , 6:213–234.		688
638			689
639			690
640	Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2022. Quantifying memorization across neural language models. <i>arXiv preprint arXiv:2202.07646</i> .		691
641			692
642			693
643			694
644	Melis Çetinçelik, Caroline F Rowland, and Tineke M Snijders. 2021. Do the eyes have it? a systematic review on the role of eye gaze in infant language development. <i>Frontiers in psychology</i> , 11:589096.		695
645			696
646			697
647			698
648	Tyler A Chang and Benjamin K Bergen. 2022. Word acquisition in neural language models. <i>Transactions of the Association for Computational Linguistics</i> , 10:1–16.		699
649			700
650			701
651			702
652	Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu Nguyen, and Yoshua Bengio. 2018. Babyai: A platform to study the sample efficiency of grounded language learning. <i>arXiv preprint arXiv:1810.08272</i> .		703
653			704
654			705
655			706
656			707
657	Alejandrina Cristia. 2023. A systematic review suggests marked differences in the prevalence of infant-directed vocalization across groups of populations. <i>Developmental Science</i> , 26(1):e13265.		708
658			709
659			710
660			711
661	Alejandrina Cristia, Emmanuel Dupoux, Michael Gerven, and Jonathan Stieglitz. 2019. Child-directed speech is infrequent in a forager-farmer population: A time allocation study. <i>Child development</i> , 90(3):759–773.		712
662			713
663			714
664			715
665			716
666	Emmanuel Dupoux. 2018. Cognitive science in the era of artificial intelligence: A roadmap for reverse-engineering the infant language-learner. <i>Cognition</i> , 173:43–59.		717
667			718
668			719
669			720
670	Linnea Evanson, Yair Lakretz, and Jean-Rémi King. 2023. Language acquisition: do children and language models follow similar learning stages? <i>arXiv preprint arXiv:2306.03586</i> .		721
671			722
672			723
673			724
674	Larry Fenson et al. 2007. Macarthur-bates communicative development inventories.		725
675			726
676	Michael C Frank, Mika Braginsky, Daniel Yurovsky, and Virginia A Marchman. 2017. Wordbank: An open repository for developmental vocabulary data. <i>Journal of child language</i> , 44(3):677–694.		727
677			728
678			729
679			730
680	Samuel J Gershman and David M Blei. 2012. A tutorial on bayesian nonparametric models. <i>Journal of Mathematical Psychology</i> , 56(1):1–12.		731
681			732
682			733
	Alex Graves, Greg Wayne, and Ivo Danihelka. 2014. Neural turing machines. <i>arXiv preprint arXiv:1410.12345</i> .		734
			735
	Yanzhu Guo, Guokan Shang, Michalis Vazirgiannis, and Chloé Clavel. 2023. The curious decline of linguistic diversity: Training language models on synthetic text. <i>arXiv preprint arXiv:2311.09807</i> .		736
			737
	Michael Hahn and Marco Baroni. 2019. Neural models learn morphological, syntactic and semantic aspects from unsegmented text. <i>arXiv preprint arXiv:1902.12345</i> .		738
			739
	Betty Hart, Todd R Risley, and John R Kirby. 1997. Meaningful differences in the everyday experience of young american children. <i>Canadian Journal of Education</i> , 22(3):323.		740
			741
	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. <i>Neural computation</i> , 9(8):1735–1780.		742
			743
	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM Computing Surveys</i> , 55(12):1–38.		744
			745
	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In <i>International Conference on Machine Learning</i> , pages 15696–15707. PMLR.		746
			747
	Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. Deduplicating training data mitigates privacy risks in language models. In <i>International Conference on Machine Learning</i> , pages 10697–10707. PMLR.		748
			749
	Jodi Kearns. 2014. Librivox: Free public domain audiobooks. <i>Reference Reviews</i> , 28(1):7–8.		750
			751
	Christo Kirov and Ryan Cotterell. 2018. Recurrent neural networks in linguistic theory: Revisiting pinker and prince (1988) and the past tense debate. <i>Transactions of the Association for Computational Linguistics</i> , 6:651–665.		752
			753
	Kushal Lakhota, Eugene Kharitonov, Wei-Ning Hsu, Yossi Adi, Adam Polyak, Benjamin Bolte, Tu-Anh Nguyen, Jade Copet, Alexei Baevski, Abdelrahman Mohamed, et al. 2021. On generative spoken language modeling from raw audio. <i>Transactions of the Association for Computational Linguistics</i> , 9:1336–1354.		754
			755
	Marvin Lavechin, Maureen de Seyssel, Marianne Métais, Florian Metz, Abdelrahman Mohamed, Hervé Bredin, Emmanuel Dupoux, and Alejandrina Cristia. 2024. Modeling early phonetic acquisition from child-centered audio data. <i>Cognition</i> , 245:105734.		756
			757
	Marvin Lavechin, Maureen De Seyssel, Hadrien Titeux, Hervé Bredin, Guillaume Wisniewski, Alejandrina		758

