# Can AI Be as Creative as Humans?

**Haonan Wang**[1]    **James Zou**[2]    **Michael Mozer**[3]    **Linjun Zhang**[4]    **Anirudh Goyal**[3]
**Alex Lamb**[5]    **Zhun Deng**[6]    **Michael Qizhe Xie**[1]    **Hannah Brown**[1]    **Kenji Kawaguchi**[1]

[1]National University of Singapore    [2]Stanford University    [3]Google DeepMind
[4]Rutgers University    [5]Microsoft Research    [6]Columbia University

Project Page: `ai-relative-creativity.github.io`

## Abstract

Creativity serves as a cornerstone for societal progress and innovation, but its assessment remains a complex and often subjective endeavor. With the rise of advanced generative AI models capable of tasks once reserved for human creativity, the study of AI's creative potential becomes imperative for its responsible development and application. This paper addresses the complexities in defining and evaluating creativity by introducing a new concept called *Relative Creativity*. Instead of trying to define creativity universally, we shift the focus to whether AI can match the creative abilities of a hypothetical human. This perspective draws inspiration from the Turing Test, expanding upon it to address the challenges and subjectivities inherent in evaluating creativity. This methodological shift facilitates a statistically quantifiable evaluation of AI's creativity, which we term *Statistical Creativity*. This approach allows for direct comparisons of AI's creative abilities with those of specific human groups. Building on this foundation, we discuss the application of statistical creativity in contemporary prompt-conditioned autoregressive models. In addition to defining and analyzing a measure of creativity, we introduce an actionable training guideline, effectively bridging the gap between theoretical quantification of creativity and practical model training. Through these multifaceted contributions, the paper establishes a cohesive, continuously evolving, and transformative framework for assessing and fostering statistical creativity in AI models.

## 1   Introduction

Creativity, usually deemed as a quintessential human trait, is not just an individual trait but a transformative force that shapes societies, catalyzing advancements in science, technology, and the arts. It forms the backbone of innovation, fuels economic growth, and social change (Amabile, 1996; Boden, 2003; Kirkpatrick, 2023). In the rapidly evolving landscape of the digital age, artificial intelligence (AI) has introduced transformative avenues for creative endeavors. Advanced generative deep learning models have not only shown proficiency in solving complex problems, such as drug and protein synthesis (Jumper et al., 2021), but they have also excelled in artistic pursuits, including composing poetry and crafting narratives (Franceschelli and Musolesi, 2023). Additionally, these models have displayed remarkable aptitude for generating novel ideas, even outpacing MBA students in terms of both quality and uniqueness of innovative product and service concepts (Terwiesch and Ulrich, 2023). As AI's generative capabilities blur the lines between human and machine-generated work, this raises the stakes for the study of creativity, especially for the creativity of AI.

Traditionally, human creativity has been extensively studied and analyzed across various disciplines, such as psychology, philosophy, and cognitive science (Amabile, 1996; Boden, 2003; Kirkpatrick, 2023). However, there is still no consensus on defining "creativity", primarily due to the subjective nature involved in the various definitions proposed in scholarly literature (Runco and Jaeger, 2012; Sawyer, 2012). Even if we consider the widely accepted definition of creativity as a blend of novelty and quality (Boden, 2003; Câmara Pereira, 2007), the inherent subjectivity of these criteria remains problematic. What is deemed novel and of quality can differ greatly across various cultures, disciplines, and time periods. For instance, a computer science paper written in iambic pentameter might be deemed highly novel but of low quality by one community, whereas another
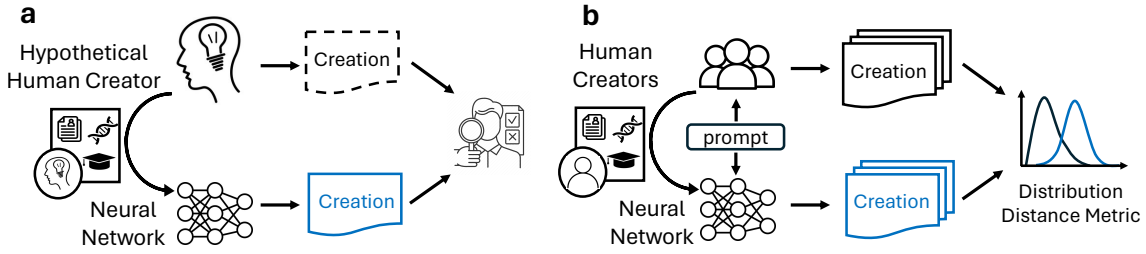
1

Figure 1: Illustration of Relative Creativity (**a**) and Statistical Creativity (**b**). (**a**): Relative Creativity assessed by AI's ability to create art indistinguishable from that of a hypothetical human, given the same biographical influences. (**b**): Statistical Creativity measured by AI's ability to generate creations to prompts that are indistinguishable from those of existing human creators, as determined by a distribution distance metric.

might view it as higher in quality but less novel. This lack of consensus on the creativity of humans hinders progress in understanding and developing creativity in AI. Additionally, in the study of AI creativity, sophisticated AI models are inherently designed to generate outputs reflecting their training data (Foster, 2022). This raises critical questions about the authenticity of their creations—whether they are genuinely original or merely repackaged elements, a point highlighted by Somepalli et al. (2023). This complexity further complicates the task of defining and analyzing creativity in AI. Against this backdrop, establishing an objective framework for evaluating AI's creativity is an essential step towards the responsible development and further evaluation of AI systems.

In this work, we aim to establish a concrete framework for exploring creativity in artificial intelligence. Diverging from traditional methods that attempt to answer the question, "Can AI be creative?"—a question mired in the ambiguous and elusive task of defining creativity absolutely—we instead pivot towards a more concrete inquiry: *"Can AI be as creative as humans?"* To formally structure this question, we introduce the concept of **Relative Creativity** (as defined in Definition 1), where an AI model is deemed as creative as a hypothetical, yet realistic, human creator if it can produce works indistinguishable from that creator, as determined by an evaluator. Building on this, we propose a method to evaluate whether and to what degree a model achieves relative creativity through **Statistical Creativity** (outlined in Theorem 1). This method involves comparing AI outputs with those of actual human creators, coupling with a theoretical guarantee to affirm that statistical creativity is a practical and valid measure for creativity. In line with recent advancements in large language model (LLM), such as GPT-4 (OpenAI, 2023) and Llama (Touvron et al., 2023), our research investigates statistical creativity in auto-regressive models that excel in next-token prediction and generation within a prompting paradigm, as detailed in Corollaries 1 and 2. Furthermore, to connect the theoretical quantification of creativity with practical methods for attaining creativity in model training, we introduce the **Statistical Creative Loss** (outlined in Corollary 3). This loss function offers a viable strategy for nurturing statistical creativity from training. The corresponding analysis also clarifies the relationship between the log-likelihood of next-token prediction and the necessary volume of training samples required for developing models capable of exhibiting statistical creativity. Overall, our framework presents a comprehensive, evolving methodology for evaluating and enhancing statistical creativity in AI, effectively bridging theoretical concepts with practical application.

## 2 Result Overview

In an effort to establish a concrete framework for exploring creativity in artificial intelligence, Section 3 introduces the concept of Relative Creativity (see Definition 1). This concept, echoing the Turing Test's comparative approach to intelligence assessment, proposes that an AI model can be deemed as creative as a hypothetical human creator if its outputs are indistinguishable from those of the creator, as judged by a designated evaluator. Applying this concept directly to AI model evaluation encounters practical challenges, mainly due to the hypothetical nature of the creator, which hinders empirical validation. To overcome this, Section 4.1 presents

Statistical Creativity (Theorem 1), a methodology that makes assessing Relative Creativity both feasible and realistic, anchoring it in the context of existing individuals' creative outputs. Further exploration in Section 4.2 delves into autoregressive models, a prevalent technique in generative models, introducing Autoregressive Statistical Creativity (Theorem 2). Moreover, Section 4.3 applies this framework to cutting-edge model paradigm, like prompt-conditioned large language models (LLMs), as explored in Brown et al. (2020). This section unveils Prompt-Contextualized Autoregressive Statistical Creativity (Corollary 2). In Section 5, we introduce the Statistical Creative Loss (Corollary 3), bridging the theoretical-practical gap and providing actionable guidelines for training creative AI models. The analysis here clarifies the relationship between the upper bound of next-token prediction log-likelihood and the required volume of training samples for achieving Statistical Creativity. Our findings highlight that the quantity of creator-creation pairs are crucial, rather than merely having a large volume of creation data. This insight renders the concept of Statistical Creativity particularly applicable in current AI frameworks based on next-token prediction (as highlighted in Remark 4). We strive for this paper not only to offer theoretical perspectives but to guide the discourse on AI creativity, advocating for relative evaluations to spur empirical research and establish a framework for assessing and augmenting the creative capabilities of AI models.

We summarize our contributions as follows.

- **Introduction of Relative Creativity:** We present the concept of Relative Creativity, a notion of evaluating the creative capabilities of artificial intelligence (AI) systems. This concept focuses on a comparative evaluation, eschewing absolute definitions in favor of a more relative understanding of creativity.

