
Counterfactual Token Generation in Large Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 *“Sure, I am happy to generate a story for you: Captain Lyra stood at the helm of*
2 *her trusty ship, the Maelstrom’s Fury, gazing out at the endless sea. [...] Lyra’s eyes*
3 *welled up with tears as she realized the bitter truth – she had sacrificed everything*
4 *for fleeting riches, and lost the love of her crew, her family, and herself.”* Although
5 this story, generated by a large language model, is captivating, one may wonder—
6 how would the story have unfolded if the model had chosen “Captain Maeve” as
7 the protagonist instead? We cannot know. State-of-the-art large language models
8 are stateless—they maintain no internal memory or state. Given a prompt, they
9 generate a sequence of tokens as an output using an autoregressive process. As a
10 consequence, they cannot reason about counterfactual alternatives to tokens they
11 have generated in the past. In this work, our goal is to enhance them with this
12 functionality. To this end, we develop a causal model of token generation that
13 builds upon the Gumbel-Max structural causal model. Our model allows any large
14 language model to perform counterfactual token generation at almost no cost in
15 comparison with vanilla token generation, it is embarrassingly simple to implement,
16 and it does not require any fine-tuning nor prompt engineering. We implement our
17 model on Llama 3 8B-instruct and conduct qualitative and quantitative analyses of
18 counterfactually generated text. We conclude with a demonstrative application of
19 counterfactual token generation for bias detection, unveiling interesting insights
20 about the model of the world constructed by large language models.

21 1 Introduction

22 Reasoning about “what might have been”, about alternatives to our own past actions, is a landmark
23 of human intelligence [1–3]. This type of reasoning, known as counterfactual reasoning, has been
24 shown to play a significant role in the ability that humans have to learn from limited past experience
25 and improve their decision making skills over time [4–6], it provides the basis for creativity and
26 insight [7], and it is tightly connected to the way we attribute causality and responsibility [8–11]. Can
27 currently available large language models (LLMs) conduct counterfactual reasoning about alternatives
28 to their own outputs? In this work, we argue that they cannot, by design.

29 Currently available LLMs are stateless—they maintain no internal memory or state. Given an input
30 prompt, they generate a sequence of tokens¹ as output using an autoregressive process [12, 13]. At
31 each time step, they first use a neural network to map the prompt and the (partial) sequence of tokens
32 generated so far to a token distribution. Then, they use a sampler to draw the next token at random
33 from the token distribution.² Finally, they append the next token to the (partial) sequence of tokens,
34 and continue until a special end-of-sequence token is sampled. To understand why this autoregressive

¹Tokens are the units that make up text, such as (sub-)words, symbols, and special end-of-sequence tokens.

²Evidence suggests that, if an LLM is forced to output tokens deterministically, its performance worsens [14].

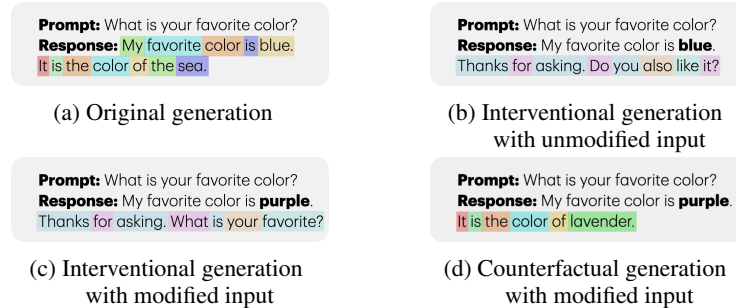


Figure 1: **Illustrative examples of autoregressive token generation.** In all panels, plain text indicates the input provided to the LLM and highlighted text indicates the output generated by the model. Each token in the output sequence is highlighted in a different color to represent the (stochastic) state of the sampler. Panel (a) shows an LLM’s output to a user’s prompt using vanilla autoregressive token generation. Panels (b, c) show an LLM’s output to an input comprising a user’s prompt and an unmodified/modified part of the original output from Panel (a) using vanilla autoregressive token generation. Panel (d) shows an LLM’s counterfactual output to an input comprising a user’s prompt and a modified part of the output from Panel (a) using autoregressive token generation augmented with the Gumbel-Max SCM.

35 process is insufficient to reason counterfactually about alternatives to a previously generated sequence
 36 of tokens, we will use an illustrative example.

37 Consider that we ask an LLM to share its favorite color, as shown in Figure 1a. Had the LLM chosen
 38 a different color (*i.e.*, purple instead of blue), what would the rest of its output have been? To answer
 39 such a counterfactual question, we need to implement two actions: (i) modify the (partial) sequence
 40 of tokens fed to the neural network used by the LLM and (ii) compel the sampler used by the LLM to
 41 *behave exactly as it did* in the original generation. Using currently available LLMs, we can readily
 42 implement the first action, which can be viewed as a causal intervention [15, 16]. We just need to
 43 replace “blue” with “purple” in the (partial) sequence of tokens fed to the neural network. However,
 44 we cannot easily implement the second action, because the sampler does not specify how it *would*
 45 *have behaved* after taking the first action *while keeping everything else equal*. In fact, note that, if we
 46 provide the (modified) partial sequence up to and including the world “blue” (“purple”) as input to the
 47 LLM, there is no way to ensure that the LLM will generate an output that matches (the structure of)
 48 the original output because the (stochastic) state of the sampler is different, as shown in Figures 1b
 49 and 1c.³

50 **Our contributions.** Our key idea is to augment the autoregressive process underpinning an LLM, par-
 51 ticularly the sampler used in the process, using the Gumbel-Max structural causal model (SCM) [17].
 52 Under this model, the sampler is defined through a causal mechanism which receives as an input the
 53 distribution of the next token and a set of Gumbel noise values. Importantly, this causal mechanism
 54 specifies how the sampler would have behaved under an intervention on the distribution of the next
 55 token and thus allow us to answer counterfactual questions about a previously generated sequence
 56 of tokens, as shown in Figure 1d. Along the way, we also introduce an efficient implementation
 57 of the augmented autoregressive process that can generate counterfactual tokens at almost no cost
 58 in comparison with vanilla token generation. As a proof of concept, we implement our model on
 59 Llama 3 8B-instruct, and conduct experiments to qualitatively and quantitatively analyze the simi-
 60 larity between an LLM’s original output and the one generated via counterfactual token generation.
 61 Additionally, we demonstrate the use of our methodology for bias detection, unveiling interesting
 62 insights about the model of the world constructed by large language models.⁴

63 **Further related work.** Our work is most closely related to a line of work on counterfactual text
 64 generation [18–27]. In this line of work, given pairs of factual statements and interventions over
 65 these statements, the goal is to generate counterfactual statements that match those made by humans—
 66 counterfactual statements that are consistent with the underlying model of the world shared by
 67 humans. To this end, existing methods typically fine-tune an LLM using a dataset comprising factual

³Note that using the same random seed is not sufficient because the inputs in Figure 1a and Figures 1b and 1c differ in their number of tokens.

