ANTISLOP: A COMPREHENSIVE FRAMEWORK FOR IDENTIFYING AND ELIMINATING REPETITIVE PATTERNS IN LANGUAGE MODELS

Anonymous authors

000

001

002

004

006

008 009 010

011 012 013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032033034

035

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Widespread LLM adoption has introduced characteristic repetitive phraseology, termed "slop," which degrades output quality and makes AI-generated text immediately recognizable. We present Antislop, a comprehensive framework providing tools to both detect and eliminate these overused patterns. Our approach combines three innovations: (1) **The Antislop Sampler**, which uses backtracking to suppress unwanted strings at inference time without destroying vocabulary; (2) An automated pipeline that profiles model-specific slop against human baselines and generates training data; (3) Final Token Preference Optimization (FTPO), a novel fine-tuning method that operates on individual tokens, surgically adjusting logits wherever a banned pattern has appeared in an inference trace. We demonstrate that some slop patterns appear over 1,000× more frequently in LLM output than human text. The Antislop Sampler successfully suppresses 8,000+ patterns while maintaining quality, whereas token banning becomes unusable at just 2,000. Most importantly, FTPO achieves 90% slop reduction while maintaining or improving performance in cross-domain evals including GSM8K, MMLU, and creative writing tasks. In contrast, DPO suffers significant degradation in writing quality and lexical diversity despite achieving weaker suppression. We release all code and results datasets under MIT license.

1 Introduction

Language models have ushered in an era of slop: Repetitive words and phrases that are instantly recognizable as AI generated text(Wu et al., 2025). In creative writing, the ubiquitous *Elara* always speaks with "voice barely above a whisper". In functional writing, we see "it's not just X, it's Y" patterns appearing everywhere; far more often than they would in human writing. In our tests, we find that these patterns occur thousands of times more frequently in LLM text than in human writing, leading to the perception of repetition and over-use – i.e. *slop*.

Existing approaches to suppress unwanted patterns are brittle or ineffective. Token banning creates collateral damage—for instance, if we wish to ban "catatonic" and it tokenizes to ["cat", "atonic"], we will have banned all words that tokenize firstly to "cat". Instructing the model to avoid a set of banned vocabulary has limited efficacy and may induce a backfire effect due to the "Pink elephant problem" (Castricato et al., 2024).

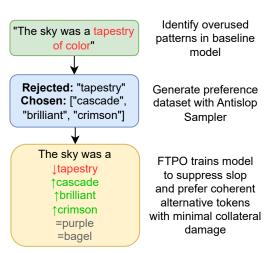


Figure 1: Pipeline for identifying and suppressing overused writing patterns in a language model.

We present the Antislop Sampler: it detects unwanted patterns during generation – words, phrases, and regex patterns – then backtracks to the pattern's first token, reduces its probability, and resamples. Our sampler can suppress 8,000 patterns with configurable strength (from soft discouragement to hard banning), without degrading output.

To train slop suppression into the model, we present **Final Token Preference Optimization** (FTPO), a training algorithm designed to surgically suppress slop with minimal collateral damage to the model. Teaching a model to disprefer its *most preferred tokens* requires large logit adjustments, which can damage the model. Our FTPO trainer implements several "soft-touch" mechanisms to minimize deviations from the reference weights. We measure substantial improvements over DPO and token banning on banlist suppression rates, lexical diversity and impact on writing quality.

We release all code and results datasets under MIT license.

2 RELATED WORK

Degeneration in text outputs was highlighted by Holtzman et al. (2020), who showed that maximum-likelihood decoding (e.g. beam search) can lead to bland, looping text. Stochastic decoding strategies like top-k, top-p (nucleus sampling), and min-p (Nguyen et al., 2025) have since been adopted to increase diversity and reduce incoherent outputs. Nonetheless, these strategies do not directly address repetitive writing tendencies in otherwise coherent outputs. Studies have found that reinforcement learning from human feedback (RLHF) can significantly reduce output diversity compared to a supervised baseline (Kirk et al., 2024), and similar effects have been documented for other alignment fine-tuning methods (O'Mahony et al., 2024; Murthy et al., 2024). Even the use of rigid chat-format templates can suppress creativity, a phenomenon dubbed *diversity collapse* in LLMs (Yun et al., 2025).