- **Integration of Subjectivity in Creative Assessment:** Relative creativity acknowledges the inherent subjectivity in creativity and incorporates the subjectivity into the comparison process, akin to how the Turing Test assesses intelligence through comparison rather than fixed definitions. The subjective nature of creativity is crystallized in the choice of the anchor population. This allows the study of AI creativity to maintain a degree of objectivity.

- **Conceptualization of Statistical Creativity:** To facilitate empirical studies, we introduce Statistical Creativity, a means for evaluating whether and to what extent an AI model achieves relative creativity by assessing its ability to replicate the creative abilities of a specific human population.

- **Application to Autoregressive Models:** Delving into autoregressive models, we propose a practical measure of statistical creativity for autoregressive models with next-token prediction. This measure can be applied to prompt-conditioned autoregressive models, ensuring its applicability to cutting-edge Large Language Models (LLMs).

- **Development of Statistical Creativity Loss:** We propose the Statistical Creativity Loss, which aligns theoretical principles with practical training guidelines for fostering creative AI, and provides insights into collecting creator-creation data and optimizing training objectives.

## 3 Relative Creativity

### 3.1 What is Relative Creativity?

Diverging from traditional methods setting an absolute threshold or checklist to determine if AI can be deemed creative, that necessitates a foray into the contentious task of defining creativity universally, we shift the study by proposing an creativity notion for AI called ***Relative Creativity***. Relative creativity sidesteps the difficulties associated with defining creativity in absolute terms, reminiscent the Turing Test's approach to intelligence evaluation (Turing, 2009). Note, the Turing Test eschews absolute definitions of intelligence, instead opting for a relative metric that contrasts machine behavior with human responses in conversational scenarios. In

the evaluation of creativity, an AI model is deemed "relatively creative" if it can generate creations that are indistinguishable from those of a hypothetical yet plausible human creator, as judged by an evaluator. In simpler terms, relative creativity posits that an AI model is as creative as a hypothetical human if its creations, when conditioned on the biographical data or characteristics of this hypothetical creator, are indistinguishable from what that human would have created. The way in which this notion of creativity is "relative" is that it depends on the individual to which the entity is being compared. For instance, an AI system might seem highly creative when compared to a non-expert human creator but may seem less so when compared to a expert designer or artist. Relative creativity distinguishes itself through acknowledges the inherently subjective facets of creativity—such as originality, divergent thinking, and problem-solving skills—these elements are integrated into the anchor selection process. Because the subjectivity of creativity is crystallized into the choice of the anchor human, relative creativity leaves the study of AI creativity to be objective. Regardless of whether the creative outputs of a specific individual or group satisfy the varied, subjective standards for what constitutes creativity, this notion reframes the debate, enabling a focused, empirical investigation into whether an AI model can replicate the creative capacities of a predefined human benchmark.

## 3.2 Notation and Formal Definition

**Preliminary Notations.** Imagine if the task is to write poems. Let $\mathcal{X}$ represent a finite set of poems, denoted by $|\mathcal{X}| < \infty$, where each individual poem in the set is given as $x \in \mathcal{X}$. Consider a generative model $q$ that receives a specific form of information $I \in \mathcal{I}$. We assume that $I$ consists of partial information pertinent to poets, such as personal background, artistic education, societal and historical Context, etc. The model outputs a creation according to the conditional distribution $q(\cdot|I)$. This format of conditional generation is widely adapted by current practical models (Rombach et al., 2021; Brown et al., 2020). let $C$ be a set containing complete information about poets, with each poet's data represented by $c \in C$. While $c$ offers a holistic view of a creator, $I[c]$ only provides a subset to prevent the AI from merely duplicating creations from $c$. Notably, $c$ includes poems linked to the given poem, whereas $I[c]$ purposely omits these. For example, $I[c]$ might include AI-generated or AI-engineer-designed synthetic data, such as imagined profiles, background details, simulated upbringing environments, cognitive patterns and experiences.

**Probability Distribution and Evaluator Function.** Define $\mathcal{D}_C$ as a probability distribution over $C$. This set, $C$, is versatile, covering all potential creators, whether they exist now, in the past, or in the future. An evaluation function $L$ determines if the generative model $q$ mirrors a particular human creator defined by the information $c$. Specifically, $L(q, c) = 0$ indicates successful emulation, while $L(q, c) = 1$ denotes a failure.

**Defining AI Creativity.** Let $\bar{c} \sim \mathcal{D}_C$ stand for a hypothetical creator, one not grounded in reality, from the creator distribution. Here, $\mathcal{M}(\bar{c}) = q(\cdot \mid I[\bar{c}])$. Given these components, we outline AI creativity as the AI model's capability to emulate a new, plausible, yet non-existent human creator by given biography of the virtual creator, evaluated by $L$. This concept is formalized in the subsequent definition of *relative creativity*:

**Definition 1** (**Relative Creativity**). An AI model, denoted as $\mathcal{M}$, achieves $\delta$-creativity (with respect to evaluator $L$ under creator distribution $\mathcal{D}_C$) if it is indistinguishable from a plausible, yet non-existent human creator to the degree where $L(\mathcal{M}(\bar{c}), \bar{c}) = 0$ with a probability of at least $1 - \delta$ for $\bar{c} \sim \mathcal{D}_C$.

**Remark 1.** Relative creativity scrutinizes the creations of the creator and the AI model's results—derived from reasoning over a human creator's biography. In doing so, we leverage the same informational foundation—the human's life history and personalized knowledge—to construct a direct comparison of creativity between the AI and the hypothetical human. Relative Creativity acknowledges the inherently subjective nature of creativity, which encompasses aspects like originality, divergent thinking, and problem-solving skills. These subjective elements are considered in the process of choosing the human benchmark for comparison.

While we can theoretically ascertain the value of $L(\mathcal{M}(\bar{c}), \bar{c})$, its practical assessment becomes intricate due to the inaccessibility of $\bar{c}$ — given that its corresponding human doesn't exist. To address this, a theoretical framework that can bridge $L(\mathcal{M}(\bar{c}), \bar{c})$ with observable real-world data is required, enabling us to empirically measure whether AI possesses the relative creativity.

# 4 Measuring Relative Creativity

In this section, we introduce a methodology to quantify AI's creative potential and establish a measurable framework. This approach evaluates whether and to what extent a model achieves Relative Creativity, basing the assessment on observable outputs of existing human creators.

## 4.1 Statistical Creativity

Consider the evaluation metric $E_0(q)$, an empirical measure designed to estimate the indistinguishability between the creative abilities of an AI model and those of human creators. The human creators are represented as $(c_i)_{i=1}^n \sim (\mathcal{D}_C)^{\otimes n}$, $c_i$ is independent and identically sampled from creator distribution $\mathcal{D}_C$. Formally, $E_0(q)$ is defined as:

$$E_0(q) = \frac{1}{n} \sum_{i=1}^n L(q(\cdot \mid I[c_i]), c_i).$$

A low $E_0(q)$ value for a generative model signifies its high resemblance to the majority of creators $c_i$. In these instances, it becomes challenging for the evaluator to distinguish between the outputs of the model $q(\cdot \mid I[c_i])$ and the human creators $c_i$.

To facilitate the adjustment from theoretical creators (as defined in Definition 1) to actual creators, whose data is readily obtainable, we present the concept of Statistical Creativity (outlined in Theorem 1). This theorem explicitly outlines the conditions under which an AI model can be classified as exhibiting $\delta$-creativity, in relation to a specific group of human creators, as assessed by evaluator $L$. Specifically, if $E_0(q) < \delta$ and we have a sufficiently large sample set for evaluation $\left(n \geq \frac{\ln(1/\delta')}{2(\delta - E_0)^2}\right)$, then the AI can be classified as $\delta$-creativity as the humans from $\mathcal{D}_C$.

**Theorem 1** (**Statistical Creativity**). *Suppose we have a positive integer $n \in \mathbb{N}_+$, and positive real numbers $\delta, t > 0$. Let $\mathcal{M}(c) = q(\cdot \mid I[c])$ be an AI model. If this model satisfies $E_0(q) < \delta$ and $n \geq \frac{\ln(1/t)}{2(\delta - E_0)^2}$, then the AI model $\mathcal{M}$ is $\delta$-creativity (w.r.t. $L$ under $\mathcal{D}_C$), with probability at least $1 - t$ over the draw of $(c_i)_{i=1}^n \sim (\mathcal{D}_C)^{\otimes n}$.*

*Proof.* The proof is presented in Appendix A. □

**Remark 2.** Moving beyond a "creative or not" categorization, statistical creativity presents a sophisticated perspective where creativity is characterized by its deviation, $\delta$, to a pre-selected existing human creators as benchmarks. Crucially, statistical creativity does not require an AI to perfectly replicate human creators. Instead, it highlights the importance of achieving a certain level of similarity that is appreciable through the lens of the evaluator, represented by $L$. When a human assumes the role of the evaluator, a creatively successful AI is expected to skillfully mimic a novel creator, as perceived from the human evaluator's perspective.