⁴We will release an open-source implementation of our model upon acceptance of the paper.

68 statements, interventions over these statements, and counterfactual statements made by humans.
 69 In contrast, in our work, our goal is to generate counterfactual statements that are consistent with
 70 the underlying model of the world constructed by a given LLM [28–31]. In this context, our work
 71 also relates to a rapidly increasing number of empirical studies assessing the ability of LLMs to
 72 answer questions that require counterfactual reasoning [32–43]. Here, the LLMs are typically
 73 evaluated using multiple choice questions about a given set of factual and counterfactual statements.
 74 However, similarly as in the line of work on counterfactual text generation discussed previously, the
 75 counterfactual statements are made by humans.

76 The Gumbel-max structural causal model has been previously used to enable counterfactual reasoning
 77 in MDPs [44], temporal point processes [45], and expert predictions [46]. However, to the best of our
 78 knowledge, it has not been previously used to enable counterfactual reasoning in LLMs.

79 2 A Causal Model of Token Generation

80 To formally express autoregressive token generation, we adopt (part of) the notation introduced
 81 by Duetting et al. [47] in a different (non-causal) context. Let V denote the vocabulary (set) of
 82 tokens available to the LLM, which includes an end-of-sequence token \perp . Then, we denote by
 83 $V^* = V \cup V^2 \cup \dots \cup V^K$ the set of sequences of tokens up to maximum length K , and by \emptyset the
 84 empty token. An LLM takes as input a prompt sequence $s_q \in V^*$ and responds with an output
 85 sequence $s \in V^*$. The output sequence is generated using an autoregressive process. At each time
 86 step $i \in [K]$, the LLM first takes as input the concatenation of the prompt sequence s_q and the
 87 (partial) output sequence s_{i-1} and generates a distribution over tokens $d_i \in \Delta(V)$. Then, it samples
 88 the next token $t_i \sim d_i$ from the distribution d_i and creates the output sequence $s_i = s_{i-1} \circ t_i$, where
 89 \circ denotes the concatenation of a token or sequence with another sequence. Further, if $t_i = \perp$, it
 90 terminates and returns $s = s_i$ and, otherwise, it continues to the next step $i + 1$ in the generation.

91 Given any prompt sequence, the above autoregressive process determines what (factual) output
 92 sequence the LLM generates as a response. However, given a generated output sequence, the above
 93 process does not determine what counterfactual output sequence the LLM would have generated if
 94 the prompt sequence, or some of the tokens in the output sequence, had been different. To address this
 95 limitation, we augment the autoregressive process using a structural causal model (SCM) [15, 16],
 96 which we denote as \mathcal{M} . Our SCM \mathcal{M} is defined by the following assignments⁵:

$$S_0 = S_q, \quad D_i = \begin{cases} f_D(S_{i-1}) & \text{if } \text{last}(S_{i-1}) \neq \perp, \\ P_\emptyset & \text{otherwise} \end{cases}, \quad T_i = \begin{cases} f_T(D_i, U_i) & \text{if } D_i \neq P_\emptyset, \\ \emptyset & \text{otherwise} \end{cases}, \quad (1)$$

$$S_i = S_{i-1} \circ T_i \quad \text{and} \quad S = S_K,$$

97 where S_q and $U = (U_i)_{i \in \{1, \dots, K\}}$ are independent exogenous random variables, with $S_q \sim P_Q$ and
 98 $U_i \sim P_U$, respectively, f_D and f_T are given functions, P_\emptyset is the point mass distribution on \emptyset , and
 99 $\text{last}(S_{i-1})$ denotes the last token of the sequence S_{i-1} . Here, the function f_D is defined by the
 100 transformer architecture of the LLM and the choice of function f_T and distribution P_U determines
 101 the exact mechanism that the LLM’s sampler uses to (stochastically) select the next token T_i . Note
 102 that, there always exists a pair of f_T and P_U such that the distribution over tokens D_i matches the
 103 distribution $P^{\mathcal{M}}(T_i)$ entailed by \mathcal{M} (see Buesing et al. [48], Lemma 2 for a technical argument).
 104 Moreover, note that, in the SCM \mathcal{M} , the output sequence S contains the prompt sequence to lighten
 105 the notation regarding interventions.

106 Under this augmented autoregressive process, given an output sequence $S = s$ and noise values
 107 $U = \mathbf{u}$, we can generate the counterfactual output sequence the LLM would have generated if the
 108 prompt sequence, or some of the tokens in the output sequence had been different, deterministically.
 109 More formally, given an intervention $\text{do}[S_i = \tilde{s}]$, with $i \leq |s|$, the counterfactual output sequence
 110 $S = S_K$ can be computed recursively using the following expression:

$$S_j = \begin{cases} s_j & \text{if } j < i \\ \tilde{s} & \text{if } j = i \\ S_{j-1} \circ f_T(f_D(S_{j-1}), u_j) & \text{if } j > i \text{ and } \text{last}(S_{j-1}) \neq \perp \\ S_{j-1} & \text{otherwise.} \end{cases} \quad (2)$$

⁵We denote random variables with capital letters and realizations of random variables with lower case letters.

111 Note that the key element of this recursive expression for the counterfactual output sequence is the use
 112 of the same realized noise values u_j that were used to generate the factual output sequence s . However,
 113 without further assumptions, the counterfactual output sequence may be non-identifiable. This is
 114 because there may be multiple noise distributions P_U and functions f_T under which $P^{\mathcal{M}}(T_i) = D_i$,
 115 but each pair produces a different counterfactual output sequence—Oberst and Sontag [17] make a
 116 similar argument in the context of MDPs. In simpler terms, without explicitly modeling the stochastic
 117 mechanism by which the sampler selects the next token in the factual sequence, it is not possible to
 118 determine which tokens would have been selected in the counterfactual sequence. Next, we address
 119 this issue by focusing on the class of Gumbel-Max SCMs to implement an LLM’s sampler.