736	Cristia, and Emmanuel Dupoux. 2022. Can statistical learning bootstrap early language acquisition? a modeling investigation.	<i>on empirical methods in natural language processing (EMNLP)</i> , pages 1532–1543.	790
737			791
738			
739	Marvin Lavechin, Yaya Sy, Hadrien Titeux, María Andrea Cruz Blandón, Okko Räsänen, Hervé Bredin, Emmanuel Dupoux, and Alejandrina Cristia. 2023. Babyslm: language-acquisition-friendly benchmark of self-supervised spoken language models. <i>arXiv preprint arXiv:2306.01506</i> .	Shannon M Pruden, Kathy Hirsh-Pasek, Roberta Michnick Golinkoff, and Elizabeth A Hennon. 2006. The birth of words: Ten-month-olds learn words through perceptual salience. <i>Child development</i> , 77(2):266–280.	792
740			793
741			794
742			795
743			796
744			
745	Angeliki Lazaridou, Alexander Peysakhovich, and Marco Baroni. 2016. Multi-agent cooperation and the emergence of (natural) language. <i>arXiv preprint arXiv:1612.07182</i> .	Yasaman Razeghi, Robert L Logan IV, Matt Gardner, and Sameer Singh. 2022. Impact of pretraining term frequencies on few-shot reasoning. <i>arXiv preprint arXiv:2202.07206</i> .	797
746			798
747			799
748			800
749	Gaël Le Godais, Tal Linzen, and Emmanuel Dupoux. 2017. Comparing character-level neural language models using a lexical decision task. In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers</i> , pages 125–130.	Alexa R Romberg and Jenny R Saffran. 2010. Statistical learning and language acquisition. <i>Wiley Interdisciplinary Reviews: Cognitive Science</i> , 1(6):906–914.	801
750			802
751			803
752			
753			804
754			805
755	Alessandro Lenci. 2018. Distributional models of word meaning. <i>Annual review of Linguistics</i> , 4:151–171.	David E Rumelhart, Geoffrey E Hinton, James L McClelland, et al. 1986. A general framework for parallel distributed processing. <i>Parallel distributed processing: Explorations in the microstructure of cognition</i> , 1(45-76):26.	806
756			807
757	Brian MacWhinney and Catherine Snow. 1985. The child language data exchange system. <i>Journal of child language</i> , 12(2):271–295.	Jenny R Saffran, Richard N Aslin, and Elissa L Newport. 1996. Statistical learning by 8-month-old infants. <i>Science</i> , 274(5294):1926–1928.	808
758			809
759			810
760	Jessica Maye, Janet F Werker, and LouAnn Gerken. 2002. Infant sensitivity to distributional information can affect phonetic discrimination. <i>Cognition</i> , 82(3):B101–B111.	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. <i>arXiv preprint arXiv:1508.07909</i> .	811
761			812
762			813
763			814
764	Jennifer K Mendoza and Caitlin M Fausey. 2021. Quantifying everyday ecologies: Principles for manual annotation of many hours of infants’ lives. <i>Frontiers in psychology</i> , 12:710636.	Ilya Sutskever, James Martens, and Geoffrey Hinton. 2011. Generating text with recurrent neural networks. In <i>Proceedings of the 28th International Conference on Machine Learning (ICML-11)</i> , pages 1017–1024.	815
765			816
766			817
767			818
768	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. <i>arXiv preprint arXiv:1301.3781</i> .	Walter JB Van Heuven, Pawel Mandera, Emmanuel Keuleers, and Marc Brysbaert. 2014. Subtlex-uk: A new and improved word frequency database for british english. <i>Quarterly journal of experimental psychology</i> , 67(6):1176–1190.	819
769			820
770			821
771			822
772	Tu Anh Nguyen, Maureen de Seyssel, Robin Algayres, Patricia Roze, Ewan Dunbar, and Emmanuel Dupoux. 2022. Are word boundaries useful for unsupervised language learning? <i>arXiv preprint arXiv:2210.02956</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	823
773			824
774			825
775			826
776			827
777	Tu Anh Nguyen, Benjamin Muller, Bokai Yu, Marta R Costa-Jussa, Maha Elbayad, Sravya Popuri, Paul-Ambroise Duquenne, Robin Algayres, Ruslan Mavlyutov, Itai Gat, et al. 2024. Spirit-lm: Interleaved spoken and written language model. <i>arXiv preprint arXiv:2402.05755</i> .	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.	828
778			829
779			830
780			831
781			832
782			833
783	Ludovica Pannitto and Aurélie Herbelot. 2020. Recurrent babbling: evaluating the acquisition of grammar from limited input data. <i>arXiv preprint arXiv:2010.04637</i> .	Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. 2019. Neural network acceptability judgments. <i>Transactions of the Association for Computational Linguistics</i> , 7:625–641.	834
784			835
785			836
786			
787	Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In <i>Proceedings of the 2014 conference</i>	Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. Byt5: Towards a token-free future with pre-trained byte-to-byte models. <i>Transactions of the Association for Computational Linguistics</i> , 10:291–306.	837
788			838
789			839

