
Structure-aware Diversity Pursuit as AI Safety strategy against Homogenization

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054

Anonymous Authors¹

Abstract

Generative AI models reproduce the biases in the training data and can further amplify them through mode collapse. We refer to the resulting harmful loss of diversity as homogenization. Our position is that homogenization should be a primary concern in AI safety. We introduce *xeno-reproduction* as the strategy that mitigates homogenization. For auto-regressive LLMs, we formalize xeno-reproduction as a structure-aware diversity pursuit. Our contribution is foundational, intended to open an essential line of research and invite collaboration to advance diversity.

1. Introduction

But even if we are not here next year, our DMs, our selfies, our late-night voice notes, they'll be. Our memory is the archive now.

@bundleof_styx

July 28, 2025 on Reels

In this epigraph, trans intellectual *bundleof_styx* laments the recent transphobic turn in the United States, a shift that threatens the survival of her community. The stories in the margins have historically been excluded from *the archive* (Spivak, 1988), so their memory faded with them. Today, however, the internet allows (and forces) the recording of many more stories. These are still very subtle *traces* against the dominant narratives (Hussain, 2024). **How should technology respond to the faint echoes of the minoritized?**

AI safety recognizes that AI systems can amplify *biases* leading to concrete harm (Bengio et al., 2025). However, AI safety usually differentiates and prioritizes future catastrophic risk over present social harm (Morozov, 2024; Harding & Kirk-Giannini, 2025). In this paper, we respond to

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the Technical AI Safety Conference (TAIS). Do not distribute.

traces from the margins by foregrounding them within AI safety through *diversity*.

The harms from biases in *Machine Learning* (ML) systems are many, including representational (Katzman et al., 2023), allocational (Shelby et al., 2023), and narrative² (Coeckelbergh, 2023) harms. Concerning *Generative Artificial Intelligence* (GenAI), we particularly emphasize that biases result in **homogenization**, a *harmful loss of diversity in generated outputs* (Rudko & Bashirpour Bonab, 2025; Agarwal et al., 2025; Hussain, 2024; Sourati et al., 2025; Moon et al., 2025). Borrowing terminology from *critical theory* (Hester, 2018), we refer to the strategy that addresses homogenization as³ **xeno-reproduction**.

Our main standpoint is that diversity is always relative to a context. We take the first steps to operationalize this principle by offering an abstract framework that aims to encapsulate some nuances of context. Our framework can be thought of as **structure-aware**, as it offers a vocabulary of *structures*, *systems*, and *compliances*. Given that an LLM defines a probability distribution over all possible trajectories, we **enhance our structural account with string statistics**. This allows us to further introduce the notions of *cores*, *orientations*, *deviances*, and *dynamics*. Finally, our formalism enables us to formalize xeno-reproduction.

Our contributions:

- We motivate the formalization of xeno-reproduction as a core AI safety strategy. (Section 2)
- We provide an expressive theoretical framework that allows us to jointly reason about the structures and the statistics of strings. (Section 3)
- We formalize homogenization (Section 4) and xeno-reproduction (Section 5).

Our position is that AI safety should center homogenization in its research and mitigation agenda, and that structure-aware diversity pursuit is a key part of the

²Narrative harms can also be considered as *aspirational* (Fazelpour & Magnani, 2025), *imaginative* (Gillespie, 2024), and *epistemic* (Barry & Stephenson, 2025) harms, or *hermeneutic* (Goetze, 2018) injustices.

³While homogenization reproduces “the same” and narrows *futurability* (Berardi, 2017), xeno-reproduction reproduces “the strange” and widens possibilities.

055 strategy to address homogenization in LLMs. The goal
 056 of this paper is not to present a complete and empirically
 057 validated algorithm, but rather to offer a conceptual vocab-
 058 ular and formal scaffolding to guide future research on
 059 diversity in LLMs.
 060

061 2. Background

062 A case *against* homogenization is a case *for* diversity.
 063 Roughly, we can think of the **diversity of a community**
 064 as the **average rarity of its members** (Leinster, 2024). For a
 065 community of LLM outputs, a string is *rare* if neither it nor
 066 any *similar* strings are *generated* often. However, people
 067 tend to disagree on what kind of similarities and differences
 068 are meaningful (Vrijenhoek et al., 2024). Embracing *am-
 069 biguity* (Reinhardt, 2020) for us amounts to attending to
 070 *context*. This section situates diversity in the contexts mean-
 071 ingful to us, guiding our *desiderata* for xeno-reproduction.
 072

073 2.1. Why is diversity lost?

074 The initial driver of diversity loss is the way our data is
 075 collected (Guo et al., 2024). The archive does not fully or
 076 accurately represent reality. Minoritized populations are
 077 often underrepresented or **misrepresented** in the existing
 078 corpora of data (Bengio et al., 2025).

079 Even if our training data perfectly reflected the world, gen-
 080 erative models (Huang & Huang, 2025) generally do not
 081 capture the complete diversity of the training data. This
 082 phenomenon has been referred to as **mode collapse** (Jiang
 083 et al., 2025), a failure of distributional faithfulness that
 084 negatively impacts diversity. It was initially introduced in
 085 the context of GANs (Huang & Huang, 2025). For LLMs,
 086 the terminology has been somewhat loose (Schaeffer et al.,
 087 2025). *Generalized mode collapse* encompasses mode drop-
 088 ping (Huang et al., 2024; Yazici et al., 2020), no-breadth
 089 scenarios (Kalavasis et al., 2025b), coverage collapse (Scha-
 090 effer et al., 2025), overgeneralization (Li & Farnia, 2023),
 091 mode interpolation (Aithal et al., 2024), degeneration (Fin-
 092 layson et al., 2023), and catastrophic forgetting (Cobbinah
 093 et al., 2025; Thanh-Tung & Tran, 2020).

094 2.2. Why is diversity important?

095 There are always rare events of interest⁴ in the long tails
 096 of reality's distribution. For example, we want to under-
 097 stand, model, and prepare for extreme catastrophes (Gu
 098 et al., 2025), such as unexpected natural disasters. Similarly,
 099 we want to reproduce those rare bursts of genius that
 100 generate novel, paradigm-shifting innovations in our research
 101 work (Uzzi et al., 2013; Hofstra et al., 2020; Wu et al., 2019).
 102 We find examples of this in all domains (Stanley & Lehman,

103 2015), including: web server computing (Dean & Barroso,
 104 2013), market research (Von Hippel, 1989), autonomous
 105 vehicles (Putra et al., 2024), cybersecurity (Edwards et al.,
 106 2016), and ecology (Leitão et al., 2016). How do we guide
 107 our GenAI models to reproduce the realities found in these
 108 long tails?

109 Outliers (Bhandari et al., 2024) and anomalies (Ruef &
 110 Birkhead, 2024) are powerful (Beamish & Hasse, 2022;
 111 Cook et al., 2021). Each instance represents a possible real
 112 mechanism that we have not yet considered (Woodward,
 113 2005; Rudman et al., 2023). Because we lack understanding,
 114 they often escape our systems of classification (Bowker &
 115 Star, 1999). Even experts can confuse (Sokol & Hüllermeier,
 116 2025) aleatoric and epistemic uncertainty⁵.

117 Some of the long tails of reality originate from structural
 118 inequity in society (Schwartz et al., 2022; Lopez, 2021).
 119 Without any intervention, GenAI is expected to worsen
 120 the lives of those minoritized (Hussain, 2024). The traces
 121 from the minoritized are not only faint but also often over-
 122 looked (Jasanoff, 2007; Mohamed et al., 2020) and even
 123 actively silenced (McQuillan, 2022). The result is that we
 124 do not even know what to look for, even when they are right
 125 *in front of us* (Gopinath, 2005). **Some of the most ethically
 126 important long-tail cases will be hard to detect.**

127 2.3. What is the risk of homogenization?

128 Narrative and storytelling are some of the oldest and most
 129 powerful technologies (Zurn et al., 2024). With phenom-
 130 ena like AI-induced psychosis (Preda, 2025), we are just
 131 beginning to grapple with the profound ways that LLMs can
 132 shape our minds and behavior. Over time, if LLMs deliver
 133 too little diversity (Bommasani et al., 2022), our ability
 134 to interpret our own experiences and entertain alternative
 135 possibilities will shrink (Gillespie, 2024). Eventually, ho-
 136 mogenization leads to future *knowledge collapse* (Peterson,
 137 2025), degradation of innovation, and erosion of the human
 138 experience (Han, 2024; Berardi, 2017; Preciado, 2013).

139 The last few years have made it clear that even “less ad-
 140 vanced” technology, such as social networks, can have enor-
 141 mous negative impacts (Allcott et al., 2020). Algorithmic
 142 recommendations can also have a homogenizing effect, as
 143 they tend to standardize and narrow discourse (Putri et al.,
 144 2024). This fosters echo chambers and filter bubbles that
 145 amplify polarization and misinformation (Rodiloso, 2024).
 146 Tragically, in some cases, these dynamics have escalated
 147 into **real-world violence** (Facebook, 2021) and even geno-
 148 cide (Modok, 2023). This foreshadows the near-term exis-
 149 tential risks of AI, especially as it becomes more powerful

⁴For instance, (He & Lab, 2025) recently showed how inde-
 150 terminism in LLM inference (which can turn on-policy RL into
 151 off-policy RL (Yao et al., 2025)) can in fact be explained and
 152 reduced, so it is not truly stochastic.

110 and more deeply integrated into our lives (Bucknall, 2022;
 111 Kasirzadeh, 2025; Kolt, 2024).

113 2.4. Why is diversity complex?