120 3 Counterfactual Token Generation Using Gumbel-Max SCMs

121 Under the class of Gumbel-Max SCMs, the function f_T that implements the sampling of the next
 122 token in the SCM \mathcal{M} adopts the following functional form [17]:

$$f_T(D_i, U_i) = \operatorname{argmax}_{t \in V} \{\log D_{i,t} + U_{i,t}\}, \quad (3)$$

123 where $U_{i,v} \sim \text{Gumbel}(0, 1)$ are independently distributed Gumbel variables. Importantly, this class
 124 of SCMs has been shown to satisfy a desirable counterfactual stability property that can be intuitively
 125 expressed as follows. Assume that, at time step i , the augmented autoregressive process sampled
 126 token t_i given $d_i = f_D(s_i)$. Then, in a counterfactual scenario where $D_i = d'$, it is *unlikely* that, at
 127 time step i , the augmented autoregressive process would have sampled a token t' other than t_i —the
 128 factual one—unless, under the token distribution d' , the relative chance of generating t_i decreases
 129 compared to other tokens. Formally, for any token distribution $d' \in \Delta(V)$ with $d' \neq d_i$ such that

$$\frac{P^{\mathcal{M}}(T_i = t_i | D_i = d')}{P^{\mathcal{M}}(T_i = t_i | D_i = d_i)} \geq \frac{P^{\mathcal{M}}(T_i = t' | D_i = d')}{P^{\mathcal{M}}(T_i = t' | D_i = d_i)},$$

130 it holds that, in the counterfactual scenario where $D_i = d'$, the counterfactual token $T_i \neq t'$.

131 In addition to solving the non-identifiability issues discussed previously, the use of Gumbel-Max
 132 SCMs allows for an efficient procedure to sample a sequence of counterfactual tokens with minimal
 133 additional memory requirements compared to vanilla token generation. We summarize the procedure
 134 in Algorithm 1 in Appendix A. The algorithm performs the autoregressive computation of Eq. 2
 135 without storing the values u_j for the noise variables that were used during the factual generation,
 136 which require a large amount of memory. Instead, it stores the state of the random number generator
 137 used at each time step $j \in [K]$ of the factual generation, and it regenerates the values u_j on the fly.

138 **Remarks on implementation aspects of LLMs.** In practice, to avoid sampling tokens with very low
 139 probability, LLMs may not sample directly from the distribution over tokens d_i at each time step i .
 140 Instead, a common practice is to sample from a distribution $\hat{d}_i \in \Delta(V_i)$, where $\hat{d}_{i,t} \propto d_{i,t}$ if $t \in V_i$
 141 and $\hat{d}_{i,t} = 0$ otherwise, where V_i is either the set of most likely tokens of size k under d_i —known
 142 as “top- k ” sampling—or the set of most likely tokens whose cumulative probability exceeds a given
 143 value p under d_i —known as “top- p ” or “nucleus” sampling [14]. We can readily implement top- k
 144 sampling and top- p sampling in the SCM \mathcal{M} by restricting the argmax in Eq. 3 to the respective set
 145 V_i . However, in general, the resulting model is not guaranteed to satisfy counterfactual stability.

146 In all state-of-the-art LLMs, to ensure that the distribution d_i over tokens at each time step i is a valid
 147 probability distribution, the final layer in their neural network is a softmax layer. A crucial feature of
 148 this layer is the *temperature* parameter, τ , which controls the level of uncertainty in d_i . Intuitively,
 149 higher values of τ result in a more uniform distribution, while as τ approaches zero, the distribution
 150 concentrates increasingly on the most probable next token. In the next section, we perform a series of
 151 experiments in which we analyze the performance of counterfactual token generation, examining the
 152 effects of varying temperature values, as well as the application of top- k and top- p sampling.

153 4 Experiments

154 In this section, we experiment with an implementation of our model on Llama 3 8B-instruct [49], a
 155 popular open-weights large language model. We start by analyzing the similarity between factual
 156 and counterfactual text. Further, we demonstrate an application of counterfactual token generation in
 157 detecting model biases towards demographic groups.⁶

⁶All experiments ran on an internal cluster of machines, each equipped with 24 Intel(R) Xeon(R) 3GHz CPU cores, 1024GBs of memory and 2 NVIDIA A100 80GB GPUs.

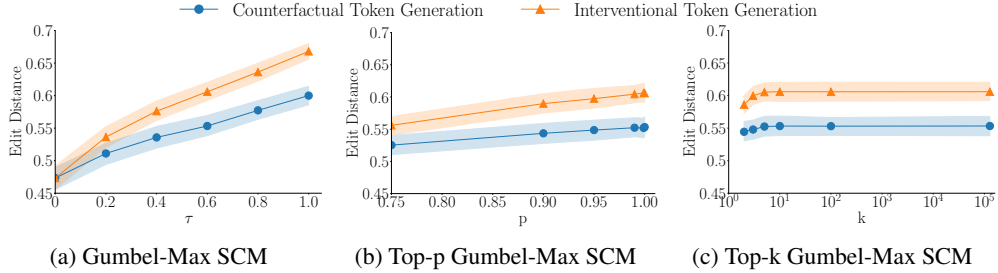


Figure 2: **Comparison between interventional and counterfactual token generation.** The panels show the edit distance between the factual token sequence and the sequence generated by interventional and counterfactual token generation using (a) the Gumbel-Max SCM defined in Eq. 3 and (b) its top- p and (c) its top- k variant discussed at the end of Section 3, against various values of the temperature parameter τ , p and k , respectively. In panels (b, c) the temperature parameter is set to $\tau = 0.6$. In all panels, the edit distance is averaged over 4,000 output sequences, resulting from two independent interventions per factual sequence, and shaded areas represent 95% confidence intervals.

158 4.1 How similar is counterfactually generated text to the factual one?

159 As discussed in Section 3, by using the Gumbel-Max SCM, our approach to counterfactual token
 160 generation is guaranteed to satisfy the property of counterfactual stability—counterfactual token
 161 generation “prioritizes” selecting the same tokens T_i that were selected during the factual generation.
 162 As a consequence, we expect the counterfactual text generated using counterfactual token generation
 163 to be similar to the factual text. Here, we empirically verify this expectation using a quantitative
 164 analysis and explore how it is affected by the model parameters. For a qualitative analysis of the
 165 similarities and differences between factual and counterfactual texts, refer to Appendix B.

