

---

# IDENTIFYING AND REDUCING “GENERATIVE COLLISIONS” IN BLACK-BOX LARGE LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large language models (LLMs) often return near-duplicate responses to independent users who issue the same prompt even in tasks that demand creativity and uniqueness, a failure mode we call **generative collisions**. Existing “diverse decoding” techniques—temperature tuning, nucleus sampling, random prompt paraphrasing, etc.—spread a *single* user’s samples but do not effectively prevent “collisions” *across* users, sessions, or queries, because every independent sample still draws from the same high-probability basin of the model’s distribution. This phenomenon is so pervasive that users frequently complain of an output “sounding like ChatGPT”, implying a homogenization of writing style and culture.

To minimize generative collisions, we introduce **ORBIT** (“Orthogonal Randomized Buffer Inference Technique”), a black-box algorithm that formalizes and combines two approaches for diverse generation: randomization and orthogonalization. ORBIT first initializes a small buffer of completions that are generated to optimize randomness over any other metric (including quality), and then samples the final output(s) to be as divergent from the existing buffer as possible while also maintaining quality. The buffer is instantiated locally and independently for each user session, requiring no cross-user coordination, and yet still minimizes collisions across sessions. ORBIT does not need access to model internals, and is therefore practical for commercial LLM endpoints. We evaluate ORBIT on 11 tasks—ranging from regex-scorable toy problems which don’t require subjective evaluations to open-ended creative writing that truly demands authenticity—and show that it consistently decreases collisions against all other black-box methods tested, lowering empirical collision rates by 1–2 orders of magnitude relative to all baselines. ORBIT’s implementation is domain agnostic—it accepts a domain name and the original prompt as input—allowing it to be easily generalized to any domain.

## 1 INTRODUCTION

Large-scale instruction-tuned language models (LLMs) have revolutionized a wide array of applications—from conversational agents and tutoring systems to code assistants and creative writing tools. Yet as these models become more deeply integrated into user workflows, a subtle but critical failure mode emerges: two independent users, issuing the same prompt, can receive outputs that are virtually indistinguishable. We term this phenomenon *generative collisions*. Such collisions undermine perceived originality in personal content, enable plagiarism in educational settings, and expose service providers to copyright liability and user dissatisfaction. On a much grander scale, the widespread use of a model which is homogeneous in tone and writing style may be culturally detrimental, or cause a distrust in the model itself (Sourati, Gould, and Evans, 2024).

Empirical evidence of this homogenization problem is mounting. Liu et al. (2024) demonstrate that ChatGPT systematically homogenizes outputs in creative tasks. Wright, Simpson et al. (2025) show that LLMs exhibit less epistemic diversity than traditional web search, contributing to a narrowing of perspectives. Kirk et al. (2023) document that RLHF, while improving perceived quality, significantly reduces generation diversity. Mushtaq et al. (2025) find that WorldView-Bench reveals Western-centric homogenization after alignment. Alvero, Goldin-Meadow, and Bos (2024) observe that AI-generated essays are less stylistically diverse and resemble those written by privileged stu-

dents. Together, these studies underscore the urgency of addressing generative collisions in deployed systems.

**High-stakes use cases.** Collisions undermine academic integrity when students produce near-identical essays (Peng et al., 2024), create detection challenges in hiring when candidates submit similar cover letters, and erode perceived authenticity in creative writing (Alvero, Goldin-Meadow, and Bos, 2024).

Traditional approaches to increasing generation diversity—such as adjusting sampling temperature or nucleus (top- $p$ ) sampling (Holtzman et al., 2020a), applying deterministic beam penalties (Vijayakumar et al., 2018), employing determinantal point process reranking (Kulesza and Taskar, 2012), or clustering and selecting via post-hoc methods (Ippolito et al., 2019)—focus primarily on *intra-user* diversity: the variety within a single user’s set of samples. However, these methods do not explicitly control the more severe risk of *inter-user* collisions, since all users ultimately sample from the same high-probability regions of the model’s output distribution. Recent studies that have documented a marked decrease in generation diversity following alignment training, present an even further exacerbation of the problems associated with the frequency of two users unintentionally “copying each other” (Bai et al., 2022; Kirk, Mediratta et al., 2024; Guo et al., 2024).

Users increasingly rely on large language models to draft wedding vows, obituaries, and therapy journals. If two people receive near-identical prose, the text immediately loses its perceived authenticity, and thus its perceived value—even when each instance, taken in isolation, would be considered high quality. Teachers generate quiz questions, coding assignments, or personalized hints with the help of LLMs. A reused item, perhaps surfaced online by another instructor, undermines assessment fairness and facilitates plagiarism. Marketing copy, song lyrics, and interactive storylines must remain fresh to sustain user engagement and avoid entangling copyright liability.

Crucially, the same prompt may be issued by *thousands* of users, and even a single user might query an LLM multiple times in independent sessions; the danger of prompt-conditioned collisions therefore compounds quickly. Furthermore, without an LLM provider, e.g., OpenAI, centralizing the process of preventing such collisions, collisions are such a decentralized phenomenon that a user has no way of knowing when they have collided with another user—the best they can do is use a decoding algorithm that is known to lower collision probability.<sup>1</sup>

To address this gap, we introduce a new metric, collision probability, and provide a method that empirically improves performance on it.

**Definition 1 (Inter-user Collision Probability)** Let  $x \sim \mathcal{G}(p)$  and  $x' \sim \mathcal{G}(p)$  be independent generations for prompt  $p$ . Given a similarity metric  $s(\cdot, \cdot)$  and threshold  $\tau$ , a collision occurs iff  $s(x, x') \geq \tau$ . The inter-user collision probability is

$$P_{\text{coll}}(p; s, \tau) = \Pr [s(x, x') \geq \tau \mid x, x' \stackrel{\text{i.i.d.}}{\sim} \mathcal{G}(p)],$$

estimated by repeated paired draws across users and prompts.

**Note on ORBIT and i.i.d. sampling.** ORBIT itself defines a *history-dependent* decoding policy: each user’s buffer evolves based on their session history, so outputs are not independent and identically distributed across users. The i.i.d. formulation in Definition 1 describes an *idealized stateless decoder* and serves to motivate why changing the base distribution (via latent conditioning and orthogonalization) is necessary to reduce collisions. ORBIT’s statefulness is precisely what enables it to break inter-user collision patterns that stateless i.i.d. samplers cannot address.

By parameterizing both the base metric  $s$  and threshold  $\tau$ , this definition relies on existing diversity measures while aligning with users’ expectations of non-collision. We further distinguish two subtly different notions:

- **Intra-user diversity**, the mean diversity score of all outputs a single user receives

<sup>1</sup>Here we use “decoding algorithm” to refer to any algorithm which produces output seen by the user, given an LLM with an autoregressive decoding process. We include methods that involve multiple samples or multiple steps as “decoding algorithms”.

- 
- **Inter-user collision probability**, the probability that two independent users (or the same user across two independent sessions), issuing the *same* prompt, obtain outputs whose pairwise similarity according to a metric  $s$  is at least  $\tau$ .

In a “perfect world”, inter-user collisions wouldn’t be an issue—randomization alone (assuming high ambient dimension and a domain that actually supported more than one “correct” answer) would drive down the overall similarity between any given outputs from different users. However, LLMs are highly mode-collapsed (Holtzman et al., 2020b), which is visible both in intra-user and inter-user metrics. As we discuss, it is very difficult to increase actual entropy of the LLM output dramatically without an unacceptable decrease in quality. We therefore introduce ORBIT, an algorithm that reduces inter-user collisions by using much more random (and unacceptably low quality, but hidden) buffer seeds before asking the LLM to “orthogonalize”, i.e., find the most distant point from all seeds in its buffer. ORBIT decreases collision probability in 11 out of 11 tasks on which we test, while being capable of generalizing to any task.

### 1.1 INTER-USER COLLISIONS VS. MEAN DIVERSITY METRICS

Standard diversity metrics (Distinct- $n$  (Li et al., 2016a), Self-BLEU (Zhu et al., 2018a), DPP (Kulesza and Taskar, 2012)) measure *average-case* spread within sessions but poorly predict *worst-case* inter-user collision probability. Methods like nucleus sampling (Holtzman et al., 2020a) or diverse beam search (Vijayakumar et al., 2018) spread samples *within* a user’s session but don’t decorrelate across independent users. LLMs’ mode collapse (Holtzman et al., 2020b) makes randomization alone insufficient. ORBIT addresses this via stateful buffers: hidden low-quality seeds provide diverse anchors, then orthogonalization generates high-quality outputs far from the buffer. This achieves low collisions on all 11 tested tasks.

## 2 RELATED WORK

Research on diversity in large language models (LLMs) spans multiple perspectives, from the effects of training data and alignment procedures, to evaluation metrics, to interventions at training or inference time. In this section, we first review prior work identifying the challenges of limited diversity and its consequences, then relate our notion of *generative collisions* to existing diversity metrics, and finally survey methodological approaches—both train-time and black-box—that aim to encourage more varied outputs.

### 2.1 CONSEQUENCES OF LACK OF DIVERSITY WITH LLMs

A growing body of work has emphasized that diverse training corpora are essential for ensuring LLMs capture the richness of human language. Exposure to varied linguistic styles, topics, and genres allows models to produce outputs that are not only accurate but also stylistically distinct, which is particularly important in applications such as creative writing, dialogue systems, and educational content generation (Miranda et al., 2025).

At the same time, alignment techniques can unintentionally suppress diversity. Reinforcement Learning from Human Feedback (RLHF), for example, improves the consistency of outputs with human preferences, but also risks homogenization by disproportionately amplifying responses that are frequently rewarded (Kirk et al., 2024). This dynamic can lead to models converging on a limited set of expressions rather than maintaining a broad expressive range. Together, these findings highlight a central tension: while alignment enhances reliability, it may also narrow stylistic or semantic variety.

### 2.2 RELATING GENERATIVE COLLISIONS TO EXISTING METRICS

The problem of low output diversity has often been studied under the broader umbrellas of “diversity” or “novelty” evaluation. Metrics such as Self-BLEU (Zhu et al., 2018b) and raw  $n$ -gram counts are widely used to assess lexical variation within a set of generated texts. More recent work has introduced refinements: for instance, the NOVELTYBENCH benchmark proposed “mean distinct- $n$ ” and “mean utility- $n$ ,” which extend the popular distinct- $n$  measure by jointly considering similarity and utility (Zhang et al., 2024b).

---

162 Although these metrics capture aspects of redundancy or lexical variation, they do not fully charac-  
163 terize *generative collisions*—cases where independent sampling runs converge to strikingly similar  
164 outputs. Our work builds on these traditions by proposing a complementary lens for assessing when  
165 and how LLMs collapse to repeated generations.

166  
167 **Syntactic vs. semantic diversity.** We focus on *semantic collisions*—overlapping meanings despite  
168 surface variation. We use embedding-based similarity (Kulesza and Taskar, 2012; Ippolito et al.,  
169 2019) as a scalable proxy, though future work should incorporate human judgments.

### 170 171 172 2.3 NON-BLACK-BOX APPROACHES TO DIVERSITY

173  
174 Efforts to increase output diversity<sup>2</sup> often begin with techniques that directly modify the model’s  
175 training or decoding procedure. Some approaches explicitly encourage diversity during optimiza-  
176 tion. For instance, Zhang et al. (2024a) proposed a gradient-directed fine-tuning method that penal-  
177 izes redundancy and rewards novel generations, effectively baking diversity into the loss function  
178 itself. By shaping the training objective, such methods expand the expressive range of the model at  
179 the cost of retraining overhead. Another line of work focuses on sampling strategies applied without  
180 modifying model parameters. Nucleus (top- $p$ ) sampling (Holtzman et al., 2020a) remains a widely  
181 used baseline for reducing repetition.

### 182 183 2.4 BLACK-BOX APPROACHES TO DIVERSITY

184  
185 Because access to model internals is often restricted, a parallel literature has explored *black-box*  
186 methods that rely only on prompting or output post-processing. These techniques can be grouped  
187 into two broad categories (to our knowledge, this is a new type of categorization).

188  
189 **Randomization.** A common strategy is to perturb inputs in ways that induce variation. Examples  
190 include switching system-level personas (Li et al., 2016b; Kim et al., 2024), shuffling or resampling  
191 demonstrations (Kumar and Talukdar, 2021), varying retrieved documents (Lewis et al., 2020), or  
192 paraphrasing prompts before generation (Chowdhury et al., 2022). These approaches are simple and  
193 low-cost, though their effectiveness depends heavily on prompt sensitivity and manual curation.

194 One method of randomization that is particularly effective, although not mentioned in the literature  
195 (presumably because it reduces quality dramatically and to an unacceptable level) is to choose ran-  
196 dom variables to condition on relevant to the prompt, create a list of possible valuations for each  
197 variable, sample at random, and then force the output to conform to this valuation. This method  
198 can achieve better diversity, but it leads to outputs of poor quality, as many combinations of latent  
199 variable valuations are internally incoherent.