843 Lili Yu, Dániel Simig, Colin Flaherty, Armen Agha-  
844 janyan, Luke Zettlemoyer, and Mike Lewis. 2024.  
845 Megabyte: Predicting million-byte sequences with  
846 multiscale transformers. *Advances in Neural Infor-*  
847 *mation Processing Systems*, 36.

848 Wang Zhu, Hexiang Hu, Jiacheng Chen, Zhiwei Deng,  
849 Vihan Jain, Eugene Ie, and Fei Sha. 2020. Babywalk:  
850 Going farther in vision-and-language navigation by  
851 taking baby steps. *arXiv preprint arXiv:2005.04625*.

months	# words	# char	# utt
4-8	3.4M	17.7M	0.3M
9-18	7.0M	36.6M	0.7M
19-28	14.7M	76.7M	1.4M
29-36	27.7M	144.9M	2.6M

Table 1: Statistics of train data

Hyperparameter	Value
Max sequence length	40960
Batch size	40960
Learning rate	0.0001
Learning rate scheduler	inverse sqrt
Warmup steps	10000
Optimizer	Adam
Adam-beta1	0.9
Adam-beta2	0.98
Dropout	0.1
LSTM hyperparameter	Value
decoder layers	3
hidden size	1024
embedding dimension	200
Transformer hyperparameter	Value
Transformer layers	3
Intermediate hidden size	2048
Attention heads	8
Attention dropout	0.1

Table 2: Language model training hyperparameters.

## A Appendix 852

### A.1 Training data Details 853

854 Table 1 shows details of training dataset. All the  
855 digits and punctuation are removed and all the char-  
856 acters are lower-cased, with special tokens inserted  
857 as word boundaries. Language model training hy-  
858 perparameters are listed in Table 2. Each model  
859 was trained on four A40 GPUs.

### A.2 Lexical diversity across different sets 860

861 We investigated linguistic diversity of different test  
862 sets. Figure 8 shows the type-token ratios of dif-  
863 ferent sets. The CHLD-directed speech is less  
864 lexically diverse compared with other sets, which  
865 corresponds to previous language acquisition study  
866 that caregivers tend to repeat same words to scaf-  
867 fold lexical learning.

868 The overall inspection of the generated data pat-  
869 terns correspond to prior observation of declin-  
870 ing lexical diversity of generated data (Guo et al.,  
871 2023). And the decreased lexical diversity might

train data	3.4M	7.0M	14.7M	27.7M
<b>child production</b>	aha	bye bye mommy	he eaten toes	it is there mom
<b>LSTM(un</b>	wounded	for woman says	smile some that	professor at once busy
<b>Trans(un)</b>	to	o he had been	him to bring	aaia with all his
<b>Prompt</b>	it tries to	one side only	smile when you	but while giving
<b>LSTM(prompted)</b>	make	of his wife	king of him	thing they must never
<b>Trans(prompted)</b>	me	horse the wilderness	will have me	the greater i

Table 3: Examples of generated sequences. The boundary marker is replaced with blank space for the ease of reading. We show the generations with the temperature = 1.0 as examples

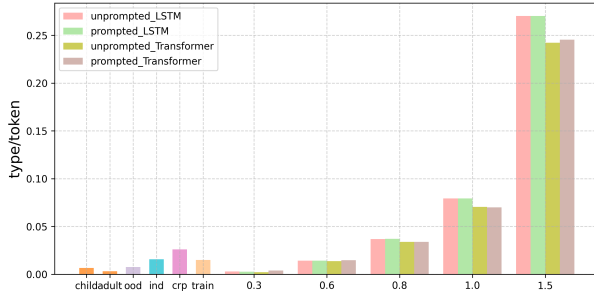


Figure 8: Type/token ratios in different datasets

stem from the large proportion of missing words, which might not necessarily be compensated by the amount of OOV words.