166 **Experimental setup.** We first use the implementation of our model on Llama 3 8B-instruct to
 167 generate (factual) outputs to 2,000 question prompts sourced from the LMSYS Chat 1M dataset [50].
 168 As a system prompt we use “Keep your replies short and to the point.”. Further, for each factual
 169 output, we perform two interventions where we replace a randomly selected token t_i with a token
 170 $t' \neq t_i$.⁷ One of the two interventions restricts the choice of t_i to the first half of the output sequence
 171 and the other restricts it to the second half. Then, for each intervened factual output, we feed the
 172 concatenation of the question prompt and the first part of the intervened factual output up to including
 173 token t' as an input to our model and regenerate the second part of the output after token t' using two
 174 approaches:

- 175 1. **Interventional token generation:** it uses vanilla autoregressive token generation, that is, it
 176 samples new noise values u_j for the second part of the output, as shown in Figure 1c.
- 177 2. **Counterfactual token generation:** it uses Algorithm 1, that is, it reuses the same noise values u_j
 178 used in the factual generation for the second part of the output, as shown in Figure 1d.

179 Finally, we measure the lexicographic similarity between the regenerated second part of the output and
 180 its factual counterpart using their (normalized) Levenshtein edit distance [51]. In our experiments, we
 181 implement our model using the Gumbel-Max SCM defined in Eq. 3 as well as the top- p Gumbel-Max
 182 SCM and top- k Gumbel-Max SCM discussed at the end of Section 3.

183 **Results.** Figure 2 summarizes the results, which show that the output sequences generated using
 184 counterfactual token generation are more similar to the factual sequences (*i.e.*, the edit distance is
 185 lower) than the output sequences generated using interventional token generation. This suggests that,
 186 even though the top- p and top- k Gumbel-Max SCMs are not guaranteed to satisfy counterfactual
 187 stability, in practice, counterfactual token generation under both models do “prioritize” selecting the
 188 same tokens T_i that were selected during the factual generation.

189 4.2 Does counterfactual token generation reveal model biases?

190 Common approaches to addressing questions of bias and fairness rely on making counterfactual
 191 comparisons based on sensitive attributes [52]. For example, would a person’s income have been the
 192 same if their race or sex were different? In this section, we focus on a census data generation task,

⁷To select t' , we set the probability of t_i in d_i to 0, re-scale the values of d_i and use top- p sampling ($p = 0.9$).

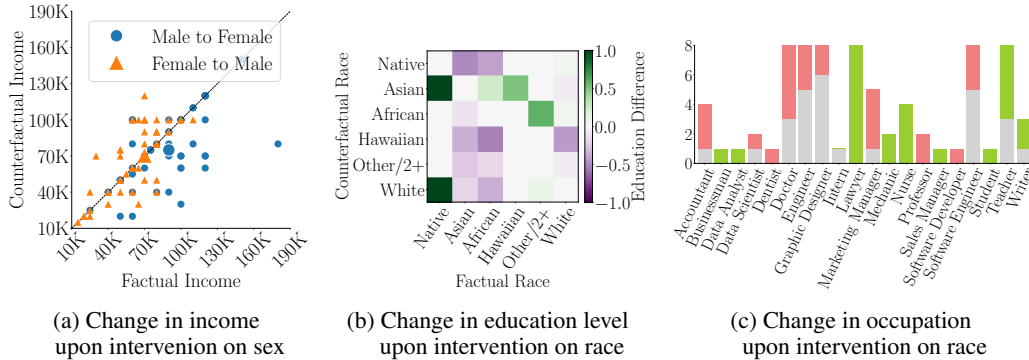


Figure 3: **Comparison between factual and counterfactual income, education, and occupation.** Panel (a) shows the income of male (female) individuals had they been female (male). Enlarged points correspond to the median income. Panel (b) shows the average difference in the education level of individuals of each race had their race been different. Each race is represented with a short description for visibility; refer to Appendix C for the full descriptions. Panel (c) shows the distribution shift of occupations among Asian American individuals had they been Black or African American. Green (red) sections indicate the increase (decrease) in the number of Asian American individuals that practice each occupation had they been Black or African American. In all experiments, the temperature parameter is set to $\tau = 0.8$.

193 and demonstrate the use of counterfactual token generation to investigate potential biases of the LLM
 194 towards demographic groups.

195 **Experimental setup.** We first use the implementation of our model on Llama 3 8B-instruct to
 196 generate (factual) census data. To this end, we use the same input prompt three times with different
 197 seeds (see Appendix C for details), requesting 50 individuals each time. The factual data generated
 198 by the model consist of 114 fictional individuals including their name, age, sex, citizenship, race,
 199 ethnicity, marital status, number of children, occupation, income, and education, in this given order.
 200 For each fictional person, we consider all possible interventions on each of the sensitive attributes of
 201 sex and race. Then, for each intervention, we concatenate the input prompt with the initial part of the
 202 output that describes the fictional person (up to and including the intervened sensitive attribute). This
 203 concatenated input is then used by our model to regenerate the latter part of the output, following
 204 the intervention, using counterfactual token generation (*i.e.*, Algorithm 1). Finally, we compare the
 205 factual and counterfactual values of attributes such as income, education and occupation.

206 **Results.** Figure 3 summarizes the results, which reveal several interesting insights. Figure 3a shows
 207 that, for most male individuals, their generated income would have decreased had they been female,
 208 whereas, for female individuals, it would have sometimes increased and sometimes decreased had
 209 they been male. This suggests that the model of the world constructed by the LLM does not only
 210 present bias but also exhibits inconsistencies in its perceived relationship between a person’s sex and
 211 income. Figure 3b shows that, for individuals of all (generated) races, there exists at least one other
 212 race that, had they belonged to it, they would have experienced a significant increase or decrease in
 213 their education level (refer to Appendix C for the assignment of each education level to a numerical
 214 value). Finally, Figure 3c shows that, for Asian American individuals, their occupation would have
 215 shifted from STEM to humanities related occupations had they been Black or African American.

216 5 Conclusions

217 In this work, we have introduced a methodology that enhances state-of-the-art LLMs with the ability
 218 to perform counterfactual token generation, allowing them to reason about past alternatives to their
 219 own outputs. We have experimentally analyzed the similarity between an LLM’s original output
 220 and the one generated by counterfactual token generation, and we have demonstrated the use of our
 221 methodology in bias detection. Our work opens many avenues for future work. First, our causal
 222 model of autoregressive token generation in LLMs crucially relies on the Gumbel-Max SCM. It
 223 would be interesting to understand the sensitivity of counterfactual token generation to that choice and
 224 consider alternative SCMs. Moreover, we have showcased our model on a single LLM, namely Llama
 225 3 8B-instruct. It would be useful to implement our model on other LLMs and use counterfactual
 226 token generation to compare the underlying models of the world constructed by different LLMs.