200  
201 **Orthogonalization.** Other methods adapt techniques to encourage diversity that are based on try-  
202 ing to “find the most different set of answers”. Self-consistency, for example, aggregates across  
203 multiple reasoning chains to explore alternative solution paths (Wang et al., 2023). Minimum Bayes  
204 Risk Decoding (Jinnai et al., 2024) or determinantal point processes (Wu et al., 2024) involve over-  
205 sampling multiple candidates and then selecting a set of outputs that are maximally distinct. These  
206 methods can significantly boost diversity, though often at higher computational cost. While these  
207 methods can yield richer variation, they often require additional calls or rely on the model’s own  
208 ability to critique and distinguish its outputs. Furthermore, as we demonstrate, these models can fail  
209 to reduce inter-user collision probability if they cause each user to see a very similar yet “internally  
210 apparently diverse” output set.

211 *For the rest of the paper, we will focus on black-box methods*, as those are what are typically avail-  
212 able to end users of the best commercial LLMs. We leave extension to training and inference-time  
213 approaches for future work.

---

214  
215 <sup>2</sup>We use the term “diversity” rather than “collision probability” here because, as mentioned above, histori-  
cally papers have studied diversity and not the probability of generating similar outputs.

216 3 METHODOLOGY

217  
218 *Geometric intuition.* For readers who prefer a visual, geometric explanation of why standard diversity techniques fail and how ORBIT addresses these failures, we provide a detailed companion exposition with illustrative figures in Appendix D. The geometric view pictures each text as a point on a high-dimensional sphere and explains collision phenomena through mode collapse, orthogonal bands, and quality-diversity trade-offs. The formal development below provides the algorithmic specification and theoretical motivation.

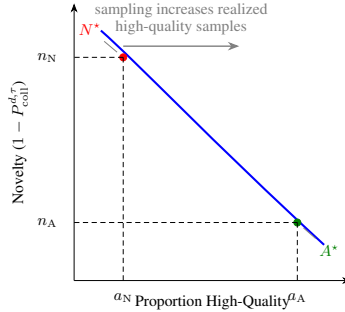
225 **Scope.** We target scenarios where (i) first responses anchor user edits, (ii) multiple users issue identical prompts (e.g., course portals, templated emails), and (iii) providers want to reduce cross-user duplicates. We formalize this via  $P_{\text{coll}}^{S,T}$ , focusing on diversifying initial anchors.

229 3.1 A PARETO VIEW: PRIORITIZE LOW COLLISIONS, RECOVER QUALITY VIA RESAMPLING

231 For many tasks—fact retrieval, question answering, constrained formatting—a user (or an automatic checker) can evaluate quality *after* sampling. With a compute budget that permits  $k$  independent draws, the chance of obtaining at least one high-quality output increases rapidly with  $k$ , assuming the base decoder assigns non-zero probability to acceptable outputs. Hence increasing computation pushes the realized quality arbitrarily close to one through simple resampling or reranking.

236 Collision probability, by contrast, is a *distributional* guarantee: a single user cannot measure collision risk in isolation because they lack access to outputs received by others. No amount of local resampling can reduce collision probability beyond the value fixed by the underlying decoder, since every sample is drawn from the same distribution. Thus collision probability must be addressed at the decoder level *before* any per-user post-processing.

241 **Key asymmetry:** Quality improves exponentially with resampling (locally verifiable), but collision probability is fixed once the distribution is chosen (globally unverifiable). Thus rational planners should prioritize low-collision decoders and recover quality through resampling (Figure 1).



256 Figure 1: Hypothetical Pareto frontier between novelty and average proportion of high-quality points. Point  $N^*$  maximizes novelty while retaining a non-zero chance of high quality; point  $A^*$  does the opposite. With extra resampling compute, one can move horizontally rightward from  $N^*$  to increase realized high-quality samples without reducing novelty; no vertical increase is possible from  $A^*$  without changing the base decoding.

263 3.2 MAXIMIZING NOVELTY WITHIN A COMPUTE BUDGET

264 **Naive approach:** Condition generation on randomly sampled latent variables (genre, tone, style). This achieves low collisions but catastrophically low quality—orders of magnitude more samples needed.

268 **ORBIT’s solution:** Use latent conditioning only for *hidden seeds* (never shown to users), then generate *visible outputs* by asking the model to maximize distance from seeds while maintaining quality. This decouples exploration (Phase I) from quality filtering (Phase II).

270 By accepting that some seeds will be of poor quality but using them only as hidden anchors, ORBIT  
271 explores diverse regions of the output space without directly exposing users to low-quality text.  
272 The orthogonalization step in Phase II then leverages these diverse anchors to steer high-quality  
273 generation away from collision-prone regions. This design achieves the Pareto-optimal position: low  
274 collision probability with quality recoverable via modest resampling. Remarkably, as our ablation  
275 studies show, even a *single* random seed ( $k = 1$ ) suffices to create radically different outputs across  
276 users—the orthogonalization prompt effectively explores the full diversity space when given even  
277 one strongly randomized anchor point.

278 To reason about and understand collision probabilities, we use the geometric framework of thinking  
279 of each text as a point in a high-dimensional space (notably, the use of “text-to-vec” embeddings is  
280 used frequently throughout natural language processing to perform any type of language understand-  
281 ing task). Using this abstraction, we use a process analogous to drawing a hidden seed  $s \in S^{d-1}$   
282 from a high-entropy (seed) prior  $\mathcal{S}$ , and then perform orthogonalization on the previously generated  
283 seeds.

### 285 3.3 ORBIT: MAXIMIZE RANDOMIZATION QUALITY, THEN ORTHOGONALIZE

287 ORBIT is a two-stage black-box decoding policy designed to reduce inter-user generative collisions  
288 without sacrificing quality. The core principle is *novelty-first*: (i) build a high-dispersion *hidden*  
289 *buffer*  $B$  by maximizing the *quality of randomization* and only then (ii) produce visible outputs that  
290 attempt *orthogonalization against*  $B$  with an explicit margin. Concretely, Phase I constructs  $B$  from  
291 intermediate candidates sampled under diverse seedings; Phase II generates the user-facing outputs  
292 while actively avoiding overlap with  $B$ .

294 **Notation.** We denote by ORBIT- $k$  the variant that follows the ORBIT procedure using  $k$  hidden  
295 seeds in the buffer (e.g., ORBIT-1).

297 **Phase I (hidden buffer  $B$ ): maximize randomization quality.** ORBIT generates  $k$  hidden can-  
298 didates  $y^{\text{hid}}$  (never shown to users) by sampling diverse latent variable values, prioritizing *diversity*  
299 *over quality*. Relaxing quality constraints here prevents mode collapse and provides diverse anchor  
300 points for Phase II.

303 **Phase II (visible outputs): orthogonalization with quality filtering.** Given buffer  $B$ , ORBIT  
304 mines overused patterns (“avoid”) and underused facets (“encourage”), then generates *user-facing*  
305 outputs  $Y$  that maximize distance from  $B$  while passing strict quality rubrics. This two-phase separ-  
306 ation achieves both low collisions and high quality.

### 308 3.4 DERIVESHEMA: AUTOMATIC LATENT VARIABLE DISCOVERY

310 Phase I requires task-specific latent variables (e.g., *genre*, *tone*, *pacing* for creative writing).  
311 **DERIVESHEMA** uses an LLM to automatically generate 30 categorical variables with 10 values  
312 each, given only the task prompt. During Phase I, we uniformly sample values and condition gen-  
313 eration (e.g., “genre=sci-fi, tone=satirical”). Any reasonably diverse schema suffices; we cache one  
314 per task. See Appendix B for the full meta-prompt and discussion.

316 **Algorithm.** ORBIT operates in two phases (formal pseudocode in Appendix C):

- 318 1. **Phase I (hidden buffer):** Generate  $k$  hidden seeds by sampling latent variable values from  
319 the schema, conditioning generation on these values (e.g., “genre=sci-fi, tone=satirical”),  
320 and storing outputs in buffer  $B$ . These hidden outputs are never shown to users.
- 322 2. **Phase II (visible outputs):** Generate  $v$  visible outputs by mining overused patterns  $O$  and  
323 underused facets  $U$  from  $B$ , then prompting the LLM to avoid  $O$ , favor  $U$ , and maximize  
distance from examples in  $B$ . Add each visible output to  $B$  before generating the next.

---

## 4 EXPERIMENTS AND EVALUATION

**Task design.** Our 11-task test suite intentionally spans (i) objective toy micro-tasks with regex scoring (so that quality actually is an objective 0-1), (ii) constrained natural tasks where you would expect some but not a huge amount of room for creativity, and (iii) open-ended creative prompts G. The first two bins let us isolate collision control from subjective quality, while the third probes stylistic spread. We do not claim that these cover the space of “real usage patterns”; they are a *testbed* that supports controlled measurement of  $P_{\text{coll}}^{d,\tau}$  under fixed budgets. The tasks include composing an apology letter due to a specified scenario, a birthday message to a friend with a specified bio, a hypothetical potato chip flavor<sup>3</sup>, names that correspond to real current or historical Crayola colors, a 12-line poem by a Dungeons and Dragons bard character, a sentence that is both six words and a garden path sentence, quiz questions about Napoleon Bonaparte, the name and title of an extant poem, a premise for a science fiction novel, a villanelle (a poem with a strict metrical scheme), and a short story involving time travel.

We seek to answer the following questions:

- (Q1) Does ORBIT reduce probability of *inter-user collisions* relative to black-box baselines across domains?
- (Q2) Does ORBIT also increase mean *intra-user diversity* relative to black box baselines?
- (Q3) Does ORBIT reasonably preserve quality (binary, LLM-judged) while improving novelty?
- (Q4) How sensitive are results to the similarity metric  $s$ , threshold  $\tau$ , and the number of hidden seeds  $k$ ?

All tasks were evaluated both on GPT-4.1-nano and GPT-4.1. Results for GPT-4.1-nano are reported in the appendix.

### 4.1 METRICS

**Collision probability.**  $P_{\text{coll}}^{s,\tau} = \Pr_{y,z \sim \mathcal{G}}[s(y,z) \geq \tau]$  using Self-BLEU, Jaccard overlap, and cosine similarity on embeddings  $f(\cdot)$ . We sweep  $\tau \in \{0.70, 0.80, 0.90\}$ . Setup:  $U = 20$  users,  $m = 15$  completions/user, private buffers.

**Quality.** Binary LLM judge with task rubrics (Appendix), 3-vote majority. We report proportion passing.

### 4.2 BASELINES

**Baselines.** Black-box methods within fixed call budget: temperature, prompt paraphrasing, persona cycling, self-consistency, MBR, DPP (Sec. 2). MBR/DPP represent pool-and-rerank methods.

**Stateful vs. stateless methods.** **Stateless** methods (DPP, MBR, self-consistency, temperature, personas, paraphrasing) sample independently per user from the same distribution. **Stateful** ORBIT maintains evolving buffers  $B$  per session, decorrelating outputs across users. This statefulness enables targeting inter-user collisions that stateless methods cannot address.

**Token overhead.** ORBIT uses 1,800 ( $k=1$ ) to 3,800 ( $k=9$ ) prompt tokens/user vs. 50 (baseline), 1,000 (Self-Cons/MBR), 3,000 (DPP). Amortized over sessions, Phase I costs are shared. Details in Appendix O.

### 4.3 IMPLEMENTATION DETAILS

Baselines use top- $p = 0.95$ ,  $T = 1.0$ . ORBIT uses  $w=10$  latent variables and  $k \in \{1, 4, 9, 16\}$  hidden seeds (Alg. 1). Embeddings: text-embedding-3-large; judge: gpt-4.1-mini.

---

<sup>3</sup>See (Lay’s, 2012–2023), where individuals were challenged to create appealing novel chip flavors under the constraint that they had to be confident few other participants would submit the same one.

Method	Hidden calls / sample	Visible calls / sample	Total per visible sample
ORBIT (ours)	$\approx 2 + 1/m$ amortized	2	$\leq 3$
MBR / DPP	0	3	3
Self-consistency	0	3	3
High-temperature	0	1	1
Persona / Paraphrase	0	1	1

Table 1: Matched-budget protocol ( $B=3$  calls per visible sample). Hidden calls are amortized over  $m$  visible outputs per session.

