

# LARGE LANGUAGE MODELS SUFFER FROM THEIR OWN OUTPUT: AN ANALYSIS OF THE SELF-CONSUMING TRAINING LOOP

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Language Models (LLM) are already widely used to generate content for a variety of online platforms. As we are not able to safely distinguish LLM-generated content from human-produced content, LLM-generated content is used to train the next generation of LLMs, giving rise to a self-consuming training loop. From the image generation domain we know that such a self-consuming training loop reduces both quality and diversity of images finally ending in a model collapse. However, it is unclear whether this alarming effect can also be observed for LLMs. Therefore, we present the first study investigating the self-consuming training loop for LLMs. Further, we propose a novel method based on logic expressions that allows us to unambiguously verify the correctness of LLM-generated content, which is difficult for natural language text. We find that the self-consuming training loop produces correct outputs, however, the output declines in its diversity depending on the proportion of the used generated data. Fresh data can slow down this decline, but not stop it. Further, we observe similar results on a real natural language dataset. Given these concerning results, we encourage researchers to study methods to negate this process.

## 1 INTRODUCTION

Transformer-based language models have received much attention in the machine learning community in recent years. Especially large language models (LLM) trained on massive amounts of data from the internet became state of the art in many benchmarks (Brown et al., 2020) and specialized conversational LLM applications like ChatGPT<sup>1</sup> have a massive influence on society already. LLMs can be used for multiple tasks ranging from code generation (Chen et al., 2021; Sobania et al., 2022; Fan et al., 2023) and automated program repair (Sobania et al., 2023) to text summarization (Yang et al., 2023) and teaching assistance (Baidoo-Anu & Ansah, 2023).

Due to their convincing generated outputs, LLMs can be used to generate a large amount of content that is posted online to coding platforms like GitHub and Stackoverflow<sup>2</sup>, social media platforms like Reddit<sup>3</sup> and other platforms on the internet. Even academic writing is already being influenced by LLM outputs (Geng & Trotta, 2024; Liang et al., 2024). Such LLM-generated text is often hard to distinguish from human-generated content (Sadasivan et al., 2023) and in turn might unwillingly be used to train the next generation of LLMs. Even paying for human-generated content might not be an option in the future, as workers at paid services like Amazon’s Mechanical Turk also use LLMs to produce content (Veselovsky et al., 2023). Consequently, a self-consuming training loop emerges in which future models are trained repeatedly on LLM-generated data from previous generations. This process was first observed for generative models in the image domain (Martínez et al., 2023; Alemohammad et al., 2023; Shumailov et al., 2024; Bertrand et al., 2023). These studies found that this self-consuming training loop leads to a decline in quality and diversity of generated images, ultimately resulting in a so called model collapse. This was also observed for repeatably fine-tuning LLMs leading to a decrease in diversity (Guo et al., 2023; Shumailov et al., 2024). However, it

<sup>1</sup><https://openai.com/blog/chatgpt>

<sup>2</sup><https://www.microsoft.com/en-us/Investor/events/FY-2023/Morgan-Stanley-TMT-Conference>

<sup>3</sup><https://www.vice.com/en/article/jg5qy8/reddit-moderators-brace-for-a-chatgpt-spam-apocalypse>

is unclear what happens with LLMs that are trained in such a self-consuming training loop from scratch, as usually done in real-world applications like ChatGPT.

Therefore, we present the first study analyzing the behavior of LLMs trained over many generations in a self-consuming training loop. We conduct experiments on a GPT-style model in different settings and measure both quality and diversity of samples from the trained model at each generation in the self-consuming training loop. The settings differ in the way a dataset is created for each generation (so called data cycles) as well as the proportion of original *real* and LLM-generated *synthetic* data samples. To better analyze the behavior of the trained models we conduct our experiments on a dataset consisting of logic expressions. In contrast to natural language, this logic expressions can be evaluated unambiguously. This allows us to analytically and accurately measure correctness and diversity of the generated samples. Furthermore, to confirm the generalizability of our results we conduct additional experiments on a natural language dataset.

We find that repeatedly training new models with synthetic data from previous models initially improves quality. However, diversity degenerates and the learned distribution inevitably collapses to a single point. In extreme cases this happens already after less than 10 generations. Additionally, we find that the speed at which diversity degenerates depends on the data cycle as well as on the proportion of real and synthetic data. Fresh real data added during the data cycle slows down but can not negate the effects of a self-consuming training loop.

In summary, our main contributions are as follows:

- The first comprehensive empirical study of the self-consuming training loop for LLMs trained from scratch,
- A novel method to unambiguously evaluate quality/correctness and diversity of LLM-generated outputs,
- An in-depth analysis and discussion of the effects and implications of this self-consuming training loop for LLMs.

Following this introduction, Sect. 2 describes the self-consuming training loop and gives an overview of related work. In Sect. 3, we describe our experimental setting. Section 4 presents the experimental results, followed by a discussion in Sect. 5 and limitations in Sect. 6. Section 7 concludes the paper.

## 2 THE SELF-CONSUMING TRAINING LOOP

Current generations of LLMs are usually trained on large amounts of unstructured data like text and code gathered from the Internet (Brown et al., 2020). Due to their generative nature and capacity those models can in turn be used to generate new *synthetic* data, often indistinguishable from the original *real* data (Sadashivan et al., 2023). This synthetic data ends up back on the internet and thus in the next large dataset, which is used to train the next generation of LLMs. These LLMs in turn produce new content and data, setting in motion a repetitive cycle in which new generations of models are trained each time with a higher proportion of synthetic data from previous generations. We call this process a self-consuming training loop, depicted in Figure 1.

More specifically, consider a dataset  $\mathcal{D}_0$  consisting of real data points  $x \in X$  sampled from the original distribution  $P_X$ . A generative model  $\mathcal{M}_t$  is trained on this original dataset  $\mathcal{D}_0$  until the data is sufficiently fitted, producing the

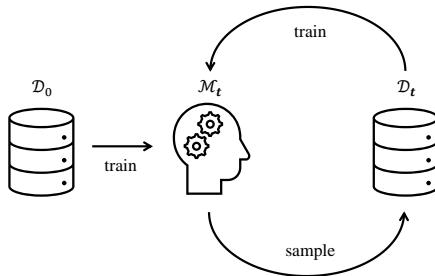


Figure 1: Self-consuming training loop: In the first generation  $t = 1$  a model  $\mathcal{M}_t$  is initially trained on a real dataset  $\mathcal{D}_0$ . From this model  $\mathcal{M}_t$  a sample  $\mathcal{S}_t$  is drawn to build a new dataset  $\mathcal{D}_t$ . The new dataset  $\mathcal{D}_t$  is in turn used to train a new model  $\mathcal{M}_{t+1}$  from scratch for the next generation  $t + 1$ . This process is repeated iteratively until the maximum number of generations  $T$  is reached.

108 first generation  $t = 1$  of generative models. In a self-consuming training loop, a new set of  $m$  syn-  
 109 thetic data points  $\mathcal{S}_t$  is now sampled from the previous generation of generative models  $\mathcal{M}_t$ . The  
 110 next generation of generative models  $\mathcal{M}_{t+1}$  is then trained from scratch on the new dataset  $\mathcal{D}_t$ .  
 111 This self-consuming training loop is repeated until the maximum number of generations  $t = T$  is  
 112 reached. Algorithm 1 in Appendix A presents a pseudo-code of this self-consuming training loop.

113 This self-consuming training loop differs from self-training to improve LLMs (Huang et al., 2022;  
 114 Wang et al., 2022; Gulcehre et al., 2023; Singh et al., 2023; Zhang et al., 2024) as the synthetic  
 115 data is not used intentionally, but rather as a consequence of LLM-generated data published on the  
 116 internet without being labeled as such.

117 **Related Work** Current research has analyzed the self-consuming training loop mostly in the context  
 118 of image generating models. Martínez et al. (2023) first studied the self-consuming training loop  
 119 for diffusion models and find that the self-consuming training loop leads to a collapse in diversity  
 120 of the generated images. Other work have given theoretical frameworks for this phenomenon and  
 121 investigate more complex data cycles in the context of image generation (Alemohammad et al.,  
 122 2023; Shumailov et al., 2024; Bertrand et al., 2023). Both Alemohammad et al. (2023) and Bertrand  
 123 et al. (2023) find that fresh real data can lead to stability within the self-consuming training loop.

124 Shumailov et al. (2024) observe a degeneration of diversity for Variational Autoencoders and Gaus-  
 125 sian Mixture Models if some of the output is used again as input. Additionally, the authors perform  
 126 experiments for iteratively fine-tuning LLMs in a self-consuming way and observe a degradation in  
 127 quality. They find that degradation is less strong than in the image generation context and that  
 128 keeping some original data helps mitigating this phenomenon. Other contemporary work also an-  
 129 alyzed repeated fine-tuning of LLMs in a self-consuming way and observed a decrease in lexical  
 130 diversity (Guo et al., 2023).

131 However, new generations of LLMs are usually trained from scratch with web scraped datasets while  
 132 fine-tuning is mainly done with curated datasets. So the analysis of the self-consuming training loop  
 133 when training from scratch is of pressing concern. To the best of our knowledge, we are the first  
 134 to study the behavior of LLMs trained in a self-consuming loop from scratch. Furthermore, self-  
 135 consuming loops in LLMs have not yet been analyzed with a method that allows to unambiguously  
 136 evaluate the quality and diversity of generated model output.