## 4.2 Measuring Statistical Creativity of Autoregressive Model

In this section, we delve into applying Statistical Creativity to autoregressive models, which are fundamental to large language models (LLMs) recognized for exhibiting a certain level of creative capabilities (Zhao et al., 2023). We specifically explore the specific form of Statistical Creativity in this context, thereby enhancing our understanding of it within next-token prediction mechanisms.

Our proposition is straightforward: if an AI can generate sequences (like poems or stories) that mirror the works of a group of human artists, it demonstrates a level of creativity comparable to that group. To quantify this assessment, we introduce the metric $E_1(q)$. This metric estimates the indistinguishability between the creative abilities of an autoregressive model and those of human creators by measuring the log-likelihood of next-token predictions across diverse pairs of creators and their creations.

$$E_1(q) = -\frac{1}{n} \sum_{i=1}^n \frac{1}{r(c_i)} \sum_{t=1}^T \log q(x_i^{(t)} \mid x_i^{(t-\omega:t-1)}, I[c_i]).$$

In this formula, the pairs $(x_i, c_i)_{i=1}^n$ are drawn independently and uniformly from $\mathcal{D}$, wherein $\mathcal{D}(x, c) = \mathcal{D}_C(c) \cdot p(x \mid c)$. The term $r(c_i) = \tau + H[p(\cdot \mid c_i)]$ is of interest, where $p(\cdot \mid c_i)$ delineates the authentic (yet unknown) distribution of creations for the creator $c_i$. And $H[p(\cdot \mid c_i)]$ signifies the entropy of $p(\cdot \mid c_i)$. The positive constant $\tau$ sets a threshold, details of which will be delved into later. Breaking it down further, the creation $x$ consists of $T$ tokens. Define $T \geq 1$ and represent $x = \{x^{(t)}\}_{t=1}^T$. The term $q(x \mid I[c_i]) = \prod_{t=1}^T q(x^{(t)} \mid x^{(t-\omega:t-1)}, I[c_i])$ is predicated on a context window size of $\omega \geq 0$. For specific conditions where $t = 1$ or $\omega = 0$, the expression simplifies to $q(x^{(t)} \mid I[c_i])$. Likewise, for $t < \omega$, the notation changes to $(t - \omega : t - 1) \triangleq (1 : t - 1)$.

**Remark 3.** It's worth noting that by setting $T = 1$, the non-auto-regressive structure is retained, preserving the metric's versatility. In addition, an intriguing aspect of $E_1(q)$ is the term $\frac{1}{r(c_i)}$. This term plays a pivotal role in ensuring that if the entropy in the creation process by a creator $c_i$ is vast, the corresponding log-likelihood receives lesser weight. This adjustment effectively captures the inherent unpredictability and diversity characterizing human creativity.

Intuitively, the evaluator $L$ is unable to discern between a synthetic AI creator $q$ and an actual creator $p$ if the KL divergence between the two remains minimal. To crystallize this concept, we frame it as the following assumption:

**Assumption 1.** *There exists a positive threshold $\tau$ such that $L(q(\cdot \mid I[c]), c) = 0$ whenever $D_{\mathrm{KL}}(p(\cdot \mid c) \parallel q(\cdot \mid I[c])) < \tau$, for any $c \in C$.*

Within this context, $\tau > 0$ stands as the threshold fulfilling the condition stipulated in Assumption 1. Furthermore, we introduce $r_{\min}$ as the least positive real number ensuring $r_{\min} \leq r(c)$ with near certainty for $c \sim \mathcal{D}_C$. The above assumption delineates the conditions under which an AI is deemed statistically creative. With the condition, the following theorem presents the measurement of statistical creativity in the next-token prediction setup.

**Theorem 2 (Autoregressive Statistical Creativity).** *Provided Assumption 1 remains valid, and for a model $\mathcal{M}$ such that its negative log likelihood $-\log q(x \mid I[c]) \leq M$ nearly always over $(x, c) \sim \mathcal{D}$, where $M \in \mathbb{R}^+$. If $n \in \mathbb{N}_+$, $\delta, t > 0$, and $\mathcal{M}$ is an AI model where $E_1 < \delta$ and $n \geq \frac{M^2 \ln(1/t)}{2 r_{\min}^2 (\delta - E_1)^2}$, then the model $\mathcal{M}$ is deemed $\delta$-creativity (with respect to $\mathcal{D}_C$), with a probability at least $1 - t$ across the sampling of $(x_i, c_i)_{i=1}^n \overset{i.i.d.}{\sim} \mathcal{D}$.*

*Proof.* The proof is presented in Appendix A. $\qquad\square$

**Remark 4.** Contrasting it with Theorem 1, Theorem 2 reveals that demonstrating creativity does not necessarily require strictly imitating each creator. Instead, it emphasizes the significance of the log-likelihood in next-token prediction across a wide range of creators. Furthermore, Theorem 2 addresses a critical question: Is it essential to have a large number of samples from each creator to demonstrate a model's statistical creativity? This inquiry is vital, considering the potential bottleneck in creative AI development due to the extensive collection of each creator's works. Theorem 2 provides insight into this issue, indicating that amassing a vast dataset for each creator is not as crucial as previously thought. What is more important, as Theorem 2 highlights, is the diversity of creator-creation pairs (i.e., $n$ in $\{(x_i, c_i)\}_{i=1}^n$). This implies that with a sufficiently large number of creators, the sample size per creator can be relatively small, thereby alleviating data collection challenges in developing creative AI models.

## 4.3 Measuring Statistical Creativity of Large Language Model

The technique of prompting stands as a potent tool for unlocking the capabilities inherent in foundation models, especially Large Language Models (LLMs) (Bommasani et al., 2022). In modeling this prompting setup, we denote $U$ as a set encompassing various contexts of creations. Each context within this set is denoted as $u$. We define $\mathcal{D}_U$ as the probability distribution over $U$. Define $\mathcal{D}_{U,C}(u, c) = \mathcal{D}_C(c) \mathcal{D}_U(u)$, where $\{u_i\}_{i=1}^n \overset{i.i.d.}{\sim} (\mathcal{D}_U)^{\otimes n}$ symbolizes a sequence of existing context prompts, assumed to be independently and

identically distributed. The notation $\bar{u} \sim \mathcal{D}_U$ corresponds to a possible new context prompt. To clarify the distinction between the variables $c$ and $u$, consider a practical example involving GPT-4 (OpenAI, 2023). Here, the creator's information, such as a biography, is represented by $c$ and serves as the system prompt, while $u$ corresponds to the user's input prompt.

To facilitate the measurement of statistical creativity in cutting-edge models, we expand our previously outlined functions to incorporate the variable $u$. Consequently, functions such as $L(\bar{q}, c)$ and $q(x|I[c])$ are updated to $L(\bar{q}, u, c)$ and $q(x|u, I[c])$, respectively. This modification is consistently applied across related definitions. For instance, $\mathcal{M}(\bar{c})$ is now redefined as $\mathcal{M}(\bar{z}) = q(\cdot|\bar{u}, I[\bar{c}])$, where $\bar{z} = (\bar{u}, \bar{c})$ represents the combined user and system prompts. Building on these modifications, we refine the definition of statistical creativity as follows:

**Definition 2.** An AI generative model, denoted as $\mathcal{M}$, is termed $\delta$-creativity (w.r.t. $L$ under $\mathcal{D}_{U,C}$), if it behaves in a manner indistinguishable from a hypothetical (yet plausible) human creator when faced with new context prompts. This is characterized by the condition $L(\mathcal{M}(\bar{z}), \bar{z}) = 0$ with a probability of at least $1 - \delta$ over the sampling of $(\bar{u}, \bar{c}) \sim \mathcal{D}_{U,C}$.

The aforementioned definition expands upon Definition 1 by integrating it into a prompting setup. Subsequently, we offer results that are analogous to those discussed in Section 4. We introduce the metric $E_2 = \frac{1}{n} \sum_{i=1}^{n} L(\mathcal{M}(z_i), z_i)$, where $(z_i)_{i=1}^{n} \sim (\mathcal{D}_{U,C})^{\otimes n}$.

**Corollary 1 (Prompt-Contextualized Statistical Creativity).** *For a given positive integer $n \in \mathbb{N}_+$, $\delta, t > 0$, and constants $\delta, t > 0$, suppose $\mathcal{M}$ is an AI model satisfying $E_2 < \delta$ and $n \geq \frac{\ln(1/t)}{2(\delta - E_2)^2}$. Then, the AI model $\mathcal{M}$ is $\delta$-creativity under $\mathcal{D}_{U,C}$ with a probability not less than $1 - t$ over the sampling of $(z_i)_{i=1}^{n} \sim (\mathcal{D}_{U,C})^{\otimes n}$.*

*Proof.* The proof is presented in Appendix A. □

**Assumption 2.** *There exists $p$ and $\tau > 0$ such that $L(q(\cdot \mid u, I[c]), u, c) = 0$ if $D_{\mathrm{KL}}(p(\cdot \mid u, c) \parallel q(\cdot \mid u, I[c])) < \tau$, applicable for any $c \in C$ and $u \in U$.*

Correspondingly, let's define $E_3 = -\frac{1}{n} \sum_{i=1}^{n} \frac{1}{r(u_i, c_i)} \sum_{t=1}^{T} \log q(x_i^{(t)} \mid x_i^{(t-\omega:t-1)}, u_i, I[c_i])$, where the triple $(x_i, u_i, c_i)_{i=1}^{n} \sim \mathcal{D}^{\otimes n}$ with $\mathcal{D}(x, u, c) = \mathcal{D}_{U,C}(u, c)p(x \mid u, c)$. Define $r(u, c) = \tau + H[p(\cdot \mid u, c)]$. Let $r_{\min}$ be a positive real number such that $r_{\min} \leq r(u, c)$ almost surely over $(u, c) \sim \mathcal{D}_{U,C}$. Then, we have the following corollary for prompt-contextualized autoregressive statistical creativity.