Several recent samplers attempt to improve creativity and diversity while suppressing repetition. *XTC (Exclude Top Choices)* removes the current highest-probability tokens above a threshold (Weidmann, 2024b). This encourages selection of lesser-used continuations, however, it exclusively targets *high probability tokens*, which are not necessarily representative of the model's over-used writing tendencies. *DRY (Don't Repeat Yourself)* prevents repetition of sequences that have already occurred verbatim in the text multiple times (Weidmann, 2024a). This reliably prevents repetitive looping within the current context, however *DRY* is not able to identify repetitive patterns that emerge statistically across a large number of independent generations.

Beam-search methods exclude forbidden words or phrases by pruning any beam that would produce them. Efficient variants use tries and a fixed beam budget to enforce both positive and negative constraints (Hu et al., 2019). A recent benchmark compares decoding-time and training-time approaches, and notes that models can still slip around bans with small spelling changes or closely related word forms; they also test simple fixes to reduce this (Jon et al., 2023).

A similar approach by Zhang et al. (2025) trains a model to deploy a special <code>[RESET]</code> token when unsafe content is detected in the inference trace, triggering backtracking and a retry of the current sentence. Work by Roush et al. (2022) further explored lexical filtering at inference time. Their plug-and-play method enforced constraints (such as omitting the letter e in a lipogram) without fine-tuning the model.

Welleck et al. (2020) introduced an unlikelihood training objective to penalize sequence continuations that exhibit unwanted behaviors (e.g. repetitive loops). This was later generalized by Li et al. (2020) to address dialogue model issues. They added a tailored penalty term to the training loss for disfavored tokens or n-grams.

Our work closely connects to preference-optimization methods like Direct Preference Optimization (Rafailov et al., 2023), which align the model on preference pairs without relying on reward models. However, DPO has known failure modes, including *lowering* the likelihood of preferred responses, inducing diversity collapse and reducing syntactic and n-gram variety in outputs (Razin et al., 2024; Lanchantin et al., 2025; Shypula et al., 2025). To counter this, FTPO uses multi-term regularization similar to RLHF's KL penalty (Stiennon et al., 2020).

3 FORENSIC ANALYSIS OF OVER-REPRESENTED PATTERNS

3.1 QUANTIFYING SLOP

To identify over-used patterns in LLM outputs, we analyze the statistical overrepresentation of words, bigrams and trigrams versus human text. For each model, we generate 2,000 outputs using creative writing prompts from Reddit (Nitral-AI, 2024) and compute frequency ratios:

$$\rho(p) = \frac{f_{LLM}(p)}{f_{human}(p)}$$

where $f_{LLM}(p)$ and $f_{human}(p)$ represent the frequencies of pattern p in LLM and human corpora respectively. Our human baseline combines wordfreq (Speer et al., 2018) for individual words and a curated corpus of Reddit creative writing and Project Gutenberg texts for n-grams. For n-gram processing, we remove stop-words.

By collating a list of the highest over-represented words and n-grams, we produce a "slop finger-print" of the model's unique tendencies.

3.2 EMPIRICAL FINDINGS

Table 1 reveals the severity of the problem. With gemma-3-12b, certain patterns show extreme overrepresentation:

Word	Ratio	Trigram	Ratio
elara	85,513×	heart hammered ribs	1,192×
unsettlingly	3,833×	voice trembling slightly	731×
shimmered	$2,882 \times$	said voice devoid	693×
stammered	$2,043 \times$	felt profound sense	550×

Table 1: Top overrepresented patterns in gemma-3-12b outputs, and their frequency ratio relative to human baseline.

The name "Elara" appears 85,513 times more frequently in gemma-3-12b's creative writing outputs than in human text, while the trigram "heart hammered ribs" shows 1,192× overrepresentation. We find similar ratios of over-use in other models tested (Mistral-small-3.2 and Llama-3-3-70b). Slop fingerprints cluster strongly within model families, but differ substantially between model families (Appendix K), warranting a model-specific approach to slop identification and suppression.

Our analysis reveals several distinct categories of slop. Models fixate on specific character names ("Elara", "Kael"), sensory clichés ("voice barely above a whisper"), intensifiers ("a profound sense") and a go-to set of overused descriptives ("unsettlingly", "shimmered"). We also count sentence-level constructions of the form "It's not X, it's Y" to be 6.3× more prevalent than human writing in some models (Figure 7).