### 138 3 EXPERIMENTAL SETUP

139 In this section, we present our logic expression dataset, which is the foundation for the verification  
 140 of the language model’s output. Additionally, we specify the natural language experiments. Fur-  
 141 thermore, we explain the data cycles we analyze in our experiments and describe the used model  
 142 architecture in detail.

#### 145 3.1 VERIFICATION WITH A LOGIC EXPRESSION DATASET

146 LLMs usually generate natural language texts and their performance is typically measured by using  
 147 similarity metrics like the BLEU score (Papineni et al., 2002), ROUGE score (Lin, 2004) and BERT  
 148 score (Zhang et al., 2019) as well as perplexity (Jurafsky & Martin, 2009). However, those metrics  
 149 rely on measuring similarity of outputs with expected reference data, which can only serve as a proxy  
 150 for quality of a language model (Callison-Burch et al., 2006; Gehrmann et al., 2023). Consequently,  
 151 we propose using a dataset consisting of logic expressions. Those expressions can be represented as  
 152 a sequence making them a good fit for language modeling. In contrast to natural language text, we  
 153 can easily and systematically evaluate the quality of logic expressions by verifying their correctness.  
 154 An expression is defined to be syntactically correct, if it can be parsed without an error. The semantic  
 155 correctness can be evaluated if an expression is either `True` or `False`. If the original dataset only  
 156 consists of `True` expressions, then a high-quality trained model should also only generate `True`  
 157 expressions.

158 We build a logic expression as a tree in a recursive way specified in the function  
 159 `GenerateLogicExpression(d)`, where  $d$  is the desired depth of a logic expression tree. If the  
 160 desired tree depth is reached when calling the function, a random Boolean is returned (either `True`  
 161 or `False`). Otherwise, a logic operator (`not`, `and`, `or`) is selected and the function is called again

162 recursively with  $d - 1$ . To build the entire original dataset we use this recursive function to sample  
 163 new random logic expression trees with a random depth between  $d_{min}$  and  $d_{max}$  until we have a  
 164 dataset with  $m$  unique expressions that evaluate in a `True` Boolean expression. The complexity  
 165 of the dataset can be easily controlled by adjusting  $d_{max}$ . For our experiments we chose a initial  
 166 dataset size of  $m = 10,000$  with expressions of minimum depth  $d_{min} = 1$  and maximum depth  
 167  $d_{max} = 5$ . Algorithm 2 in Appendix B describes the generation of our logic expression dataset in  
 168 pseudo-code.

169 The resulting logic expression trees can then be saved as strings, e.g., `not ( True and`  
 170 `False )`. This allows us to evaluate them in Python using the `eval()` function to test whether  
 171 they are correct or trigger an error and whether they evaluate to a `True` or `False` Boolean. Ad-  
 172 ditionally, we can encode those strings to a sequence of tokens (`True`, `False`, `not`, `and`,  
 173 `or`, `(`, `)`, `<eos>`, where `<eos>` is a stop token and indicates the end of an expression) and  
 174 use these sequences to train a language model. Examples of logical expressions can be found in  
 175 Appendix E.1.

### 177 3.2 NATURAL LANGUAGE EXPERIMENTS

178 While using the logic expression dataset enables us to design a controlled experimental setting in  
 179 which we can unambiguously measure both quality and diversity of a trained model, this dataset  
 180 might not fully cover the characteristics of real textual data. Therefore, we conduct additional  
 181 experiments using the *tiny Shakespeare dataset* (Karpathy, 2015), a collection of 40,000 lines of  
 182 Shakespeare text with over 1,1 million characters. Since there is no way to analytically investigate  
 183 the correctness of the generated text for this dataset, we focus on the diversity in this experiments.  
 184 Examples of the original dataset as well as generated samples within the self-consuming training  
 185 loop can be found in Appendix E.2.

### 187 3.3 MEASURING THE DIVERSITY OF THE MODEL’S OUTPUT

188 A generative model does not only need to generate correct outputs but also a diverse set of outputs.  
 189 To measure the diversity of a sample from a LLM in the logic expression experiments we use the  
 190 Levenshtein diversity (Beijering et al., 2008; Wittenberg et al., 2023) in our experiments. This metric  
 191 calculates the pairwise Levenshtein distance between each expression in the sample normalized by  
 192 the number of tokens of the longer expression of each pairwise comparison and averaged over the  
 193 number of pairwise comparisons. If the average pairwise normalized Levenshtein distance is close  
 194 to zero the diversity within a sample is low. A value closer to one suggests a stronger diversity within  
 195 a sample. The Levenshtein distance itself is a metric for measuring the distance between two strings  
 196 and is defined as the number of edits (deletion, addition, substitution) required to change one string  
 197 to another (Levenshtein et al., 1966). Normally, this edit distance is calculated on a character basis,  
 198 however, we calculate it on the tokens from our encoding of the logic expressions to accommodate  
 199 for the different length each token has in character representation.  
 200

201 To measure the diversity of a sample in the natural language experiments we employ the  $n$ -gram  
 202 diversity score (Meister et al., 2023; Li et al., 2023; Padmakumar & He, 2024) which measures  
 203 the ratio between unique and total  $n$ -grams in a text. A smaller value indicates less diversity. We  
 204 also investigate the compression rate (Shaib et al., 2024) and the vocabulary size of the generated  
 205 samples.

206 We provide further details for all diversity metrics in Appendix C.

### 208 3.4 DATA CYCLES

209 We define a *data cycle* as the way a dataset  $\mathcal{D}_t$  is constructed in generation  $t$  from the original data  
 210  $\mathcal{D}_0$ , potential fresh data  $\mathcal{F}_t$  from the original distribution, and the generated data from the current  
 211 and previous generations  $\mathcal{S}_{1..t}$ . This way of constructing a new dataset may vary and different  
 212 data cycles are possible. In the image generation domain, previous work suggests that the self-  
 213 consuming training loop is influenced by the data cycle being used (Alemohammad et al., 2023).  
 214 Therefore, inspired by their work we conduct our experiments with the following four different data  
 215 cycles (also depicted in Fig. 2):

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

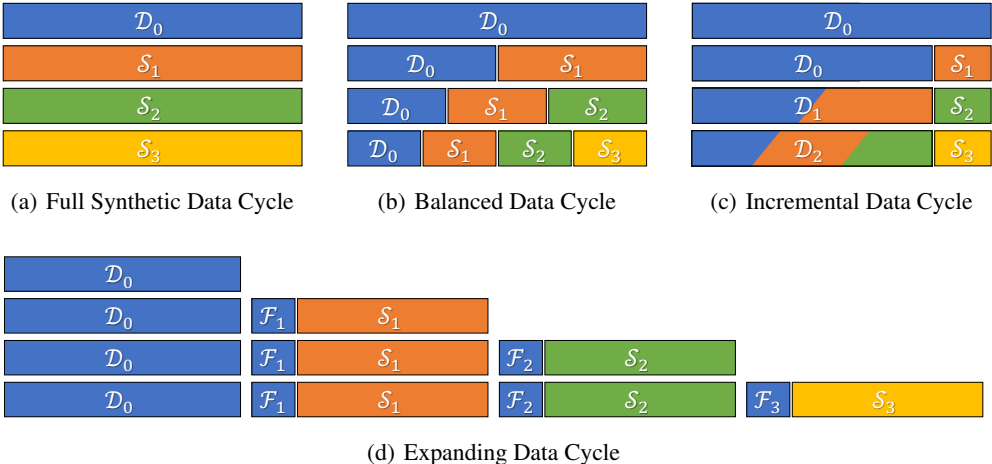


Figure 2: Each sub figure represents a different data cycle exemplary for four generations. The first row in each sub figure is the original dataset  $\mathcal{D}_0$ . The second, third and fourth row depict the dataset  $\mathcal{D}_1$ ,  $\mathcal{D}_2$ , and  $\mathcal{D}_3$  respectively.

**Full Synthetic Data Cycle:** In the most extreme case of a self-consuming training loop a new model  $\mathcal{M}_t$  is only trained on the generated data from the last generation so that  $\mathcal{D}_{t-1} = \mathcal{S}_{t-1}$ . We call this a *full synthetic* data cycle. Normally, new datasets would still contain the original data when they are collected in practice. However, we can use this data cycle to study the most extreme changes in behavior of models from generation to generation.

**Balanced Data Cycle:** We refer to the second data cycle as the *balanced* data cycle. In this data cycle, we construct the new dataset  $\mathcal{D}_t$  from equal parts of all the previous samples  $\mathcal{S}_{1..t}$  and the original dataset  $\mathcal{D}_0$  so that every previous generation contributes  $m * \frac{1}{t+1}$  logic expressions to the new dataset  $\mathcal{D}_t$  of size  $m$ .

**Incremental Data Cycle:** The third data cycle we study is called the *incremental* data cycle. In this data cycle, the new dataset is created by a portion  $(1 - \lambda)$  of the last dataset  $\mathcal{D}_{t-1}$  and a portion  $\lambda$  of new sampled data  $\mathcal{S}_t$  so that  $\mathcal{D}_t = (1 - \lambda) * \mathcal{D}_{t-1} + \lambda * \mathcal{S}_t$  while holding the size  $m$  of the dataset constant each generation. A  $\lambda = 1$  would result in a full synthetic data cycle. We chose  $\lambda = 0.1$  for our initial experiments and later used different values of  $\lambda$  to study the influence of this parameter on diversity.