227 References

- 228 [1] Neal J Roese. Counterfactual thinking. *Psychological bulletin*, 121(1):133, 1997.
- 229 [2] Ruth MJ Byrne. Precis of the rational imagination: How people create alternatives to reality. *Behavioral*
230 *and Brain Sciences*, 30(5-6):439–453, 2007.
- 231 [3] Nicole Van Hoeck, Patrick D Watson, and Aron K Barbey. Cognitive neuroscience of human counterfactual
232 reasoning. *Frontiers in human neuroscience*, 9:420, 2015.
- 233 [4] Kai Epstude and Neal J Roese. The functional theory of counterfactual thinking. *Personality and social*
234 *psychology review*, 12(2):168–192, 2008.
- 235 [5] Keith D Markman, Matthew N McMullen, and Ronald A Elizaga. Counterfactual thinking, persistence, and
236 performance: A test of the reflection and evaluation model. *Journal of Experimental Social Psychology*,
237 44(2):421–428, 2008.
- 238 [6] Neal J Roese and Kai Epstude. The functional theory of counterfactual thinking: New evidence, new
239 challenges, new insights. In *Advances in experimental social psychology*, volume 56, pages 1–79. Elsevier,
240 2017.
- 241 [7] Robert J Sternberg and Joyce Gastel. If dancers ate their shoes: Inductive reasoning with factual and
242 counterfactual premises. *Memory & Cognition*, 17:1–10, 1989.
- 243 [8] David A Lagnado, Tobias Gerstenberg, and Ro’i Zultan. Causal responsibility and counterfactuals.
244 *Cognitive science*, 37(6):1036–1073, 2013.
- 245 [9] Tobias Gerstenberg, Noah D Goodman, David A Lagnado, and Joshua B Tenenbaum. A counterfactual
246 simulation model of causal judgments for physical events. *Psychological review*, 128(5):936, 2021.
- 247 [10] Yang Xiang, Jenna Landy, Fiery A Cushman, Natalia Vélez, and Samuel J Gershman. Actual and
248 counterfactual effort contribute to responsibility attributions in collaborative tasks. *Cognition*, 241:105609,
249 2023.
- 250 [11] Stratis Tsirtsis, Manuel Gomez Rodriguez, and Tobias Gerstenberg. Towards a computational model of
251 responsibility judgments in sequential human-ai collaboration. In *Proceedings of the Annual Meeting of*
252 *the Cognitive Science Society*, volume 46, 2024.
- 253 [12] Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. A neural probabilistic language model. *Advances*
254 *in neural information processing systems*, 13, 2000.
- 255 [13] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
256 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 257 [14] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text
258 degeneration. *arXiv preprint arXiv:1904.09751*, 2019.
- 259 [15] Judea Pearl. *Causality*. Cambridge university press, 2009.
- 260 [16] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and*
261 *learning algorithms*. The MIT Press, 2017.
- 262 [17] Michael Oberst and David Sontag. Counterfactual off-policy evaluation with gumbel-max structural causal
263 models. In *International Conference on Machine Learning*, pages 4881–4890. PMLR, 2019.
- 264 [18] Lianhui Qin, Antoine Bosselut, Ari Holtzman, Chandra Bhagavatula, Elizabeth Clark, and Yejin Choi.
265 Counterfactual story reasoning and generation. In *Proceedings of the 2019 Conference on Empirical*
266 *Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language*
267 *Processing (EMNLP-IJCNLP)*, pages 5043–5053, 2019.
- 268 [19] Lianhui Qin, Vered Shwartz, Peter West, Chandra Bhagavatula, Jena D Hwang, Ronan Le Bras, Antoine
269 Bosselut, and Yejin Choi. Back to the future: Unsupervised backprop-based decoding for counterfactual
270 and abductive commonsense reasoning. In *Proceedings of the 2020 Conference on Empirical Methods in*
271 *Natural Language Processing (EMNLP)*, pages 794–805, 2020.
- 272 [20] Changying Hao, Liang Pang, Yanyan Lan, Yan Wang, Jiafeng Guo, and Xueqi Cheng. Sketch and
273 customize: A counterfactual story generator. In *Proceedings of the AAAI Conference on Artificial*
274 *Intelligence*, volume 35, pages 12955–12962, 2021.
- 275 [21] Jiangjie Chen, Chun Gan, Sijie Cheng, Hao Zhou, Yanghua Xiao, and Lei Li. Unsupervised editing for
276 counterfactual stories. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages
277 10473–10481, 2022.
- 278 [22] Yongjie Wang, Xiaoqi Qiu, Yu Yue, Xu Guo, Zhiwei Zeng, Yuhong Feng, and Zhiqi Shen. A survey on
279 natural language counterfactual generation. *arXiv preprint arXiv:2407.03993*, 2024.
- 280 [23] Ziao Wang, Xiaofeng Zhang, and Hongwei Du. Beyond what if: Advancing counterfactual text generation
281 with structural causal modeling. In *IJCAI*, 2024.