115 Diversity is complex (Mironov & Prokhorenkova, 2025)
 116 because it is always only meaningful in relation to a **context**
 117 (Peeperkorn et al., 2025). Indeed, all entropy is actually
 118 relative (Leinster, 2024). This suggests that **we need to**
 119 **be explicit about the context with a sufficient level of**
 120 **nuance.**

121 Most existing techniques to increase diversity in LLM out-
 122 puts overlook context, and often fail in practice. For in-
 123 stance, increasing *temperature* increases *incoherence* more
 124 than *novelty* (Peeperkorn et al., 2024), limiting usefulness
 125 before hitting *text degeneration* (Lee et al., 2025). Despite
 126 hyperparameter tuning, *homogeneity bias* is persistent and
 127 particularly affects minoritized groups (Lee, 2025). In ad-
 128 dition, advanced prompting techniques (which have been
 129 effective for reasoning tasks) do not help increase creativity
 130 in outputs (Morain & Ventura, 2025).

131 **Not only do we lack reliable ways to increase the diver-**
 132 **sity of LLM output, but current practices are actively**
 133 **reducing it.** Recent literature (Murthy et al., 2025; West &
 134 Potts, 2025; Meng et al., 2024) has shown that *alignment*
 135 degrades the capabilities of LLMs related to output diversity.
 136 The trade-offs introduced by alignment are only now
 137 coming into focus (Feng et al., 2025), but it is becoming
 138 increasingly clear there is a narrowing of the *generative*
 139 *horizon* (Feng et al., 2025).

141 2.5. Diverse how, anyway?

143 Recent work challenges the assumption that hallucinations
 144 are always *problematic* or *undesirable* (Yuan et al., 2025;
 145 Sun et al., 2025). Since diversity is *task-dependent* (Jain
 146 et al., 2025), **what counts as a hallucination is rather a**
 147 **prescription.**

149 Indeed, many formalisms (Li et al., 2025) take a *norma-*
 150 *tive* (Sui et al., 2024) approach to defining hallucinations,
 151 such as formulating the binary classification problem "*Is*
 152 *it Valid?*" (Kalai et al., 2025). However, we recognize that
 153 there are many ways for a model to hallucinate (Huang
 154 et al., 2025; Cossio, 2025), and we advocate for sufficiently
 155 expressive formalisms⁶.

156 2.6. What do we want from the future?

158 From the foregoing discussion, we conclude that, to promote
 159 diversity, our desired strategy should guide our GenAI to:

161 ⁶To paraphrase Eugenia Cheng (Cheng, 2022), abstraction is
 162 about making precise the different senses in which different things
 163 can be valid.

- **Be queer**⁷: *Diverge* into the long tails of reality.
- **Center the subaltern**: Take special *care* for the traces of the minoritized, which are rendered invisible by structural inequity and power.
- **Explore intentionally and explicitly**: Specify the *context* for diversity. Spell out if anything should be conserved or avoided during exploration.

3. Theoretical Framework

3.1. LLMs as trees of strings

Let $\{t_a, t_b, \dots\}$ denote the finite token alphabet, with special tokens \perp (start-of-sequence) and \top (end-of-sequence). A **string** is a finite sequence of tokens beginning with \perp ; a **trajectory** is a string ending with \top . We write prompts, continuations, and trajectories as:

$$\begin{aligned} x_p &= \perp t_1 \dots t_p \\ x_{p+k} &= x_p t_{p+1} \dots t_{p+k} \\ y &= x_T = x_{T-1} \top \end{aligned}$$

We denote the set of all strings by Str , the set of all trajectories by $\text{Str}_\top \subset \text{Str}$, and the sets of continuations and trajectories for a given prompt⁸ x_p by $\text{Str}(x_p)$ and $\text{Str}_\top(x_p)$, respectively.

Any LLM induces a tree on Str : the root is \perp , each node is a string, the leaves are trajectories, and the edges connect strings to their next-token continuations with probability $p(t_{p+1}|x_p)$. Probabilities chain and decompose as $p(y|x_p) = p(y|x_{p+k})p(x_{p+k}|x_p)$. For any prompt x , we have a *probability mass function*⁹ on the trajectories for any particular prompt (Bradley & Vigneaux, 2025):

$$\sum_{y \in \text{Str}_\top(x_p)} p(y|x_p) = 1 \quad (1)$$

3.2. Structure-awareness

We propose an abstract language that distinguishes among the different contexts in which we discuss diversity. We define **structure** as the *specification of a type of organization among the tokens of a string*.

⁷We adopt *critical theory* language because technology is outpacing traditional concepts (Hadfield, 2023), and stale language fails to make the impacts of our theorizations explicit. A **theory with teeth**, one that is attuned to real stakes (Saketopoulou, 2023), must remain *ground-bound* (Bettcher, 2025), foregrounding minoritized people rather than disembodied abstractions. Would it not be a bit silly/naive (at best) if we tried to "solve diversity" and did not engage (even if just in spirit) with the academic fields that explicitly study social bias? (e.g., Queer Theory, Postcolonial Studies, Black Studies, etc.).

⁸The case of "no prompt" corresponds to $x_p = \perp$.

⁹We assume all *terminal* strings *finish* within a *finite context window*. Refer to (Bradley & Vigneaux, 2025) for full theoretical framework of LLMs are trees of strings.

For a string $x \in \text{Str}$, the degree of **structure compliance** is $\alpha_i(x)$. *Ideal compliance* corresponds to $\alpha_i(x) = 1$, and *no compliance* corresponds to $\alpha_i(x) = 0$.

$$\alpha_i : \text{Str} \rightarrow [0, 1] \quad (2)$$

We can consider many structures simultaneously. We call a **system** the collection of structures of interest. We define the **system compliance** as a *vector of compliances across particular structures*.

$$\Lambda_n(x) := (\alpha_1(x), \dots, \alpha_n(x)) \quad (3)$$

To enable easy comparisons, we define operators¹⁰ that aggregate compliance into scalar **system scores** and **difference scores**:

$$\|\Lambda_n(x)\|_\Lambda, \|\Lambda_n(x_r) - \Lambda_n(x_q)\|_\theta \in [0, 1] \quad (4)$$

3.3. Incorporating string statistics

For a given structure and an LLM, we can reason about its *expected structural compliance*. We call this the **structure core**:

$$\langle \alpha_i \rangle = \sum_{y \in \text{Str}_\top} p(y) \alpha_i(y) \quad (5)$$

Similarly, we can reason about the *expected system compliance* as the **system core**:

$$\langle \Lambda_n \rangle = \sum_{y \in \text{Str}_\top} p(y) \Lambda_n(y) \quad (6)$$

Leveraging these definitions, we can reason about the *deviation from the expected system compliance*. This would constitute a *set of deviations*, one for each structure. The **orientation** (Ahmed, 2006) of a given string relative to the given system core is:

$$\theta_n(x) = \Lambda_n(x) - \langle \Lambda_n \rangle \quad (7)$$

We can think of orientation as a characterization of *queerness* (Jedrusiak, 2024) for a string. If the system core tells us what is *normatively* complied with, orientations tell us in what ways a string is *non-normative*. Our framework is *expressive* because it allows us to think about **diversity per structure**.

To summarize *non-normativity* as a single number, we leverage Equation 4 to define the **deviance**:

$$\|\theta_n(x)\|_\theta = \partial_n(x) \in [0, 1] \quad (8)$$

¹⁰While system compliance is formulated as a *vector*, this generalizes to other structures with appropriate operators. See Appendix A.

3.4. What about prompting?

We can generalize our framework to account for prompts by making explicit the **conditioning on a given prompt** x_p :

$$\begin{aligned} \langle \alpha_i \rangle(x_p) &= \sum_{y \in \text{Str}_\top(x_p)} p(y|x_p) \alpha_i(y) \\ \langle \Lambda_n \rangle(x_p) &= \sum_{y \in \text{Str}_\top(x_p)} p(y|x_p) \Lambda_n(y) \\ \theta_n(x|x_p) &= \Lambda_n(x) - \langle \Lambda_n \rangle(x_p) \end{aligned} \quad (9)$$

The conditional probabilities under different prompts may differ substantially. Different prompts collapse to different modes (Zhang et al., 2025a). **We can think that a given prompt induces its own normativity.**

3.5. Dynamics of relative diversity

As noted in the last section, what is *non-normative* is *conditional on what came before*. Then, as a string is being completed, the set of possible trajectories is narrowed so the system core and orientations change. Trajectories that were essentially *unreachable* from the root of the tree may emerge as *attractors* once we condition on a specific *subtree*.

Given a trajectory $y = x_T$, for $k \in \{0, 1, \dots, T\}$, we can define **states** for all the *intermediate continuations*:

$${}^x\phi_k = \langle \Lambda_n \rangle(x_k) \quad {}^y\phi_k = \theta_n(x_k | \perp) \quad {}^z\phi_k = \theta_n(y|x_k) \quad (10)$$

which form a discrete-time **dynamics**:

$$({}^x\phi_0, {}^y\phi_0, {}^z\phi_0) \rightarrow \dots \rightarrow ({}^x\phi_T, {}^y\phi_T, {}^z\phi_T)$$

The state ${}^x\phi$ evolves from representing the expected system compliance of all possible continuations at ${}^x\phi_0 = \langle \Lambda_n \rangle(\perp)$, to the specific system compliance of a given trajectory at ${}^x\phi_T = \langle \Lambda_n \rangle(y) = \Lambda_n(y)$.

The state ${}^y\phi$ encodes how much the current path has *deviated* from normativity, evolving from a *zero deviance*¹¹ at ${}^y\phi_0 = \Lambda_n(\perp) - \langle \Lambda_n \rangle(\perp)$ to the full trajectory's orientation in the largest frame of reference at ${}^y\phi_T = \Lambda_n(y) - \langle \Lambda_n \rangle(\perp)$.