**Expanding Data Cycle:** The first three data cycles create a dataset of equal size in each generation. In practice, however, the dataset would grow each generation by adding new generated data to the already existing dataset. Additionally, fresh data samples  $\mathcal{F}_t$  from the original distribution (e.g. human generated) would also be added at each generation. Consequently, the last data cycle we study is an *expanding* data cycle. At each generation the new dataset is created so that  $\mathcal{D}_t = \mathcal{D}_{t-1} + (1 - \lambda) * \mathcal{F}_t + \lambda * \mathcal{S}_t$ , where  $\lambda$  is the portion of generated data we add each generation and  $(1 - \lambda)$  is the portion of fresh real data added at each generation. For this data cycle, we chose  $\lambda = 0.9$  for our initial experiments and later used different values of  $\lambda$  to study the influence of fresh data in an expanding data cycle on diversity.

### 3.5 MODEL ARCHITECTURE AND TRAINING

We employ a GPT-style LLM using the open-source implementation *nanoGPT*<sup>4</sup> (MIT license) in our experiments. The model accepts a context of up to 256 tokens and consists of 6 attention layers with 6 attention heads each and an embedding dimensionality of 384, resulting in roughly 10.6 million parameters.

<sup>4</sup><https://github.com/karpathy/nanoGPT>

270 During training we use a batch size of 64 and a dropout rate of 0.2, training for 5000 iterations  
 271 minimizing cross entropy loss, starting with a learning rate of  $10^{-3}$  decaying to  $10^{-4}$ , to achieve a  
 272 sufficient fitting of the training data. We split the dataset in 90% training and 10% validation data.  
 273 During training we calculate the validation error every 250 iterations and use the model parameters  
 274 with the lowest validation error as our final model  $M_t$  for each generation  $t$  in the self-consuming  
 275 training loop.

276 We use the trained model at each generation to sample 10,000 logic expressions. We also performed  
 277 experiments for larger amounts of logical expressions (20,000, 30,000, & 40,000) and observed  
 278 no impact on the results for varying dataset sizes (see Appendix D.4). During sampling, we auto-  
 279 regressively generate new tokens with temperature 0.8 and feed them back into the model until a stop  
 280 token  $\langle \text{eos} \rangle$  is sampled up to a maximum of 200 tokens per expression. We run the self-consuming  
 281 training loop for  $T = 50$  generations in each experiment.

282 We use the same model and setup to run the self-consuming training loop for our natural language  
 283 experiments but tokenize the text on a character level resulting in a vocabulary size of 65. During  
 284 sampling, we auto-regressively generate text of 1000 tokens in length until our sample  $S_t$  is of  
 285 equal size as the original dataset  $\mathcal{D}_0$  (approximately 1.1 million tokens per generation). For the  
 286 expanding data cycle in the natural language experiments we only use a subset of  $\mathcal{D}_0$  in order to  
 287 have enough fresh original data for 50 generations. Consequently, we generate less tokens each  
 288 generation (approximately 186,000 tokens per generation).

289 All experiments are performed on a workstation using consumer NVIDIA graphics cards.

## 291 4 RESULTS

294 In this section, we present our experimental results in terms of correctness and diversity of model  
 295 outputs trained in a self-consuming training loop.

### 297 4.1 CORRECTNESS OF GENERATED CONTENT

299 We first study the correctness of the expressions within a generated sample  $S_t$  from a model  $\mathcal{M}_t$  at  
 300 each generation  $t$ . As described in Sect. 3.1, we consider an expression to be syntactically correct if  
 301 it can be parsed without error. Since the original dataset  $\mathcal{D}_0$  only consists of `True` expressions, we  
 302 consider a generated expression to be semantically correct if it also evaluates to `True`. A semanti-  
 303 cally correct expression is also syntactically correct.

304 Figure 3 displays the composition of samples  $S_t$  generated during a self-consuming training loop.  
 305 Each subplot presents the results for a different data cycle: a) full synthetic, b) balanced, c) incre-  
 306 mental, and d) expanding. Every bar in a subplot displays the composition of  $S_t$  generated from  
 307 model  $\mathcal{M}_t$  at generation  $t$ , except for the first bar in each subplot which displays  $\mathcal{D}_0$ . The green  
 308 portion of a bar indicates the number of syntactically correct expressions that evaluate to `True` in  
 309 that sample. The yellow part of a bar represents the number of syntactically correct expressions  
 310 that evaluate to `False`, and the red part of a bar displays the number of syntactically incorrect  
 311 expressions that result in an error when being parsed.

312 Overall, we see that only very few expressions are syntactically incorrect. This indicates that the  
 313 models are sufficiently trained and can correctly learn the syntactic rules of a logic expression. We  
 314 see a drop of around 20% in the number of semantically correct expressions from the original data  
 315  $\mathcal{D}_0$  to the first sample  $S_1$  in every data cycle. Interestingly, the number of semantically correct  
 316 expressions increases afterwards over the course of the self-consuming training loop. The speed  
 317 at which this number increases depends on the data cycle. We can see the fastest increase for the  
 318 full synthetic data cycle in which the whole sample consists of `True` expressions by generation  
 319  $t = 6$ . The incremental data cycle takes consistently longer but also nearly reaches this point by  
 320 generation  $t = 36$ . While the balanced data cycle has an initially steeper increase in semantically  
 321 correct expressions, it does not reach the point of a completely semantically correct sample by the  
 322 end of 50 generations but comes very close to it. Lastly, the expanding data cycle shows this trend  
 323 as well. Overall, we see that the self-consuming training loop seems to help with generating more  
 semantically correct expressions in our experiments and the rate at which this happens differs by  
 data cycle.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

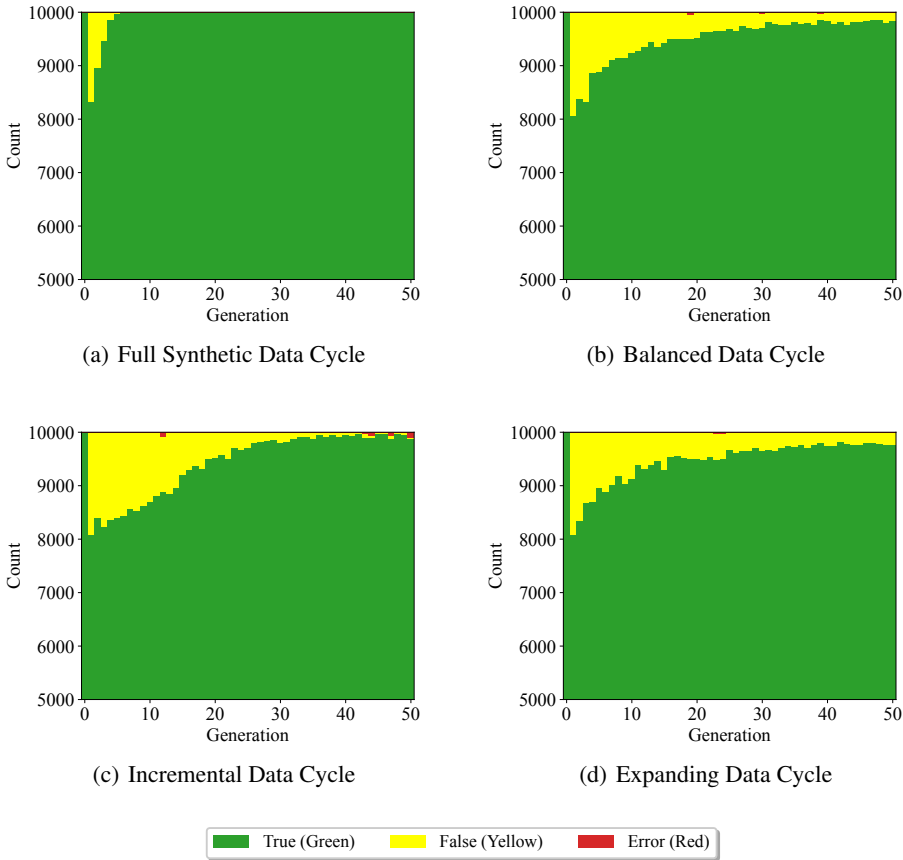


Figure 3: The composition of each sample  $\mathcal{S}_t$  from model  $\mathcal{M}_t$  at generation  $t$  (the first bar displays the composition of  $\mathcal{D}_0$ ) with regards to the number of syntactically and semantically correct expressions. Green indicates syntactically correct expressions that evaluate to `True`, yellow indicates syntactically correct expressions that evaluate to `False`, and red indicates syntactically incorrect expressions that result in an error when being parsed. Each subplot displays the results for a different data cycle.

#### 4.2 DIVERSITY OF GENERATED CONTENT

While we can see an increase in correctness over generations in our experiments, contemporary work in the image generation domain suggests, that this comes with a loss of diversity (Martínez et al., 2023; Alemohammad et al., 2023; Bertrand et al., 2023). Therefore, in this section we study the diversity within a sample  $\mathcal{S}_t$  of model  $\mathcal{M}_t$  for each generation  $t$ .

Figure 4 displays the average pairwise normalized Levenshtein distance over generations of a self-consuming training loop for different data cycles. For each data cycle, we observe a decrease in diversity over the course of generations, with the degree of decrease varying depending on the data cycle. For the full synthetic data cycle we see a steep decline in diversity with a collapse into a single point by generation  $t = 39$ . The incremental data cycle is initially stable, but decreases in diversity from the 10th generation onwards, with a decrease of 68% in diversity by the end of the 50 generations. The balanced and expanding data cycle also decrease in diversity, however, this decrease is way slower, with a 30% decrease in diversity for the balanced data cycle and a 22% decrease for the expanding data cycle. While not all data cycles fully collapse in diversity by generation 50, ultimately, we expect all of them to eventually reach zero diversity if the self-consuming training loop is run for enough generations.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

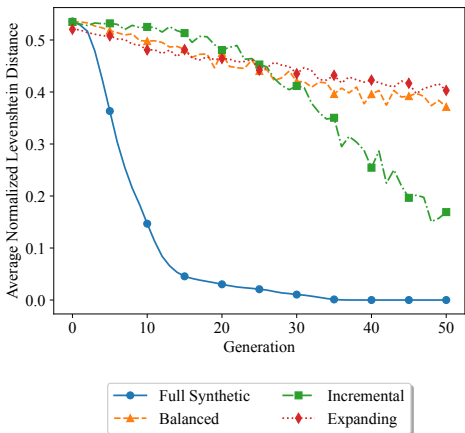


Figure 4: Average pairwise normalized Levenshtein distance over generations of a self-consuming training loop for different data cycles. Line markers added every 5 generations for better display.