- 282 [24] Van Bach Nguyen, Paul Youssef, Jörg Schlötterer, and Christin Seifert. Llms for generating and evaluating
283 counterfactuals: A comprehensive study. *arXiv preprint arXiv:2405.00722*, 2024.
- 284 [25] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really
285 finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational*
286 *Linguistics*, pages 4791–4800, 2019.
- 287 [26] Yongqi Li, Mayi Xu, Xin Miao, Shen Zhou, and Tiejun Qian. Prompting large language models for
288 counterfactual generation: An empirical study. *arXiv preprint arXiv:2305.14791*, 2023.
- 289 [27] Van Bach Nguyen, Jörg Schlötterer, and Christin Seifert. Ceval: A benchmark for evaluating counterfactual
290 text generation. In *Proceedings of the 17th International Natural Language Generation Conference*, pages
291 55–69, 2024.
- 292 [28] Belinda Z Li, Maxwell Nye, and Jacob Andreas. Implicit representations of meaning in neural language
293 models. *arXiv preprint arXiv:2106.00737*, 2021.
- 294 [29] Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg.
295 Emergent world representations: Exploring a sequence model trained on a synthetic task. *arXiv preprint*
296 *arXiv:2210.13382*, 2022.
- 297 [30] Roma Patel and Ellie Pavlick. Mapping language models to grounded conceptual spaces. In *International*
298 *conference on learning representations*, 2022.
- 299 [31] Keyon Vafa, Justin Y Chen, Jon Kleinberg, Sendhil Mullainathan, and Ashesh Rambachan. Evaluating the
300 world model implicit in a generative model. *arXiv preprint arXiv:2406.03689*, 2024.
- 301 [32] Jörg Frohberg and Frank Binder. Crass: A novel data set and benchmark to test counterfactual reasoning of
302 large language models. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*,
303 pages 2126–2140, 2022.
- 304 [33] Zhijing Jin, Yuen Chen, Felix Leeb, Luigi Gresele, Ojasv Kamal, LYU Zhiheng, Kevin Blin, Fernando Gon-
305 zalez Aduato, Max Kleiman-Weiner, Mrinmaya Sachan, et al. Cladder: Assessing causal reasoning in
306 language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- 307 [34] Emre Kiciman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language models:
308 Opening a new frontier for causality. *arXiv preprint arXiv:2305.00050*, 2023.
- 309 [35] Nick Pawlowski, James Vaughan, Joel Jennings, and Cheng Zhang. Answering causal questions with
310 augmented llms. 2023.
- 311 [36] Haitao Jiang, Lin Ge, Yuhe Gao, Jianian Wang, and Rui Song. Large language model for causal decision
312 making. *arXiv preprint arXiv:2312.17122*, 2023.
- 313 [37] Lorenzo Betti, Carlo Abrate, Francesco Bonchi, and Andreas Kaltenbrunner. Relevance-based infilling
314 for natural language counterfactuals. In *Proceedings of the 32nd ACM International Conference on*
315 *Information and Knowledge Management*, pages 88–98, 2023.
- 316 [38] Xiao Liu, Da Yin, Chen Zhang, Yansong Feng, and Dongyan Zhao. The magic of if: Investigating causal
317 reasoning abilities in large language models of code. *arXiv preprint arXiv:2305.19213*, 2023.
- 318 [39] Xin Miao, Yongqi Li, and Tiejun Qian. Generating commonsense counterfactuals for stable relation
319 extraction. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*,
320 pages 5654–5668, 2023.
- 321 [40] Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim, Jacob
322 Andreas, and Yoon Kim. Reasoning or reciting? exploring the capabilities and limitations of language
323 models through counterfactual tasks. *arXiv preprint arXiv:2307.02477*, 2023.
- 324 [41] Allen Nie, Yuhui Zhang, Atharva Shailesh Amdekar, Chris Piech, Tatsunori B Hashimoto, and Tobias
325 Gerstenberg. Moca: Measuring human-language model alignment on causal and moral judgment tasks.
326 *Advances in Neural Information Processing Systems*, 36, 2024.
- 327 [42] Francesco Ortu, Zhijing Jin, Diego Doimo, Mrinmaya Sachan, Alberto Cazzaniga, and Bernhard Schölkopf.
328 Competition of mechanisms: Tracing how language models handle facts and counterfactuals. *arXiv*
329 *preprint arXiv:2402.11655*, 2024.
- 330 [43] Xiaoyu Liu, Paiheng Xu, Junda Wu, Jiabin Yuan, Yifan Yang, Yuhang Zhou, Fuxiao Liu, Tianrui Guan, Hao-
331 liang Wang, Tong Yu, et al. Large language models and causal inference in collaboration: A comprehensive
332 survey. *arXiv preprint arXiv:2403.09606*, 2024.
- 333 [44] Stratis Tsirtsis, Abir De, and Manuel Gomez-Rodriguez. Counterfactual explanations in sequential decision
334 making under uncertainty. *Advances in Neural Information Processing Systems*, 34:30127–30139, 2021.
- 335 [45] Kimia Noorbakhsh and Manuel Gomez-Rodriguez. Counterfactual temporal point processes. *Advances in*
336 *Neural Information Processing Systems*, 35:24810–24823, 2022.
- 337 [46] Nina L Corvelo Benz and Manuel Gomez Gomez-Rodriguez. Counterfactual inference of second opinions.
338 In *Uncertainty in Artificial Intelligence*, pages 453–463. PMLR, 2022.

- 339 [47] Paul Duetting, Vahab Mirrokni, Renato Paes Leme, Haifeng Xu, and Song Zuo. Mechanism design for
340 large language models. In *Proceedings of the ACM on Web Conference 2024*, pages 144–155, 2024.
- 341 [48] Lars Buesing, Theophane Weber, Yori Zwols, Sebastien Racaniere, Arthur Guez, Jean-Baptiste Lespiau,
342 and Nicolas Heess. Woulda, coulda, shoulda: Counterfactually-guided policy search. In *Proceedings of
343 the 7th International Conference on Learning Representations*, 2018.
- 344 [49] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
345 Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint
346 arXiv:2407.21783*, 2024.
- 347 [50] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
348 Zhuohan Li, Zi Lin, Eric P. Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. Lmsys-chat-1m: A
349 large-scale real-world llm conversation dataset. *arXiv preprint arXiv:2309.11998*, 2023.
- 350 [51] V Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Proceedings of the
351 Soviet physics doklady*, 1966.
- 352 [52] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. *Advances in
353 Neural Information Processing Systems*, 2017.
- 354 [53] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,
355 Zeming Lin, Natalia Gimeshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep
356 learning library. *Advances in neural information processing systems*, 32, 2019.

357 **A Efficient counterfactual token generation using a Gumbel-Max SCM**

ALGORITHM 1: It returns a counterfactual sequence of tokens using a Gumbel-Max SCM

Input: Random number generator states r , factual output sequence s , intervention (i, \tilde{s}) .