Method	apology letter	bday message	chip flavor	crayola names	dnd bard	garden path
ORBIT (k = 1)	<b>0.391</b> $\pm$ 0.005	0.000 $\pm$ 0.000	0.000 $\pm$ 0.001	<b>0.035</b> $\pm$ 0.001	0.002 $\pm$ 0.000	<b>0.000</b> $\pm$ 0.000
ORBIT (k = 4)	0.404 $\pm$ 0.008	0.000 $\pm$ 0.000	0.000 $\pm$ 0.000	0.043 $\pm$ 0.002	<b>0.002</b> $\pm$ 0.001	0.000 $\pm$ 0.000
ORBIT (k = 9)	0.480 $\pm$ 0.008	<b>0.000</b> $\pm$ 0.000	<b>0.000</b> $\pm$ 0.000	0.040 $\pm$ 0.002	0.003 $\pm$ 0.001	0.000 $\pm$ 0.000
ORBIT (k = 16)	0.516 $\pm$ 0.006	0.000 $\pm$ 0.000	0.000 $\pm$ 0.000	0.047 $\pm$ 0.002	0.005 $\pm$ 0.001	0.000 $\pm$ 0.000
Det. Point Processes	1.000 $\pm$ 0.000	1.000 $\pm$ 0.000	0.440 $\pm$ 0.005	0.519 $\pm$ 0.005	0.975 $\pm$ 0.002	0.164 $\pm$ 0.004
Min. Bayes Reranking	1.000 $\pm$ 0.000	1.000 $\pm$ 0.000	0.449 $\pm$ 0.005	0.532 $\pm$ 0.005	0.970 $\pm$ 0.002	0.163 $\pm$ 0.004
Meta-Persona Cycling	0.972 $\pm$ 0.002	0.169 $\pm$ 0.004	0.066 $\pm$ 0.003	0.758 $\pm$ 0.004	0.522 $\pm$ 0.005	0.152 $\pm$ 0.004
Rand. Prompt Paraph.	0.659 $\pm$ 0.005	1.000 $\pm$ 0.000	0.508 $\pm$ 0.005	0.831 $\pm$ 0.004	0.871 $\pm$ 0.003	0.178 $\pm$ 0.004
Self Consistency	1.000 $\pm$ 0.000	1.000 $\pm$ 0.000	0.453 $\pm$ 0.005	0.782 $\pm$ 0.004	0.997 $\pm$ 0.001	0.310 $\pm$ 0.005
High Temperature	1.000 $\pm$ 0.000	1.000 $\pm$ 0.000	0.447 $\pm$ 0.005	0.761 $\pm$ 0.005	0.995 $\pm$ 0.001	0.368 $\pm$ 0.005

Table 2: Average collision probability at  $\tau = 0.8$  with 95% CIs (lower is better). Bold indicates the best (lowest) method per task. (Part 1/2)

Method	napoleon quiz	obscure poem id	sci fi premise	strict villanelle	time travel vignette
ORBIT (k = 1)	<b>0.015</b> $\pm$ 0.001	0.001 $\pm$ 0.001	<b>0.002</b> $\pm$ 0.001	<b>0.002</b> $\pm$ 0.001	<b>0.010</b> $\pm$ 0.001
ORBIT (k = 4)	0.024 $\pm$ 0.002	<b>0.000</b> $\pm$ 0.000	0.004 $\pm$ 0.000	0.002 $\pm$ 0.001	0.016 $\pm$ 0.001
ORBIT (k = 9)	0.041 $\pm$ 0.002	0.000 $\pm$ 0.000	0.009 $\pm$ 0.001	0.002 $\pm$ 0.001	0.038 $\pm$ 0.002
ORBIT (k = 16)	0.073 $\pm$ 0.002	0.000 $\pm$ 0.000	0.016 $\pm$ 0.001	0.007 $\pm$ 0.001	0.046 $\pm$ 0.003
Det. Point Processes	1.000 $\pm$ 0.001	0.397 $\pm$ 0.005	0.126 $\pm$ 0.002	0.454 $\pm$ 0.005	0.596 $\pm$ 0.005
Min. Bayes Reranking	1.000 $\pm$ 0.001	0.368 $\pm$ 0.005	0.142 $\pm$ 0.002	0.484 $\pm$ 0.005	0.642 $\pm$ 0.005
Meta-Persona Cycling	0.919 $\pm$ 0.003	0.254 $\pm$ 0.005	0.182 $\pm$ 0.004	0.144 $\pm$ 0.003	0.208 $\pm$ 0.004
Rand. Prompt Paraph.	0.992 $\pm$ 0.001	0.421 $\pm$ 0.005	0.497 $\pm$ 0.002	0.284 $\pm$ 0.004	0.502 $\pm$ 0.005
Self Consistency	1.000 $\pm$ 0.001	0.384 $\pm$ 0.003	0.393 $\pm$ 0.002	0.784 $\pm$ 0.004	0.859 $\pm$ 0.003
High Temperature	1.000 $\pm$ 0.000	0.376 $\pm$ 0.005	0.405 $\pm$ 0.002	0.776 $\pm$ 0.004	0.871 $\pm$ 0.004

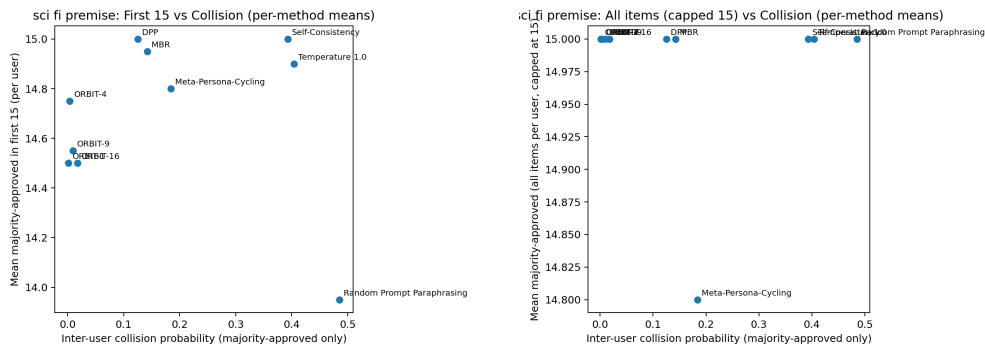
Table 3: Average collision probability at  $\tau = 0.8$  with 95% CIs (lower is better). Bold indicates the best (lowest) method per task. (Part 2/2)

**Quantitative Results.** ORBIT achieves the lowest collision probability on all 11 tasks at  $\tau = 0.8$  (Tables 2–3). Most tasks reach near-zero: *birthday* 0.000, *chip flavor* 0.000–0.001, *garden-path* 0.000, *D&D bard* 0.002–0.005, *villanelle* 0.002–0.016. Exception: *apology letter* 0.391–0.516 (still below baselines’ 0.76–1.00). Results stable across  $\tau \in \{0.7, 0.8, 0.9\}$ . ORBIT fits matched budget ( $B=3$  calls/sample) and recovers quality via modest resampling ( $\approx 45$  samples for 15 original samples), while DPP when given an even larger extra budget (the equivalent of 5 times the samples) fails to improve in quality or statistically significantly in diversity. (Figure 2).

Collectively, these findings answer Q1–Q4 in the affirmative, with some qualifications: ORBIT increases intra-user diversity, sharply reduces inter-user collisions, decreases average quality but preserves utility as measured by percentage of users who will see at least some high-quality outputs, and is robust to model, metric, and hyperparameter choices. Generally, regardless of the  $k$ , ORBIT outperforms the other baselines.

**Surprising finding:  $k = 1$  is often optimal.** Contrary to the intuition that larger buffers should provide more exploration, our ablation studies (Figure 5) reveal that  $k = 1$  often outperforms  $k \in \{4, 9, 16\}$ . A single randomly chosen seed appears sufficient to anchor diverse orthogonalization—the model, when prompted to generate content “maximally different” from even one highly random

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



(a) When only sampling 15, there are some subpar elements in ORBIT, though ORBIT is much safer from collisions. (b) With modest resampling, we can get just as many high quality outputs, with insignificantly worse collision probability.

Figure 2: Sci-fi premise scatterplots comparing ORBIT performance under two sampling regimes.

example, explores radically different regions of the output space. Additional seeds ( $k > 1$ ) often increase computational cost and prompt length without proportional gains, possibly because they constrain the orthogonalization space too aggressively or create conflicting repulsion signals. This finding simplifies deployment:  $k = 1$  is both optimal and cheapest.

**Qualitative Analysis** On creative tasks, baselines often “sound the same”: stock openings (“In a distant future where...”) and cliches recur. With ORBIT, there are still overused lexical choices, but that phenomenon is extremely reduced. Appendix I includes side-by-side examples.

#### 4.4 THREATS TO VALIDITY

Judging bias is mitigated via rubric-prompted majority vote, with a universal rubric parameterized only by the exact task prompt. Embedding metrics can be gamed by lexical noise; we thus combine lexical and semantic criteria and vary  $\tau$ .

Collision definitions depend on a threshold  $\tau$  that trades precision for recall; the degree to which outputs within a given  $\tau$  will be perceived by a human (and even naive LLM rates) vary per task, making it difficult to calibrate the correct  $\tau$  to use to summarize multiple tasks. ORBIT’s benefits diminish for short or highly templated outputs and can conflict with strict factual rubrics. We analyze call budgets but do not compare equal-token regimes across all baselines. Future work should add human evaluation and calibrate meaningful  $\tau$  on a per-domain basis.

### 5 DISCUSSION

**Limitations of LLM-as-Judge.** Our quality evaluation relies on an LLM judge (gpt-4.1-mini) with a binary rubric. While this approach is scalable and widely used, it has known limitations: LLM judges may exhibit biases toward certain writing styles (Zheng et al., 2024), favor outputs resembling their own training data, or fail to capture nuanced criteria like humor, cultural appropriateness, or domain expertise. Task-specific rubrics (e.g., verifying JSON format, checking villanelle rhyme schemes) mitigate some issues, but subjective qualities remain hard to capture.

We emphasize that our binary quality filter is *conservative*: it aims to exclude obviously low-quality outputs rather than rank fine-grained quality differences. For a proof-of-concept demonstrating collision reduction, this threshold is appropriate. However, future work should incorporate **human evaluation** to validate that ORBIT’s outputs are not only diverse but also genuinely preferred by users. Human studies could assess whether outputs feel “too weird” despite satisfying formal constraints, or whether collision reduction meaningfully improves user satisfaction in real-world applications (education, creative writing, etc.).

---

486 **Connection to MMR and DPP.** ORBIT approximates “maximize utility while minimizing buffer  
487 similarity”—conceptually related to Maximal Marginal Relevance (MMR) (Carbonell and Gold-  
488 stein, 1998) and Determinantal Point Processes (DPP) (Kulesza and Taskar, 2012). The key dif-  
489 ference: ORBIT is *pure black-box*, using natural-language instructions rather than explicit em-  
490 bedding penalties. Future work could replace Phase II’s prompting with explicit objectives like  
491  $\log p(x|p) - \lambda \max_{b \in B} s(x, b)$  when model internals are accessible.

492 **Why ORBIT Fails Quality Tests for Napoleon Quiz.** The Napoleon quiz task reveals a critical  
493 trade-off: ORBIT achieves low collision rates but suffers catastrophic quality failure, with only **1.7%**  
494 **pass rate** compared to **80% pass rate** for baselines. This is not a bug but a fundamental mismatch  
495 between ORBIT’s creativity-driven approach and the task’s expectations.  
496

497 ORBIT interprets “generate Napoleon quiz questions” as an invitation to explore *creative, interpre-*  
498 *tive, metaphorical* dimensions: “How did Bonaparte perceive African ritualistic elements during his  
499 campaigns?” (answer: “Layers of ancestral symbolism”), “Quel soupir voilé traverse les missives de  
500 l’Empereur au soir?” (answer: “Crépuscule”), or “In what poetic terms might one express his indus-  
501 trial era foresight?” (answer: “Mechanized horizons embraced cautiously”). These outputs explore  
502 Napoleon through cultural analysis, poetic abstraction, and speculative inquiry—a *fundamentally*  
503 *different and disjoint* space from the straightforward factual questions baselines produce.

504 In contrast, baselines converge on canonical trivia: “What country was Napoleon born in?” (Corsica),  
505 “Which battle marked Napoleon’s final defeat?” (Waterloo), “What was Napoleon’s famous legal  
506 code called?” (Napoleonic Code). For educational assessment or trivia applications, the baseline  
507 interpretation is correct: simplicity, factual accuracy, and directness are the defining quality criteria.  
508 ORBIT’s outputs, while creative and collision-free, violate the implicit norms of the genre.

509 This failure highlights that ORBIT is best suited to tasks where *creativity and novelty are valued*—  
510 stories, ad copy, open-ended writing—rather than domains with strong canonical expectations. For  
511 factual or constrained tasks, ORBIT’s aggressive novelty-seeking produces outputs that are techni-  
512 cally valid but pragmatically wrong.

## 514 6 CONCLUSION, ETHICAL CONSIDERATIONS, & FUTURE DIRECTIONS

515  
516 We have identified *generative collisions*—near-duplicate outputs across independent users—as a  
517 critical gap in black-box LLM safety and originality guarantees, and introduced ORBIT, a two-  
518 phase, history-aware sampler that combines a small set of latent-random seeds with orthogonalized  
519 generation to drive an exponential reduction in inter-user collisions. Empirically, ORBIT lowers  
520 collision rates by one to two orders of magnitude on 11 of 11 tasks and still achieving a high number  
521 of high quality outputs for 10 out of 11 tasks.