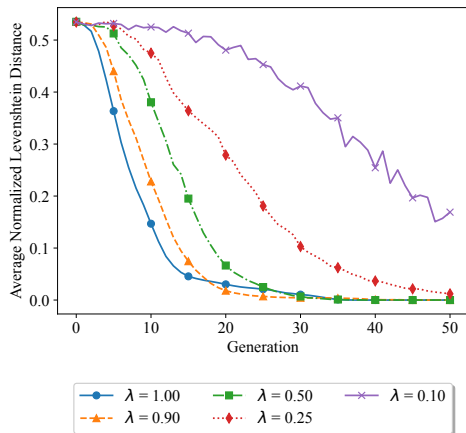


Figure 5: Average pairwise normalized Levenshtein distance over the course of a self-consuming training loop for the expanding data cycle for different portions  $\lambda$  of generated data and  $(1 - \lambda)$  of fresh data. Line markers added every 5 generations for better display.

To get further insight into the decline in diversity, we also inspect the number of unique expressions sampled per generation for the full synthetic data cycle. We find that the number of unique expressions is stable at first while the diversity already declines and then rapidly decreases within 5 generations to very few unique expressions (see Appendix D.1).

To better understand the influence of the amount of generated data in each generation on the diversity, we also study the incremental data cycle in more detail with different portions  $\lambda$  of new sampled data. Figure 5 displays the average pairwise normalized Levenshtein distance over generations in a self-consuming training loop for the incremental data cycle with different values of  $\lambda$ . We observe, that the loss in diversity is stronger for a higher share of generated data added in each generation of the self-consuming training loop. Even with as little as  $\lambda = 0.25$ , we reach the complete loss of diversity within 50 generations. Only for our experimental run with  $\lambda = 0.10$  the diversity does not drop to zero, but to a value around 0.17 after 50 generations. However, we see that the speed at which diversity decreases does not scale linear with the amount of generated data. Instead we observe a more exponential increase in speed where larger amounts of generated data lead to an even quicker decay in diversity.

Lastly, we study the expanding data cycle in more detail with regard to the proportion between generated data and fresh data added to the training data at each generation. Figure 6 displays the average pairwise normalized Levenshtein distance over generations in a self-consuming training loop for the expanding data cycle with different portions  $\lambda$  of generated data and  $(1 - \lambda)$  fresh real data. We observe, that while each configuration declines in diversity, the rate of this decline depends on the proportion between generated and fresh data. The more fresh data is added at each generation, the slower diversity decreases. With no fresh data or only 1% of fresh data the diversity decreases by approximately 36%. If the portion of fresh data is 25% we only observe a decrease in diversity of approximately 11%. This indicates that fresh data cannot stop the self-consuming training loop but enough fresh data can slow down the rate of the decline in diversity.

### 4.3 DIVERSITY OF GENERATED NATURAL LANGUAGE

To extend our results past the domain of logic expression we also study the self-consuming training loop for a real textual dataset. We first investigate the diversity of generated content. Figure 7 displays the  $n$ -gram diversity over generations in a self-consuming training loop for different data



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

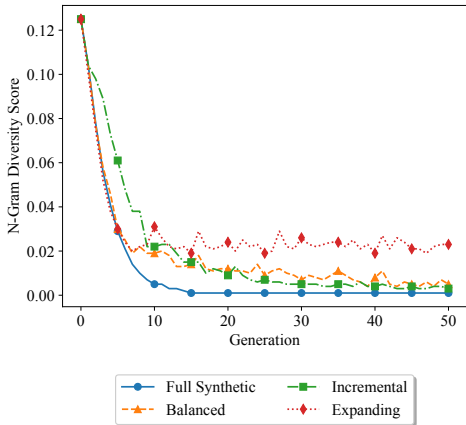
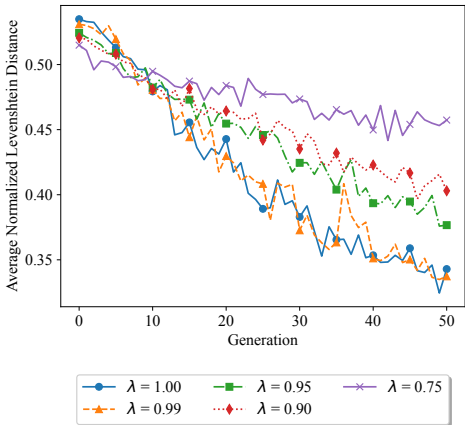


Figure 6: Average pairwise normalized Levenshtein distance over the course of a self-consuming training loop for the expanding data cycle for different portions  $\lambda$  of generated data and  $(1 - \lambda)$  of fresh data. Line markers added every 5 generations for better display.

Figure 7: N-gram diversity score over the course of a self-consuming training loop for different data cycles for textual data. The sample size is truncated to equal length for better comparability. Line markers added every 5 generations for better display.

cycles. As diversity scores for text are sensitive to text length, we truncated the samples to be of equal size as suggested in (Shaib et al., 2024), to better compare the data cycles with each other. Similar to our results in the logic expression domain, we observe a decline in diversity for all data cycles. As expected, this is most severe for the full synthetic data cycle dropping to zero  $n$ -gram diversity by generation 15, followed by the incremental and balanced data cycle. The expanding data cycle retains the most diversity but also declines strongly over the course of the self-consuming training loop. We also investigate the compression rate and vocabulary size and find similar effects (see Appendix D.5). Overall, this confirms our findings from the logic expression dataset in terms of diversity also for textual data.

While we cannot evaluate the correctness of the generated automatically, we manually inspect text samples for a qualitative assessment. We find that while the models produce some hallucinated words in first generations the quality of the generated text is subjectively mostly correct and Shakespeare like. However, the stark decline in diversity quickly leads to repetitive outputs and makes the models unusable, further confirming the effects of the self-consuming training loop. More details and examples of generated text can be found in Appendix E.2.

## 5 DISCUSSION

In our experiments we found that iteratively re-training LLMs from scratch with self generated data might help with correctness of model output. This is sometimes already done in practice where the output of LLMs is used to improve their performance (Huang et al., 2022; Wang et al., 2022; Gulcehre et al., 2023; Singh et al., 2023; Zhang et al., 2024) or generate new training data (Yu et al., 2023). In this self-training setting the origin of data is known and researchers can carefully filter and adjust the data. However, the self-consuming training loop is an emerging trend where the data source is not known and LLM-generated content from the internet is unwillingly used to train new models. This trend is concerning as we also found that this self-consuming training loop eventually leads to a drastic loss of diversity of a model’s output. Our results confirm the findings on the self-consuming training loop in the field of image generation (Martínez et al., 2023; Alemohammad et al., 2023; Shumailov et al., 2024; Bertrand et al., 2023), further highlighting the importance of this issue for generative models and LLMs in particular. These LLMs have an even greater impact on society than image generation models, leading to more LLM-generated data being produced and

486 potentially having a greater impact on society as a whole, as many people already rely on good  
487 model performance.

488  
489 Consequently, as we expect that a self-consuming training loop will harm model performance in the  
490 future, researchers and practitioners should carefully choose their data when training their models  
491 and test those models sufficiently for diversity. Additionally, the machine learning community at  
492 large needs to find ways to deal with this problem as generated text data is already present in new  
493 internet scraped datasets. While our results indicate that fresh real data can slow down the self-  
494 consuming training loop, other studies suggest that this is not an option as we will run out of real  
495 data eventually and have to rely on generated data for future models (Villalobos et al., 2022). Addi-  
496 tionally, LLM-generated data has already found its way onto the internet, even infecting academic  
497 publications (Geng & Trotta, 2024; Liang et al., 2024), and LLM-generated data may soon overtake  
real data on the internet.

498 Furthermore, this LLM-generated data is hard to differentiate from real data (Kreps et al., 2022;  
499 Sadasivan et al., 2023; Huschens et al., 2023), making it difficult to curate datasets scraped from the  
500 web. This makes the self-consuming training loop inevitable in the future. Therefore, we advise to  
501 further study ways to maintain diversity of generative model outputs like quality diversity methods  
502 (Pugh et al., 2015; Fontaine & Nikolaidis, 2021).

## 503 504 6 LIMITATIONS

505  
506 One cause for limitations of our work are restricted computational resources. First, we ran our  
507 experiments with a smaller GPT-style model than for example GPT-4. Even though we expect that  
508 larger models with billions of parameters still suffer from the same loss of diversity if trained in a  
509 self-consuming training loop, we note that the behavior may vary for those larger models. However,  
510 a study for models of those size is not feasible with current compute options nor is it responsible in  
511 terms of the expected CO<sub>2</sub> emissions (Strubell et al., 2019). Nevertheless, we performed additional  
512 experiments for the full synthetic data cycle with varying model sizes in Appendix D.3 and still  
513 observe the loss in diversity across smaller and larger models. Therefore, we expect that our results  
514 can be generalized to models of different size.