**Corollary 2 (Prompt-Contextualized Autoregressive Statistical Creativity).** *Let Assumption 2 hold. Let $q$ be given such that $-\log q(x \mid u, I[c]) \leq M$ almost surely over $(x, u, c) \sim \mathcal{D}$ for some $M > 0$. Let $n \in \mathbb{N}_+$, $\delta, t > 0$, and $\mathcal{M}$ be an AI model such that $E_3 < \delta$ and $n \geq \frac{M^2 \ln(1/t)}{2r_{\min}^2(\delta - E_3)^2}$. Then, the AI model $\mathcal{M}$ is $\delta$-creativity, with probability at least $1 - t$ over the draw of $(x_i, u_i, c_i)_{i=1}^{n} \sim \mathcal{D}^{\otimes n}$.*

*Proof.* The proof is presented in Appendix A. □

The corollaries presented above provide an expansion of our prior theorems, illustrating the ways in which our creativity framework can be extended and reimagined with the introduction of a prompting setup. This expansion not only enhances our framework but also lays the groundwork for a deeper understanding of the behavior of large language models.

# 5 Achieving Statistical Creativity through Training

The discussions in the preceding sections have focused on establishing framework to determine whether trained models exhibit statistical creativity. This leads to a pivotal question: How can we steer a model towards achieving creativity? This section aims to crystallize and elucidate such an approach.

We revisit and scrutinize the terms previously introduced in the context of statistical creativity, particularly focusing on terms like $E_2$ and $E_3$. The derivable nature of these terms suggests a pathway towards defining *Statistical Creativity Loss*, Equation (1). This concept encapsulates the essence of our approach, linking

the theoretical underpinnings of statistical creativity with practical measures for evaluating and enhancing AI models' creative capabilities.

$$L_{\bar{z}}[q] = L(q(\cdot \mid \bar{u}, I[\bar{c}]), \bar{z}) \tag{1}$$

For the notation, we follow the setup of Section 4.3, where $\bar{z} = (\bar{u}, \bar{c})$ encapsulating the combined user and system prompts. We consider a learning algorithm $\mathcal{A}$ and denote its output as $q_S = \mathcal{A}(S)$, given the sample set $S$. This is established under the condition that $-\log q_S(x \mid u, I[c]) \leq M$ almost surely over $v \sim \mathcal{D}$, where $S = (v_i)_{i=1}^n \sim \mathcal{D}^{\otimes n}$. Besides, for the following discussion, we define $\psi(q, v) = -\frac{1}{r(z)} \log q(x \mid u, I[c])$ where $v = (x, u, c)$.

A pivotal aspect of our discussion is centered on the concept of generalization in deep learning (Kawaguchi et al., 2022b). We introduce a placeholder notation $Q(t)$ to denote for potential overfitting of the model to its training dataset. The choice to employ a placeholder is strategic, given that the exact metrics and quantifications of overfitting—captured by $Q(t)$—are at the forefront of ongoing research in the realm of generalization bounds. Note, this modular approach allows for easy integration of future advancements in the study of generalization into our framework, by simply updating the placeholder, $Q(t)$. For clarity, we will also present some potential manifestations of $Q(t)$ based on the prevailing understanding from the study of deep learning generalization. Concretely, for any $t > 0$, we define $Q(t)$ such that with probability at least $1 - t$ over an draw of $S = (v_i)_{i=1}^n \sim \mathcal{D}^{\otimes n}$, the following inequality holds:

$$\mathbb{E}_v[\psi(q_S, v)] - \frac{1}{m} \sum_{i=1}^m \psi(q_S, v_i) \leq \tilde{\mathcal{O}}\left(\sqrt{\frac{Q(t)}{n}}\right).$$

For deep neural networks, several foundational works have explored their complexity. For instance, with the Rademacher complexity bound (Bartlett and Mendelson, 2002; Mohri et al., 2012; Shalev-Shwartz and Ben-David, 2014), the size-independent sample complexity of neural networks (Golowich et al., 2018) provides

$$Q(t) = B^2 \rho \prod_{j=1}^{\rho} M_F^2(j) + \ln(1/t),$$

where $\rho$ is the number of layers of the network, $B$ is the upper bound on the Euclidean norm of the network, and $M_F(j)$ is the upper bound on the Frobenius norm of the weight matrix at $j$-th layer.

In addition, with the recent result in information theory (Kawaguchi et al., 2023), we have,

$$Q(t) = \min_{j \in \{1, \ldots, \rho\}} I(X; Z_j|Y) + I(S; \Theta_j) + \ln(1/t),$$

where $I(X; Z|Y)$ is the conditional mutual information between the network input $X$ and the output $Z_j$ of $j$-th layer, conditioned on the target label $Y$. Here, $I(S; \phi_j)$ is the mutual information between $S$ and the set of weight parameters from the first layer to $j$-th layer.

Besides, the sample complexity based on robustness (Xu and Mannor, 2012; Kawaguchi et al., 2022a) indicates,

$$Q(t) = c(S)^2 + \mathcal{N}((2\sqrt{n})^{-1}, \mathcal{V}, \kappa) + \ln(1/t),$$

where $\mathcal{V}$ designates a chosen compact space of $v$ relative to metric $\kappa$. The constant $c(S)$ is Lipschitz and defined as: $|\psi(q_S, v) - \psi(q_S, v')| \leq c(S)\kappa(v, v')$ for any pair $(v, v') \in \mathcal{V}$. Notably, $\mathcal{N}((2\sqrt{n})^{-1}, \mathcal{V}, \kappa)$ is the $\epsilon$-covering number of $\mathcal{V}$ (as detailed in Definition 1 of Xu and Mannor, 2012).

Moreover, for any model configuration, if we chose $\mathcal{A} : \mathcal{S} \to \mathcal{Q}$ with $|\mathcal{Q}| < \infty$, an use of concentration inequalities (e.g., Kawaguchi et al., 2022b) provides

$$Q(t) = \ln\left(\frac{|\mathcal{Q}|}{t}\right).$$

Different from the traditional discussions on model generalization, which often focus on metrics like classification accuracy, Corollary 3 studies the creative performance of models trained using the Statistical Creativity Loss, Equation (1). This shift in focus underscores our commitment to exploring and understanding the unique aspects of creativity in AI models.

**Corollary 3.** *Let Assumption 2 hold. Then, with probability at least $1 - \delta$ over an draw of $S = (v_i)_{i=1}^n \sim \mathcal{D}^{\otimes n}$ and $\bar{z} \sim \mathcal{D}_{U,C}$,*

$$L_{\bar{z}}[q_S] < 2\delta^{-1}E[q_S] + \mathcal{O}\left(\sqrt{\frac{\delta^{-2}Q(\delta/2)}{n}}\right) \leq \frac{2}{\delta r_{\min}}\tilde{E}[q_S] + \mathcal{O}\left(\sqrt{\frac{\delta^{-2}Q(\delta/2)}{n}}\right). \tag{2}$$

*Proof.* The proof is presented in Appendix A. □

**Remark 5.** Corollary 3 establishes that $L_{\bar{z}}[q_S]$ is constrained by $2\delta^{-1}E[q_S] + \mathcal{O}(\sqrt{\frac{Q(\delta/2)}{n}})$ (alternatively, $\frac{2}{\delta r_{\min}}\tilde{E}[q_S] + \mathcal{O}(\sqrt{\frac{\delta^{-2}Q(\delta/2)}{n}})$), where the first term minimizes with the training objective $E[q_S]$ (or $\tilde{E}[q_S]$), and the latter by increasing the sample size $n$. This corollary underscores that to achieve statistical creativity, one can minimize training objectives like $E[q_S]$ or $\tilde{E}[q_S]$ while simultaneously augmenting the training sample size $n$.

# 6 Related Works

In this section, we provide a concise overview of previous definitions on creativity and applications in language and vision.