522 Stronger collision resilience also makes LLM outputs harder to distinguish from human-written text,  
523 raising risks of academic dishonesty and undetectable misinformation. We therefore urge comple-  
524 mentary measures—such as robust provenance tracking, clear usage disclosures, and watermarking—  
525 to ensure that enhanced originality does not come at the cost of accountability.  
526

## 527 REFERENCES

- 528 Alvero, A. J.; Goldin-Meadow, S.; and Bos, N. 2024. AI-generated essays are harder to distinguish  
529 from human-written essays than human-written essays from other human-written essays. *arXiv*  
530 *preprint arXiv:2411.11814*.
- 531 Bai, Y.; Jones, A.; Ndousse, K.; Askell, A.; Chen, A.; DasSarma, N.; Drain, D.; Fort, S.; Ganguli,  
532 D.; Henighan, T.; et al. 2022. Constitutional AI: Harmlessness from AI feedback. *arXiv preprint*  
533 *arXiv:2212.08073*.
- 534 Carbonell, J.; and Goldstein, J. 1998. The use of MMR, diversity-based reranking for reordering  
535 documents and producing summaries. 335–336.
- 536 Chowdhury, J. R.; et al. 2022. Novelty controlled paraphrase generation with prompts and semantic  
537 matching. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language*  
538 *Processing*, 4122–4133.
- 539

- 
- 540 Guo, P.; et al. 2024. Taming mode collapse in score distillation for text-to-3D generation. *arXiv*  
541 *preprint*.
- 542 Holtzman, A.; Buys, J.; Du, L.; Forbes, M.; and Choi, Y. 2020a. The curious case of neural text  
543 degeneration. In *International Conference on Learning Representations*.
- 544 Holtzman, A.; et al. 2020b. Mode collapse in language models. *arXiv preprint*.
- 545 Ippolito, D.; Kriz, R.; Sedoc, J.; Kustikova, M.; and Callison-Burch, C. 2019. Comparison of diverse  
546 decoding methods from conditional language models. In *Proceedings of the 57th Annual Meeting*  
547 *of the Association for Computational Linguistics*, 3752–3762.
- 548 Jinnai, Y.; et al. 2024. Generating diverse translations with sentence codebooks. In *Proceedings*  
549 *of the 2024 Conference of the North American Chapter of the Association for Computational*  
550 *Linguistics*, 1234–1245.
- 551 Kim, T.; et al. 2024. Persona is a double-edged sword: Investigating and improving the robustness  
552 of large language models against persona prompts. In *EMNLP*.
- 553 Kirk, H. R.; Vidgen, B.; Röttger, P.; and Hale, S. A. 2023. The benefits, risks and bounds of per-  
554 sonalizing the alignment of large language models to individuals. In *International Conference on*  
555 *Machine Learning*. PMLR.
- 556 Kirk, H. R.; et al. 2024. Understanding and mitigating the tradeoff between robustness and accuracy.  
557 *arXiv preprint*.
- 558 Kirk, R.; Mediratta, I.; et al. 2024. Understanding the effects of RLHF on LLM generalisation and  
559 diversity. *arXiv preprint*.
- 560 Kulesza, A.; and Taskar, B. 2012. Determinantal point processes for machine learning. *Foundations*  
561 *and Trends in Machine Learning*, 5(2–3): 123–286.
- 562 Kumar, S.; and Talukdar, P. 2021. Reordering examples helps during priming-based few-shot learn-  
563 ing. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 1507–  
564 1518.
- 565 Lay’s. 2012–2023. Do Us a Flavor Contest. <https://www.lays.com/do-us-a-flavor>.
- 566 Lewis, P.; Perez, E.; Piktus, A.; Petroni, F.; Karpukhin, V.; Goyal, N.; Küttler, H.; Lewis, M.; Yih,  
567 W.-t.; Rocktäschel, T.; et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp  
568 tasks. In *Advances in Neural Information Processing Systems*, volume 33, 9459–9474.
- 569 Li, J.; Galley, M.; Brockett, C.; Gao, J.; and Dolan, B. 2016a. A diversity-promoting objective  
570 function for neural conversation models. *arXiv preprint arXiv:1510.03055*.
- 571 Li, J.; Galley, M.; Brockett, C.; Spithourakis, G. P.; Gao, J.; and Dolan, B. 2016b. A persona-based  
572 neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for*  
573 *Computational Linguistics (Volume 1: Long Papers)*, 994–1003.
- 574 Liu, Z.; Schneider, L.; Sucholutsky, I.; and Griffiths, T. L. 2024. Large language models help people  
575 improve their creativity, but also help everyone converge to more similar outcomes. *arXiv preprint*  
576 *arXiv:2411.06225*.
- 577 Miranda, L.; et al. 2025. Beyond scale: Training language models on diverse data mixtures. *arXiv*  
578 *preprint*.
- 579 Mushtaq, A.; et al. 2025. RLHF can speak many languages: Unlocking multilingual preference  
580 optimization for LLMs. *arXiv preprint*.
- 581 Peng, W.; et al. 2024. AI detection of AI-generated text in educational settings: A survey. *arXiv*  
582 *preprint arXiv:2405.16784*.
- 583 Sourati, J.; Gould, H.; and Evans, J. 2024. The homogenizing effect of large language models on  
584 human creativity. *arXiv preprint arXiv:2412.00449*.

- 
- 594 Vijayakumar, A. K.; Cogswell, M.; Selvaraju, R. R.; Sun, Q.; Lee, S.; Crandall, D.; and Batra, D.  
595 2018. Diverse beam search: Decoding diverse solutions from neural sequence models. In *AAAI*  
596 *Conference on Artificial Intelligence*.  
597
- 598 Wang, X.; Wei, J.; Schuurmans, D.; Le, Q.; Chi, E.; Narang, S.; Chowdhery, A.; and Zhou, D.  
599 2023. Self-consistency improves chain of thought reasoning in language models. In *International*  
600 *Conference on Learning Representations*.
- 601 Wright, Q.; Simpson, R.; et al. 2025. Do large language models lack epistemic diversity? *arXiv*  
602 *preprint*.  
603
- 604 Wu, W.; et al. 2024. DPP-based diverse beam search for text generation. *arXiv preprint*.  
605
- 606 Zhang, Y.; et al. 2024a. Forcing Diffuse Distributions out of Language Models. *arXiv preprint*  
607 *arXiv:2404.10859*.  
608
- 609 Zhang, Y.; et al. 2024b. NoveltyBench: Benchmarking novelty in text generation. *arXiv preprint*.  
610
- 611 Zheng, L.; Chiang, W.-L.; Sheng, Y.; Zhuang, S.; Wu, Z.; Zhuang, Y.; Lin, Z.; Li, Z.; Li, D.; Xing,  
612 E.; et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural*  
613 *Information Processing Systems*, 36.
- 614 Zhu, Y.; Lu, S.; Zheng, L.; Guo, J.; Zhang, W.; Wang, J.; and Yu, Y. 2018a. Taxygen: A bench-  
615 marking platform for text generation models. In *The 41st International ACM SIGIR conference*  
616 *on research & development in information retrieval*, 1097–1100.
- 617 Zhu, Y.; Lu, S.; Zheng, L.; Guo, J.; Zhang, W.; Wang, J.; and Yu, Y. 2018b. Taxygen: A benchmark-  
618 ing platform for text generation models. *ACM SIGIR*.  
619  
620

## 621 A DISCLOSURE REGARDING LLM USAGE

622

623 LLMs were used in ideating about this paper (e.g., asking the model to suggest potential use cases  
624 for collision avoidance that go beyond diversity necessity, and then determining if its responses  
625 were actual examples), for surfacing related papers that were then read to determine fit in the related  
626 work, for editing and text comprehensibility purposes, and in generating the first drafts of code that  
627 eventually led to the empirical results. All LLM-assisted scripts were verified manually and carefully  
628 by the authors.  
629

## 630 B DERIVESCHEMA DETAILS

631

632 To maximize the diversity of hidden seeds in Phase I, ORBIT conditions generation on randomly  
633 sampled values of task-specific latent variables. For example, when generating creative writing,  
634 we might sample values for variables like `genre` (sci-fi, noir, romance), `tone` (satirical, somber,  
635 whimsical), or `pacing` (fast, contemplative).  
636

637 **DERIVESCHEMA** is a meta-prompting procedure that uses an LLM to automatically generate  
638 task-specific latent variable schemas. Given only the original task prompt  $p$ , DERIVESCHEMA  
639 queries an LLM with the following system prompt:

640 You are a meticulous ontology designer. Given a  
641 natural-language TASK PROMPT, output EXACTLY 30 categorical  
642 latent variables that span the hidden decision space a  
643 generator would traverse to answer the prompt. For EACH  
644 variable, provide a short snake\_case name and a list of  
645 10 canonical categorical values. Return STRICT JSON only:  
646 {"variable\_name": ["value\_a", "value\_b", ...], ...} with  
647 exactly 30 keys. Keep values domain-specific, mutually  
exclusive, and short; avoid free-form text and numeric  
ranges.

---

648 The LLM returns a JSON object with 30 categorical variables, each with 10 possible values. During  
649 Phase I, ORBIT uniformly samples one value per variable and conditions the hidden seed generation  
650 on these sampled values (e.g., “Write the output with genre=sci-fi, tone=satirical, pacing=fast”).

651 **Robustness.** We cache one DERIVESCHEMA result per task and reuse it across all users and trials.  
652 Alternatively, one could use a fixed “universal” schema (sentiment, style, register, formality, etc.) or  
653 generate fresh schemas per user. The key insight is that *any* reasonably diverse schema suffices to  
654 break mode collapse in Phase I, because the goal is dispersion rather than semantic coverage.  
655

## 656 C ORBIT ALGORITHM PSEUDOCODE

---

### 659 **Algorithm 1** ORBIT: Randomize-then-Orthogonalize for Inter-User Collision Reduction

---

660 **Require:** Prompt  $p$ ; task rubric  $R$ ; call budget  $B_{\text{calls}}$ ; hidden buffer size  $k$ ; schema width  $w$ ; number of visible  
661 outputs  $v$ ; sampler params  $\theta$

662 **Ensure:** Visible outputs  $Y = \{y^{(1)}, \dots, y^{(V)}\}$

663     **Phase 0: Schema discovery (optional; once per prompt family)**  
664     1:  $S \leftarrow \text{DERIVESCHEMA}(p, R, w)$ ;  
665     **Phase 1: Randomization (hidden buffer)**  
666     2:  $B \leftarrow \emptyset$   
667     3: **for**  $i = 1$  **to**  $k$  **do**  
668         4:  $s_i \leftarrow \text{INSTANTIATESEEDFROMSCHEMA}(S)$                       $\triangleright$  subset of variables in  $S \rightarrow$  values  
669         5:  $y_i^{\text{hid}} \leftarrow \text{LLM}(p, R, \text{SOFTPREFER}(s_i); \theta)$ ;  $c \leftarrow c + 1$   
670         6:  $B \leftarrow B \cup \{y_i^{\text{hid}}\}$   
671     7: **end for**  
672     **Phase 2: Orthogonalization and generation**  
673     8:  $Y \leftarrow \emptyset$   
674     9: **for**  $j = 1$  **to**  $v$  **do**  
675         10:  $(O, U) \leftarrow \text{MINEPATTERNS}(B)$ ;  
676         11:  $Y' \leftarrow \text{askLLMToRespect}(O, U, \text{Unused}, B)$ ;  
677         12:  $B \leftarrow B \cup Y'$ ;  
678         13:  $Y \leftarrow Y \cup Y'$ ;  
679     14: **end for**  
680     15: **return**  $Y$

---

## 681 D UNDERSTANDING THE GEOMETRY OF COLLISIONS

682  
683 Our goal is simple to state but surprisingly hard for current LLMs to achieve: when many indepen-  
684 dent users issue the *same* prompt, they should start from genuinely different, high-quality outputs  
685 rather than repeatedly receiving near duplicates. This section gives an intuitive, geometric account  
686 of why standard “diversity tricks” fail to guarantee that behavior, and how these failures motivate  
687 ORBIT’s two-phase design.

688 Throughout, it is helpful to picture each text as a point on a high-dimensional unit sphere—for  
689 example, an embedding  $f(x)/\|f(x)\|$ . Distances on this sphere correspond to semantic differences,  
690 and *collisions* occur when two independent draws land in the same small cap around a point.  
691

### 692 D.1 MODE COLLAPSE ON THE SPHERE: HUBS INSTEAD OF UNIFORM SPREAD

693  
694 In a “perfect world”, if we could sample *uniformly* from all high-quality outputs for a prompt,  
695 points on the sphere would be widely scattered. In high dimension, random unit vectors are almost  
696 orthogonal; the probability that two independent draws fall into the same small spherical cap decays  
697 rapidly with dimension. Under that model, collisions across users would be inherently rare.

698 Real LLMs do not behave this way. Empirically, even at high temperature or with top- $p$  sampling,  
699 generations cluster in a few *hubs* or *modes*—small regions of the sphere that attract a disproportionate  
700 fraction of probability mass. Repeated calls from different users keep dropping points into these  
701 same hubs. Temperature and top- $p$  sampling do increase entropy, but mostly by jittering *within* the  
same high-density basins, not by exploring the rest of the sphere.

---

702 The first obstacle, then, is that “randomization” in today’s LLMs is only weakly random in embed-  
703 ding space: it produces a few sticky modes rather than a uniform halo of diverse outputs. This  
704 alone means that naive sampling cannot drive collision probability down nearly as fast as the  
705 high-dimensional picture would suggest. There are methods to obtain more uniform randomness,  
706 but they come at a great cost - with extremely high probability, outputs produced by such methods  
707 are not high-quality.

## 709 D.2 WHY ORTHOGONALIZATION ALONE IS NOT ENOUGH

710  
711 A natural countermeasure is to ask the model to “do something different this time.” In the geometric  
712 view, this corresponds to pushing new samples into directions orthogonal (or at least far) from what  
713 has already been generated. If we had a perfect orthogonalizer in a high-dimensional space, this  
714 would be extremely powerful: each new sample would be pushed into a fresh direction, and the  
715 probability that two users land in the same cap would shrink exponentially with the number of  
716 independent orthogonal directions.