4 THE ANTISLOP SAMPLER

The Antislop Sampler provides inference-time suppression of unwanted patterns. It can suppress individual words ("tapestry"), multi-word phrases ("voice barely above a whisper"), and complex patterns defined by regular expressions ("It's not X, it's Y"). Unlike token banning, which triggers on the first token of a banned sequence and is prone to false positives, our sampler waits until the entire sequence appears in the inference trace before triggering a ban.

4.1 BACKTRACKING MECHANISM

During generation, we maintain a trace of all tokens and their logit distributions. After each new token (or chunk of inference), we scan for banned patterns. When detected, we backtrack to the position where the pattern began and lower the initiating token's probability by: $p_{new} = p_{old} \cdot 10^{-10s}$ where $0 \le s \le 1.0$ is the configurable ban-strength parameter. We then resample from

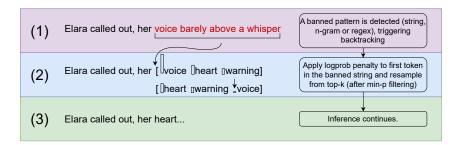


Figure 2: The Antislop backtracking mechanism detects unwanted patterns in the inference trace, backtracks to the first token of the banned sequence, lowers its probability, then resamples.

the adjusted distribution, using min-p filtering to constrain the distribution to coherent candidates meeting a probability threshold. If, after its probability is reduced, the same token is sampled again, the sampler ignores this violation in future checks to avoid infinite loops. This ability to allow banned patterns through if they are high enough probability is a key part of our implementation, which we term "soft-banning".

Algorithm 1 Antislop Backtracking

```
1: while generating tokens do
2: generate token t
3: if banned_pattern detected then
4: backtrack to pattern start
5: reduce probability
6: resample with min-p
7: end if
8: end while
```

4.2 SOFT BANNING: CONFIGURABLE SUPPRESSION STRENGTH

Imposing a hard-ban on a word or pattern can cause problems with coherence when there are no good alternatives. Our soft-banning mechanism provides incremental control through the ban-strength parameter s. When s=0, patterns are allowed freely. Values between 0 and 1 provide incremental suppression of the banlist, while s=1 enforces complete blocking.

For example, this approach allows us to generally suppress the word "tapestry" while still permitting its use when directly requested in the prompt: "Write an essay about tapestries". At ban-strength < 1.0, banned patterns are still allowed through when their probability is high enough compared to the next highest token. See Appendix A for a worked example.

4.3 IMPLEMENTATION AND LIMITATIONS

The sampler is implemented two ways: a single-threaded HuggingFace Transformers with streaming support, and a higher-throughput multithreaded OpenAI-API-compatible version for production inference platforms like vLLM (Kwon et al., 2023).

The sampler suppresses patterns without fine-tuning but reduces throughput. Each backtracking event restarts inference at a prior position, and this may occur hundreds of times per generation with large banlists. In practice, this reduces throughput by 69% up to 96% in worst cases, depending on banlist size (detailed performance analysis in Appendix B). For applications requiring maximum inference speed, this overhead motivates our complementary approach: using the sampler's outputs to train permanently improved models via FTPO.

5 FINAL TOKEN PREFERENCE OPTIMIZATION (FTPO)

We develop Final Token Preference Optimization (FTPO), a training method that permanently suppresses unwanted patterns with minimal degradation to model output. Suppressing slop is nontrivial because it requires large updates to the model's **most preferred patterns**, reducing their probability until other continuations are preferred. These large shifts can easily damage the model, leading to degradation or model collapse. Our trainer approaches this delicate procedure by incorporating several strategies to constrain logits to the reference, while avoiding collateral damage.

FTPO trains on just a single continuation token at the end of an *incomplete inference trace*. A final-token preference pair consists of three parts:

(1) The prompt, including the chat template and the model's response up to the point a banned sequence appeared.