## 6.1 Definitions of Creativity

Research into creativity has spanned several decades, with the topic remaining a central point of interest in psychology, cognitive science, and philosophy. While there has been an abundance of research on the subject, there is no universally agreed upon definition. Indeed, literature showcases an array of over a hundred proposed definitions, underscoring the multifaceted essence of creativity (Aleinikov et al., 2000; Treffinger, 1996; Boden, 2003; Elgammal et al., 2017). Within the field of computational creativity, a commonly accepted framework for assessing creativity uses the "four P's": namely, *person* (the creator of a work), *press* (the environmental context for a work), *process* (how a work is created), and *product* (the work itself) (Jordanous, 2016). Boden (2003) presents a triadic criterion for gauging machine creativity, highlighting artifacts or ideas that are "new, surprising, and valuable." Building upon this foundation, Boden (2003) identifies three nuanced forms of creativity: Combinatorial, Exploratory, and Transformational. This classification melds the *process* and *product* dimensions of the four P's in order to focus on the production of surprising or novel outputs. The previous works consider *absolute* creativity at the level of each creation given each environment without mathematical concreteness. We depart from previous work here. Instead of considering *absolute* creativity at the creation level, we consider the *relative* creativity at the meta-level of parallel universes consisting of humans and AI algorithms. This allows us to avoid the question of what is creative in an absolute sense, and thus enables us to rigorously define objectives and desired properties via concrete mathematical formula instead of English descriptions.

## 6.2 Applications in Vision and Language

The field of creative image generative models has seen significant growth, raising questions regarding machines' ability to produce creative art. Hertzmann (2018) delved into this, highlighting intersections between computer graphics and artistic innovation. Xu et al. (2012a) introduced creative 3D modeling that aligns with user preferences while ensuring variety. Generative Adversarial Networks (GANs) (Goodfellow et al., 2014), used by Elgammal et al. (2017), drive the creation of distinctive artistic styles by maximizing deviations from known styles. Sbai et al. (2018) further this deviation by encouraging models to differ from training set styles. The perspective of creative generation as composition is evident in works by Ge et al. (2021) and Ranaweera (2016), emphasizing the integration of detailed elements. Vinker et al. (2023) build on this by fragmenting personalized concepts into visual elements for innovative reassembly, enriching creative output.

Parallel to advancements in vision, the evolution of language models has spurred inquiries into optimizing data usage to bolster their adaptability across various domains, tasks, and languages (Gururangan et al., 2020; Devlin et al., 2019; Conneau et al., 2020). Researchers have shown an inclination to deploy language models in deciphering nuances in human communication (Schwartz et al., 2013; Wu et al., 2022). This insight is further harnessed to refine classification models (Hovy, 2015; Flek, 2020). With the increasing popularity of generative models, there has also been an interest in controllable text generation, in which a model's output must satisfy constraints such as politeness (Saha et al., 2022; Sennrich et al., 2016), sentiment (Liu et al., 2021; Dathathri et al., 2019; He et al., 2020), or other stylistic constraints. Finally, text style transfer (TST), in which the goal is to transform the style of an input text to a set goal style, has also become a popular task. Style may refer to a range of different text and author specific features including politeness (Madaan et al., 2020), formality (Rao and Tetreault, 2018; Briakou et al., 2021), simplicity (Zhu et al., 2010; van den Bercken et al., 2019; Weng et al., 2019; Cao et al., 2020), author (Xu et al., 2012b; Carlson et al., 2018), author gender (Prabhumoye et al., 2018), and more (Jin et al., 2022). While all these applications seek to apply elements of creativity in generative models, none directly defines creativity or seeks to directly optimize for it. Instead, the focus is on improving models' ability to excel at predefined tasks as proxies for creativity. In contrast, our research pivots on establishing a theoretical foundation for creativity. This framework of creativity naturally includes preceding insights on diversity and the quality of generation. We anticipate that our contributions will lay the groundwork for future endeavors, guiding the enhancement of model creativity.

# 7    Conclusion

In this paper, we have embarked on an in-depth exploration of creativity within the realm of artificial intelligence. Our journey began with the introduction of Relative Creativity, a paradigm-shifting concept that redefines creativity assessment in AI systems from a comparative standpoint, moving away from absolute standards. This approach not only recognizes the inherent subjectivity in creative processes but also cleverly integrates it, drawing inspiration from the Turing Test's comparative method of assessing intelligence. Progressing further, we introduced the concept of Statistical Creativity, which is instrumental in bridging the gap between theoretical constructs and empirical evaluation. By focusing on whether AI can mirror the creative output of specific human groups, this concept enables the quantifiable assessment of AI creativity and enhances the practical applicability of our theoretical framework. A significant stride was made with the application of these concepts to autoregressive models. We developed a practical measure for assessing statistical creativity in these models, particularly with next-token prediction. This measure's adaptability to the prompting paradigm ensures its relevance and applicability to contemporary AI models, illustrating our commitment to keeping pace with technological advancements. Moreover, the introduction of Statistical Creativity Loss marks a cornerstone in our theoretical development. This concept provides invaluable insights for optimizing training objectives and dataset sizes, serving as a guidepost for achieving statistical creativity within a robust theoretical framework. In conclusion, our contributions lay a solid foundation for future explorations into AI creativity. We have not only advanced the theoretical discourse but also provided practical tools and methodologies for assessing and enhancing AI's creative potential. Our work stands as a testament to the dynamic nature of AI research, and we aspire for it to act as a guiding beacon in the ongoing discourse on AI creativity, encouraging relative evaluations and empirical research. This paper, therefore, serves as both a theoretical contribution and a practical guide for those seeking to harness and understand the creative capacities of AI models.

**Limitations.** This paper primarily focuses on conceptualizing Relative Creativity and proposing a methodology to assess the extent to which an AI model achieves this type of creativity. The development of a comprehensive dataset and the establishment of a benchmark for creativity evaluation are forthcoming endeavors. Additionally, we will test current advanced AI models to evaluate their creative performance. As we progress, we anticipate integrating our methods into the existing evaluation benchmark toolkit, align with the evolving landscape of AI technologies and assessment.

# References

Aleinikov, A. G., Kackmeister, S., and Koenig, R. (2000). *Creating Creativity: 101 Definitions (what Webster Never Told You)*. Alden B. Dow Creativity Center Press, Midland, MI.

Amabile, T. M. (1996). *Creativity and innovation in organizations*, volume 5. Harvard Business School Boston.

Bartlett, P. L. and Mendelson, S. (2002). Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482.

Boden, M. A. (2003). *The Creative Mind: Myths and Mechanisms*. Routledge, London, UK.

Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N., Chen, A., Creel, K., Davis, J. Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K., Goodman, N., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D. E., Hong, J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P. W., Krass, M., Krishna, R., Kuditipudi, R., Kumar, A., Ladhak, F., Lee, M., Lee, T., Leskovec, J., Levent, I., Li, X. L., Li, X., Ma, T., Malik, A., Manning, C. D., Mirchandani, S., Mitchell, E., Munyikwa, Z., Nair, S., Narayan, A., Narayanan, D., Newman, B., Nie, A., Niebles, J. C., Nilforoshan, H., Nyarko, J., Ogut, G., Orr, L., Papadimitriou, I., Park, J. S., Piech, C., Portelance, E., Potts, C., Raghunathan, A., Reich, R., Ren, H., Rong, F., Roohani, Y., Ruiz, C., Ryan, J., Ré, C., Sadigh, D., Sagawa, S., Santhanam, K., Shih, A., Srinivasan, K., Tamkin, A., Taori, R., Thomas, A. W., Tramèr, F., Wang, R. E., Wang, W., Wu, B., Wu, J., Wu, Y., Xie, S. M., Yasunaga, M., You, J., Zaharia, M., Zhang, M., Zhang, T., Zhang, X., Zhang, Y., Zheng, L., Zhou, K., and Liang, P. (2022). On the opportunities and risks of foundation models.

Briakou, E., Lu, D., Zhang, K., and Tetreault, J. (2021). Olá, bonjour, salve! XFORMAL: A benchmark for multilingual formality style transfer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3199–3216, Online. Association for Computational Linguistics.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Câmara Pereira, F. (2007). *Creativity and artificial intelligence: a conceptual blending approach*. Mouton de Gruyter.

Cao, Y., Shui, R., Pan, L., Kan, M.-Y., Liu, Z., and Chua, T.-S. (2020). Expertise style transfer: A new task towards better communication between experts and laymen. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1061–1071, Online. Association for Computational Linguistics.

Carlson, K., Riddell, A., and Rockmore, D. (2018). Evaluating prose style transfer with the bible. *Royal Society open science*, 5(10):171920.

Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., and Stoyanov, V. (2020). Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Dathathri, S., Madotto, A., Lan, J., Hung, J., Frank, E., Molino, P., Yosinski, J., and Liu, R. (2019). Plug and play language models: A simple approach to controlled text generation. *arXiv preprint arXiv:1912.02164*.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Elgammal, A., Liu, B., Elhoseiny, M., and Mazzone, M. (2017). CAN: Creative Adversarial Networks, Generating "Art" by Learning About Styles and Deviating from Style Norms. arXiv:1706.07068 [cs].

Flek, L. (2020). Returning the N to NLP: Towards Contextually Personalized Classification Models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7828–7838, Online. Association for Computational Linguistics.

Foster, D. (2022). *Generative deep learning*. " O'Reilly Media, Inc.".

Franceschelli, G. and Musolesi, M. (2023). On the creativity of large language models. *arXiv preprint arXiv:2304.00008*.