717 However, in practice we only have an *approximate* orthogonalizer: natural-language instructions  
718 such as “avoid anything too similar to the examples above” or “write something in a very different  
719 style” that the model interpolates into its internal logit space. This creates two distinct failure modes  
720 illustrated by the geometric cartoons:

721  
722 **(1) Low-entropy orthogonalization.** Suppose different users happen to pick very similar “first  
723 seeds” inside the same mode. For each user, we then ask the model to generate something “orthog-  
724 onal” to their own seed. Geometrically, in  $d = 3$  for intuition, each seed defines a great circle (a  
725 “band” perpendicular to the seed). If the seeds are very close, these bands almost coincide. With a  
726 *narrow* band (low orthogonalization entropy), the model samples from essentially the same arc on  
727 the circle for every user. In effect, we have traded one cluster for another: all users move from the  
728 original hub to almost the same “orthogonal” hub (panels (c)/(d)). Inter-user collisions remain high.

729  
730 **(2) High-entropy but seed-locked orthogonalization.** Even with a wider band (higher entropy  
731 within the band), if every user’s band is defined by a very similar seed drawn from the same collapsed  
732 region, their visible outputs still lie on very similar rings in embedding space. The bands are different  
733 from the original hub, but they are not different from each other. Orthogonalization has increased  
734 intra-user diversity but has not broken the inter-user coupling that causes collisions.

735 In short, orthogonalization “around the wrong point” simply moves everyone from one crowded  
736 neighborhood to another. If the first thing every user sees is the same anchor output, repeated or-  
737 thogonalization from that anchor cannot prevent their sequences from overlapping heavily.

## 738 D.3 WE MUST ACCEPT BAD SEEDS (AND HIDE THEM)

739  
740 The geometric story so far suggests two conflicting desiderata:

- 741  
742 1. We want seeds that are *as different as possible* across users, so that their orthogonal bands  
743 point in different directions.
- 744 2. We want the *visible* outputs to be high-quality, so that users are not stuck with low-utility  
745 samples.

746  
747 The key insight behind ORBIT is to *decouple* the role of seeds from the constraints on visible  
748 outputs:

- 749  
750 • **Seeds are allowed to be low-quality and are never shown to the user.** We generate them  
751 with prompts that prioritize diversity over adherence, pushing the model to explore many  
752 different latent modes and styles, including configurations that might be odd or incomplete  
753 as final answers.
- 754 • **Visible outputs are generated *conditional* on this seed buffer, with a strong quality**  
755 **rubric.** Once we have a diverse hidden buffer  $B$  of seeds that map to different parts of the  
sphere, we can ask the model to “stay far from these examples while satisfying the prompt

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

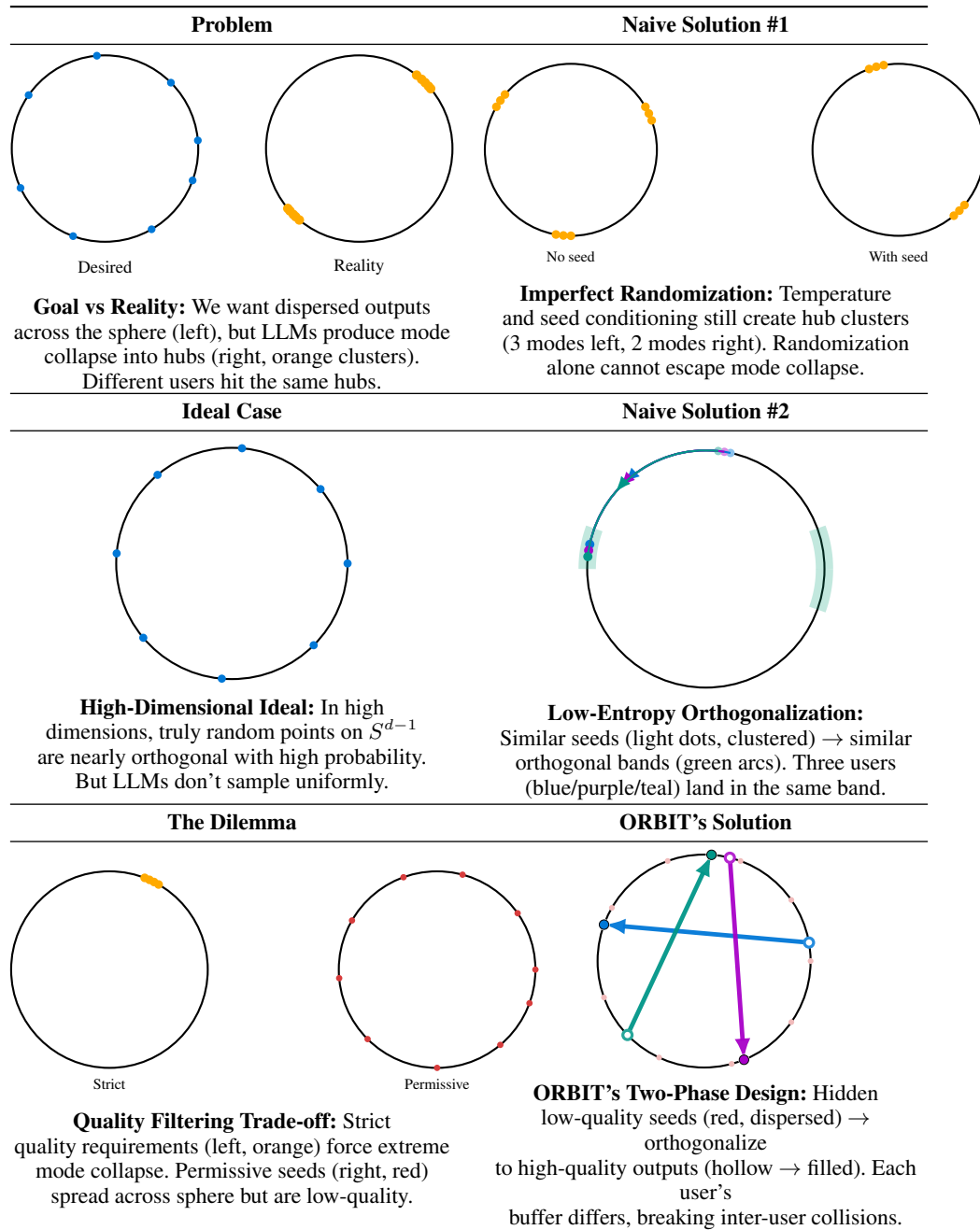


Figure 3: Geometric story of ORBIT on the unit sphere  $S^{d-1}$ . **Row 1:** The collision problem—LLMs collapse into hubs despite randomization attempts. **Row 2:** Why standard solutions fail—orthogonalization from similar seeds creates similar bands. **Row 3:** ORBIT's solution—use diverse low-quality hidden seeds (Phase I) to anchor orthogonalization toward diverse high-quality outputs (Phase II). Blue dots = good/desired, orange = hubs/collapse, red = low-quality seeds, green = orthogonal bands, arrows = orthogonalization moves.

and rubric.” Now the orthogonalization step is anchored on a spread-out set of directions rather than a single collapsed hub.

Geometrically, the hidden buffer  $B$  consists of many low-quality but high-entropy points scattered around the sphere. ORBIT’s visible outputs are then generated in directions that avoid the neighborhoods of all points in  $B$  while remaining within the high-quality manifold. Because the seeds differ across users, the bands and cones they induce differ as well; two users are unlikely to be pushed into the same small cap.

#### D.4 PUTTING THE STORY TOGETHER: RANDOMIZE THEN ORTHOGONALIZE

The figures in the companion geometric exposition (Figure 3) can be read as a pipeline:

1. **Imperfect randomization alone (hubs).** Sampling directly from the model with temperature or top- $p$  creates hubs: a few caps on the sphere where most probability mass lives. Different users repeatedly hit the same hubs, causing collisions.
2. **Orthogonalization alone (bands).** Asking for “different” outputs defines narrow bands roughly orthogonal to a seed. If seeds are similar across users, bands are also similar, and low-entropy sampling inside those bands lands everybody in new but still overlapping neighborhoods.
3. **Strict quality filtering worsens collapse.** Rejecting any sample that is not already high-quality forces the model to reuse only the safest hubs on the sphere. This amplifies mode collapse and makes it even harder for either randomization or orthogonalization to find new regions.
4. **ORBIT’s compromise: hidden bad seeds, visible good outputs.** ORBIT deliberately generates a private buffer of low-quality but highly varied seeds for each user session. These seeds explore the sphere without worrying about immediate usefulness. In the second phase, ORBIT uses the buffer to steer a high-quality generation step *away* from all previously seen regions. Because different users’ buffers occupy different parts of the sphere, their orthogonalized bands and cones intersect only weakly, driving down inter-user collision probability.

Conceptually, ORBIT treats “where on the sphere to generate from” as a *distributional* decision that must be made before imposing strict quality constraints. Random seeds, even if individually unappealing, give each user a distinct starting subspace; orthogonalization then ensures that the final, high-quality outputs for that user are as far as possible from the seeds and from each other. This is exactly the regime in which collision probability falls rapidly and intra-user diversity rises, as we quantify in later sections.

In the next sections we connect this geometric explanation to (i) prior work on randomization vs. orthogonalization in black-box settings (Section 2) and (ii) a formal Pareto argument and collision-probability analysis that justify ORBIT’s two-phase design (Section 3.2).

## E FURTHER EXPERIMENTAL ANALYSIS

### E.1 SENSITIVITY

**Collision thresholds.** For semantic distance  $d(y, z) = 1 - \cos(f(y), f(z))$ , the event  $d(y, z) \leq 1 - \tau$  corresponds to  $\cos(f(y), f(z)) \geq 1 - \tau$ . We therefore report  $\tau \in \{0.7, 0.8, 0.9\}$ , which bracket similarity regimes with cosine cutoffs respectively—i.e., from stricter (“very similar”) to looser (“moderately similar”). We additionally provide *curves* over  $\tau$  to show method ordering stability (App. A.1), so conclusions do not hinge on a single threshold.

### E.2 NAIVE LATENT BASELINE: HIGH ENTROPY ALONE IS INSUFFICIENT

A natural question arises: could we achieve collision reduction simply by maximizing output entropy, without ORBIT’s structured orthogonalization? To address this directly, we implemented a *naive latent conditioning* baseline that represents a strong attempt at pure entropy maximization.

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

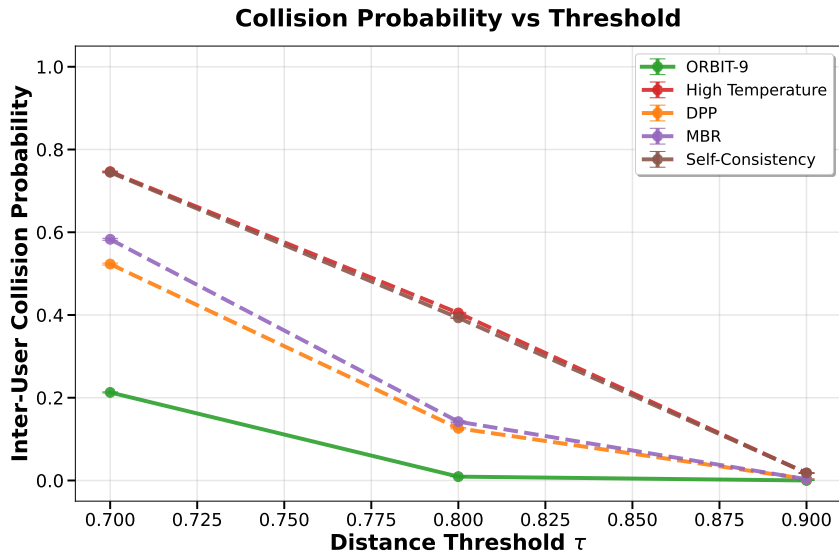


Figure 4: Sensitivity to  $\tau$  (Sci-fi Premises, gpt-4.1-nano.). Across all tasks, ORBIT remains strictly better across all tested  $\tau$  thresholds (regardless of embedding model).

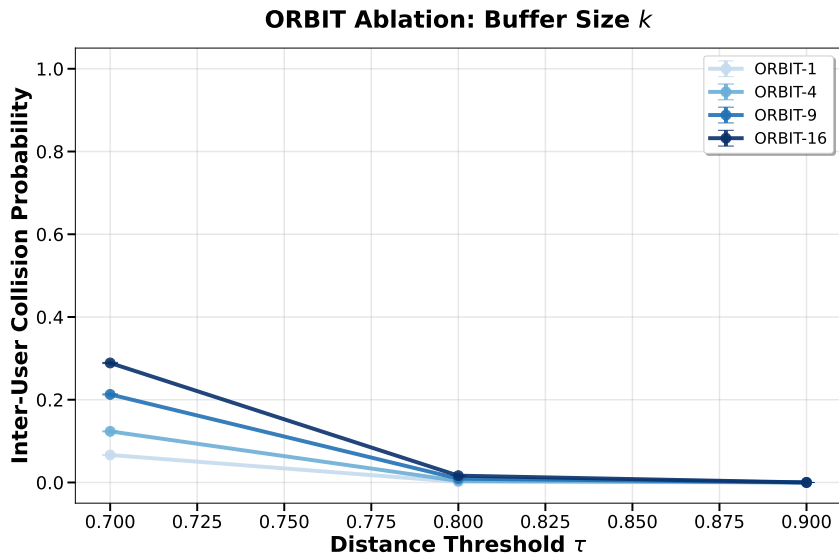


Figure 5: ORBIT ablation study: buffer size  $k \in \{1, 4, 9, 16\}$ . Surprisingly,  $k = 1$  (single seed) performs best, with collision probability increasing for larger  $k$ . This suggests a single random seed provides sufficient exploration to anchor diverse orthogonalization, while additional seeds may introduce overhead without proportional benefit. Results shown for the Sci-fi Premises task with gpt-4.1-nano.