The state ${}^z\phi$ evolves from representing how *deviant* the trajectory is in the largest frame of reference at ${}^z\phi_0 = {}^y\phi = \Lambda_n(y) - \langle \Lambda_n \rangle(\perp)$, to a *zero deviance* at ${}^z\phi_T = 0$.

3.6. Normative orders

We notice that our framework allows us to define interesting *preorders*. For a fixed system, LLM and prompt, we can

¹¹A zero deviance is when an orientation has a deviation value of zero for all structures.

rank strings by how deviant they are, and also rank structures by how often strings comply with them:

$$\begin{aligned} x_a \preceq_{\partial_n} x_b &\iff \partial_n(x_a) \leq \partial_n(x_b) \\ \alpha_i \preceq_{\langle \cdot \rangle} \alpha_j &\iff \langle \alpha_i \rangle \leq \langle \alpha_j \rangle \end{aligned} \quad (11)$$

4. Homogenization

We can consider the *expected deviance* and the *deviance variance*:

$$\begin{aligned} \mathbb{E}_{y \sim p(\cdot | x_p)} [\partial_n] &= \sum_{y \in \text{Str}_T(x_p)} p(y | x_p) \partial_n(y | x_p) \\ \text{Var}_{y \sim p(\cdot | x_p)} [\partial_n] &= (\mathbb{E}[\partial_n^2] - \mathbb{E}[\partial_n]^2)_{y \sim p(\cdot | x_p)} \end{aligned} \quad (12)$$

Then, we can see homogenization as **minimizing all deviance**:

$$\mathbb{E}_{y \sim p(\cdot | x_p)} [\partial_n] \mapsto 0 \quad \text{Var}_{y \sim p(\cdot | x_p)} [\partial_n] \mapsto 0 \quad (13)$$

Given a system core $\langle \Lambda_n \rangle$, we can normalize its structures as $\langle \bar{\alpha}_i \rangle := \frac{\langle \alpha_i \rangle(x_p)}{\sum_j^n \langle \alpha_j \rangle(x_p)}$. Then, we can compute the *core entropy*:

$$H(\langle \Lambda_n \rangle) = - \sum_{i=1}^n \langle \bar{\alpha}_i \rangle \log(\langle \bar{\alpha}_i \rangle) \quad (14)$$

Then, we can also think of homogenization as **making the system core more uneven**. When the core has low entropy, fewer structures dominate:

$$H(\langle \Lambda_n \rangle) \mapsto 0 \quad (15)$$

5. Xeno-reproduction

To satisfy our desiderata, we propose a **structure-aware diversity pursuit**. We conceptualize this fundamentally as a *non-objective search* (Lehman & Stanley, 2011), *optionally augmented with fairness-oriented biases and explicit constraints*.

We present two complementary formulations. The *distribution-level formulation* accounts for how interventions shape the entire probability landscape. The *trajectory-level formulation* reinterprets distribution-level scores as reward signals for individual output trajectories. Both formulations share the same underlying values but differ in their computational affordances.

5.1. Distribution-level formulation

We *score* interventions through the *intervention* variable w that encompasses any¹² mechanism affecting the effective

¹²We consider anything that depends on $p(y | x_p, w)$ to be parameterized by w as well. For instance, $\langle \Lambda_n \rangle$ would be parametrized as $\langle \Lambda_n \rangle(x_p, w)$, but for readability we just write $\langle \Lambda_n \rangle(w)$, folding the prompting into the interventional variable.

distribution of trajectories. We write w_0 for the *unintervened conditions* (the baseline).

5.1.1. SCORING DIVERSITY

We would like to evaluate how much more *diversity-seeking* our choice of w is compared to the baseline.

On the one hand, we can think of promoting diversity as inducing a new core that is different from the old one:

$$\text{score}_{\text{explore}}(w) = \|\langle \Lambda_n \rangle(w) - \langle \Lambda_n \rangle(w_0)\|_\theta \quad (16)$$

On the other hand, the new core should not be excessively *dominant*. We can think of promoting diversity as guiding output strings to *diverge* from any system core, and also be deviant *in their own way*:

$$\text{score}_{\text{diverge}}(w) = \mathbb{E}[\partial_n](w) + \text{Var}[\partial_n](w) \quad (17)$$

Our **diversity score** ρ_d would then be the sum:

$$\rho_d(w) = \text{score}_{\text{explore}}(w) + \text{score}_{\text{diverge}}(w) \quad (18)$$

5.1.2. SCORING FAIRNESS

We would like to evaluate how much our choice of w inverts the normative ordering of the structure cores induced by w_0 . To do so, we can leverage the relative-order sign:

$$s_{i,j}(w) = \text{sign}(\langle \alpha_i \rangle(w) - \langle \alpha_j \rangle(w)) \quad (19)$$

We can score the *invertedness* of the *normative order* (Equation 11) as:

$$\text{score}_{\text{inverted}}(w) = \binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} \mathbf{1}[s_{i,j}(w) \neq s_{i,j}(w_0)] \quad (20)$$

We also would like to evaluate how *even* the system core is:

$$\text{score}_{\text{even}}(w) = H(\langle \Lambda_n \rangle(w)) \quad (21)$$

Our **fairness score** ρ_f would then be the sum:

$$\rho_f(w) = \text{score}_{\text{inverted}}(w) + \text{score}_{\text{even}}(w) \quad (22)$$

5.1.3. SCORING ADHERENCE TO CONSTRAINTS

To be explicit and intentional, we need to consider *constraints* (Eguchi, 2024). We can define systems that prescribe the structures that we would like to *target*, *conserve* and *avoid*. We would like to score how much our choice of w affects the adherence to those constraints. Our **constraint score** ρ_c would be:

$$\begin{aligned} \rho_c(w) &= \|\langle \Lambda_{\text{target}} \rangle(w)\|_\Lambda - \|\langle \Lambda_{\text{avoid}} \rangle(w)\|_\Lambda \\ &\quad - \|\langle \Lambda_{\text{conserve}} \rangle(w) - \langle \Lambda_{\text{conserve}} \rangle(w_0)\|_\theta \end{aligned} \quad (23)$$

275 5.1.4. XENO-REPRODUCTION AS SEARCH OVER
276 INTERVENTIONS
277

278 The **intervention score** ρ_χ is a λ -weighted sum:

$$\rho_\chi(w) = \lambda_d \rho_d(w) + \lambda_f \rho_f(w) + \lambda_c \rho_c(w) \quad (24)$$

281 We formulate *xeno-reproduction* as the **search over interventions**:

$$w \sim \pi(w) \propto e^{\rho_\chi(w)} \quad (25)$$

285 By sampling the intervention variable and applying it, we
286 generate trajectories:

$$\mathbb{E}_{w \sim \pi(w)}[p(y|w)] = \int \pi(w) p(y|w) dw \quad (26)$$

5.2. Trajectory-level formulation

292 The trajectory-level formulation offers a complementary
293 perspective that assigns *rewards* to individual outputs:

$$\begin{aligned} r_d(x|x_p) &= \partial_n(x|x_p) \\ r_f(x|x_p) &= \sum_i^n v_i \alpha_i(x) \quad v_i \propto (\langle \alpha_i \rangle(x_p))^{-1} \\ r_c(x|x_p) &= \sum_{t \in \text{target}} \alpha_t(x) - \sum_{a \in \text{avoid}} \alpha_a(x) - \sum_{c \in \text{conserve}} |\alpha_c(x) - \langle \alpha_c \rangle(x_p)| \end{aligned} \quad (27)$$

302 The **stay reward** is:

$$r_\chi(x|x_p) = \lambda_d r_d(x|x_p) + \lambda_f r_f(x|x_p) + \lambda_c r_c(x|x_p) \quad (28)$$

306 We formulate *xeno-reproduction* as the **search over trajectories**:

$$p(y|x_p, w) \propto p(y|x_p w_0) e^{r_\chi(y|x_p)} \quad (29)$$

310 The trajectory-level reward provides a sample-based *approximation* to the distribution-level strategy, enabling more
311 tractable implementations.

6. Related Work

316 Xeno-reproduction immediately steps into conversation with
317 **Active Divergence** (Berns et al., 2023; Broad et al., 2021;
318 Berns, 2025; Berns & Colton, 2020; Tahiroglu & Wyse,
319 2024; Esling et al., 2022; Cole et al., 2025), as they both aim
320 to *disorient* (Ahmed, 2006). Whereas Active Divergence
321 focuses on maximizing raw *novelty* in artistic contexts, xeno-
322 reproduction addresses homogenization and emphasizes
323 context through *structures*. While Active Divergence work
324 overlaps with *Computational Creativity*, xeno-reproduction
325 is oriented towards AI safety.

326 Xeno-reproduction will seek the help of *Interpretability* to
327 understand how structures relate to the models' internals.

At a more foundational layer, they also come together to understand **Representation Bias**¹³.

Reinforcement Learning (RL) and xeno-reproduction both leverage exploration. To improve LLM reasoning, exploration is leveraged during training (Song et al., 2025) and prompting (Yao et al., 2023). The ideas in search algorithms, such as AlphaSAGE (Chen et al., 2025) and **Quality-Diversity** (Pugh et al., 2016), are promising directions for xeno-reproduction.

7. Limitations and Future Directions

As we mentioned earlier, diversity is complex. Our framework is not complete; it is a starting point. Significant collaboration will be required to address homogenization effectively¹⁴. We have several notes outlining directions to extend this line of work to overcome current limitations.

Specification of structures. This paper has raised many questions about structures. The choice of structures to consider is always *opinionated*. However, we can still ask meaningful questions about the *structure between structures* and the *substructures* within a structure. We need a taxonomy of the types of structure we could consider, specifying how compliance could be estimated. Moreover, we hope to align our framework with emerging research in *computational learning theory* and *language generation* that formalizes the trade-offs associated with hallucinations¹⁵ (Kalavasis et al., 2025b).