Ge, S., Goswami, V., Zitnick, L., and Parikh, D. (2021). Creative sketch generation. In *International Conference on Learning Representations*.

Golowich, N., Rakhlin, A., and Shamir, O. (2018). Size-independent sample complexity of neural networks. In *Conference On Learning Theory*, pages 297–299. PMLR.

Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. *Advances in neural information processing systems*, 27.

Gururangan, S., Marasović, A., Swayamdipta, S., Lo, K., Beltagy, I., Downey, D., and Smith, N. A. (2020). Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360.

He, J., Wang, X., Neubig, G., and Berg-Kirkpatrick, T. (2020). A probabilistic formulation of unsupervised text style transfer. *arXiv preprint arXiv:2002.03912*.

Hertzmann, A. (2018). Can computers create art? In *Arts*, volume 7, page 18. MDPI.

Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, pages 13–30.

Hovy, D. (2015). Demographic Factors Improve Classification Performance. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 752–762, Beijing, China. Association for Computational Linguistics.

Jin, D., Jin, Z., Hu, Z., Vechtomova, O., and Mihalcea, R. (2022). Deep Learning for Text Style Transfer: A Survey. *Computational Linguistics*, 48(1):155–205.

Jordanous, A. (2016). Four PPPPerspectives on Computational Creativity in Theory and in Practice. *Connection Science*, 28(2):294–216.

Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A., Potapenko, A., et al. (2021). Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589.

Kawaguchi, K., Deng, Z., Ji, X., and Huang, J. (2023). How does information bottleneck help deep learning? In *International Conference on Machine Learning (ICML)*.

Kawaguchi, K., Deng, Z., Luh, K., and Huang, J. (2022a). Robustness implies generalization via data-dependent generalization bounds. In *International Conference on Machine Learning*, pages 10866–10894. PMLR.

Kawaguchi, K., Kaelbling, L. P., and Bengio, Y. (2022b). Generalization in deep learning. *Mathematical Aspects of Deep Learning*, pages 112–148.

Kirkpatrick, K. (2023). Can ai demonstrate creativity? *Communications of the ACM*, 66(2):21–23.

Liu, R., Jia, C., Wei, J., Xu, G., Wang, L., and Vosoughi, S. (2021). Mitigating political bias in language models through reinforced calibration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14857–14866.

Madaan, A., Setlur, A., Parekh, T., Poczos, B., Neubig, G., Yang, Y., Salakhutdinov, R., Black, A. W., and Prabhumoye, S. (2020). Politeness transfer: A tag and generate approach. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1869–1881, Online. Association for Computational Linguistics.

Mohri, M., Rostamizadeh, A., and Talwalkar, A. (2012). *Foundations of machine learning*. MIT press.

OpenAI (2023). Gpt-4 technical report.

Prabhumoye, S., Tsvetkov, Y., Salakhutdinov, R., and Black, A. W. (2018). Style transfer through back-translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 866–876, Melbourne, Australia. Association for Computational Linguistics.

Ranaweera, W. L. (2016). Exquimo: An exquisite corpse tool for co-creative 3d shape modeling.

Rao, S. and Tetreault, J. (2018). Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 129–140, New Orleans, Louisiana. Association for Computational Linguistics.

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. (2021). High-resolution image synthesis with latent diffusion models.

Runco, M. A. and Jaeger, G. J. (2012). The standard definition of creativity. *Creativity research journal*, 24(1):92–96.

Saha, P., Singh, K., Kumar, A., Mathew, B., and Mukherjee, A. (2022). Countergedi: A controllable approach to generate polite, detoxified and emotional counterspeech. In Raedt, L. D., editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 5157–5163. International Joint Conferences on Artificial Intelligence Organization. AI for Good.

Sawyer, K. (2012). Extending sociocultural theory to group creativity. *Vocations and Learning*, 5(1):59–75.

Sbai, O., Elhoseiny, M., Bordes, A., LeCun, Y., and Couprie, C. (2018). Design: Design inspiration from generative networks. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0.

Schwartz, H. A., Eichstaedt, J. C., Kern, M. L., Dziurzynski, L., Ramones, S. M., Agrawal, M., Shah, A., Kosinski, M., Stillwell, D., Seligman, M. E. P., and Ungar, L. H. (2013). Personality, Gender, and Age in the Language of Social Media: The Open-Vocabulary Approach. *PLOS ONE*, 8(9):e73791. Publisher: Public Library of Science.

Sennrich, R., Haddow, B., and Birch, A. (2016). Controlling politeness in neural machine translation via side constraints. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 35–40, San Diego, California. Association for Computational Linguistics.

Shalev-Shwartz, S. and Ben-David, S. (2014). *Understanding machine learning: From theory to algorithms*. Cambridge university press.

Somepalli, G., Singla, V., Goldblum, M., Geiping, J., and Goldstein, T. (2023). Diffusion art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6048–6058.

Terwiesch, C. and Ulrich, K. (2023). M.b.a. students vs. chatgpt: Who comes up with more innovative ideas? *The Wall Street Journal*.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. (2023). Llama: Open and efficient foundation language models.

Treffinger, D. J. (1996). *Creativity, Creative Thinking, and Critical Thinking: In Search of Definitions*. Center for Creative Learning, Sarasota, FL.

Turing, A. M. (2009). *Computing machinery and intelligence*. Springer.

van den Bercken, L., Sips, R.-J., and Lofi, C. (2019). Evaluating neural text simplification in the medical domain. In *WWW'19 The World Wide Web Conference (WWW)*, pages 3286–3292, United States. Association for Computing Machinery (ACM). WWW 2019 : The Web Conference 2019, 30 years of the web, WWW'19 ; Conference date: 13-05-2019 Through 17-05-2019.

Vinker, Y., Voynov, A., Cohen-Or, D., and Shamir, A. (2023). Concept decomposition for visual exploration and inspiration. *arXiv preprint arXiv:2305.18203*.

Weng, W.-H., Chung, Y.-A., and Szolovits, P. (2019). Unsupervised clinical language translation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '19, page 3121–3131, New York, NY, USA. Association for Computing Machinery.

Wu, Y., Suchanek, F., Vasilescu, I., Lamel, L., and Adda-Decker, M. (2022). Using a knowledge base to automatically annotate speech corpora and to identify sociolinguistic variation. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, Marseille, France. European Language Resources Association.

Xu, H. and Mannor, S. (2012). Robustness and generalization. *Machine learning*, 86(3):391–423.

Xu, K., Zhang, H., Cohen-Or, D., and Chen, B. (2012a). Fit and diverse: Set evolution for inspiring 3d shape galleries. *ACM Transactions on Graphics (TOG)*, 31(4):1–10.

Xu, W., Ritter, A., Dolan, W. B., Grishman, R., and Cherry, C. (2012b). Paraphrasing for style. In *Proceedings of COLING 2012*, pages 2899–2914.

Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., Du, Y., Yang, C., Chen, Y., Chen, Z., Jiang, J., Ren, R., Li, Y., Tang, X., Liu, Z., Liu, P., Nie, J., and rong Wen, J. (2023). A survey of large language models. *ArXiv*, abs/2303.18223.

Zhu, Z., Bernhard, D., and Gurevych, I. (2010). A monolingual tree-based translation model for sentence simplification. In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pages 1353–1361, Beijing, China. Coling 2010 Organizing Committee.

# Appendix A  Proofs

## A.1  Proof of Theorem 1

*Proof.* Since $L(q(\cdot \mid I[c]), c) \geq 0$, by using Markov's inequality, it holds that with probability at least $1 - \delta$ over the draw of $\bar{c} \sim \mathcal{D}$,

$$L(q(\cdot \mid I[\bar{c}]), \bar{c}) < \delta^{-1} \mathbb{E}_{c \sim \mathcal{D}_C}[L(q(\cdot \mid I[c]), c)].$$

Since $L(q(\cdot \mid I[c]), c) \in \{0, 1\} \subset [0, 1]$, via Hoeffding's inequality (Hoeffding, 1963), it holds that with probability at least $1 - t$ over the draw of $(c_i)_{i=1}^n \sim (\mathcal{D}_C)^{\otimes n}$,

$$\mathbb{E}_{c \sim \mathcal{D}_C}[L(q(\cdot \mid I[c]), c)] < \frac{1}{n} \sum_{i=1}^n L(q(\cdot \mid I[c_i]), c_i) + \sqrt{\frac{\ln(1/t)}{2n}}$$