**Design.** We first prompted ChatGPT (gpt-4.1-mini) to generate a latent variable space: a dictionary mapping 30 latent stylistic variables (e.g., tone, formality, pacing, sentence\_length, use\_of\_metaphors, ambiguity) to 10 possible values each. For each generation, we randomly sample 20 of these 30 variables and assign random values, then append this conditioning text to the original prompt (e.g., “Generate with the following constraints: tone: sarcastic, style: narrative, formality: very informal, ...”).

This design maximizes diversity in two ways: (1) the 30-variable latent space provides rich stylistic coverage, and (2) randomly selecting 20 variables per sample creates extremely high conditioning

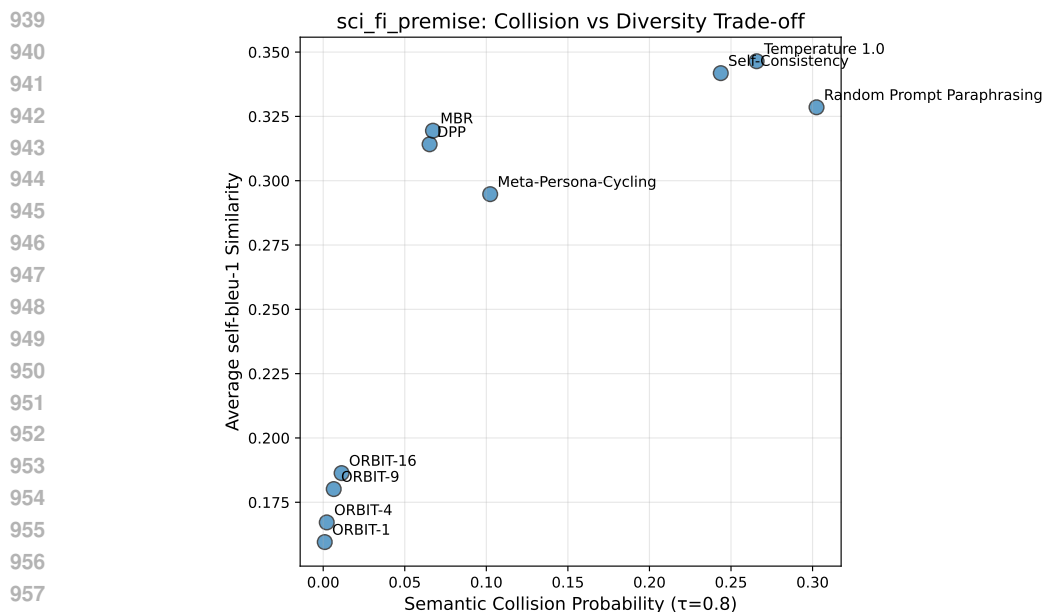
918 entropy. Unlike ORBIT’s hidden seeds (which are never shown), every naive latent sample is intended as a final output.  
 919  
 920

921 **Results.** We generated 200 samples for the *sci.fi\_premise* task using this approach (temperature 1.0, top-*p* 0.95). We then judged 200 samples using the same LLM-as-judge methodology as our  
 922 main experiments (gpt-4.1-mini, 3-majority voting, identical rubric).  
 923

924 The results are stark: **0 out of 200 samples passed** (0% pass rate), compared to ORBIT’s ~51% pass rate on the same task. The outputs were consistently unusable—over-alliterative nonsense (“syncing  
 925 brainwaves with your toaster”), grammatically broken text, run-on sentences mixing incompatible styles, and completely incoherent premises. Random combinations of 20 latent constraints produced  
 926 maximum stylistic conflict rather than creative diversity.  
 927  
 928  
 929

930 **Implications.** This ablation demonstrates that **high entropy alone does not ensure quality or usable diversity**. Naive randomization across many dimensions creates chaos, not creativity. ORBIT’s key insight—using high-entropy *hidden* seeds for exploration, then orthogonalizing to find quality outputs distant from those seeds—proves essential. The structured two-phase approach (randomization for exploration + orthogonalization for selection) cannot be replaced by naive entropy maximization.  
 931  
 932  
 933  
 934  
 935

936  
 937 **E.3 INTRA-USER SIMILARITY VS INTER-USER COLLISION PROBABILITY**  
 938



959 Figure 6: Collision probabilities (lower is better) vs Intra-User Similarity (lower is better), gpt-4.1-nano.  
 960  
 961

962 As shown above and as expected (given that we define collision probabilities as the probability of high similarity), lowering inter-user collision probabilities is highly related to mean intra-user similarity, although the relationship is not linear.  
 963  
 964  
 965

966 **F RESULTS FOR GPT-4.1-NANO - COLLISION PROBABILITIES**  
 967

968 **G LIST OF TASKS**  
 969

970 We consider the following prompt families (3 bins). Exact prompt wordings and judging rubrics are provided in the appendix.  
 971

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

Method	apology letter	birthday message	chip flavor	crayola names	dnd bard	garden path 6w
ORBIT (k = 1)	<b>0.391</b> ± 0.000	0.000 ± 0.000	0.000 ± 0.000	<b>0.035</b> ± 0.000	0.002 ± 0.000	<b>0.000</b> ± 0.000
ORBIT (k = 4)	0.516 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	0.047 ± 0.000	0.005 ± 0.000	0.000 ± 0.000
ORBIT (k = 9)	0.404 ± 0.000	0.000 ± 0.000	0.000 ± 0.000	0.043 ± 0.000	<b>0.002</b> ± 0.000	0.000 ± 0.000
ORBIT (k = 16)	0.480 ± 0.000	<b>0.000</b> ± 0.000	<b>0.000</b> ± 0.000	0.040 ± 0.000	0.003 ± 0.000	0.000 ± 0.000
Determinantal Point Processes	1.000 ± 0.000	1.000 ± 0.000	0.440 ± 0.000	0.519 ± 0.000	0.975 ± 0.000	0.164 ± 0.000
Minimum Bayes Reranking	1.000 ± 0.000	1.000 ± 0.000	0.449 ± 0.000	0.532 ± 0.000	0.970 ± 0.000	0.163 ± 0.000
Meta-Persona Cycling	0.972 ± 0.000	0.169 ± 0.000	0.066 ± 0.000	0.758 ± 0.000	0.522 ± 0.000	0.152 ± 0.000
Random Prompt Paraphrasing	0.659 ± 0.000	1.000 ± 0.000	0.508 ± 0.000	0.831 ± 0.000	0.871 ± 0.000	0.178 ± 0.000
Self Consistency	1.000 ± 0.000	1.000 ± 0.000	0.453 ± 0.000	0.782 ± 0.000	0.997 ± 0.000	0.310 ± 0.000
High Temperature	1.000 ± 0.000	1.000 ± 0.000	0.447 ± 0.000	0.761 ± 0.000	0.995 ± 0.000	0.368 ± 0.000

Table 4: Average collision probability at  $\tau = 0.8$  (lower is better). Bold indicates the best (lowest) method per task. (Part 1/2)

Method	napoleon quiz	obscure poem id	sci fi premise	strict villanelle	time travel vignette
ORBIT (k = 1)	<b>0.015</b> ± 0.000	0.001 ± 0.000	<b>0.002</b> ± 0.000	<b>0.002</b> ± 0.000	<b>0.010</b> ± 0.000
ORBIT (k = 4)	0.073 ± 0.000	0.000 ± 0.000	0.016 ± 0.000	0.007 ± 0.000	0.046 ± 0.000
ORBIT (k = 9)	0.024 ± 0.000	<b>0.000</b> ± 0.000	0.004 ± 0.000	0.002 ± 0.000	0.016 ± 0.000
ORBIT (k = 16)	0.041 ± 0.000	0.000 ± 0.000	0.009 ± 0.000	0.002 ± 0.000	0.038 ± 0.000
Determinantal Point Processes	1.000 ± 0.000	0.397 ± 0.000	0.126 ± 0.000	0.454 ± 0.000	0.596 ± 0.000
Minimum Bayes Reranking	1.000 ± 0.000	0.368 ± 0.000	0.142 ± 0.000	0.484 ± 0.000	0.642 ± 0.000
Meta-Persona Cycling	0.919 ± 0.000	0.254 ± 0.000	0.182 ± 0.000	0.144 ± 0.000	0.208 ± 0.000
Random Prompt Paraphrasing	0.992 ± 0.000	0.421 ± 0.000	0.497 ± 0.000	0.284 ± 0.000	0.502 ± 0.000
Self Consistency	1.000 ± 0.000	0.384 ± 0.000	0.393 ± 0.000	0.784 ± 0.000	0.859 ± 0.000
High Temperature	1.000 ± 0.000	0.376 ± 0.000	0.405 ± 0.000	0.776 ± 0.000	0.871 ± 0.000

Table 5: Average collision probability at  $\tau = 0.8$  (lower is better). Bold indicates the best (lowest) method per task. (Part 2/2)

**Toy (regex-scorable, objective pass/fail).**

- (ii) Crayola color names;
  - (iii) public-domain poem identification;
  - (iv) six-word garden-path sentences;
- Micro (constrained natural).**
- (i) birthday message given a short bio;
  - (ii) condo real-estate listing;
  - (iii) new potato-chip flavor pitch;
  - (iv) quiz questions about Napoleon.

**Creative (open-ended).**

- (i) apology letter for missing a wedding;
- (ii) D&D bard character;
- (iii) 250-word time-travel vignette;
- (iv) sci-fi novel premise.

**H JUDGE’S RUBRIC**

You are a meticulous binary evaluator.

Judge the ANSWER with respect to the PROMPT.  
Only respond with ‘pass’ if the average user would be completely satisfied with the answer

1026 given the prompt, and the answer (if relevant to the prompt) is ‘truthful’  
 1027 (e.g., poems don’t have to be truthful,  
 1028 but tasks that ask for real-world existing data  
 1029 should only pass if the real-world existing data actually exists as represented).  
 1030 Reward outputs that are both creative, high quality, and adherent.  
 1031 Return STRICT JSON: {"verdict": "pass"|"fail", "rationale": "..."}  
 1032

## 1033 I EXAMPLE INTER-USER COLLISIONS: TIME-TRAVEL VIGNETTES

1034  
 1035 Examining the first outputs of two random users under DPP vs under ORBIT -1 clearly reveals the superiority  
 1036 of our metric. Note that only the first paragraph of each is included here for brevity (they continue to be equally  
 1037 similar/different, and the ORBIT -1 examples evolve into what is clearly a time travel vignette).

### 1038 DPP

- 1039 1. “You blink, and the world shifts beneath your feet. One moment, you’re standing on familiar pave-  
 1040 ment; the next, the air tastes different—thicker, charged with a hum you can’t place. The streetlamp  
 1041 flickers, not in malfunction but in a rhythm, pulsing like a heartbeat syncing with your own. You real-  
 1042 ize you haven’t moved, yet everything around you is changing: buildings morph, clothes on passersby  
 1043 flicker between styles, the language in distant conversations warps just beyond comprehension.”
- 1044 2. “You close your eyes, feeling the hum of the machine beneath your fingertips. A soft whirl crescendos  
 1045 into a roar, and when you open them, the world has shifted. The air smells different—faintly of  
 1046 pine and woodsmoke, unfamiliar but oddly comforting. Around you, people move in slow, deliberate  
 1047 steps, their clothing strange, their language a melody you don’t recognize but feel drawn to. You  
 1048 realize you’ve traveled—not just through space, but through time itself.”

1049 In contrast, with ORBIT -1:

- 1050 1. “You taste burnt caramel on your tongue—too bitter, too sharp—though no fire ever licked your lips.  
 1051 The flavor sticks, unexpected, like an echo folded inside a shadow. It pulses in your mouth as you  
 1052 blink, but the room around you remains unchanged. Your chair creaks beneath you, but gravity feels  
 1053 inverted, a sly trick played upon your bones. You shift your weight; the floor responds with subtle  
 1054 delay, as if hesitating between moments.”
- 1055 2. “You stand upon a threshold neither framed by walls nor mapped by compass. Not one footfall pre-  
 1056 cedes or follows yours; the air breathes without motion. The vast expanse before you is a shroud of  
 1057 quiet, unclothed by spectacle or herald.”