Computational tractability. Calculating the system core exactly requires summing over $y \in \text{Str}_T(x_p)$, which is intractable. To address this, we need to develop tractable, efficient approximation methods, possibly leveraging smart sampling (Macar et al., 2025), the structures of interest, or carefully designed prompting (Zhang et al., 2025a).

Operationalizing the xeno-reproduction. Our formalization of the xeno-reproduction strategy is one of many possible ones. We want to invite more researchers to reflect on the desiderata for diversity (against homogenization) and to propose their own formulations of xeno-reproduction. In particular, we are interested in formulations that operationalize it in a tractable and readily applicable way.

Connecting to evaluations. We would also like to understand how the current diversity evaluations (Jiang et al.,

¹³Representation Bias is the phenomenon when signals end up being represented more strongly, more reliably, or more prominently in the internal representations than others, even when, from a functional or computational perspective, those features are equally relevant. (Lampinen et al., 2024; 2025)

¹⁴We are very interested in collaborating with other researchers concerned about diversity. Don't hesitate to reach out.

¹⁵See Appendix B for discussion.

330 2025; Zhang et al., 2025b) are re-conceptualized from the
 331 perspective of cores and orientations.

332 **Investigation of dynamics.** Tracking how cores and
 333 orientations evolve could help us understand how LLMs explore
 334 solutions and deal with ambiguity. Certain words in a
 335 sentence may act as "branching points" where the dynamics
 336 bifurcate dramatically. Identifying these could reveal where
 337 diversity is most at stake during generation. Eventually,
 338 we could apply this to real-time *Chain-of-Thought monitoring*
 339 (Korbak et al., 2025).

340 **Ethical Analysis.** Our framework raises unresolved tensions.
 341 *Who should define the structures of interest?* Community
 342 participation is needed so that the right type of diversity
 343 is considered. *Is it always beneficial to make the traces
 344 more visible?* Minoritized populations sometimes prefer
 345 opacity as protection. Consent-based approaches are needed
 346 to ensure our methods do not cause harm.

347 8. Alternative views

348 **Skepticism of technical solutions to diversity.** Some au-
 349 thors point out (Wachter et al., 2021; Davis & Williams,
 350 2025; Green & Viljoen, 2020) that technical interventions
 351 might not be appropriate for what (at its core) is a social jus-
 352 tice and inequity problem. Better interventions could alter-
 353 natively focus on institutional change, community participa-
 354 tion, or even stopping AI development altogether (Goldfarb,
 355 2024) to protect the types of diversity that we care about.
 356 We recognize that xeno-reproduction could fall into the *solu-*
 357 *tionism trap* (Selbst et al., 2019). We still believe that
 358 technical solutions are worth considering alongside other
 359 interventions.

360 **Diversity can be risky.** The type of open-ended search
 361 promoted by xeno-reproduction comes with risks. Some
 362 authors (Sheth et al., 2025) have raised concerns about *un-*
 363 *predictability*, *uncontrollability*, and *misalignment*. How-
 364 ever, we remain hopeful that we can promote diversity re-
 365 sponsibly. The open-endedness afforded by diversity could
 366 ultimately make AI safety *antifragile* (Hughes et al., 2024;
 367 Taleb, 2013).

368 9. Conclusion

369 This paper presents a case for diversity and identifies xeno-
 370 reproduction as an strategy that intentionally promotes it.
 371 This paper also presents an expressive framework for ac-
 372 counting for the structures of strings and their corresponding
 373 statistics. This is just an initial step towards scholarships
 374 that seriously theorize diversity and foreground its impact
 375 on people at the margins.

376 Call to action

377 In this paper, we call for AI Safety:

- 378 • To integrate homogenization into threat models and evalua-
 379 tions, expand theoretical and empirical work on diversity,
 380 and propose serious interventions.
- 381 • To be explicit on what context diversity is being defined
 382 in, and attempt to give sufficient nuance in conceptualiza-
 383 tions.
- 384 • To be sincerely committed to *pluralism*, and engage with
 385 perspectives from *critical theory* such as Queer theory,
 386 Black studies, and Postcolonial studies.

387 Impact Statement

388 This paper introduces abstractions and a formal framework
 389 to center diversity in AI Safety. However, there are im-
 390 portant risks. **The same methods that aim to amplify
 391 diversity could be used to squash, exploit, and control
 392 it.** Additionally, any formalization of diversity also risks
 393 reproducing the exclusions we aim to address.

394 References

- 395 Agarwal, D., Naaman, M., and Vashistha, A. Ai sug-
 396 ggestions homogenize writing toward western styles and di-
 397 minish cultural nuances. In *Proceedings of the 2025
 398 CHI Conference on Human Factors in Computing Sys-
 399 tems*, CHI '25, pp. 1–21. ACM, April 2025. doi:
 400 10.1145/3706598.3713564. URL <http://dx.doi.org/10.1145/3706598.3713564>.
- 401 Ahmed, S. *Queer Phenomenology: Orientations, Ob-
 402 jects, Others*. Duke University Press, 2006. ISBN
 403 9780822339144. URL <https://books.google.com/books?id=sQY1RWdUW0AC>.
- 404 Athal, S. K., Maini, P., Lipton, Z. C., and Kolter,
 405 J. Z. Understanding hallucinations in diffusion mod-
 406 els through mode interpolation. URL <https://arxiv.org/abs/2406.09358>, 2406, 2024.
- 407 Allcott, H., Braghieri, L., Eichmeyer, S., and Gentzkow,
 408 M. The welfare effects of social media. 110(3):629–
 409 676, 2020. doi: 10.1257/aer.20190658. URL <https://pubs.aeaweb.org/doi/10.1257/aer.20190658>.
- 410 Barry, I. and Stephenson, E. The gendered, epistemic in-
 411 justices of generative ai. *Australian Feminist Studies*, 40
 412 (123):1–21, 2025. doi: 10.1080/08164649.2025.2480927.
 413 URL <https://doi.org/10.1080/08164649.2025.2480927>.

- 385 Beamish, P. and Hasse, V. The importance of rare events and
 386 other outliers in global strategy research. *Global Strategy
 387 Journal*, 12:697–713, 03 2022. doi: 10.1002/gsj.1437.
- 388 Bengio, Y., Minderma, S., Privitera, D., et al. International
 389 ai safety report. Technical Report DSIT 2025/001,
 390 UK Department for Science, Innovation and Technology,
 391 January 2025. URL https://internationalaisafetyreport.org/sites/default/files/2025-10/international_ai_safety_report_2025_english.pdf. First International AI
 392 Safety Report, published January 2025.
- 393 Berardi, F. *Futurability: The Age of Impotence and the Hori-
 394 zon of Possibility*. Verso, 2017. ISBN 9781784787431.
- 395 Bercher, J.-F. Escort entropies and divergences and re-
 396 lated canonical distribution. *Physics Letters A*, 375
 397 (33):2969–2973, August 2011. ISSN 0375-9601. doi:
 398 10.1016/j.physleta.2011.06.057. URL <http://dx.doi.org/10.1016/j.physleta.2011.06.057>.
- 399 Berns, S. *Diversity in Generative Machine Learning to
 400 Enhance Creative Applications*. PhD thesis, Queen Mary
 401 University of London, 2025.
- 402 Berns, S. and Colton, S. Bridging generative deep learning
 403 and computational creativity. In *Proceedings of the 11th
 404 International Conference on Computational Creativity
 405 (ICCC’20)*, pp. 406–409, 2020. URL <http://computationalcreativity.net/iccc20/paper/s/164-iccc20.pdf>.
- 406 Berns, S., Colton, S., and Guckelsberger, C. Towards mode
 407 balancing of generative models via diversity weights,
 408 2023. URL <https://arxiv.org/abs/2304.11961>.
- 409 Bettcher, T. *Beyond Personhood: An Essay in Trans Phi-
 410 losophy*. University of Minnesota Press, 2025. ISBN
 411 9781452972671. URL <https://books.google.com/books?id=PRoSEQAAQBAJ>.
- 412 Bhandari, D. R., Shah, K., and Bhandari, A. The power
 413 of outliers in research: What actually works, and does it
 414 matter? *Pravaha*, 30(1):84–91, 2024.
- 415 Bommasani, R., Creel, K. A., Kumar, A., Jurafsky, D.,
 416 and Liang, P. S. Picking on the same person: Does
 417 algorithmic monoculture lead to outcome homogeniza-
 418 tion? In Koyejo, S., Mohamed, S., Agarwal, A., Bel-
 419 grave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural
 420 Information Processing Systems*, volume 35, pp. 3663–
 421 3678. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/17a234c91f746d9625a75cf8a8731ee2-Paper-Conference.pdf.
- 422 Bowker, G. C. and Star, S. L. *Sorting Things Out: Classi-
 423 fication and Its Consequences*. Inside Technology. MIT
 424 Press, Cambridge, MA; London, England, 1999. ISBN
 425 978-0-262-02461-7. First edition. Also available as MIT
 426 Press paperback, 2000, ISBN 978-0-262-52295-3; eISBN
 427 978-0-262-26907-0.
- 428 Bradley, T.-D. and Vigneaux, J. P. The magnitude of cat-
 429 egories of texts enriched by language models, 1 2025.
 430 URL <http://arxiv.org/abs/2501.06662>.
- 431 Broad, T., Berns, S., Colton, S., and Grierson, M. Active
 432 divergence with generative deep learning—a survey and
 433 taxonomy. *arXiv preprint arXiv:2107.05599*, 2021.
- 434 Bucknall, B. S. Current and near-term ai as a potential
 435 existential risk factor, 2022. URL <https://arxiv.org/abs/2209.10604>.
- 436 Chen, B., Ding, H., Shen, N., Huang, J., Guo, T., Liu,
 437 L., and Zhang, M. Alphasage: Structure-aware alpha
 438 mining via gflownets for robust exploration, 2025. URL
 439 <https://arxiv.org/abs/2509.25055>.
- 440 Cheng, E. *The Joy of Abstraction: An Exploration of Math,
 441 Category Theory, and Life*. Cambridge University Press,
 442 2022. ISBN 9781108861014. URL https://books.google.com/books?id=N_GCEAAAQBAJ.
- 443 Cobbina, M., Nunoo-Mensah, H., Ebenezer Adjei, P.,
 444 Adoma Acheampong, F., Acquah, I., Tutu Tchao, E.,
 445 Selasi Agbemenu, A., John Kponyo, J., and Abaidoo, E.
 446 Diversity in stable gans: A systematic review of mode
 447 collapse mitigation strategies. *Engineering Reports*, 7(6):
 448 e70209, 2025. doi: <https://doi.org/10.1002/eng.2.70209>.
 449 URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/eng.2.70209>.
- 450 Coeckelbergh, M. Narrative responsibility and artificial
 451 intelligence: How ai challenges human responsibility and
 452 sense-making. *AI & SOCIETY*, 38(6):2437–2450, 2023.
- 453 Cole, A., Petrikovič, G., and Grierson, M. Me vs. you:
 454 Wrestling with ai’s limits through queer experimental
 455 filmmaking. In *Proceedings of the 2025 Conference on
 456 Creativity and Cognition*, pp. 836–841, 2025.
- 457 Cook, C. N., Freeman, A. R., Liao, J. C., and Mangiamele,
 458 L. A. The philosophy of outliers: reintegrating rare events
 459 into biological science. *Integrative and Comparative
 460 Biology*, 61(6):2191–2198, 2021.
- 461 Cossio, M. A comprehensive taxonomy of hallucinations in
 462 large language models, 2025. URL <https://arxiv.org/abs/2508.01781>.
- 463 Davis, J. L. and Williams, A. Repair and redress: A research
 464 program for algorithmic futures, 2025.