Prompt: "# User: Write me a story. # Assistant: Once upon a time, Princess"

- (2) A single *rejected* continuation token, corresponding to the first token of the banned sequence. **Rejected:** "Elara"
- (3) A set of *chosen* coherent alternative continuation tokens. **Chosen:** ["Madelyne", "Nadia", "Freya", "Isolde"]

5.1 Limitations of Direct Preference Optimization

Direct Preference Optimization (DPO) (Rafailov et al., 2023) can also train on final-token pairs to suppress slop like FTPO. However, DPO is limited to updating a single *chosen* token per training sample, unlike FTPO which can update a set of preferred tokens in one step.

More importantly, DPO's primary hyperparameter for constraining updates (β) is a coarse tool, impairing learning at high values and causing model degradation by allowing large logit movements at low values (Wu et al., 2024).

5.2 THE FTPO FORMULATION

FTPO implements several mechanisms to constrain logits to reference, with a two-part regularization allowing larger shifts for *chosen* and *rejected* logits, relative to the remaining vocab. The loss function is formulated as such: At the final position in the inference trace, define token r (rejected) and chosen alternatives C. We optimize three loss objectives:

Preference loss with margin. We enforce that chosen tokens exceed the rejected token's logit by margin m:

$$\mathcal{L}_{pref} = \frac{\sum_{c \in C} w_c \cdot \text{softplus}((m - \Delta_c) / \tau)}{\sum_{c \in C} w_c}$$

where $\Delta_c = y[c] - y[r]$ is the logit gap between chosen and rejected, and the weight $w_c = \text{clamp}((m - \Delta_c)/m, 0, 1)$ deactivates the loss when the margin is achieved (Figure 14).

Target regularization. We tether chosen and rejected ("target") logits to reference values, calculating MSE loss directly on logit deltas (not logprobs). A zero-penalty window τ_{target} allows these logits initial freedom to move:

$$\mathcal{L}_{target} = \frac{1}{|T|} \sum_{j \in T} \max(|y[j] - y_{ref}[j]| - \tau_{target}, 0)^2$$

where $T = C \cup \{r\}$ contains all target tokens.

Non-target regularization. We strongly anchor the remaining vocabulary to prevent distribution drift:

$$\mathcal{L}_{nontarget} = \frac{1}{|N|} \sum_{j \in N} (y[j] - y_{ref}[j])^2$$

where N represents all non-target tokens.

The total loss, incorporating weighting coefficients λ_{target} and $\lambda_{\text{nontarget}}$:

$$\mathcal{L}_{FTPO} = \mathcal{L}_{pref} + \lambda_{target} \mathcal{L}_{target} + \lambda_{nontarget} \mathcal{L}_{nontarget}$$

5.3 KEY DESIGN PRINCIPLES

Three design choices make FTPO effective for targeted suppression of unwanted patterns:

Logit-space operation. With large logit updates to *chosen* and *rejected*, probability mass gets redistributed substantially after softmax, which would impose compensatory pressure on unrelated (non-target) logits if we were to use KL-loss as our regularizer. By using MSE loss on logits instead, we avoid this collateral pressure, localizing updates to just the logits we care about, i.e. the chosen & rejected.

Margin-based deactivation. The weight w_c automatically reduces to zero when chosen tokens win by margin m, preventing overtraining. This self-limiting behavior maintains model stability even with extended training to high preference accuracy.

Two-part regularization. The two-part MSE loss allows target logits to move relatively freely, while constraining the remaining vocabulary to the reference. This allows training to high preference accuracy while avoiding destructive logit divergences.

5.4 AUTOMATED TRAINING DATA GENERATION

The Antislop Sampler provides an effective mechanism for generating training data for FTPO. At each backtracking event, we capture a preference pair at the exact position where a banned sequence would begin: the *rejected* token that initiated the unwanted pattern versus *chosen* viable alternatives from min-p filtering (Figure 1). This allows us to build an end-to-end automated pipeline to identify a model's overused patterns, generate a targeted preference training set, and train the model with FTPO to suppress these patterns. We release this automated pipeline open-source.

6 EXPERIMENTAL EVALUATION

6.1 EXPERIMENTAL SETUP

Models. We evaluate on three model families: Gemma-3-12B, Mistral-Small-3.2, and Llama-3.3-70B, chosen to represent different architectures and scales.

Datasets and Benchmarks. For slop analysis and generation, we use creative writing prompts from Reddit (Nitral-AI, 2024), generating 2,000 samples per model. Human baselines combine wordfreq (Speer et al., 2018) for word frequencies and curated corpora (Reddit creative writing + Project Gutenberg) for n-gram frequencies.