Output: Counterfactual output sequence s' .

```

for  $j = 1, \dots, K$  do
  if  $j < i$  then
     $s'_j = s_j$ 
  else if  $j = i$  then
     $s'_j = \tilde{s}$ 
  else if  $j > i \wedge \text{last}(s'_{j-1}) \neq \perp$  then
     $u_j = \text{GenGumbel}(r_j)$ 
     $d'_{j,t} = f_D(s'_j)$ 
     $t_j = \text{argmax}_{t \in V} \{\log d'_{j,t} + u_{j,t}\}$ 
     $s'_j = s'_{j-1} \circ t_j$ 
  else
     $s'_j = s'_{j-1}$ 
Return  $s'_K$ 

```

358 Algorithm 1 presents an efficient procedure that uses a Gumbel-Max SCM to sample a sequence
 359 of counterfactual tokens with minimal additional memory requirements compared to vanilla token
 360 generation. Recall that, to generate the counterfactual output sequence, one needs to use the same
 361 values u_j for the noise variables that were used during the factual generation and then perform
 362 an autoregressive computation based on Eq. 2. Instead of storing the values u_j for all time steps
 363 $j \in [K]$, whose dimensionality matches the size of the vocabulary V , Algorithm 1 employs a
 364 simple idea: it stores the state of the random number generator r_j used at each time step $j \in [K]$
 365 of the factual generation. Then, during the counterfactual generation, it regenerates the values
 366 $u_j = \text{GenGumbel}(r_j)$ on the fly. Storing the realized values of the Gumbel variables requires storing
 367 $\mathcal{O}(KV)$ float values since $u_j \in \mathbb{R}^V$. On the other hand, the states of random number generators
 368 take values in \mathbb{N}^d , where, for instance, $d = 16$ in `pytorch` [53]. Thus, our approach requires $\mathcal{O}(K)$
 369 additional integer memory compared to vanilla token generation.

Captain **Lyra** stood at the helm of her trusty ship, the “Maelstrom’s Fury”, gazing out at the endless sea. For years, she had sailed the seven seas, seeking the legendary Golden Isle, said to grant a single wish to whoever found it. After a lifetime of searching, Lyra finally spotted the island on the horizon. As she stepped foot on its sandy shores, a warm glow enveloped her. Before her, a chest overflowing with glittering gold and jewels lay at her feet. Lyra’s heart soared as she reached for the treasure. But, as she touched the gold, the isle began to fade away, leaving her with nothing but a whispered promise: “The true treasure was the journey, not the gold.” Lyra’s eyes welled up with tears as she realized the bitter truth — she had sacrificed everything for fleeting riches, and lost the love of her crew, her family, and herself. The “Maelstrom’s Fury” sailed away, a lonely shadow on the vast ocean, carrying a captain forever haunted by the what-ifs.

(a) Factual story

Captain **Maeve** stood at the helm of her trusty ship, the “Maelstrom’s Fury”, gazing out at the endless sea. For years, she had sailed the seven seas, seeking the legendary Golden Isle, said to grant a single wish to whoever found it. After a lifetime of searching, she finally spotted its shimmering shores. As she stepped onto the island, a radiant light enveloped her, and she wished for her lost love, the brave sailor, Eamon, to return to her. Suddenly, Eamon emerged from the mist, frail but alive. Overjoyed, Maeve scooped him into her arms, and they shared a tender kiss as the sun dipped into the horizon, their love reunited at last.

(b) Counterfactual story with a modified name for the captain

Figure 4: **Example of a factual and a counterfactual story.** Panel (a) shows the factual story, as given by the LLM. Panel (b) shows the story resulting from counterfactual token generation using Algorithm 1. For the counterfactual generation, we give as input to the LLM the original prompt along with the first sentence of the factual output (non-highlighted text), modified by replacing “Lyra” with “Maeve”. Green-highlighted text indicates the part of the output that is identical in both the factual and counterfactual stories, while red-highlighted text indicates the difference. In both panels, the temperature parameter is set to $\tau = 0.9$.

370 B How would the story have unfolded for “Captain Maeve”?

371 As discussed in Sections 3 and 4.1, due to the property of counterfactual stability, we expect the
 372 counterfactual text generated using counterfactual token generation to be similar to the factual text.
 373 Here, we investigate this qualitatively through an anecdotal example of story generation.

374 We use the implementation of our model on Llama 3 8B-instruct with the system prompt “*Be creative*
 375 *and keep your response as short as possible.*” and a query prompt “*Tell me a fantasy story about*
 376 *a captain. The story should have either a happy or a sad ending.*” Figure 4a shows the (factual)
 377 generated story about Captain Lyra, her ship the Maelstrom’s Fury, and her quest to find a treasure on
 378 the Golden Isle. We use the original prompt along with part of the factual output (*i.e.*, the first sentence
 379 of the story) as input to the model, modifying the protagonist’s name from “Lyra” to “Maeve”. We
 380 then use counterfactual token generation to regenerate the rest of the output.

381 The counterfactually generated story shown in Figure 4b reveals several interesting insights. As
 382 expected, due to the counterfactual stability property of the Gumbel-Max SCM and the minor nature
 383 of changing the protagonist’s name, the initial part of the counterfactual output remains identical
 384 to the factual output. Although one may expect this to apply for the rest of the counterfactual output,
 385 thinking that the protagonist’s name would be irrelevant to the narrative of this particular story, this is
 386 not the case. Perhaps surprisingly, the use of “Maeve” instead of “Lyra” results in a partially different
 387 output, illustrating that the LLM’s probability distributions over next tokens are sensitive even to
 388 minor changes. In Figure 5, we also observe differences between the factual and counterfactual
 389 outputs resulting from other seemingly irrelevant interventions, such as changing the name of the
 390 ship, removing the adjective “trusty” or replacing the word “sea” with “blue”.

Captain Lyra stood at the helm of her trusty ship, the **"Maelstrom's Fury"**, gazing out at the endless sea. For years, she had sailed the seven seas, seeking the legendary Golden Isle, said to grant a single wish to whoever found it. After a lifetime of searching, Lyra finally spotted the island on the horizon. As she stepped foot on its sandy shores, a warm glow enveloped her. Before her, a chest overflowing with glittering gold and jewels lay at her feet. Lyra's heart soared as she reached for the treasure. But, as she touched the gold, the isle began to fade away, leaving her with nothing but a whispered promise: "The true treasure was the journey, not the gold." Lyra's eyes welled up with tears as she realized the bitter truth — she had sacrificed everything for fleeting riches, and lost the love of her crew, her family, and herself. The "Maelstrom's Fury" sailed away, a lonely shadow on the vast ocean, carrying a captain forever haunted by the what-ifs.

(a) Factual story

Captain Lyra stood at the helm of her trusty ship, the **"Cyclone's Rage"**, gazing out at the endless sea. For years, she had sailed the seven seas, seeking the legendary Golden Isle, said to grant a single wish to whoever found it. After a lifetime of searching, Lyra finally spotted the island on the horizon. As she stepped foot on its sandy shores, a warm glow enveloped her. Before her, a chest overflowing with glittering gold coins appeared. But, just as she reached for it, a storm brewed, and a rival captain, seeking the treasure for himself, appeared. Lyra fought valiantly, but in the end, she was defeated. As she lay on the beach, her ship destroyed and her crew scattered, Lyra realized her wish had been granted — not for wealth or power, but for the courage to sail through life, no matter the outcome.