Thus, with probability at least $1 - t$ over the draw of $(c_i)_{i=1}^n \sim (\mathcal{D}_C)^{\otimes n}$, it holds that with probability at least $1 - \delta$ over the draw of $\bar{c} \sim \mathcal{D}$,

$$L(q(\cdot \mid I[\bar{c}]), \bar{c}) < \delta^{-1} E_0 + \delta^{-1} \sqrt{\frac{\ln(1/t)}{2n}},$$

where $E_0 = \frac{1}{n} \sum_{i=1}^n L(q(\cdot \mid I[c_i]), c_i)$. Here, we have that if $n \geq \frac{\ln(1/t)}{2(\delta - E_0)^2}$ and $\delta - E_0 > 0$,

$$\delta^{-1} E_0 + \delta^{-1} \sqrt{\frac{\ln(1/t)}{2n}} \leq 1.$$

Thus, if $n \geq \frac{\ln(1/t)}{2(\delta - E_0)^2}$ and $\delta - E_0 > 0$, with probability at least $1 - t$ over the draw of $(c_i)_{i=1}^n \sim (\mathcal{D}_C)^{\otimes n}$, it holds that with probability at least $1 - \delta$ over the draw of $\bar{c} \sim \mathcal{D}$,

$$L(q(\cdot \mid I[\bar{c}]), \bar{c}) < 1.$$

Since $L(q(\cdot \mid I[\bar{c}]), \bar{c}) \in \{0, 1\}$, then $L(q(\cdot \mid I[\bar{c}]), \bar{c}) = 0$. □

## A.2  Proof of Theorem 2

*Proof.* From Assumption 1,

$$\mathbb{E}_{c \sim \mathcal{D}_C}[L(q(\cdot \mid I[c]), c)] \leq \mathbb{E}_{c \sim \mathcal{D}_C}[\mathbb{1}\{D_{\mathrm{KL}}(p(\cdot \mid c) \parallel q(\cdot \mid I[c])) \geq \tau\}].$$

We bound the right-hand side of this inequality by using a function similar to hinge loss with a slope as follow: for any $r(c) > 0$,

$$\mathbb{E}_{c \sim \mathcal{D}_C}[L(q(\cdot \mid I[c]), c)] = \mathbb{E}_{c \sim \mathcal{D}_C}[\mathbb{1}\{\tau - D_{\mathrm{KL}}(p(\cdot \mid c) \parallel q(\cdot \mid I[c])) \leq 0\}]$$
$$\leq \mathbb{E}_{c \sim \mathcal{D}_C}\left[\max\left(0, 1 - \frac{1}{r(c)}(\tau - D_{\mathrm{KL}}(p(\cdot \mid c) \parallel q(\cdot \mid I[c])))\right)\right].$$

Here, from the definition of the KL divergence,

$$D_{\mathrm{KL}}(p(\cdot \mid c) \parallel q(\cdot \mid I[c])) = \mathbb{E}_{x \sim p(\cdot \mid c)}\left[\log \frac{p(x \mid c)}{q(x \mid I[c])}\right] = \mathbb{E}_{x \sim p(\cdot \mid c)}[\log p(x \mid c)] - \mathbb{E}_{x \sim p(\cdot \mid c)}[\log q(x \mid I[c])].$$

Substituting this into the above inequality,

$$\mathbb{E}_{c \sim \mathcal{D}_C}[L(q(\cdot \mid I[c]), c)] \leq \mathbb{E}_{c \sim \mathcal{D}_C}\left[\max\left(0, 1 - \frac{\tau}{r(c)} + \frac{D_{\mathrm{KL}}(p(\cdot \mid c) \parallel q(\cdot \mid I[c]))}{r(c)}\right)\right]$$

$$= \mathbb{E}_{c \sim \mathcal{D}_C} \left[ \max \left( 0, A + B \right) \right],$$

where $A = 1 - \frac{\tau - \mathbb{E}_{x \sim p(\cdot|c)} [\log p(x|c)]}{r(c)}$ and $B = -\frac{\mathbb{E}_{x \sim p(\cdot|c)} [\log q(x|I[c])]}{r(c)}$. Here, we notice that $B \geq 0$ because $r(c) > 0$ and $-\log q(x \mid I[c]) \geq 0$ (here we use the fact that $q(\cdot \mid I[c])$ is the probability mass function instead of density). By using the fact of $B \geq 0$ and the definition of maximum operator, we expand the cases of different values of $A$ and $B$ as

$$\max \left( 0, A + B \right) = \begin{cases} A + B & \text{if } A + B \geq 0 \\ 0 & \text{if } A + B < 0 \end{cases}$$

$$= \begin{cases} A + B & \text{if } A \geq 0 \\ A + B & \text{if } -B \leq A < 0 \\ 0 & \text{if } A < -B \end{cases}$$

$$\leq \begin{cases} A + B & \text{if } A \geq 0 \\ B & \text{if } -B \leq A < 0 \\ B & \text{if } A < -B \end{cases}$$

where the last line follows from the fact that $A + B \leq B$ if $-B \leq A < 0$ and $0 \leq B$. Since the last two cases output the same value $B$ whenever $A < 0$, we can simplify this expression as

$$\max \left( 0, A + B \right) \leq \begin{cases} A + B & \text{if } A \geq 0 \\ B & \text{if } A < 0 \end{cases}$$

$$= \max(0, A) + B.$$

Substituting this into the above inequality of $\mathbb{E}_{c \sim \mathcal{D}_C} [L(q(\cdot \mid I[c]), c)] \leq \mathbb{E}_{c \sim \mathcal{D}_C} \left[ \max \left( 0, A + B \right) \right]$,

$$\mathbb{E}_{c \sim \mathcal{D}_C} [L(q(\cdot \mid I[c]), c)] \leq \mathbb{E}_{c \sim \mathcal{D}_C} [\max(0, A)] + \mathbb{E}_{c \sim \mathcal{D}_C} [B].$$

We now set $r(c) > 0$ to remove the 1st term on the right-hand side of this inequality as

$$A \leq 0 \iff 1 - \frac{\tau - \mathbb{E}_{x \sim p(\cdot|c)} [\log p(x \mid c)]}{r(c)} \leq 0$$

$$\iff r(c) \leq \tau - \mathbb{E}_{x \sim p(\cdot|c)} [\log p(x \mid c)].$$

Here, notice that $-\mathbb{E}_{x \sim p(\cdot|c)} [\log p(x \mid c)] = H[p(\cdot \mid c)]$, which is the entropy of $p(\cdot \mid c)$. Thus, by setting $r(c) = \tau + H[p(\cdot \mid c)]$, we have that $\mathbb{E}_{c \sim \mathcal{D}_C} [\max(0, A)] = 0$ and thus,

$$\mathbb{E}_{c \sim \mathcal{D}_C} [L(q(\cdot \mid I[c]), c)] \leq \mathbb{E}_{c \sim \mathcal{D}_C} [B] = \mathbb{E}_{c \sim \mathcal{D}_C} \left[ -\frac{1}{r(c)} \mathbb{E}_{x \sim p(\cdot|c)} [\log q(x \mid I[c])] \right]$$

$$= \mathbb{E}_{c \sim \mathcal{D}_C} \mathbb{E}_{x \sim p(\cdot|c)} \left[ -\frac{1}{r(c)} \log q(x \mid I[c]) \right]$$

$$= \mathbb{E}_{(x,c) \sim \mathcal{D}} \left[ -\frac{1}{r(c)} \log q(x \mid I[c]) \right],$$

where $\mathcal{D}(x, c) = \mathcal{D}_C(c) p(x \mid c)$. Here, since $-\log q(x \mid I[c]) \in [0, M]$ and $\frac{1}{r(c)} \in [0, \frac{1}{r_{\min}}]$ almost surely, we have that $-\frac{1}{r(c)} \log q(x \mid I[c]) \in [0, \frac{M}{r_{\min}}]$ almost surely. (The discussion regarding the upper bound of the negative log likelihood, $-\log q$, is detailed in Appendix B.2.) Thus, by using Hoeffding's inequality (Hoeffding, 1963), it holds that with probability at least $1 - t$ over the draw of $(x_i, c_i)_{i=1}^n \sim \mathcal{D}^{\otimes n}$,

$$\mathbb{E}_{(x,c) \sim \mathcal{D}} \left[ -\frac{1}{r(c)} \log q(x \mid I[c]) \right] < -\frac{1}{n} \sum_{i=1}^n \frac{1}{r(c_i)} \log q(x_i \mid I[c_i]) + \frac{M}{r_{\min}} \sqrt{\frac{\ln(1/t)}{2n}}.$$

16

Combining with above two inequalities,

$$\mathbb{E}_{c \sim \mathcal{D}_C}[L(q(\cdot \mid I[c]), c)] < -\frac{1}{n} \sum_{i=1}^{n} \frac{1}{r(c_i)} \log q(x_i \mid I[c_i]) + \frac{M}{r_{\min}} \sqrt{\frac{\ln(1/t)}{2n}}.$$

Since the log of the products is the sum of logs, using $q(x \mid I[c]) = \prod_{t=1}^{T} q(x^{(t)} \mid x^{(t-\omega:t-1)}, I[c])$,

$$\mathbb{E}_{c \sim \mathcal{D}_C}[L(q(\cdot \mid I[c]), c)] < -\frac{1}{n} \sum_{i=1}^{n} \frac{1}{r(c_i)} \sum_{t=1}^{T} \log q(x_i^{(t)} \mid x_i^{(t-\omega:t-1)}, I[c_i]) + \frac{M}{r_{\min}} \sqrt{\frac{\ln(1/t)}{2n}}.$$