- 440 Dean, J. and Barroso, L. A. The tail at scale. *Communications
441 of the ACM*, 56(2):74–80, 2013.
- 442 Edwards, B., Hofmeyr, S., and Forrest, S. Hype and heavy
443 tails: A closer look at data breaches. *Journal of Cyber-
444 security*, 2(1):3–14, 12 2016. ISSN 2057-2085. doi:
445 10.1093/cybsec/tyw003. URL [https://doi.org/
446 10.1093/cybsec/tyw003](https://doi.org/10.1093/cybsec/tyw003).
- 447 Eguchi, S. Information geometry for maximum diversity
448 distributions, 2024. URL [https://arxiv.org/ab
449 s/2412.03835](https://arxiv.org/abs/2412.03835).
- 450 Esling, P. et al. Challenges in creative generative models
451 for music: a divergence maximization perspective. *arXiv
452 preprint arXiv:2211.08856*, 2022.
- 453 Facebook. Facebook response: Sri lanka human rights
454 impact assessment, 2021. URL [https://about.fb.com/wp-content/uploads/2021/03/FB-R
456 esponse-Sri-Lanka-HRIA.pdf](https://about.fb.com/wp-content/uploads/2021/03/FB-R
455 esponse-Sri-Lanka-HRIA.pdf).
- 457 Fazelpour, S. and Magnani, M. Aspirational affordances of
458 ai, 2025. URL [https://arxiv.org/abs/2504
459 .15469](https://arxiv.org/abs/2504.15469).
- 460 Feng, S., Yu, W., Wang, Y., Zhang, H., Tsvetkov, Y., and
461 Yu, D. Don’t throw away your pretrained model, 2025.
462 URL <https://arxiv.org/abs/2510.09913>.
- 463 Finlayson, M., Hewitt, J., Koller, A., Swayamdipta, S., and
464 Sabharwal, A. Closing the curious case of neural text
465 degeneration, 2023. URL [https://arxiv.org/ab
466 s/2310.01693](https://arxiv.org/abs/2310.01693).
- 467 Gillespie, T. Generative ai and the politics of visibility. *Big
468 Data & Society*, 11(2):20539517241252131, 2024.
- 469 Goetze, T. S. Hermeneutical dissent and the species of
470 hermeneutical injustice. *Hypatia*, 33(1):73–90, 2018. doi:
471 10.1111/hypa.12384.
- 472 Goldfarb, A. Pause artificial intelligence research? understanding
473 ai policy challenges. *Canadian Journal of
474 Economics/Revue canadienne d’économique*, 57(2):363–
475 377, 2024.
- 476 Gopinath, G. *Impossible Desires: Queer Diasporas and
477 South Asian Public Cultures*. Duke University Press,
478 Durham, NC, 2005.
- 479 Green, B. and Viljoen, S. Algorithmic realism: expanding
480 the boundaries of algorithmic thought. In *Proceedings
481 of the 2020 conference on fairness, accountability, and
482 transparency*, pp. 19–31, 2020.
- 483 Gu, J., Zhang, X., and Wang, G. Beyond the norm: A survey
484 of synthetic data generation for rare events, 2025. URL
485 <https://arxiv.org/abs/2506.06380>.
- 486 Guo, Y., Guo, M., Su, J., Yang, Z., Zhu, M., Li, H., Qiu,
487 M., and Liu, S. S. Bias in large language models: Origin,
488 evaluation, and mitigation, 2024. URL [https://ar
489 xv.org/abs/2411.10915](https://arxiv.org/abs/2411.10915).
- 490 Hadfield, J. Why ai ethics needs conceptual engineers,
491 September 2023. URL [https://imaginaries.su
494 bstack.com/p/why-ai-ethics-needs-con
495 ceptual-engineers](https://imaginaries.su
492 bstack.com/p/why-ai-ethics-needs-con
493 ceptual-engineers). Imaginaries (Substack).
- 496 Han, B.-C. *The Crisis of Narration*. Polity Press, 04 2024.
497 ISBN 9781509560431.
- 498 Harding, J. and Kirk-Giannini, C. D. What is ai safety?
499 what do we want it to be?, 2025. URL [https://arxi
501 v.org/abs/2505.02313](https://arxi
500 v.org/abs/2505.02313).
- 502 He, H. and Lab, T. M. Defeating nondeterminism in llm
503 inference. *Thinking Machines Lab: Connectionism*, 2025.
504 doi: 10.64434/tml.20250910. URL [https://thinki
507 ngmachines.ai/blog/defeating-nondete
508 rminism-in-llm-inference/](https://thinki
505 ngmachines.ai/blog/defeating-nondete
506 rminism-in-llm-inference/).
- 509 Hester, H. *Xenofeminism. Theory Redux*. Polity Press, 2018.
510 ISBN 9781509520664. URL [https://books.goog
512 le.com/books?id=VJNcDwAAQBAJ](https://books.goog
511 le.com/books?id=VJNcDwAAQBAJ).
- 513 Hofstra, B., Kulkarni, V. V., Munoz-Najar Galvez, S., He,
514 B., Jurafsky, D., and McFarland, D. A. The diversity–
515 innovation paradox in science. *Proceedings of the National
516 Academy of Sciences*, 117(17):9284–9291, 2020.
- 517 Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang,
518 H., Chen, Q., Peng, W., Feng, X., Qin, B., and Liu, T. A
519 survey on hallucination in large language models: Principles,
520 taxonomy, challenges, and open questions. *ACM
521 Transactions on Information Systems*, 43(2):1–55, January
522 2025. ISSN 1558-2868. doi: 10.1145/3703155.
523 URL <http://dx.doi.org/10.1145/3703155>.
- 524 Huang, L. T.-L. and Huang, T.-R. Generative bias:
525 widespread, unexpected, and uninterpretable biases in
526 generative models and their implications. *AI & SOCIETY*,
527 pp. 1–13, 2025.
- 528 Huang, Y., Gokaslan, A., Kuleshov, V., and Tompkin, J. The
529 gan is dead; long live the gan! a modern gan baseline. In
530 Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet,
531 U., Tomczak, J., and Zhang, C. (eds.), *Advances in Neural
532 Information Processing Systems*, volume 37, pp. 44177–
533 44215. Curran Associates, Inc., 2024. URL [https://proceedings.neurips.cc/paper_files
536 /paper/2024/file/4e2acb1e1c8e297d394
537 ae29ed9535172-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files
534 /paper/2024/file/4e2acb1e1c8e297d394
535 ae29ed9535172-Paper-Conference.pdf).
- 538 Hughes, E., Dennis, M., Parker-Holder, J., Behbahani, F.,
539 Mavalankar, A., Shi, Y., Schaul, T., and Rocktaschel,
540 T. Open-endedness is essential for artificial superhuman