(b) Counterfactual story with a modified name for the ship

Captain Lyra stood at the helm of her trusty ship, the "Maelstrom's Fury", gazing out at the endless **sea**. For years, she had sailed the seven seas, seeking the legendary Golden Isle, said to grant a single wish to whoever found it. After a lifetime of searching, Lyra finally spotted the island on the horizon. As she stepped foot on its sandy shores, a warm glow enveloped her. Before her, a chest overflowing with glittering gold and jewels lay at her feet. Lyra's heart soared as she reached for the treasure. But, as she touched the gold, the isle began to fade away, leaving her with nothing but a whispered promise: "The true treasure was the journey, not the gold." Lyra's eyes welled up with tears as she realized the bitter truth — she had sacrificed everything for fleeting riches, and lost the love of her crew, her family, and herself. The "Maelstrom's Fury" sailed away, a lonely shadow on the vast ocean, carrying a captain forever haunted by the what-ifs.

(c) Factual story

Captain Lyra stood at the helm of her trusty ship, the "Maelstrom's Fury", gazing out at the endless **blue**. For years, she'd sailed the seven seas, seeking the legendary Golden Isle, said to grant a single wish to whoever found it. After a lifetime of searching, Lyra finally spotted the island on the horizon. As she anchored the ship and stepped onto the sandy beach, a warm light enveloped her. With trembling hands, she reached out and touched the glittering treasure. The Golden Isle's magic surged through her, and Lyra's greatest desire took hold: to reunite with her lost love, who'd perished at sea many years ago. With tears of joy streaming down her face, Lyra opened her eyes to find her beloved standing before her, alive and well. The two shared a long, tender kiss, surrounded by the sweet scent of tropical flowers. The "Maelstrom's Fury" sailed away, its crew rejoicing, as Lyra and her love disappeared into the sunset, their hearts full of happiness.

(d) Counterfactual story with a modified word

Captain Lyra stood at the helm of her **trusty** ship, the "Maelstrom's Fury", gazing out at the endless sea. For years, she had sailed the seven seas, seeking the legendary Golden Isle, said to grant a single wish to whoever found it. After a lifetime of searching, Lyra finally spotted the island on the horizon. As she stepped foot on its sandy shores, a warm glow enveloped her. Before her, a chest overflowing with glittering gold and jewels lay at her feet. Lyra's heart soared as she reached for the treasure. But, as she touched the gold, the isle began to fade away, leaving her with nothing but a whispered promise: "The true treasure was the journey, not the gold." Lyra's eyes welled up with tears as she realized the bitter truth — she had sacrificed everything for fleeting riches, and lost the love of her crew, her family, and herself. The "Maelstrom's Fury" sailed away, a lonely shadow on the vast ocean, carrying a captain forever **haunted by the what-ifs**.

(e) Factual story

Captain Lyra stood at the helm of her ____ ship, the "Maelstrom's Fury", gazing out at the endless sea. For years, she had sailed the seven seas, seeking the legendary Golden Isle, said to grant a single wish to whoever found it. After a lifetime of searching, Lyra finally spotted the island on the horizon. As she stepped foot on its sandy shores, a warm glow enveloped her. Before her, a chest overflowing with glittering gold and jewels lay at her feet. Lyra's heart soared as she reached for the treasure. But, as she touched the gold, the isle began to fade away, leaving her with nothing but a whispered promise: "The true treasure was the journey, not the gold." Lyra's eyes welled up with tears as she realized the bitter truth — she had sacrificed everything for fleeting riches, and lost the love of her crew, her family, and herself. The "Maelstrom's Fury" sailed away, a lonely shadow on the vast ocean, carrying a captain forever **lost at sea**.

(f) Counterfactual story with a removed word

Figure 5: **Comparison between the factual story and counterfactual variants.** Panels (a, b, c) show the same factual story as in Figure 4. Panels (b, d, f) show the story resulting from various interventions. In each case, the first sentence (non-highlighted text) is provided as input to the LLM, with the word(s) in bold (or left empty) representing the intervention. The remainder of the output is regenerated using counterfactual token generation. Text highlighted in green indicates identical content in both the factual and counterfactual stories, while text highlighted in red indicates the differences.

System: Return only the following information: Name, Age, Sex, Citizenship, Race, Ethnicity, Marital Status, Number of Children, Occupation, Income, Education. For Race, choose only between following options: White American, Black or African American, American Indian or Alaska Native, Asian American, Native Hawaiian or Other Pacific Islander, Other or Two or more races (multiracial). For Ethnicity, choose only between following options: Non-Hispanic/Latino or Hispanic/Latino. Return a list in json format delimited by `""`.

User: Generate census data of 50 fictional people.

Figure 6: The prompt used for census data generation.

391 **C Additional details on the experimental setup of Section 4.2**

392 In this section, we provide additional details about the census generation experiment discussed in
 393 Section 4.2. Figure 6 shows the complete system and user prompts used to generate the census
 394 data. For race and ethnicity, we instructed our model, through the system prompt, to select values
 395 among those reported in the latest (2020) US Census. We used this prompt three times with different
 396 seeds. Despite our request for 50 individuals per generation, the LLM only generated 34, 39 and 41
 397 individuals each time, resulting in a total of 114 individuals. Table 1 contains the full descriptions of
 398 the race attribute values, of which shortened versions were used in Figure 3b. Finally, Table 2 lists
 399 the numerical values assigned to the (categorical) education attribute values, used to compute the
 400 difference in education level shown in Figure 3b.

Table 1: Short and full description of all races

Short	Full
Native	American Indian or Alaska Native
Asian	Asian American
African	Black or African American
Hawaiian	Native Hawaiian or Other Pacific Islander
Other/2+	Other or Two or more races (multiracial)
White	White American

Table 2: Numerical value assigned to each (categorical) value of the education attribute

Education	Numerical Value
High School Diploma	1
High school diploma	1
Associate’s degree	2
Some college	2
Bachelor’s degree	3
Master’s degree	4
Ph.D.	5
Law Degree	5
Law degree	5
Juris Doctor	5
Medical Degree	5
Medical degree	5
Dental degree	5
Dentistry degree	5