515  
516 Second, the main results of this work consist of only one run per experimental configuration. With  
517 more runs per experiment the results would be more robust, but conducting more runs for our exper-  
518 iments was simply not feasible with our available computing resources. As all our experiments point  
519 towards the same direction, however, we do not believe that this limitation is crucial. Nevertheless,  
520 to further mitigate this limitation we performed additional experiments for the full synthetic data  
521 cycle with fewer generations over multiple runs and with different initializations in Appendix D.2.  
These results show that the trend of a decline in diversity still holds true over multiple runs.

522  
523 Lastly, datasets used to train LLMs cover a wide range of different data sources with different  
524 characteristics. Therefore, our results might be slightly different when performed with a real web  
525 scraped corpora. However, the experiments with textual data show that our results generalize past  
526 logic expressions.

## 527 528 7 CONCLUSION

529  
530 In this paper we studied the behavior of LLMs trained in a self-consuming training loop and found  
531 that iteratively training LLMs from scratch with self generated data can initially help with correct-  
532 ness of model outputs. However, this comes at a price, as the diversity of generated data eventually  
533 degenerates and collapses into a single point. The rate of this decline in diversity depends on the  
534 data cycle creating the training data at each generation. Furthermore, we found that fresh data can  
535 slow down this process and preserve diversity of a model output for longer. This should encourage  
536 researches to carefully monitor their datasets and models to mitigate such a harmful self-consuming  
537 training loop in the future.

538 In future work, we will study how fine-tuning can influence the self-consuming training loop for  
539 LLMs. Additionally, we plan to further investigate how the effects of a self-consuming training loop  
can be slowed down or completely negated.

## REFERENCES

- 540  
541  
542 Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein  
543 Babaei, Daniel LeJeune, Ali Siahkoohi, and Richard G Baraniuk. Self-consuming generative  
544 models go mad. *arXiv preprint arXiv:2307.01850*, 2023.
- 545 David Baidoo-Anu and Leticia Owusu Ansah. Education in the era of generative artificial intelli-  
546 gence (ai): Understanding the potential benefits of chatgpt in promoting teaching and learning.  
547 *Journal of AI*, 7(1):52–62, 2023.
- 548 Karin Beijering, Charlotte Gooskens, and Wilbert Heeringa. Predicting intelligibility and perceived  
549 linguistic distance by means of the levenshtein algorithm. *Linguistics in the Netherlands*, 25(1):  
550 13–24, 2008.
- 551  
552 Quentin Bertrand, Avishek Joey Bose, Alexandre Duplessis, Marco Jiralerspong, and Gauthier  
553 Gidel. On the stability of iterative retraining of generative models on their own data. *arXiv*  
554 *preprint arXiv:2310.00429*, 2023.
- 555 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
556 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
557 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 558  
559 Chris Callison-Burch, Miles Osborne, and Philipp Koehn. Re-evaluating the role of bleu in machine  
560 translation research. In *11th conference of the european chapter of the association for computa-*  
561 *tional linguistics*, pp. 249–256, 2006.
- 562 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
563 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
564 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 565  
566 Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and  
567 Jie M Zhang. Large language models for software engineering: Survey and open problems. *arXiv*  
568 *preprint arXiv:2310.03533*, 2023.
- 569 Matthew Fontaine and Stefanos Nikolaidis. Differentiable quality diversity. *Advances in Neural*  
570 *Information Processing Systems*, 34:10040–10052, 2021.
- 571  
572 Sebastian Gehrmann, Elizabeth Clark, and Thibault Sellam. Repairing the cracked foundation: A  
573 survey of obstacles in evaluation practices for generated text. *Journal of Artificial Intelligence*  
574 *Research*, 77:103–166, 2023.
- 575  
576 Mingmeng Geng and Roberto Trotta. Is chatgpt transforming academics’ writing style? *arXiv*  
577 *preprint arXiv:2404.08627*, 2024.
- 578 Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek  
579 Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. Reinforced self-training  
580 (rest) for language modeling. *arXiv preprint arXiv:2308.08998*, 2023.
- 581  
582 Yanzhu Guo, Guokan Shang, Michalis Vazirgiannis, and Chloé Clavel. The curious decline of lin-  
583 guistic diversity: Training language models on synthetic text. *arXiv preprint arXiv:2311.09807*,  
584 2023.
- 585 Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei  
586 Han. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.
- 587  
588 Martin Huschens, Martin Briesch, Dominik Sobania, and Franz Rothlauf. Do you trust chatgpt?–  
589 perceived credibility of human and ai-generated content. *arXiv preprint arXiv:2309.02524*, 2023.
- 590 Dan Jurafsky and James H Martin. *Speech and language processing: An introduction to natural*  
591 *language processing, computational linguistics, and speech recognition*. Prentice Hall, Pearson  
592 Education, Upper Saddle River, NJ, USA, 2009.
- 593 Andrej Karpathy. char-rnn. <https://github.com/karpathy/char-rnn>, 2015.

- 594 Sarah Kreps, R Miles McCain, and Miles Brundage. All the news that’s fit to fabricate: Ai-generated  
595 text as a tool of media misinformation. *Journal of experimental political science*, 9(1):104–117,  
596 2022.
- 597 Vladimir I Levenshtein et al. Binary codes capable of correcting deletions, insertions, and reversals.  
598 In *Soviet physics doklady*, volume 10, pp. 707–710. Soviet Union, 1966.
- 600 Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori B Hashimoto, Luke  
601 Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization.  
602 In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*  
603 (*Volume 1: Long Papers*), pp. 12286–12312, 2023.
- 604 Weixin Liang, Yaohui Zhang, Zhengxuan Wu, Haley Lepp, Wenlong Ji, Xuandong Zhao, Hancheng  
605 Cao, Sheng Liu, Siyu He, Zhi Huang, et al. Mapping the increasing use of llms in scientific  
606 papers. *arXiv preprint arXiv:2404.01268*, 2024.
- 607 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization*  
608 *branches out*, pp. 74–81, 2004.
- 609 Gonzalo Martínez, Lauren Watson, Pedro Reviriego, José Alberto Hernández, Marc Juárez, and Rik  
610 Sarkar. Towards understanding the interplay of generative artificial intelligence and the internet.  
611 *arXiv preprint arXiv:2306.06130*, 2023.
- 612 Clara Meister, Tiago Pimentel, Gian Wiher, and Ryan Cotterell. Locally typical sampling. *Transac-*  
613 *tions of the Association for Computational Linguistics*, 11:102–121, 2023.
- 614 Vishakh Padmakumar and He He. Does writing with language models reduce content diversity? In  
615 *The Twelfth International Conference on Learning Representations*, 2024.
- 616 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic  
617 evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association*  
618 *for Computational Linguistics*, pp. 311–318, 2002.
- 619 Justin K Pugh, Lisa B Soros, Paul A Szerlip, and Kenneth O Stanley. Confronting the challenge  
620 of quality diversity. In *Proceedings of the 2015 Annual Conference on Genetic and Evolutionary*  
621 *Computation*, pp. 967–974, 2015.
- 622 Vinu Sankar Sadasivan, Aounon Kumar, Sriram Balasubramanian, Wenxiao Wang, and Soheil Feizi.  
623 Can ai-generated text be reliably detected? *arXiv preprint arXiv:2303.11156*, 2023.
- 624 Chantal Shaib, Joe Barrow, Jiuding Sun, Alexa F Siu, Byron C Wallace, and Ani Nenkova. Stan-  
625 dardizing the measurement of text diversity: A tool and a comparative analysis of scores. *arXiv*  
626 *preprint arXiv:2403.00553*, 2024.
- 627 Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal.  
628 Ai models collapse when trained on recursively generated data. *Nature*, 631(8022):755–759,  
629 2024.
- 630 Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Peter J Liu, James  
631 Harrison, Jaehoon Lee, Kelvin Xu, Aaron Parisi, et al. Beyond human data: Scaling self-training  
632 for problem-solving with language models. *arXiv preprint arXiv:2312.06585*, 2023.
- 633 Dominik Sobania, Martin Briesch, and Franz Rothlauf. Choose your programming copilot: A com-  
634 parison of the program synthesis performance of github copilot and genetic programming. In  
635 *Proceedings of the genetic and evolutionary computation conference*, pp. 1019–1027, 2022.
- 636 Dominik Sobania, Martin Briesch, Carol Hanna, and Justyna Petke. An analysis of the automatic  
637 bug fixing performance of chatgpt. In *2023 IEEE/ACM International Workshop on Automated*  
638 *Program Repair (APR)*, pp. 23–30, Los Alamitos, CA, USA, 2023. IEEE Computer Society.
- 639 Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep  
640 learning in nlp. *arXiv preprint arXiv:1906.02243*, 2019.