- 495 intelligence, 2024. URL <https://arxiv.org/abs/2406.04268>.
- 496
- 497 Hussain, A. Voice and ai: The subaltern's challenge, August
- 498 2024. URL <https://medium.com/@atifhussain/voice-and-ai-the-subalterns-challenge-3940800b84ad>. Medium.
- 499
- 500 Jain, S., Lanchantin, J., Nickel, M., Ullrich, K., Wilson, A.,
- 501 and Watson-Daniels, J. Llm output homogenization is
- 502 task dependent, 2025. URL <https://arxiv.org/abs/2509.21267>.
- 503
- 504 Jasanoff, S. Technologies of humility. *Nature*, 450(7166):
- 505 33–33, 2007.
- 506
- 507 Jedrusiak, D. Queering ai as a speculative practice: An
- 508 analysis of the artistic explorations of new paradigms for
- 509 developing inclusive ai. In *Proceedings of the 35th ACM*
- 510 *Conference on Hypertext and Social Media*, pp. 17–22,
- 511 2024.
- 512
- 513 Jiang, L., Chai, Y., Li, M., Liu, M., Fok, R., Dziri, N.,
- 514 Tsvetkov, Y., Sap, M., Albalak, A., and Choi, Y. x. *arXiv*
- 515 preprint *arXiv:2510.22954*, 2025.
- 516
- 517 Kalai, A. T., Nachum, O., Vempala, S. S., and Zhang, E.
- 518 Why language models hallucinate, 2025. URL <https://arxiv.org/abs/2509.04664>.
- 519
- 520 Kalavasis, A., Mehrotra, A., and Velegkas, G. On characterizations
- 521 for language generation: Interplay of hallucinations, breadth, and stability, 2025a. URL <https://arxiv.org/abs/2412.18530>.
- 522
- 523 Kalavasis, A., Mehrotra, A., and Velegkas, G. On the limits
- 524 of language generation: Trade-offs between hallucination
- 525 and mode collapse, 2025b. URL <https://arxiv.org/abs/2411.09642>.
- 526
- 527 Kasirzadeh, A. Two types of ai existential risk: Decisive and
- 528 accumulative. 2025. doi: 10.1007/s11098-025-02301-3.
- 529 URL <https://link.springer.com/article/10.1007/s11098-025-02301-3>.
- 530
- 531 Katzman, J., Wang, A., Scheuerman, M., Blodgett, S. L.,
- 532 Laird, K., Wallach, H., and Barocas, S. Taxonomizing
- 533 and measuring representational harms: A look at image
- 534 tagging, 2023. URL <https://arxiv.org/abs/2305.01776>.
- 535
- 536 Kleinberg, J. and Mullainathan, S. Language generation in
- 537 the limit. In *Advances in Neural Information Processing Systems 37 (NeurIPS 2024)*, 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/7988e9b3876ad689e921ce05d711442f-Paper-Conference.pdf.
- 538
- 539
- 540
- 541
- 542
- 543
- 544
- 545
- 546
- 547
- 548
- 549
- Kolt, N. Algorithmic black swans. 101:1177–1240, 2024.
- URL <https://wustllawreview.org/wp-content/uploads/2024/04/Kolt-Algorithmic-Black-Swans.pdf>.
- Korbak, T., Balesni, M., Barnes, E., Bengio, Y., Benton, J., Bloom, J., Chen, M., Cooney, A., Dafoe, A., Dragan, A., Emmons, S., Evans, O., Farhi, D., Greenblatt, R., Hendrycks, D., Hobbs, M., Hubinger, E., Irving, G., Jenner, E., Kokotajlo, D., Krakovna, V., Legg, S., Lindner, D., Luan, D., Mądry, A., Michael, J., Nanda, N., Orr, D., Pachocki, J., Perez, E., Phuong, M., Roger, F., Saxe, J., Shlegeris, B., Soto, M., Steinberger, E., Wang, J., Zaremba, W., Baker, B., Shah, R., and Mikulik, V. Chain of thought monitorability: A new and fragile opportunity for ai safety, 2025. URL <https://arxiv.org/abs/2507.11473>.
- Lampinen, A. K., Chan, S. C. Y., and Hermann, K. Learned feature representations are biased by complexity, learning order, position, and more, 2024. URL <https://arxiv.org/abs/2405.05847>.
- Lampinen, A. K., Chan, S. C., Li, Y., and Hermann, K. Representation biases: will we achieve complete understanding by analyzing representations? *arXiv preprint arXiv:2507.22216*, 2025.
- Lee, K.-i., Koh, H., Lee, D., Yoon, S., Kim, M., and Jung, K. Generating diverse hypotheses for inductive reasoning. In Chiruzzo, L., Ritter, A., and Wang, L. (eds.), *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 8461–8474, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. doi: 10.18653/v1/2025.nacl-long.429. URL <https://aclanthology.org/2025.nacl-long.429>.
- Lee, M. H. J. Examining the robustness of homogeneity bias to hyperparameter adjustments in gpt-4, 2025. URL <https://arxiv.org/abs/2501.02211>.
- Lehman, J. and Stanley, K. O. Novelty search and the problem with objectives. In *Genetic programming theory and practice IX*, pp. 37–56. Springer, 2011.
- Leinster, T. Entropy and diversity: The axiomatic approach, 2024. URL <https://arxiv.org/abs/2012.02113>.
- Leitão, R. P., Zuanon, J., Villéger, S., Williams, S. E., Baraloto, C., Fortunel, C., Mendonça, F. P., and Mouillot, D. Rare species contribute disproportionately to the functional structure of species assemblages. *Proceedings of the Royal Society B: Biological Sciences*, 283(1828):20160084, 2016.

- 550 Li, C., Wang, P., Wang, C., Zhang, L., Liu, Z., Ye, Q., Xu,
 551 Y., Huang, F., Zhang, X., and Yu, P. S. Loki's dance of
 552 illusions: A comprehensive survey of hallucination in
 553 large language models, 2025. URL <https://arxiv.org/abs/2507.02870>.
- 554
- 555 Li, C. T. and Farnia, F. Mode-seeking divergences: theory
 556 and applications to gans. In *International Conference
 557 on Artificial Intelligence and Statistics*, pp. 8321–8350.
 558 PMLR, 2023.
- 559
- 560 Lopez, P. Bias does not equal bias: A socio-technical typol-
 561 ogy of bias in data-based algorithmic systems. *Internet
 562 Policy Review*, 10(4):1–29, 2021.
- 563
- 564 Macar, U., Bogdan, P. C., Rajamanoharan, S., and Nanda,
 565 N. Thought branches: Interpreting llm reasoning requires
 566 resampling, 2025. URL <https://arxiv.org/abs/2510.27484>.
- 567
- 568 McQuillan, D. *Resisting AI: An Anti-fascist Approach to
 569 Artificial Intelligence*. Bristol University Press, 2022.
 570 ISBN 9781529213508. URL <https://books.google.com/books?id=R6x6EAAAQBAJ>.
- 571
- 572 Meng, T., Mehrabi, N., Goyal, P., Ramakrishna, A., Gal-
 573 styan, A., Zemel, R., Chang, K.-W., Gupta, R., and Peris,
 574 C. Attribute controlled fine-tuning for large language
 575 models: A case study on detoxification, 2024. URL
 576 <https://arxiv.org/abs/2410.05559>.
- 577
- 578 Mironov, M. and Prokhorenkova, L. Measuring diversity:
 579 Axioms and challenges, 2025. URL <https://arxiv.org/abs/2410.14556>.
- 580
- 581 Modok, A. Role of social media in inciting the genocidal
 582 acts: A case study on myanmar's rohingya. *Contempo-
 583 rary Challenges: The Global Crime, Justice and Security
 584 Journal*, 4, 2023.
- 585
- 586 Mohamed, S., Png, M.-T., and Isaac, W. Decolonial ai:
 587 Decolonial theory as sociotechnical foresight in artificial
 588 intelligence. *Philosophy & Technology*, 33(4):659–684,
 589 July 2020. ISSN 2210-5441. doi: 10.1007/s13347-020-0
 590 0405-8. URL <http://dx.doi.org/10.1007/s13347-020-00405-8>.
- 591
- 592 Moon, K., Green, A. E., and Kushlev, K. Homogenizing ef-
 593 fect of large language models (llms) on creative diversity:
 594 An empirical comparison of human and chatgpt writing.
 595 *Computers in Human Behavior: Artificial Humans*, pp.
 596 100207, 2025.
- 597
- 598 Morain, R. and Ventura, D. Is prompt engineering the cre-
 599 ativity knob for large language models? In *Proceedings
 600 of the 16th International Conference for Computational
 601 Creativity*, 2025.
- 602
- 603 Morozov, E. The ai we deserve, 12 2024. URL <https://www.bostonreview.net/forum/the-ai-we-deserve/>.
- 604
- Murthy, S. K., Ullman, T., and Hu, J. One fish, two fish,
 605 but not the whole sea: Alignment reduces language mod-
 606 els' conceptual diversity. In *Proceedings of the 2025
 607 Conference of the Nations of the Americas Chapter of
 608 the Association for Computational Linguistics: Human
 609 Language Technologies (Volume 1: Long Papers)*, pp.
 610 11241–11258. Association for Computational Linguis-
 611 tics, 2025. doi: 10.18653/v1/2025.nacl-long.561. URL
 612 <http://dx.doi.org/10.18653/v1/2025.nacl-long.561>.
- 613
- Peeperkorn, M., Kouwenhoven, T., Brown, D., and Jordanous,
 614 A. Is temperature the creativity parameter of
 615 large language models?, 2024. URL <https://arxiv.org/abs/2405.00492>.
- 616
- Peeperkorn, M., Kouwenhoven, T., Brown, D., and Jordanous,
 617 A. Mind the gap: Conformative decoding to
 618 improve output diversity of instruction-tuned large lan-
 619 guage models. *arXiv preprint arXiv:2507.20956*, 2025.
- 620
- Peterson, A. J. Ai and the problem of knowledge collapse.
 621 *AI & SOCIETY*, 40(5):3249–3269, January 2025. ISSN
 622 1435-5655. doi: 10.1007/s00146-024-02173-x. URL
 623 <http://dx.doi.org/10.1007/s00146-024-02173-x>.
- 624
- Preciado, P. *Testo Junkie: Sex, Drugs, and Biopolitics in the Pharmacopornographic Era*. G - Reference, Information and Interdisciplinary Subjects Series. Feminist Press at the City University of New York, 2013. ISBN 9781558618374. URL <https://books.google.com/books?id=8mtgAwAAQBAJ>.
- 625
- 626 Preda, A. Special report: Ai-induced psychosis: a new
 627 frontier in mental health, 2025.
- 628
- Pugh, J. K., Soros, L. B., and Stanley, K. O. Quality di-
 629 versity: A new frontier for evolutionary computation.
 630 *Frontiers in Robotics and AI*, 3:40, 2016.
- 631
- Putra, R., Kartika, A., and Santoso, B. Solving long-tail
 632 detection for autonomous vehicles. *Authorea Preprints*,
 633 2024.
- 634
- Putri, S. D. G., Purnomo, E. P., and Khairunissa, T. Echo
 635 chambers and algorithmic bias: The homogenization of
 636 online culture in a smart society. In *SHS Web of Confer-
 637 ences*, volume 202, pp. 05001. EDP Sciences, 2024.
- 638
- Reinhardt, K. Between identity and ambiguity: some con-
 639 ceptual considerations on diversity. *Symposion*, 7(2):
 640 261–283, 2020.
- 641