Since $L(q(\cdot \mid I[c]), c) \geq 0$, by using Markov's inequality, it holds that with probability at least $1 - \delta$ over the draw of $\bar{c} \sim \mathcal{D}$, $L(q(\cdot \mid I[\bar{c}]), \bar{c}) < \delta^{-1} \mathbb{E}_{c \sim \mathcal{D}_C}[L(q(\cdot \mid I[c]), c)]$. Thus, with probability at least $1 - t$ over the draw of $(x_i, c_i)_{i=1}^{n} \sim \mathcal{D}^{\otimes n}$, it holds that with probability at least $1 - \delta$ over the draw of $\bar{c} \sim \mathcal{D}$,

$$L(q(\cdot \mid I[\bar{c}]), \bar{c}) < \delta^{-1} E_1 + \frac{\delta^{-1} M}{r_{\min}} \sqrt{\frac{\ln(1/t)}{2n}}.$$

where $E_1 = -\frac{1}{n} \sum_{i=1}^{n} \frac{1}{r(c_i)} \sum_{t=1}^{T} \log q(x_i^{(t)} \mid x_i^{(t-\omega:t-1)}, I[c_i])$. Here, we have that if $n \geq \frac{M^2 \ln(1/t)}{2r_{\min}^2(\delta - E_1)^2}$ and $\delta - E_1 > 0$,

$$\delta^{-1} E_1 + \frac{\delta^{-1} M}{r_{\min}} \sqrt{\frac{\ln(1/t)}{2n}} \leq 1.$$

Thus, if $n \geq \frac{M^2 \ln(1/t)}{2r_{\min}^2(\delta - E_1)^2}$ and $\delta - E_1 > 0$, with probability at least $1 - t$ over the draw of $(x_i, c_i)_{i=1}^{n} \sim \mathcal{D}^{\otimes n}$, it holds that with probability at least $1 - \delta$ over the draw of $\bar{c} \sim \mathcal{D}$,

$$L(q(\cdot \mid I[\bar{c}]), \bar{c}) < 1.$$

Since $L(q(\cdot \mid I[\bar{c}]), \bar{c}) \in \{0, 1\}$, then $L(q(\cdot \mid I[\bar{c}]), \bar{c}) = 0$.

$\square$

## A.3   Proof of Corollary 1

*Proof.* Following all the proof steps of Theorem 1 while replacing $c$ and $\mathcal{D}_C$ by $z = (u, c)$ and $\mathcal{D}_{U,C}$, we have the following. If $E_2 < \delta$ and $n \geq \frac{\ln(1/t)}{2(\delta - E_2)^2}$, with probability at least $1 - t$ over the draw of $(u_i, c_i)_{i=1}^{n} \sim (\mathcal{D}_{U,C})^{\otimes n}$, it holds that with probability at least $1 - \delta$ over the draw of $(\bar{u}, \bar{c}) \sim \mathcal{D}_{U,C}$,

$$L(\mathcal{M}(\bar{z}), \bar{z}) < \delta^{-1} \mathbb{E}_{z \sim \mathcal{D}_{U,C}}[L(\mathcal{M}(z), z)] < \frac{\delta^{-1}}{n} \sum_{i=1}^{n} L(\mathcal{M}(z_i), z_i) + \delta^{-1} \sqrt{\frac{\ln(1/t)}{2n}} \leq 1.$$

This implies the desired statement.

$\square$

## A.4   Proof of Corollary 2

*Proof.* Following all the proof steps of Theorem 2 while replacing $z$ and $\mathcal{D}_C$ by $z = (u, c)$ and $\mathcal{D}_{U,C}$, we have the following. By setting $r(z) = \tau + H[p(\cdot \mid z)]$, if $E_3 < \delta$ and $n \geq \frac{M^2 \ln(1/t)}{2r_{\min}^2(\delta - E_3)^2}$, with probability at least $1 - t$ over the draw of $(x_i, u_i, c_i)_{i=1}^{n} \sim \mathcal{D}^{\otimes n}$, it holds that with probability at least $1 - \delta$ over the draw of $(\bar{u}, \bar{c}) \sim \mathcal{D}_{U,C}$,

$$L(\mathcal{M}(\bar{z}), \bar{z}) < \delta^{-1} \mathbb{E}_{z \sim \mathcal{D}_{U,C}}[L(\mathcal{M}(z), z)]$$
$$\leq \delta^{-1} \mathbb{E}_{z \sim \mathcal{D}_{U,C}} \left[ \max\left(0, 1 - \frac{1}{r(z)}(\tau - D_{\mathrm{KL}}(p(\cdot \mid z) \| \mathcal{M}(z))) \right) \right]$$

$$\leq \delta^{-1} \mathbb{E}_{(x,u,c) \sim \mathcal{D}} \left[ -\frac{1}{r(z)} \log q(x \mid u, I[c]) \right],$$

$$\leq -\frac{\delta^{-1}}{n} \sum_{i=1}^{n} \frac{1}{r(z_i)} \log q(x_i \mid u_i, I[c_i]) + \frac{\delta^{-1} M}{r_{\min}} \sqrt{\frac{\ln(1/t)}{2n}}$$

$$= -\frac{\delta^{-1}}{n} \sum_{i=1}^{n} \frac{1}{r(z_i)} \sum_{t=1}^{T} \log q(x_i^{(t)} \mid x_i^{(t-\omega:t-1)}, u_i, I[c_i]) + \frac{\delta^{-1} M}{r_{\min}} \sqrt{\frac{\ln(1/t)}{2n}} \leq 1.$$

This implies the desired statement. $\qquad \square$

## A.5 Proof of Corollary 3

*Proof.* Following all the proof steps of Theorem 2 while replacing $z$ and $\mathcal{D}_C$ by $z = (u, c)$ and $\mathcal{D}_{U,C}$, we have the following. Since $r(z) = \tau + H[p(\cdot \mid z)]$, we have that with probability at least $1 - \delta$ over the draw of $\bar{z} = (\bar{u}, \bar{c}) \sim \mathcal{D}_{U,C}$,

$$L_{\bar{z}}[q_S] < \delta^{-1} \mathbb{E}_{(u,c) \sim \mathcal{D}_{U,C}} [L(q_S(\cdot \mid u, I[c]), u, c)]$$

$$\leq \delta^{-1} \mathbb{E}_{(u,c) \sim \mathcal{D}_{U,C}} \left[ \max \left( 0, 1 - \frac{1}{r(u,c)} (\tau - D_{\mathrm{KL}}(p(\cdot \mid u, c) \| q_S(\cdot \mid u, I[c]))) \right) \right]$$

$$\leq \delta^{-1} \mathbb{E}_{(x,u,c) \sim \mathcal{D}} \left[ -\frac{1}{r(z)} \log q_S(x \mid u, I[c]) \right].$$

By using the definition of $Q(t)$ and taking the union bound, we have that with probability at least $1 - \delta$ over an draw of $S = (v_i)_{i=1}^{n} \sim \mathcal{D}^{\otimes n}$ and $\bar{z} \sim \mathcal{D}_{U,C}$,

$$L_{\bar{z}}[q_S] < -\frac{2\delta^{-1}}{n} \sum_{i=1}^{n} \frac{1}{r(z_i)} \log q_S(x_i \mid u_i, I[c_i]) + \mathcal{O}\left( \sqrt{\frac{Q(\delta/2)}{\delta^2 n}} \right)$$

$$= -\frac{2\delta^{-1}}{n} \sum_{i=1}^{n} \frac{1}{r(z_i)} \sum_{t=1}^{T} \log q_S(x_i^{(t)} \mid x_i^{(t-\omega:t-1)}, u_i, I[c_i]) + \mathcal{O}\left( \sqrt{\frac{Q(\delta/2)}{\delta^2 n}} \right).$$

This proves the first inequality in the desired statement. The second inequality in the desired statement follows from the fact that $\frac{1}{r(z_i)} \leq \frac{1}{r_{\min}}$ almost surely. $\qquad \square$

# Appendix B   Additional Discussions

## B.1   On the use of context prompting

The concept of relative creativity, as defined in Definition 1, does not explicitly incorporate the conditioning on a specific context prompt. Nevertheless, all derived results seamlessly adapt to such context conditioning. This adaptation is achieved by substituting the original probability space with a conditional probability space, predicated on a given context. Given that a conditional probability space retains the fundamental properties of a probability space, the application of this theory to the conditional probability scenario necessitates no alterations. Consequently, in Sections 4.2 and 4.3, we expand the statistical creativity framework to encompass the context prompting scenario. The findings, as presented in Corollary 1 and 2, evaluate statistical creativity within the setting of model performance under context prompting.

## B.2   The upper bound on the negative log likelihood

In Theorem 2, the condition of $-\log q(x \mid I[c]) \leq M$ is easily satisfiable, e.g., by using softmax over choices or any distribution that puts non-zero probability over all $x \in \mathcal{X}$. Even when we have $q$ that does not satisfy this condition, we can easily modify $q$ to $q + \epsilon$ with a normalization for some small $\epsilon > 0$: since a normalized version of $q + \epsilon$ puts non-zero probability over all $x \in \mathcal{X}$, this satisfies the condition.