- 648 Veniamin Veselovsky, Manoel Horta Ribeiro, and Robert West. Artificial artificial artificial intelli-  
649 gence: Crowd workers widely use large language models for text production tasks. *arXiv preprint*  
650 *arXiv:2306.07899*, 2023.
- 651 Pablo Villalobos, Jaime Sevilla, Lennart Heim, Tamay Besiroglu, Marius Hobbhahn, and Anson Ho.  
652 Will we run out of data? an analysis of the limits of scaling datasets in machine learning. *arXiv*  
653 *preprint arXiv:2211.04325*, 2022.
- 654 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and  
655 Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions.  
656 *arXiv preprint arXiv:2212.10560*, 2022.
- 657 David Wittenberg, Franz Rothlauf, and Christian Gagné. Denoising autoencoder genetic program-  
658 ming: strategies to control exploration and exploitation in search. *Genetic Programming and*  
659 *Evolvable Machines*, 24(2):17, 2023.
- 660 Xianjun Yang, Yan Li, Xinlu Zhang, Haifeng Chen, and Wei Cheng. Exploring the limits of chatgpt  
661 for query or aspect-based text summarization. *arXiv preprint arXiv:2302.08081*, 2023.
- 662 Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander Ratner, Ranjay Krishna, Jiaming Shen,  
663 and Chao Zhang. Large language model as attributed training data generator: A tale of diversity  
664 and bias. *arXiv preprint arXiv:2306.15895*, 2023.
- 665 Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. Rest-mcts\*: Llm  
666 self-training via process reward guided tree search. *arXiv preprint arXiv:2406.03816*, 2024.
- 667 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluat-  
668 ing text generation with bert. *arXiv preprint arXiv:1904.09675*, 2019.

## 673 A PSEUDOCODE: SELF-CONSUMING TRAINING LOOP

674 Algorithm 1 describes the self-consuming training loop in detail.

---

### 675 **Algorithm 1** Self-Consuming Training Loop

---

676 **Input:** Number of generations  $T$ , data cycle

677 **Initialize:** Sample  $\mathcal{D}_0$  from  $P_X$

678 **for**  $t = 1$  **to**  $t = T$  **do**

679     Train model  $\mathcal{M}_t$  with  $\mathcal{D}_{t-1}$

680     Sample  $\mathcal{S}_t$  from  $\mathcal{M}_t$

681     Create  $\mathcal{D}_t$  from  $\mathcal{D}_0$  and  $\mathcal{S}_{1..t}$  according to data cycle

682 **end for**

683 **Output:** Trained models  $\mathcal{M}_{1..T}$ , samples  $\mathcal{S}_{1..T}$ , datasets  $\mathcal{D}_{0..T}$

---

## B PSEUDOCODE: GENERATION OF THE LOGIC EXPRESSION DATASET

Algorithm 2 describes the generation of our logic expression dataset in detail.

---

### Algorithm 2 Generate Logic Expression Dataset

---

**Input:** Number of unique expressions  $m$ , minimum tree depth  $d_{min}$ , maximum tree depth  $d_{max}$

**Initialize:**  $\mathcal{D}_0 = \{\}$

```

function GenerateLogicExpression( $d$ ) :
if  $d = 0$  then
  return RandomBoolean()
else
   $op \leftarrow$  RandomLogicOperator()
  if  $op = \text{'not'}$  then
    return  $op +$  GenerateLogicExpression( $d - 1$ )
  else
     $leftExpr \leftarrow$  GenerateLogicExpression( $d - 1$ )
     $rightExpr \leftarrow$  GenerateLogicExpression( $d - 1$ )
    return  $leftExpr + op + rightExpr$ 
  end if
end if
end function

```

```

while  $|\mathcal{D}_0| < m$  do
   $i \leftarrow$  RandomIntegerInRange( $d_{min}, d_{max}$ )
   $expr \leftarrow$  GenerateLogicExpression( $d = i$ )
  if EvaluateBooleanExpression( $expr$ ) = True then
     $\mathcal{D}_0 \leftarrow \mathcal{D}_0 \cup \{expr\}$ 
  end if
end while

```

**Output:**  $\mathcal{D}_0$

---

## C DIVERSITY METRICS

In this section we provide more detail on the diversity metrics used for the experiments using logic expression as well as textual data.

The **average normalized Levenshtein distance** (Levenshtein diversity) (Beijering et al., 2008; Wittenberg et al., 2023) used in our logic expression experiments is defined as:

$$\text{Average Normalized Levenshtein Distance} = \frac{1}{|Q|} \sum_{(A,B) \in Q} \frac{d_L(A, B)}{\max(\text{len}(A), \text{len}(B))}$$

where  $Q$  is the set of unique pairs of logic expressions within a sample  $\mathcal{S}_t$  and  $d_L(A, B)$  is the Levenshtein distance between two logic expressions  $A$  and  $B$ .

The  **$n$ -gram diversity score** (Meister et al., 2023; Li et al., 2023; Padmakumar & He, 2024) used to investigate the diversity of textual data is defined as the number of unique  $n$ -grams divided by the total number of  $n$ -grams in a sample  $\mathcal{S}_t$ :

$$\text{N-Gram Diversity Score} = \sum_{n=1}^4 \frac{\# \text{ unique } n\text{-grams in } \mathcal{S}_t}{\# n\text{-grams in } \mathcal{S}_t}$$

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

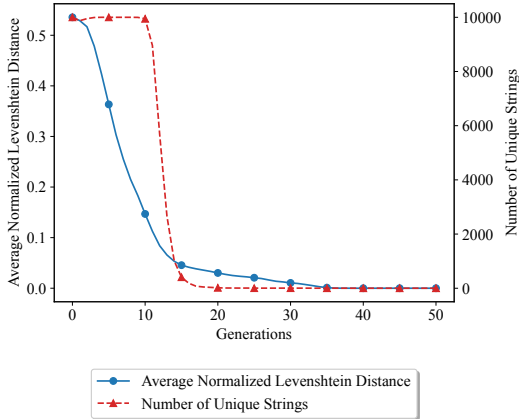


Figure 8: Average pairwise normalized Levenshtein distance (left y-axis) and number of unique expressions in a sample (right y-axis) over generations of a self-consuming training loop for the full synthetic data cycle. Line markers added every 5 generations for better display.

The **compression rate** as a measure for diversity of a text (Shaib et al., 2024) is defined as the size of a sample  $\mathcal{S}_t$  divided by the size of the same sample compressed by a compression algorithm (e.g. gZip):

$$\text{Compression Rate} = \frac{\text{size of } \mathcal{S}_t}{\text{compressed size of } \mathcal{S}_t}$$

## D ADDITIONAL EXPERIMENTAL RESULTS

### D.1 UNIQUE EXPRESSIONS GENERATED

To better understand the decline of diversity, we inspect the number of unique expressions sampled at each generation. We focus on the full synthetic data cycle, as this data cycle fully collapses to a single point within 50 generations.

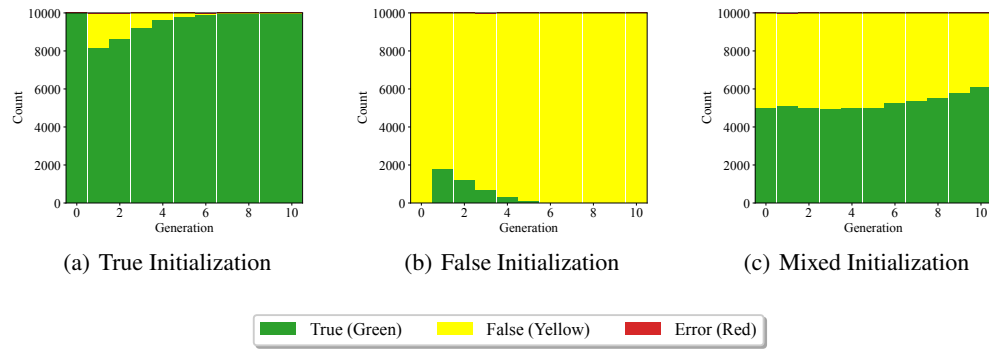
Figure 8 displays the number of unique expressions (right y-axis) generated at each generation of a self-consuming training loop in comparison to the diversity measured as the average pairwise normalized Levenshtein distance (left y-axis). We observe that in the beginning the number of unique expressions stays stable while the diversity is already in a steep decline. Only by generation  $t = 10$  the number of unique individuals starts to decline rapidly and by generation  $t = 15$  only around 400 unique expressions are sampled while the decline in diversity slows down. In generation  $t = 21$  the number of unique individuals is below 10 and in generation  $t = 39$  only a single individual is sampled at each generation. Consequently, the diversity is collapsed to zero.

This shows that it is not enough to only track unique outputs of a model, but that diversity needs to be tracked as well. By the time a decrease in unique results is noticeable, the diversity already decreased significantly.

### D.2 IMPACT OF INITIALIZATION

So far we only investigated the self-consuming training loop starting from an initial dataset consisting only of `True` expressions. To check whether our findings also hold for different initializations of data across multiple runs, we conduct further experiments. Specifically, we analyze the full synthetic data cycle with the original dataset  $\mathcal{D}_0$  only consisting of (1) `True` expressions, (2) `False` expressions, and (3) an equal mixture of both. The results from Sect. 4 indicate that the effect of the self-consuming training loop for the full synthetic data cycle are already clearly visible in the

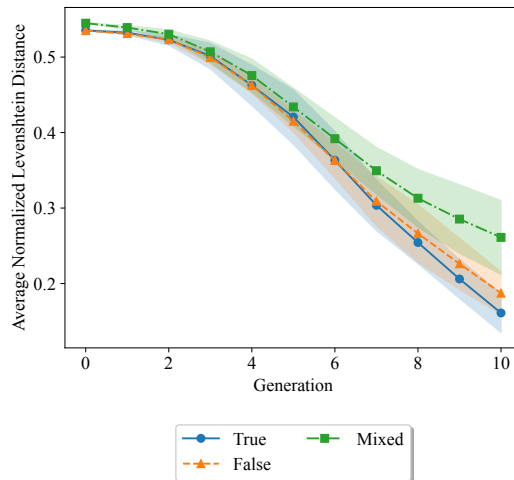
810  
811  
812  
813  
814  
815  
816  
817  
818  
819



820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831

Figure 9: The composition of each sample  $\mathcal{S}_t$  from model  $\mathcal{M}_t$  at generation  $t$  (the first bar displays the composition of  $\mathcal{D}_0$ ) with regards to the number of syntactically and semantically correct expressions averaged over 5 runs. Green indicates syntactically correct expressions that evaluate to True, yellow indicates syntactically correct expressions that evaluate to False, and red indicates syntactically incorrect expressions that result in an error when being parsed. Each subplot displays the results for the full synthetic data cycle with different initializations (True, False, Mixed).