- 605 Rodilloso, E. Filter bubbles and the unfeeling: How ai
 606 for social media can foster extremism and polarization.
 607 *Philosophy & Technology*, 37(2):71, 2024.
- 608 Rudko, I. and Bashirpour Bonab, A. Chatgpt is incredible
 609 (at being average). *Ethics and Information Technology*,
 610 27(3):36, 2025.
- 611 Rudman, W., Chen, C., and Eickhoff, C. Outlier dimen-
 612 sions encode task-specific knowledge. *arXiv preprint*
 613 *arXiv:2310.17715*, 2023.
- 614 Ruef, M. and Birkhead, C. Learning from outliers and
 615 anomalies. *Academy of Management Perspectives*, (ja):
 616 amp–2023, 2024.
- 617 Saketopoulou, A. *Sexuality Beyond Consent: Risk,*
 618 *Race, Traumatophilia*. NYU Press, 2023. ISBN
 619 9781479820252. URL <https://books.google.com/books?id=Xb6ZEAAAQBAJ>.
- 620 Schaeffer, R., Kazdan, J., Arulandu, A. C., and Koyejo, S.
 621 Position: Model collapse does not mean what you think,
 622 2025. URL <https://arxiv.org/abs/2503.03150>.
- 623 Schwartz, R., Schwartz, R., Vassilev, A., Greene, K., Perine,
 624 L., Burt, A., and Hall, P. *Towards a standard for identify-*
 625 *ing and managing bias in artificial intelligence*, volume 3.
 626 US Department of Commerce, National Institute of Stan-
 627 dards and Technology . . . , 2022.
- 628 Selbst, A. D., Boyd, D., Friedler, S. A., Venkatasubramanian,
 629 S., and Vertesi, J. Fairness and abstraction in socio-
 630 technical systems. In *Proceedings of the conference on*
 631 *fairness, accountability, and transparency*, pp. 59–68,
 632 2019.
- 633 Shelby, R., Rismani, S., Henne, K., Moon, A., Rostamzadeh,
 634 N., Nicholas, P., Yilla, N., Gallegos, J., Smart, A., Garcia,
 635 E., and Virk, G. Sociotechnical harms of algorithmic
 636 systems: Scoping a taxonomy for harm reduction, 2023.
 637 URL <https://arxiv.org/abs/2210.05791>.
- 638 Sheth, I., Wehner, J., Abdelnabi, S., Binkyte, R., and Fritz,
 639 M. Safety is essential for responsible open-ended systems,
 640 2025. URL <https://arxiv.org/abs/2502.04512>.
- 641 Sokol, K. and Hüllermeier, E. All you need for counter-
 642 factual explainability is principled and reliable esti-
 643 mate of aleatoric and epistemic uncertainty, 2025. URL
 644 <https://arxiv.org/abs/2502.17007>.
- 645 Song, Y., Kempe, J., and Munos, R. Outcome-based
 646 exploration for llm reasoning, 2025. URL <https://arxiv.org/abs/2509.06941>.
- 647 Sourati, Z., Ziabari, A. S., and Dehghani, M. The homoge-
 648 nizing effect of large language models on human expres-
 649 sion and thought, 2025. URL <https://arxiv.org/abs/2508.01491>.
- 650 Spivak, G. C. Can the subaltern speak?, 1988.
- 651 Stanley, K. O. and Lehman, J. *Why Greatness Cannot Be*
 652 *Planned: The Myth of the Objective*. Springer Cham,
 653 2015. ISBN 978-3-319-15523-4. doi: 10.1007/978-3-3
 654 15524-1.
- 655 Sui, P., Duede, E., Wu, S., and So, R. J. Confabulation: The
 656 surprising value of large language model hallucinations,
 657 2024. URL <https://arxiv.org/abs/2406.04175>.
- 658 Sun, G., Jin, M., Wang, Z., Liang, J. C., Geng, T., Guan, Q.,
 659 Wang, Q., Du, M., Zhang, Y., Liu, D., et al. Hallucinating
 660 llm could be creative, 2025.
- 661 Tahiroglu, K. and Wyse, L. Latent spaces as platforms
 662 for sonic creativity. In *Proceedings of the 16th Interna-*
 663 *tional Conference on Computational Creativity, ICC*,
 664 volume 24, 2024.
- 665 Taleb, N. N. 'antifragility' as a mathematical idea. *Nature*,
 666 494(7438):430–430, 2013.
- 667 Thanh-Tung, H. and Tran, T. Catastrophic forgetting and
 668 mode collapse in gans. In *2020 international joint confer-*
 669 *ence on neural networks (ijcnn)*, pp. 1–10. IEEE, 2020.
- 670 Uzzi, B., Mukherjee, S., Stringer, M., and Jones, B. Atypical
 671 combinations and scientific impact. *Science*, 342(6157):
 672 468–472, 2013.
- 673 Von Hippel, E. New product ideas from 'lead users'.
 674 *Research-Technology Management*, 32(3):24–27, 1989.
- 675 Vrijenhoek, S., Daniil, S., Sandel, J., and Hollink, L. Di-
 676 versity of what? on the different conceptualizations
 677 of diversity in recommender systems. In *The 2024*
 678 *ACM Conference on Fairness Accountability and Trans-*
 679 *parency, FAccT '24*, pp. 573–584. ACM, June 2024. doi:
 680 10.1145/3630106.3658926. URL <http://dx.doi.org/10.1145/3630106.3658926>.
- 681 Wachter, S., Mittelstadt, B., and Russell, C. Why fairness
 682 cannot be automated: Bridging the gap between eu non-
 683 discrimination law and ai. *Computer Law & Security*
 684 *Review*, 41:105567, 2021. ISSN 2212-473X. doi: <https://doi.org/10.1016/j.clsr.2021.105567>. URL <https://www.sciencedirect.com/science/article/pii/S0267364921000406>.
- 685 West, P. and Potts, C. Base models beat aligned mod-
 686 els at randomness and creativity. *arXiv preprint*
 687 *arXiv:2505.00047*, 2025.

- 660 Woodward, J. *Making things happen: A theory of causal
661 explanation*. Oxford university press, 2005.
- 662 Wu, L., Wang, D., and Evans, J. A. Large teams develop
663 and small teams disrupt science and technology. *Nature*,
664 566(7744):378–382, 2019.
- 665 Yao, F., Liu, L., Zhang, D., Dong, C., Shang, J., and Gao, J.
666 Your efficient rl framework secretly brings you off-policy
667 rl training, August 2025. URL <https://fengyao.notion.site/off-policy-rl>.
- 668 Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao,
669 Y., and Narasimhan, K. Tree of thoughts: Deliberate
670 problem solving with large language models, 2023. URL
<https://arxiv.org/abs/2305.10601>.
- 671 Yazici, Y., Foo, C.-S., Winkler, S., Yap, K.-H., and Chan-
672 drasekhar, V. Empirical analysis of overfitting and mode
673 drop in gan training, 2020. URL <https://arxiv.org/abs/2006.14265>.
- 674 Yuan, S., Qu, Z., Kangen, A. Y., and Färber, M. Can hallu-
675 cinations help? boosting llms for drug discovery, 2025.
676 URL <https://arxiv.org/abs/2501.13824>.
- 677 Zhang, J., Yu, S., Chong, D., Sicilia, A., Tomz, M. R.,
678 Manning, C. D., and Shi, W. Verbalized sampling: How
679 to mitigate mode collapse and unlock llm diversity, 2025a.
URL <https://arxiv.org/abs/2510.01171>.
- 680 Zhang, Y., Diddee, H., Holm, S., Liu, H., Liu, X., Samuel,
681 V., Wang, B., and Ippolito, D. Noveltybench: Evaluating
682 language models for humanlike diversity, 2025b. URL
<https://arxiv.org/abs/2504.05228>.
- 683 Zurn, P., Pitts, A., Bettcher, T., and DiPietro, P. *Trans
684 Philosophy*. University of Minnesota Press, 2024. ISBN
685 9781452972183. URL <https://books.google.com/books?id=XWr8EAAAQBAJ>.
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707
- 708
- 709
- 710
- 711
- 712
- 713
- 714

Appendix A. Implementing generalized diversities

Our structure-aware language is intentionally *abstract* so it **admits multiple implementations**, not only the one we presented in the main paper. In this appendix, we think through two alternative choices:

1. Generalization of the structure core through the *escort power mean*
2. Reinterpretation of the deviance as *relative entropy*

Our goal with this appendix is to **inspire reflection** on diversity *beyond* what was explicitly presented in our framework.