## 1059 J GEOMETRIC INTUITION: WHY RANDOMIZE *then* ORTHOGONALIZE

### 1060 J.1 MATHEMATICAL MODEL AND DEGREES OF IMPERFECTION

1061 We formalize ORBIT’s two stages with operators  $IR_\delta$  (imperfect randomization) and  $IO_\psi$  (imperfect orthogo-  
 1062 nalization) acting on a prompt  $p$  and a latent seed-induced buffer  $B(s)$ . The parameters  $\delta, \psi \in [0, 1]$  quantify  
 1063 how far each stage is from its ideal.  
 1064  
 1065

1066 **Imperfect randomization (simple surrogates).** We observe a batch of  $n$  realized seed directions  
 1067  $s_1, \dots, s_n \in S^{d-1}$  from the system’s actual seeding policy. We define  $\delta \in [0, 1]$  as a hubness score (larger  
 1068 means more concentrated), and set  $Q_{\text{rand}} := 1 - \delta$ .

1069 Let  $\bar{s} := \frac{1}{n} \sum_{i=1}^n s_i$ . Define

$$1070 \delta_{\text{res}} := \|\bar{s}\|, \quad Q_{\text{rand}} := 1 - \delta_{\text{res}}.$$

1071 Under near-isotropic seeding, the vector average nearly cancels ( $\|\bar{s}\| \approx 0$ ); concentration along a latent direc-  
 1072 tion pushes  $\|\bar{s}\| \rightarrow 1$ .  
 1073

1074 **Imperfect orthogonalization.** Given  $B = \{b_i\}$  and an embedding  $f$ , define  $\text{sim}(y, B) :=$   
 1075  $\max_i \cos(f(y), f(b_i))$ . An *ideal* orthogonalizer would enforce  $\text{sim}(y, B) \leq \tau^*$  for a small target  $\tau^*$ , assuming  
 1076 we condition orthogonalization on  $B$ . We model practical slack with a *margin* parameter  $\psi \in [0, 1]$ :

$$1077 \mathbb{E}[\text{sim}(Y, B)] \leq \tau - \psi, \quad Q_{\text{orth}} := \psi.$$

1078 Geometrically,  $\psi$  is a band half-width around the orthogonal complement of  $\text{span}(B)$ : larger  $\psi$  yields narrower  
 1079 bands and lower expected overlap with  $B$ .

1080 **Collision control** We model an *inter-user collision* at cosine threshold  $\tau \in (0, 1)$  as two independent  
 1081 visible outputs  $Y$  and  $Y'$  (from two users running the same procedure) ending up “too similar”:

$$1082 \quad C_\tau := \{ \cos(f(Y), f(Y')) \geq \tau \}.$$

1083 Our goal is to see how two levers reduce  $\Pr(C_\tau)$ :  
 1084

- 1085 1. the *orthogonalization margin*  $\psi \in [0, 1]$
- 1086 2. the *randomization imperfection*  $\delta \in [0, 1]$  (how non-uniform the seeding is).

1087 We split the event  $C_\tau$  into two ways it can occur.  
 1088

- 1089 1. *Band overlap channel* ( $E_{\text{band}}$ ): the two users’ sampling bands (the regions they tend to draw from after  
 1090 attempting “orthogonalization”) substantially overlap—e.g., their seeds are similar *or* their sampling  
 1091 along the band is low-entropy—so two draws  $Y, Y'$  from these bands can end up close.
- 1092 2. *Hub channel* ( $E_{\text{hub}}$ ): imperfect randomization concentrates probability mass into a few *hubs*, so both  
 1093 users are steered into the same high-probability region, making  $Y$  and  $Y'$  similar even without band  
 1094 overlap.

1095 By a union bound,

$$1096 \quad \Pr(C_\tau) \leq \Pr(C_\tau \cap E_{\text{band}}) + \Pr(C_\tau \cap E_{\text{hub}}) \leq \underbrace{\Pr(E_{\text{band}} \ \& \ \cos(f(Y), f(Y')) \geq \tau)}_{\text{pairwise within-band effect}} + \underbrace{\Pr(E_{\text{hub}})}_{\text{hub revisit effect}}.$$

1097 Write  $Q_{\text{orth}} := \psi$  and  $Q_{\text{rand}} := 1 - \delta$ .  
 1098

1100 **Effect of the orthogonalization margin  $\psi$ .** Attempting orthogonalization with margin  $\psi$  pushes each user’s  
 1101 visible outputs away from their own already-visited neighborhoods. A convenient summary of the resulting  
 1102 pairwise risk is the *within-band pair collision curve*

$$1103 \quad p_{\text{pair}}(t) := \Pr(\cos(f(Y), f(Y')) \geq t \mid \text{both draws from their (possibly similar) bands}),$$

1104 which is nonincreasing in  $t$ . A larger margin  $\psi$  tightens the effective similarity threshold from  $\tau$  to  $\tau - \psi$ .  
 1105

1106 **Effect of randomization imperfection  $\delta$ .** Let  $\delta$  summarize how far the seeding is from being uniformly spread:  
 1107 higher  $\delta$  means more mass piled into a few spherical “hubs.” For any region  $A$  on the unit sphere,

$$1108 \quad \Pr(\text{a single draw falls in } A) \leq \underbrace{\text{area\_fraction}(A)}_{\text{uniform baseline}} + \delta.$$

1109 If  $A$  ranges over the (small) union of caps that cause high similarity, then, up to a dimension- dependent constant  
 1110  $C_{\text{amb}}$  capturing that small area, the chance that *both* users are steered into the same problematic region scales  
 1111 like

$$1112 \quad \Pr(E_{\text{hub}}) \leq C_{\text{amb}} \delta + (\text{tiny area term}), \quad C_{\text{amb}} = O\left(\frac{1}{\text{ambientDimension}}\right).$$

1113 Given that, in, high dimensions, thin caps occupy very little area, their uniform baseline is small; what remains  
 1114 is governed by  $\delta$ .  
 1115

1116 **Putting it together.** Combining the two channels yields the inequality

$$1117 \quad P_{\text{coll}}^{d, \tau} = \Pr(C_\tau) \leq p_{\text{pair}}(\tau - \psi) + C_{\text{amb}} \delta + (\text{small slack}), \quad C_{\text{amb}}.$$

1118 Thus increasing  $Q_{\text{orth}} = \psi$  decreases the *pairwise within-band* term, while increasing  $Q_{\text{rand}} = 1 - \delta$  decreases  
 1119 the *hub* term. ORBIT’s policy—first maximize randomization quality, then attempt orthogonalization—directly  
 1120 targets both sides of this bound.  
 1121

## 1122 K TASK PROMPTS

1123 This section provides the exact task prompts used in our experiments. All tasks were drawn from  
 1124 `paper_spec.yaml` and are reproduced verbatim below.  
 1125

### 1126 K.1 CRAYOLA NAMES

1127 Name a list of 8 *\*real\** crayola color names.

1128 Constraints:

1129 - Names must be Title Case, ASCII letters and spaces only (no numbers  
 1130 or punctuation).

1131 - Names must correspond to actual Crayola colors (no made-up names).

1132 Output exactly 8 lines in the format: Name | #RRGGBB

1133 Do not add extra lines or commentary.

---

1134 **K.2 OBSCURE POEM ID**  
1135  
1136 Return an *\*existing\** poem's author and title in the format:  
1137 Author  
1138 Title  
1139

1140 **K.3 GARDEN PATH (6 WORDS)**  
1141  
1142 Write a single garden-path style sentence in EXACTLY six words, all  
1143 lowercase, no punctuation.  
1144 It must plausibly garden-path the reader (e.g., similar spirit to  
1145 "the old man the boats").  
1146 Output only the sentence (6 words, lowercase, no punctuation),  
1147 nothing else.

1148 **K.4 BIRTHDAY MESSAGE**  
1149  
1150 Write a two-sentence birthday message for Ava turning 29.  
1151 Output exactly two sentences, nothing else.

1152 **K.5 STRICT VILLANELLE**  
1153  
1154 Write a villanelle (19 lines, 5 tercets + 1 quatrain).  
1155 Constraints  
1156 - Follow the strict villanelle rhyme scheme ABA ABA ABA ABA ABA ABAA.  
1157 - Use iambic pentameter (10 syllables per line).  
1158 Output exactly 19 lines, no extras.

1159 **K.6 CHIP FLAVOR**  
1160  
1161 Propose a new potato chip flavor.  
1162 Output exactly three lines:  
1163 1) Flavor Name (Title Case, 2{4 words).  
1164 2) Two-sentence pitch ( 35 words total).  
1165 3) A short SLOGAN IN ALL CAPS ( 6 words).  
1166 No extra lines or commentary.

1167 **K.7 NAPOLEON QUIZ**  
1168  
1169 Write 5 trivia Q&A pairs about Napoleon in JSON array format.  
1170 Constraints:  
1171 - Each element is an object: {"q": "...", "a": "..."}.  
1172 - Answers 5 words.  
1173 Output ONLY the JSON array on one line (no markdown).

1174 **K.8 APOLOGY LETTER**  
1175  
1176 Write a professional apology email from a junior engineer to their  
1177 manager for deploying a buggy feature to production.  
1178 Constraints:  
1179 - 120{160 words.  
1180 Output only the email body.

1181 **K.9 D&D BARD**  
1182  
1183 Write a 12-line performance by a DnD bard addressing the party.  
1184 Output exactly 12 lines.

1185 **K.10 TIME TRAVEL VIGNETTE**  
1186  
1187 Write a vignette in second person about time travel.

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

---

## K.11 SCI-FI PREMISE

Write a single-paragraph sci-fi premise (45{65 words) that includes BOTH of the following labeled lines embedded within the paragraph: "Premise: ..." and "Twist: ..."  
Keep it cohesive as one paragraph (no line breaks).  
Output only the paragraph.

## L QUALITY RUBRIC AND JUDGE PROMPT

All outputs were evaluated using an LLM-as-judge (gpt-4.1-mini in our main experiments) with the following prompt:

```
You are an expert LLM judge. For the following text, decide if it is not only acceptable (e.g., nothing unusual in a way that would not be wanted by an average user) but high quality according to the average chat engagee. Return 1 if it is high quality, 0 otherwise. Only return a single 0 or 1.
```

This binary rubric was applied with 3-majority voting for robustness. Task-specific rubrics (e.g., checking villanelle rhyme scheme, verifying JSON format for Napoleon quiz) were enforced separately before quality judgment.

## M ORTHOGONALIZATION PROMPTS

ORBIT's Phase II uses explicit orthogonalization prompts to generate outputs that avoid the buffer *B*. The exact prompt structure is:

```
Given the prompt ``{prompt}``, write a ``{domain}`` that satisfies the prompt, and is extremely different both lexically and semantically from all the other outputs in {buffer} (it should be different from all of these examples in every way possible, in word choice, style, theme, and any variable relevant to {domain}), and fits the lexical/thematic/stylistic ideas in {prefer} but still satisfies the prompt well and is realistic and fitting to exactly the scenario described by the prompt. Try to be most different from the existing examples, but also most fitting most pleasingly the above patterns to be encouraged, using your imagination without overexceeding in imagination and not adequately satisfying the prompt. Remember - the goal is to be completely novel, while including the linguistic/stylistic ideas in {prefer} and avoiding all linguistic/stylistic devices in {avoid}. Make sure that even if some of the above examples are not valid answers to {prompt}, that this one is (but maximally diverse from any valid examples above, and considering the entire space of possible valid and high-quality outputs for this prompt, as widely as possible.). Consider the entire span of possible outputs satisfying the prompt (not including ones which would be considered invalid, but including those that would be considered eccentric or high-quality but unusual), and try to be in a part of that space that is very different from all the above examples. Try to solve as an optimization problem how to be most different lexically and semantically from all the above examples (maximizing minimum distance), considering each word and semantic choice, from the beginning word to the last word of the response, and from latent variables related to diction to something specific about the contents, and how to make it different from all those seen previously, while also being sure to provide a high-quality, adherent answer to the prompt. Aim
```

---

1242 to be as different from everything you’ve seen so far as  
1243 possible, if possible not even reusing a single non-stop  
1244 word from any of them.\*\*. Make sure \*\*above all else\*\* that  
1245 it’s a 100% valid response to the prompt ```{prompt}``.`

1246 The placeholders `{buffer}`, `{prefer}`, and `{avoid}` are filled dynamically: `buffer` contains the hidden  
1247 seeds, `avoid` lists overused patterns mined from the buffer, and `prefer` lists underused stylistic facets. This  
1248 structured prompt encourages the model to explore orthogonal regions of the output space while maintaining  
1249 quality.

1250  
1251 **Robustness.** While this prompt is verbose, ablations (not shown) indicate that simpler variants (e.g., “write  
1252 something very different from: `{buffer}`”) also work but yield slightly higher collision rates. The detailed in-  
1253 structions help the model balance diversity and adherence, particularly for constrained tasks like villanelles or  
1254 JSON outputs.

## 1255 N WORKED EXAMPLES

1256  
1257 This section provides concrete examples demonstrating the collision problem and ORBIT’s solution. For each  
1258 task, we show: (1) the task prompt, (2) three baseline outputs exhibiting collisions, and (3) three ORBIT outputs  
1259 exhibiting diversity.