### A.3 Fitted vocabulary growth curves

The figures below show the fitted sigmoid curves across different consitions

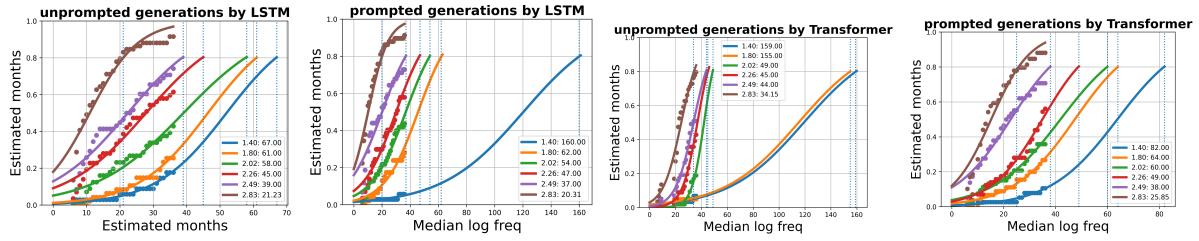


Figure 9: Fitted sigmoid curves of models in temperature of 0.3.

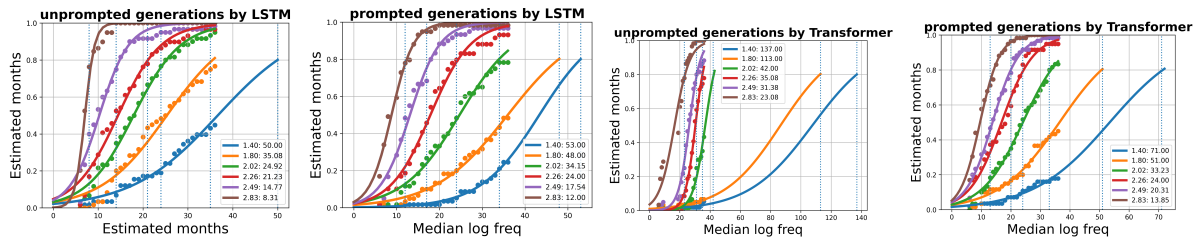


Figure 10: Fitted sigmoid curves of models in temperature of 0.6.

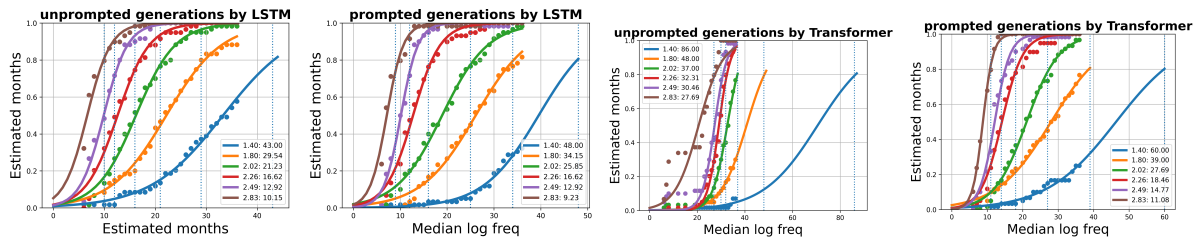


Figure 11: Fitted sigmoid curves of models in temperature of 0.8.

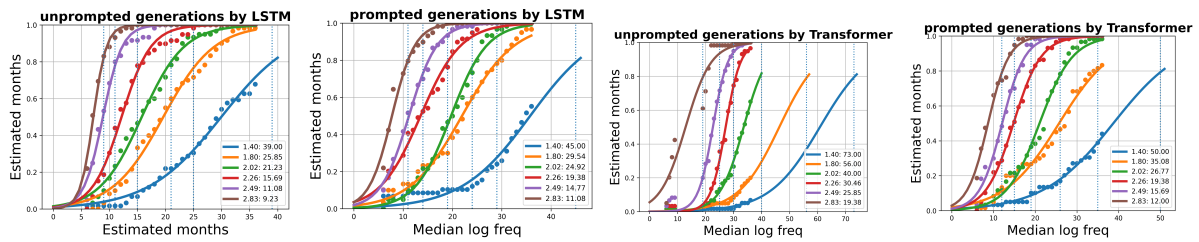


Figure 12: Fitted sigmoid curves of models in temperature of 1.0.

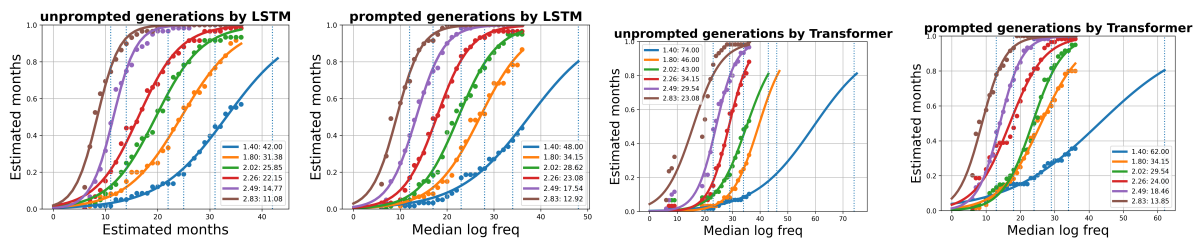


Figure 13: Fitted sigmoid curves of models in temperature of 1.5.