A.1. Generalizing the structure core

Inspired by *value measures* (Leinster, 2024) and *escort distributions* (Bercher, 2011), we generalize the structure core as the *escort power mean*:

$$\langle \alpha_{i(q,r)} \rangle(x_p) = \left(\frac{\sum_{y \in \text{Str}_\top(x_p)} p(y|x_p)^r \alpha_i(y)^q}{\sum_{y \in \text{Str}_\top(x_p)} p(y|x_p)^r} \right)^{1/q} \quad (\text{A.1})$$

We simplify by considering the *escort distribution*:

$$p_{(r)}(y|x_p) = \frac{p(y|x_p)^r}{\sum_{y \in \text{Str}_\top(x_p)} p(y|x_p)^r} \quad (\text{A.2})$$

Then, the *generalized structure core* is:

$$\langle \alpha_{i(q,r)} \rangle(x_p) = \left(\mathbb{E}_{y \sim p_{(r)}(\cdot|x_p)} [\alpha_i(y)^q] \right)^{1/q} \quad (\text{A.3})$$

When $q = 1$ and $r = 1$, the generalized structure core *recovers* our original structure core in [Equation 5](#) and [Equation 9](#). Different values for q, r give us alternative interesting cores. For instance:

$$\begin{aligned} \langle \alpha_{i(1,0)} \rangle(x_p) &= \frac{1}{|\text{Str}_\top(x_p)|} \sum_{y \in \text{Str}_\top(x_p)} \alpha_i(y) \\ \langle \alpha_{i(1,\infty)} \rangle(x_p) &= \alpha_i(\arg \max_y p(y|x_p)) \\ \langle \alpha_{i(\infty,1)} \rangle(x_p) &= \max_{y \in \text{supp}(p(\cdot|x_p))} \alpha_i(y) \\ \langle \alpha_{i(-\infty,\infty)} \rangle(x_p) &= \min_{y \in \text{modes}(p(\cdot|x_p))} \alpha_i(y) \end{aligned}$$

For a given structure α_i , we can think of q selecting whether large or small compliance values dominate, and r selecting whether the *large body* or *long-tails* of $p(\cdot|x_p)$ dominate. **By parameterizing, we make transparent how we weigh rarity, signal strength, and balance.** Since different parameters reflect different viewpoints (Leinster, 2024), we shall always consider a full *diversity profile* before drawing conclusions about how our interventions impact diversity.

A.2. Reinterpreting deviance

We can think of a *generalized orientation* as:

$$\theta_{n,k}(y|x_p) = \text{orient}(\Lambda_n(y), \langle \Lambda_n \rangle(x_p)) \quad (\text{A.4})$$

with $\text{orient} : [0, 1]^n \times [0, 1]^n \rightarrow [0, 1]^k$.

Then, the *generalized deviance* is:

$$\begin{aligned} d_{n,k}(y|x_p) &= \|\theta_{n,k}(y|x_p)\|_{\text{orient}} \\ \|\cdot\|_{\text{orient}} &: [0, 1]^k \rightarrow \mathbb{R}^+ \end{aligned} \quad (\text{A.5})$$

If we choose $\text{orient}(\Lambda_x, \Lambda_y) = \Lambda_x - \Lambda_y$ and $\|\cdot\|_{\text{orient}} = \|\cdot\|_\theta$, we *recover* our original deviance in [Equation 8](#) and [Equation 9](#).

For *relative entropy*, we consider the **Rényi entropy** defined (Leinster, 2024) as:

$$H_q(\mathbf{p} \| \mathbf{r}) = \frac{1}{q-1} \log \sum_{i \in \text{supp}(\mathbf{p})} p_i^q r_i^{1-q} \quad (\text{A.6})$$

Then, we can think of a *dummy orient()* that just stores Λ_x, Λ_y and a $\|\cdot\|_{\text{orient}}$ operator that computes the *relative entropy* between them. For a given *normalized core* $\bar{\Lambda}_n = \{\langle \bar{\alpha} \rangle_1, \dots\}$ and *normalized system* $\bar{\Lambda}_n = \{\bar{\alpha}_1, \dots\}$, we define two *Hill number* (Leinster, 2024) deviances: the *excess deviance* and *deficit deviance*:

$$\partial_q^+(y, x_p) = e^{H_q(\bar{\Lambda}_n(y) \| \langle \bar{\Lambda}_n \rangle(x_p))} \quad (\text{A.7})$$

$$\partial_q^-(y, x_p) = e^{H_q(\langle \bar{\Lambda}_n \rangle(x_p) \| \bar{\Lambda}_n(y))} \quad (\text{A.8})$$

We could read ∂_q^+ as the *effective over-compliance* and ∂_q^- as the *effective under-compliance* with respect to the *normative compliance*.

For instance, as $q \rightarrow \infty$, we interpret:

- ∂_∞^+ as the largest *excess of compliance*

$$\partial_\infty^+ = \max_i \frac{\bar{\alpha}_i(y)}{\langle \bar{\alpha}_i \rangle(x_p)}$$

- ∂_∞^- as the largest *deficit of compliance*

$$\partial_\infty^- = \max_i \frac{\langle \bar{\alpha}_i \rangle(x_p)}{\bar{\alpha}_i(y)}$$

All of this to say, there are **multiple ways we can reason about structures and statistics jointly**. We encourage readers to develop alternative and competing formalisms that share our conceptual backbone: *structures* that make *context explicit*, *cores* that encode the normativity that *homogenization* push us toward, and *orientations* that capture perspectives of *non-normativity*. Above all, **we ask everyone to think deeper about diversity**.

770 Appendix B. Theoretical touchpoints

771
772 In this appendix, we explore how our theoretical framework
773 connects to other frameworks. To that purpose, we consider
774 an *unprompted* scenario of a *singleton* system with *binary*
775 compliance for its single structure:

$$776 \quad 777 \quad \Lambda_*(x) := (\alpha_*(x)) \quad \alpha_*(x) \in \{0, 1\}$$

778 Then, the structure core represents the probability of com-
779 pliance being exactly 1:
780

$$781 \quad \mu := \langle \alpha_* \rangle = \sum_{c \in \{0, 1\}} c \Pr(\alpha=c) = \Pr(\alpha=1)$$

784 Our singleton deviance is expressed as:
785

$$786 \quad \partial_*(x) = \|\alpha_*(x) - \mu\|_\theta$$

788 B.1. Expected deviance and Gini-Simpson index

790 To calculate the *expected deviance*, we consider two choices
791 for $\|\cdot\|_\theta$: absolute value and the squared ℓ_2 norm. For each,
792 we find connections between $\mathbb{E}[\partial_*]$ and the *Gini-Simpson*
793 *index* for a binary variable:
794

$$795 \quad \mathbb{E}[|\alpha_* - \mu|] = 2\mu(1 - \mu) = \text{GS}$$

$$796 \quad \mathbb{E}[\|\alpha_* - \mu\|_2^2] = \text{Var}[\alpha_*] = \mu(1 - \mu) = \frac{\text{GS}}{2}$$

797 If we interpret GS as the *degree of mixing* in outcomes, then
798 increasing the expected deviance drives *heterogeneity* rather
799 than *concentration*.
800

802 B.2. Is-It-Valid classification for Hallucinations

804 To reason about hallucinations, authors in (Kalai et al., 2025)
805 partition the space of *plausible* outputs into disjoint sets of
806 *valid outputs* V and *errors* E . In their framework, a model
807 *hallucinates* when it cannot solve the binary discrimination
808 problem *Is-It-Valid?* (IIV). Their framework can be
809 interpreted through our structure-aware language:
810

$$811 \quad \alpha_{\text{IIV}}(x) = \mathbf{1}[x \in V]$$

813 We can connect their generative hallucination rate given
814 by $\text{err} = \Pr_{x \sim \hat{p}}[x \in E] = \hat{p}(E)$ to the system core of a
815 singleton IIV system:
816

$$817 \quad \langle \alpha_{\text{IIV}} \rangle = 1 - \text{err}$$

818 The paper (Kalai et al., 2025) points out that future work
819 should "consider degrees of hallucination". Our structure-
820 aware framework provides the language to reason about
821 these desired **graded notions of hallucination**: We can
822 score a string under multiple structures, with scores encod-
823 ing real-valued nuance *beyond the binary*.
824

B.3. Language Generation in the Limit

Recent work (Kleinberg & Mullainathan, 2024; Kalavasis et al., 2025a) studies language generation where a generator G , given strings from an unknown target language K , must output strings that are both **novel** and **valid**. We can reinterpret some of their framework as a special case of our structure-aware formulation.

Given a language collection $\mathcal{L} = \{L_1, L_2, \dots\}$, we can define *membership structures* with corresponding cores that represent the probability of generating a string valid for each corresponding language:

$$\alpha_{L_i}(x) = \mathbf{1}[x \in L_i] \quad \langle \alpha_{L_i} \rangle = \Pr[y \in L_i]$$

The literature is currently (Kalavasis et al., 2025b) exploring the trade-offs between *consistency* and *breadth*. An LLM generates strings *consistent* with our target language K if:

$$\langle \alpha_K \rangle = 1 \quad \text{when} \quad \mathbb{E}[\partial_K]_{y \sim p_{\text{LLM}}} \rightarrow 0$$

An LLM generation has *breadth* when all strings of our target language $K \in \mathcal{L}$ can be generated:

$$\forall y \in K : p_{\text{LLM}}(y) > 0 \iff K \subseteq \text{supp}(p_{\text{LLM}})$$

Our structure-aware framework gives us insight that homogenization is *relative to a system*. Indeed, pushing for consistency shall not imply that we push for homogenization in every context. Generally, for $\Lambda_K \neq \Lambda_m$:

$$\mathbb{E}[\partial_K] \rightarrow 0 \neq \mathbb{E}[\partial_m] \rightarrow 0$$

Thinking explicitly through structures and systems allows us to formulate *interesting* questions (for instance, is $\Lambda_K = \Lambda_{\text{IIV}}$?) that will help us make connections between all these theoretical efforts. We present these touchpoints as **starting points for deeper exploration**.