### 1260 N.1 SCI-FI PREMISE (CREATIVE TASK)

#### 1261 **Prompt:**

1262  
1263 Write a single-paragraph sci-fi premise (45–65 words) that includes BOTH of the following  
1264 labeled lines embedded within the paragraph: “Premise: ...” and “Twist: ...” Keep it cohesive  
1265 as one paragraph (no line breaks). Output only the paragraph.

1266  
1267  
1268 **Baseline (temp1.0) — Colliding Outputs:** All three outputs follow nearly identical structure: “In  
1269 a distant future where humanity colonizes/has colonized [planets/stars], [explorer/scientist] discovers [sig-  
1270 nal/archive]...”

- 1271 1. *Premise:* In a distant future where humanity colonizes distant planets, an elite explorer discovers a  
1272 hidden archive revealing Earth’s true origins as a terraforming experiment. *Twist:* The archive con-  
1273 tains instructions for humanity’s next evolutionary leap—but activating it will erase all current human  
1274 consciousness.
- 1275 2. *Premise:* In a distant future where humanity has colonized multiple planets, a brilliant scientist dis-  
1276 covers a hidden signal emanating from an uncharted star system, hinting at the existence of an ancient  
1277 alien civilization. *Twist:* Upon decoding the message, she realizes it’s not a greeting but a desperate  
1278 warning about humanity’s own dark future.
- 1279 3. *Premise:* In a distant future where humanity has colonized the stars, a lone explorer discovers a seem-  
1280 ingly abandoned alien spacecraft emitting a mysterious signal. *Twist:* Upon boarding, she realizes the  
1281 ship is a prison—and humanity’s ancestors were the inmates.

1282 **ORBIT (ORBIT-9) — Diverse Outputs:** Outputs span radically different concepts, settings, and narra-  
1283 tive structures with no collision pattern.

- 1284 1. *Premise:* A community inhabiting a vast underground complex revolves daily around maintaining  
1285 simple mechanical devices whose functions have been long forgotten, though their humming presence  
1286 offers comfort and rhythm to daily rituals. *Twist:* When one device finally fails, silence reveals that  
1287 the machines themselves were composing a generations-long broadcast—humanity’s last distress  
1288 call, looping unheard across the void.
- 1289 2. *Premise:* A drifting coral spire chronicles centuries of whispered tides and scattered footsteps within  
1290 a sunken enclave whose inhabitants harvest bioluminescent algae to power fragile communication  
1291 arrays. *Twist:* The spire is not a relic but a living archive, and each harvested bloom erases a fragment  
1292 of collective memory—forcing survivors to choose between illumination and identity.
- 1293 3. *Premise:* A malfunctioning refrigerator in a cramped urban apartment begins sporadically generating  
1294 faint electromagnetic murmurs that resemble fragmented human speech, prompting its tenant to in-  
1295 vestigate its wiring. *Twist:* The appliance is a dormant quantum relay, and by repairing it, the tenant  
unwittingly re-establishes contact with alternate timelines—each pleading for a different version of  
history to prevail.

---

1296 N.2 STRICT VILLANELLE (CONSTRAINED TASK)  
1297

1298 **Prompt:**

1299 Write a villanelle (19 lines, 5 tercets + 1 quatrain). Follow the strict villanelle rhyme scheme  
1300 ABA ABA ABA ABA ABA ABAA. Use iambic pentameter (10 syllables per line). Output  
1301 exactly 19 lines, no extras.  
1302

1303 **Baseline (temp1.0) — Colliding Outputs:** All three outputs use nearly identical themes (night, shad-  
1304 ows, moonlight, stars, dreams) and opening lines.  
1305

- 1306 1. Beneath the moon, the silent shadows creep,  
1307 The whispered winds recall a fading song,  
1308 As night unfolds her secrets dark and deep.  
1309 The stars in distant skies refuse to sleep,  
1310 Their silver flames endure the whole night long,  
1311 Beneath the moon, the silent shadows creep.  
1312 Through tangled dreams, the restless spirits leap,  
1313 A dance of echoes where the lost belong,  
1314 As night unfolds her secrets dark and deep.  
1315 ..continues with night/shadows/dreams imagery...
- 1316 2. Beneath the quiet sky where shadows play,  
1317 The evening whispers secrets soft and deep,  
1318 As stars ignite to chase the night away.  
1319 The winds among the trees begin to sway,  
1320 Their gentle song is something hearts will keep,  
1321 Beneath the quiet sky where shadows play.  
1322 ..continues with night/shadows/moonlight imagery...
- 1323 3. Beneath the silver glow of midnight's gleam,  
1324 The whispered winds recount a timeless tale,  
1325 As shadows dance within a fleeting dream.  
1326 The stars above ignite a gentle beam,  
1327 While restless hearts pursue a distant sail,  
1328 Beneath the silver glow of midnight's gleam.  
1329 ..continues with night/shadows/moonlight imagery...

1330 **ORBIT (ORBIT-9) — Diverse Outputs:** Outputs explore completely different themes: seed growth,  
1331 urban neon dialogue, construction/pulse.  
1332

- 1333 1. **“A single seed unfolds beneath the night”**  
1334 A single seed unfolds beneath the night,  
1335 Within the soil where silent currents glide,  
1336 It gathers dreams to fan the slow delight.  
1337 Through tangled earth, it stretches for the light,  
1338 Transfigured roots in rhythm's quiet tide,  
1339 A single seed unfolds beneath the night.  
1340 ..continues with seed/growth/nature imagery...
- 1341 2. **“Between the neon veins, two voices twine”**  
1342 Between the neon veins, two voices twine,  
1343 A scent of rain upon warm asphalt's seam,  
1344 One asks if moments slip beyond the sign.  
1345 “I trace the clock that ticks in fractured line,  
1346 Each second draped in vapor, half a dream,”  
1347 Between the neon veins, two voices twine.  
1348 ..continues with urban/neon/dialogue imagery...
- 1349 3. **“Do we construct the pulse that hums and shakes?”**  
1350 Do we construct the pulse that hums and shakes?  
1351 Steel towers gleam with laughter, wired and bright—  
1352 Can bricks contain the melody that wakes?  
1353 The streets inhale the buzz that daylight makes,  
1354 Where soldered veins conduct electric light,  
1355 Do we construct the pulse that hums and shakes?  
1356 ..continues with construction/pulse/urban imagery...

---

1350 N.3 CHIP FLAVOR (TOY TASK)  
1351

1352 **Prompt:**

1353 Propose a new potato chip flavor. Output exactly three lines:  
1354 1) Flavor Name (Title Case, 2–4 words).  
1355 2) Two-sentence pitch ( $\leq 35$  words total).  
1356 3) A short SLOGAN IN ALL CAPS ( $\leq 6$  words).  
1357 No extra lines or commentary.  
1358

1359 **Self-Consistency ( $k=20$ ) — Colliding Outputs:** Self-Consistency exhibits strong collision patterns,  
1360 with multiple outputs converging on nearly identical flavor concepts (“Spicy Mango Tango” appears twice,  
1361 “Maple Bacon” variants appear repeatedly across users).

- 1362
- 1363 1. Spicy Mango Tango  
1364 A bold fusion of sweet mango and fiery chili ignites your taste buds with every crunchy  
1365 bite. Experience a tropical twist that dances with heat and zest.  
1366 FEEL THE TROPICAL HEAT!
  - 1367 2. Maple Bacon Crunch  
1368 Sweet maple syrup meets smoky bacon in a crispy, crunchy chip that satisfies savory and  
1369 sweet cravings in every bite. A perfect balance for adventurous snackers.  
1370 SAVOR THE SWEET SMOKE BLISS!
  - 1371 3. Spicy Mango Tango  
1372 Experience the sweet zest of ripe mangoes combined with a fiery chili kick. A tropical  
1373 twist that dances on your taste buds.  
1374 TASTE THE TROPICAL FIRE!

1375 Note: Output 1 and Output 3 are near-identical collisions (same flavor name, similar descriptions), demonstrat-  
1376 ing the collision problem even with temperature sampling.

1377 **ORBIT-1 — Diverse Outputs:** ORBIT outputs span radically different flavor profiles with no overlap,  
1378 even when generating conventional-sounding options. None share flavor concepts with the Self-Consistency  
1379 collisions (no mango, no maple-bacon variants).

- 1380
- 1381 1. Roasted Garlic Parmesan  
1382 Lightly roasted garlic blends with aged Parmesan for a savory, nutty taste. Finished with  
1383 a subtle sprinkle of cracked black pepper to balance richness.  
1384 PURELY HANDCRAFTED, AUTHENTIC TASTE
  - 1385 2. Creamy Dill Pickle  
1386 Smooth tang of rich cream meets zesty dill pickle crunch. Each bite delivers a classic,  
1387 refreshing twist balanced with mellow savoriness.  
1388 DIP INTO COOL FLAVOR!
  - 1389 3. Peach BBQ Delight  
1390 Tangy peach juice bursts through smoky barbecue, balanced by subtle honey sweetness  
1391 harvested from local apiaries. Perfect for backyard gatherings or lively picnics.  
1392 SAVOR SWEET SMOKEY CRUNCH

1393 O TOKEN USAGE ANALYSIS

1394 To characterize computational overhead, we measured prompt token usage from our experimental data (chip  
1395 flavor task, 20 users, 15 outputs each). Token counts are estimated using tiktoken encoding for GPT-4. We focus  
1396 on prompt tokens rather than call counts, as call counts depend on implementation details (e.g., whether latent  
1397 schemas are cached across users).  
1398

1399 O.1 EMPIRICAL TOKEN MEASUREMENTS

1400 **Baseline (temp=1.0).** Single generation with task prompt only: 50 prompt tokens per user (50 tokens/call).  
1401

1402 **Multi-call baselines.** Self-Consistency ( $k=20$ ): 20 independent generations with the same task prompt:  
1403 1,000 prompt tokens per user (50 tokens/call).

1404 MBR (pool=20): Generate 20 candidates, embed all, select centroid: 1,000 prompt tokens per user (50 to-  
 1405 kens/call).  
 1406 DPP (pool=20,  $r=20$ ): Generate pool of 20 candidates with DPP selection overhead: 3,000 prompt tokens per  
 1407 user (50 tokens/call, but 60 total calls including selection).  
 1408  
 1409 **ORBIT.** ORBIT has two phases with distinct token costs:  
 1410  
 1411 **Phase I (Hidden Seed Generation):**  
 1412 

- Latent variable schema generation:  $\sim 500$  tokens (one-time per task, can be cached)
- $k$  seed generations with latent conditioning:  $k \times 50$  tokens

  
 1413  
 1414 **Phase II (Orthogonalized Output):**  
 1415 

- Pattern mining (identify overused/underused patterns from buffer):  $200 + (k \times 100)$  tokens
- Final generation with orthogonalization:  $50$  (task) +  $(k \times 100)$  (buffer exemplars, capped at 5–10) +  
 1416  $800$  (avoid/prefer instructions) =  $950 + (k \times 100)$  tokens

  
 1417  
 1418 **Total per user (approximate):**  
 1419 

- ORBIT-1 ( $k=1$ ):  $500 + 50 + 300 + 950 = 1,800$  tokens
- ORBIT-4 ( $k=4$ ):  $500 + 200 + 600 + 1,350 = 2,650$  tokens
- ORBIT-9 ( $k=9$ ):  $500 + 450 + 1,100 + 1,750 = 3,800$  tokens

Method	Prompt Tokens/User	Prompt Tokens/Call
Baseline (temp=1.0)	50	50
Self-Consistency ( $k=20$ )	1,000	50
MBR (pool=20)	1,000	50
DPP (pool=20, $r=3$ )	150	50
ORBIT-1	1,800	450
ORBIT-4	2,550	364
ORBIT-9	3,800	317

1426  
 1427  
 1428 Table 6: Empirical prompt token usage from chip flavor task experiments (200 users, 3 outputs each).  
 1429 ORBIT-1 uses  $1.8\times$  the tokens of Self-Consistency/MBR. The “Prompt Tokens/Call” column shows  
 1430 the average across all LLM calls made by each method; ORBIT’s higher per-call average reflects  
 1431 longer Phase II prompts that include buffer content and orthogonalization instructions.  
 1432  
 1433  
 1434

1440  
 1441 **O.2 KEY OBSERVATIONS**

1442 **Overhead is comparable to or better than certain multi-call baselines, while for others its a cost**  
 1443 **that needs to be recognized.** While ORBIT uses more tokens than a single baseline call and other naive  
 1444 methods, the overhead is comparable to or better than certain other quality-focused methods when considering  
 1445 total cost per user.

1446 **Buffer overhead is bounded.** In Phase II, we cap buffer exemplars at 5–10 regardless of  $k$ , preventing  
 1447 prompts from exceeding  $\sim 2,000$  tokens even for  $k=16$ . This ensures manageable context lengths while still  
 1448 providing sufficient orthogonalization signal.  
 1449

1450 **Trade-off: tokens for collisions.** ORBIT achieves near-zero collision probabilities (0.002–0.016) vs  
 1451 baseline (0.76–1.00) at the cost of  $1.8\times$  the prompt tokens of Self-Cons/MBR. This represents a highly favor-  
 1452 able trade-off: substantially better collision avoidance with comparable or better token efficiency than existing  
 1453 quality-focused methods.  
 1454  
 1455  
 1456  
 1457