832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847



848  
849  
850  
851

Figure 10: Average pairwise normalized Levenshtein distance over generations of a self-consuming training loop for the full synthetic data cycle with different initializations. Displayed is the mean and standard deviation at each generation for each initialization over 5 runs. Line markers added every 5 generations for better display.

852  
853

854 first 10 generations. Therefore, we conduct experiments with only 10 generations. With this saved  
855 computing time, we conduct 5 independent runs for each initialization.

856  
857  
858  
859  
860  
861  
862  
863

Figure 9 displays the composition of samples  $\mathcal{S}_t$  generated during a self-consuming training loop. Each subplot presents the composition for a different initialization for the full synthetic data cycle averaged over 5 independent runs: a) only True expressions, b) only False expressions, and c) half True expressions and half False expressions. Every bar in a subplot displays the composition of  $\mathcal{S}_t$  generated from model  $\mathcal{M}_t$  at generation  $t$ , except for the first bar in each subplot which displays  $\mathcal{D}_0$ . The green portion of a bar indicates the number of syntactically correct expressions that evaluate to True in that sample. The yellow part of a bar represents the number of syntactically correct expressions that evaluate to False, and the red part of a bar displays the number of syntactically incorrect expressions that result in an error when being parsed.



864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

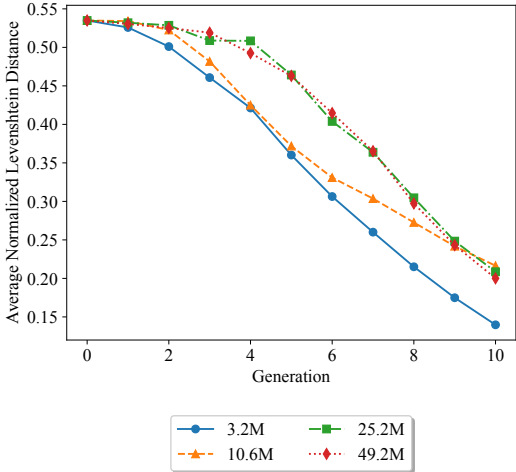


Figure 11: Average pairwise normalized Levenshtein distance over generations of a self-consuming training loop for the full synthetic data cycle with different model sizes. Line markers added every 5 generations for better display.

Similar to our results in Sect. 4.1, we observe an initial drop of `True` expressions by 20% in favor of `False` expressions for the initialization with only `True` expressions in the first generation. Afterwards, the number of `True` expressions in a sample increases again over the course of generations. For the initialization with only `False` expressions, we observe a similar effect, only in the opposite direction. The number of `False` expression first drops by around 20% and then increases again until the whole sample consists of only `False` expressions. When initialized with an equal amount of `True` and `False` expressions, we can observe that the proportion of both expression types first stays stable and then by generation  $t = 6$  starts to slightly shift towards `True` expressions. Upon further inspection of the 5 runs, we find that one run quickly shifts towards only `True` expressions, one run shifts quickly to only `False` expression, one run shifts slowly to `True` expressions, and two runs maintain equal proportions between `True` and `False` expressions across the 10 generations.

Figure 10 displays the diversity measured in the average pairwise normalized Levenshtein distance for the three different initializations. Displayed is the mean and standard deviation of five runs over the number of generations. We observe that the diversity declines for all three types of initialization, further confirming our results from Sect. 4.2. The initializations with only `True` and only `False` expressions decline at around the same speed. The initialization with equal proportions of `True` and `False` expressions declines a bit slower but shows the same effect.

Overall these findings indicate that the effects of the self-consuming training loop apply to different initializations across multiple runs, partially mitigating the limitations mentioned in Sect. 6.

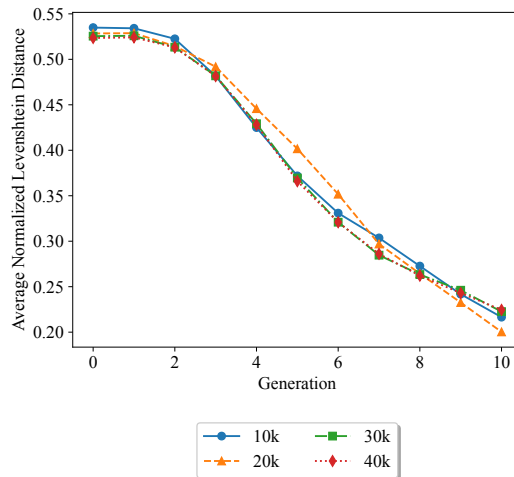
### D.3 IMPACT OF MODEL SIZE

To investigate if our results can be generalized to different model sizes and architectures we repeated our experiments for the full synthetic data cycle for 10 generations with multiple different model sizes. We conducted experiments for models with parameters ranging from 3.2 million to 49.2 million. Table 1 describes the details of those models.

Figure 11 displays the diversity measured in the average pairwise normalized Levenshtein distance for different model sizes. Similar to our results in Sect. 4.2 we observe a steady decline in diversity over generations regardless of model size. While the decrease in diversity happens quicker for the smaller model it is also clearly present in the larger models. Therefore, we argue that the effects of a self-consuming training loop is also present for different model sizes.

918  
919 Table 1: Details of the model architectures for the additional experiments. Displayed is the total  
920 number of parameters in millions, the number of layers, number of attention heads per layer, and the  
921 embedding dimensionality.

922 #PARAMETERS	#LAYERS	#ATTENTION HEADS	EMBEDDING DIMENSION
924 3.2M	4	4	256
925 10.6M	6	6	384
926 25.2M	8	8	512
927 49.2M	10	10	640



947 Figure 12: Average pairwise normalized Levenshtein distance over generations of a self-consuming  
948 training loop for the full synthetic data cycle with different dataset sizes. Line markers added every  
949 5 generations for better display.

#### 951 D.4 IMPACT OF DATASET SIZE

952 While we used 10,000 logic expressions for our main experiments, we conducted additional ex-  
953 periments for larger amounts of logical expressions (20,000, 30,000, and 40,000) to analyze the  
954 impact of varying amounts of training data on the self-consuming training loop. Specifically, we  
955 conducted experiments for the full synthetic data cycle over 10 generations.

956 Figure 12 displays the diversity measured in the average pairwise normalized Levenshtein distance  
957 for different dataset sizes. We observe that for all dataset sizes the diversity decreases over the course  
958 of generations. This indicates that the effect of the self-consuming training loop occurs regardless  
959 of dataset size.

#### 962 D.5 ADDITIONAL DIVERSITY METRICS FOR THE NATURAL LANGUAGE DATASET

963 In addition to our results in section 4.3 we also inspect the compression rate (Shaib et al., 2024) and  
964 the vocabulary size of samples obtained in the self-consuming training loop for textual data. Once  
965 again we truncate the data samples to be of equal size for better comparability between the different  
966 data cycles.

967 Figure 13 displays the compression rate over generations of a self-consuming training loop. A higher  
968 compression rate indicates a less diverse sample. Similar to our results in the logic expression do-  
969 main we observe the highest increase in compression rate and therefore decrease in diversity for the  
970 full synthetic data cycle. The incremental data cycle also has a stark increase in compression rate  
971 followed by the balanced data cycle. The expanding data cycle also shows an increase in compres-

972  
 973  
 974  
 975  
 976  
 977  
 978  
 979  
 980  
 981  
 982  
 983  
 984  
 985  
 986  
 987  
 988  
 989  
 990  
 991  
 992  
 993  
 994  
 995  
 996  
 997  
 998  
 999  
 1000  
 1001  
 1002  
 1003  
 1004  
 1005  
 1006  
 1007  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015  
 1016  
 1017  
 1018  
 1019  
 1020  
 1021  
 1022  
 1023  
 1024  
 1025

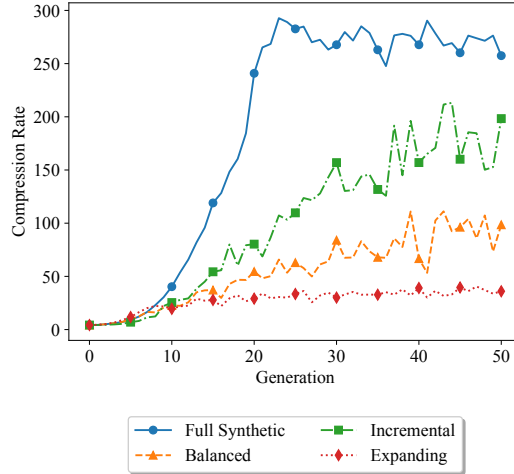


Figure 13: Compression rate over generations of a self-consuming training loop for different data cycles for textual data. A higher compression rate indicates less diversity. The sample size is truncated to equal length for better comparability. Line markers added every 5 generations for better display.

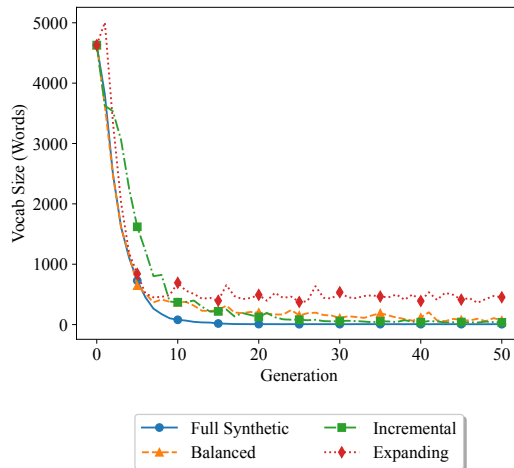


Figure 14: Vocabulary size over generations of a self-consuming training loop for different data cycles for textual data. The vocabulary size is defined as the number of unique words in a sample. The sample size is truncated to equal length for better comparability. Line markers added every 5 generations for better display.

sion rate, however, not as severe as for the other data cycles. This indicates that fresh data can slow down the effects of the self-consuming training as observed for the logic expression dataset.

Figure 14 displays the vocabulary size measured as unique words within a sample. Similar to the previous results we observe a decline in vocabulary size for all data cycles. While the full synthetic data cycle fully collapses in vocabulary size, closely followed by the balanced, and incremental data cycle, the expanding data cycle stabilizes around 500 words. However, this is still a drop in vocabulary size by 95%. Once again, this confirms a strong decline in diversity for all data cycles.

## E EXAMPLE OUTPUTS

### E.1 EXAMPLES OF LOGIC EXPRESSIONS

Table 2 provides examples of logic expressions that evaluate to `True` or `False`, as well as syntactically incorrect logic expressions (`Error`). In the erroneous example the logic expression is syntactically incorrect because of a missing parentheses. Additionally, all generated logic expressions from our experiments are provided in the supplementary material.<sup>5</sup>

Table 2: Examples of logic expressions that evaluate to `True`, `False`, or cannot be evaluated (`Error`).

TYPE	EXPRESSION
True	( not ( ( ( True and True ) and ( not True ) ) or ( ( False or False ) and ( not True ) ) ) ) <eos>
False	( ( ( ( True and False ) and ( True or False ) ) or ( not ( False or True ) ) ) and ( ( ( True and True ) and ( True or False ) ) or ( ( True or True ) and ( True or True ) ) ) ) <eos>
Error	( ( ( False or True ) and ( not True ) or ( ( not True ) or ( not False ) ) ) ) <eos>

### E.2 EXAMPLES OF NATURAL LANGUAGE

Figure 15-18 provide example outputs of models trained in a self-consuming training loop on the natural language dataset for all data cycles. Displayed is an example of the initial input and examples for generations 1, 10, 25, and 50. The examples are truncated for better display. All generated text samples are also provided in the supplementary material.

We manually inspect the text samples for the first few generations and notice a few hallucinated words that are not present in the original corpus, indicating a drop in correctness. However, overall the generated text is subjectively mostly correct and Shakespeare like with a diverse set of samples.

By generation 10 the diversity is already declining for all data cycles. This is most severe for the full synthetic data cycle which produces nearly the same output every time. The incremental data cycle is also very repetitive text with repeating words in half of the time. In the other half (not shown in the example), the output is just a repetition of a role without text (e.g. GLOUCESTER:\n\n GLOUCESTER:\n\n GLOUCESTER:\n\n). The other two data cycles (balanced and expanding) show a strong repetition of words.

By generation 25 the full synthetic data cycle is fully converged producing the same outputs for the rest of the self-consuming training loop. The incremental data cycle also shows heavy repetition or just samples a role repeatedly and by generation 50 only very few roles are sampled repeatedly without text besides a very few exceptions. The other two data cycles continue the trend of producing repetitive words and  $n$ -grams and by generation 50 only small diverse text samples are generated before the same word or  $n$ -gram is produced repeatedly.

Overall, while one could argue that the text is more *correct* (fewer hallucinated words) at the end of a self-consuming training loop, the stark decrease in diversity and the generation of meaningless text makes the language model unusable.

<sup>5</sup><https://figshare.com/s/b09d5bdb3b216998d330> - Anonymized for review

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

**Full Synthetic Data Cycle**

**Initial:**  
First Citizen:  
Before we proceed any further, hear me speak.

All:  
Speak, speak.

First Citizen:  
You are all resolved rather to die than to famish?

-----

**Generation 1:**  
YORK:  
I think you where you have seen to be envious death.

LORD STANLEY:  
No more, my lord: let me speak with men and my soul  
Your proof interior marriage, with it end.

WARWICK:  
Fair lords, and Warwick, where she lives you lies?

-----

**Generation 10:**  
WARWICK:  
What with his wisdom?

WARWICK:  
Why, what with his soul?

WARWICK:  
What with his soul?

-----

**Generation 25:**  
WARWICK:  
What with his wisdom?

WARWICK:  
What with his wisdom?

WARWICK:  
What with his wisdom?

-----

**Generation 50:**  
WARWICK:  
What with his wisdom?

WARWICK:  
What with his wisdom?

WARWICK:  
What with his wisdom?

Figure 15: Example model outputs of the full synthetic data cycle for the natural language experiments over generations of a self-consuming training loop.

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

**Balanced Data Cycle**

**Initial:**  
First Citizen:  
Before we proceed any further, hear me speak.

All:  
Speak, speak.

First Citizen:  
You are all resolved rather to die than to famish?

-----

**Generation 1:**  
ANGELO:  
Anon I would have no more for your better.

ANGELO:  
I will be so not better, whom you flow.

ISABELLA:  
I would he that would see your behalf.

-----

**Generation 10:**  
CORIOLANUS:  
You will not speak to the people of the prince of the people,  
Be so she is a man in the prince of the prince,  
And the traitor of the prince of the prince,  
And the prince of the prince of the prince,  
And the prince of the prince of the prince of the prince of the prince,  
And the seat of the prince of the state of the prince of the prince,  
And so she shall be so she is not so shall not speak to the prince of the prince of the prince,  
And the seat of the seat of the state of the prince of the state of the state,

-----

**Generation 25:**  
And give me dead, I have then,  
And I have done to see your holy profits of York,  
And so shall be so stand to the state of the man in the man in the story of the state,  
And so shall be spent for the seat of the seat of the first me of the prince,  
And the state of the state of the state of the state,  
And the state of the state of the state of the state of the state,  
And the state of the state of the state of the state,  
And the state of the state of the state of the state of the state,  
And the state of the state of the state of the state,

-----

**Generation 50:**  
Here's noble by so dead.

DUKE VINCENTIO:  
Here's a man in the prince of the prince, that the prince of the prince of the prince,  
And the prince of the prince of the prince of the prince,  
And the prince of the prince of the prince of the prince,  
And the prince of the seat of the prince of

And the state of the state of the state of the state, ... (*continues*)

Figure 16: Example model outputs of the balanced data cycle for the natural language experiments over generations of a self-consuming training loop.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

**Incremental Data Cycle**

**Initial:**  
First Citizen:  
Before we proceed any further, hear me speak.

All:  
Speak, speak.

First Citizen:  
You are all resolved rather to die than to famish?

-----

**Generation 1:**  
TYBALT:  
I'll stand my lord.

Marshal:  
Tell her her and be solemn'd and soldiers,  
For her wrong so understrange for Lancaster.

DUKE VINCENTIO:  
Go you, good good lady.

-----

**Generation 10:**  
LADY CAPULET:  
I will not say the world,  
And thou hast not said to be so for the world,  
And with the world.

LADY CAPULET:  
I have said to shame the world,  
And thou hast not said to show the world,  
And thou shalt not say the world in the world,  
That I have not said to shame the world.

-----

**Generation 25:**  
LADY ANNE:  
I will not say the world.

LADY ANNE:  
I will not say the world.

LADY ANNE:  
I will not say the world,  
And thou shalt not say the world in the world,  
And thou hast not said to the world,

-----

**Generation 50:**  
DUKE OF YORK:

DUKE OF YORK:

DUKE OF YORK:

Figure 17: Example model outputs of the incremental data cycle for the natural language experiments over generations of a self-consuming training loop.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

**Expanding Data Cycle**

**Initial:**  
First Citizen:  
Before we proceed any further, hear me speak.

All:  
Speak, speak.

First Citizen:  
You are all resolved rather to die than to famish?

-----

**Generation 1:**  
SICINIUS:  
The enemy noble are and true,  
And like it to life he putinion.

VOLUMNIA:  
So, let not for us.

CORIOLANUS:  
No, I wise the take tower?

-----

**Generation 10:**  
LADY ANNE:  
Would good we have my most for my soul, but your shouls  
The can is sould distrother particularish,  
And the such that the take tribunes,  
Where the name and the vail the people and to still grave the truest roper in the people,  
Which the port the people that who the hath the shall should to the people the people the  
people  
The people of the people and to the people tribuness of the people tribuness, the people,

-----

**Generation 25:**  
COMINIUS:  
Say would be medd to the common of the give the stords  
The people and the people the people tribuness, and the people tribuness  
Of the people tribuness of the country's noble the people,  
Which the people and the stirly tribunes of the people tribuness of the people and to the  
people tribuness of the people, and the people,  
Which the people and the people the people trick of the people the people tribuness,  
Who was in the people and the people the people and the people and the people the people  
and the people, the people, and the people,

-----

**Generation 50:**  
LADY ANNE:  
The gentleman is good more gone?

GLOUCESTER:  
You have been my lord, I will be possess to the gods gods!

GLOUCESTER:  
What the grace the people the people the people the people the people the people the people the people  
the people the people the people the people the people the people the people . . . (*continues*)

Figure 18: Example model outputs of the expanding data cycle for the natural language experiments over generations of a self-consuming training loop.