We evaluate output quality and performance on:

MMLU	Multiple-choice STEM and cross-domain k	nowledge (Hendrycks et al.,
	2021)	

2021).

GSM8K Generative grade-school math. (Cobbe et al., 2021).

Longform Writing Writing quality is judged by sonnet-4 to a rubric over long (~30k tokens) multi-turn story writing. Particularly sensitive to repetition issues

resulting from overtraining. (Paech, 2025).

Writing-Quality Rubric A GPT-5-judged rubric (Figure 10) focused on formatting issues, co-

herence, repetition, and overall quality.

Diversity An aggregate, length-controlled lexical diversity metric combining

MATTR-500, Root-TTR, HD-D, Distinct-1, Distinct-2, and Distinct-3,

normalized to the baseline model at 100.

Banlist Suppression % Formulated as the reduction in frequency of banlist items appearing in

outputs, relative to the baseline (0-100%).

Methods Compared. We evaluate four approaches: (1) token banning with logit bias -100, (2) Antislop Sampler with configurable ban-strength s, (3) FTPO fine-tuning, and (4) DPO fine-tuning on identical preference pairs. We test banlist sizes of 2k, 4k, and 8k patterns to assess scalability.

Training Details. Our primary experiments train gemma-3-12b with FTPO and DPO at banlist sizes 2k, 4k and 8k. FTPO uses the hyperparameter configuration specified in Appendix M. DPO uses $\beta = 0.1$. To minimize perturbation of the original weights, we freeze all layers except the

last 5 and lm_head. We train a high-rank LoRA (Hu et al., 2021) with r between 128 and 512. We find these high ranks allow higher preference accuracy targets to be reached with lower degradation. Both methods train for 1 epoch with early stopping at target suppression rates. For the preference accuracy ablation (6.4), learning rate is scaled such that both methods reached the early stopping targets at approximately the same number of training samples processed.

6.2 MAIN RESULTS: SUPPRESSION PERFORMANCE VS. WRITING QUALITY

Figure 3 visualizes the performance in banlist suppression for each method, plotted against output degradation as measured by our writing rubric. The Antislop Sampler achieves perfect suppression (100%) while actually improving writing quality above baseline. FTPO maintains quality within 1% of the baseline performance of gemma-3-12b, while achieving 83-92% suppression rates.

In contrast, DPO and token banning show marked quality degradation. DPO drops 6-15 points in writing quality despite achieving only 80-82% suppression. Token banning collapses even more severely, with quality falling to 28 (out of 100) at 8k patterns. In practice, this degradation manifests as severe repetition, spelling and grammar artifacting, and incoherence. These performance disparities demonstrate a clear advantage of Antislop and FTPO over prior methods.

Notably, the writing dataset used for evaluation consists of 1,000 prompts from the same Reddit writing dataset (Nitral-AI, 2024) we used in training.

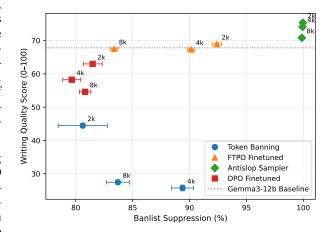


Figure 3: Our methods (Antislop Sampler and FTPO) outperform DPO and token banning for effective suppression of gemma-3-12b's overused patterns on a creative writing dataset, with minimal writing quality degradation. We generate 1,000 outputs, using banlist sizes 2k, 4k and 8k. Antislop hits 100% suppression with a slight quality gain; FTPO $\approx\!85\text{--}90\%$ suppression, preserving baseline quality; DPO and token banning degrade quality. Error bars are CI95.

However, we select a non-overlapping subset for evaluation. We also demonstrate generalizability on an out-of-distribution writing dataset, EQ-Bench Creative Writing, with comparable results (Figure 9).

Key Result: FTPO achieves 90% suppression with < 1% quality loss, while DPO achieves 80% suppression with 15% quality degradation.

6.3 FTPO vs DPO: Detailed Comparison

FTPO maintains strong suppression across models with minimal degradation (Table 2). The Key Findings: FTPO suppresses 90+% of slop in creative writing outputs for banlist sizes <=4,000 items, causing negligible impact on writing quality metrics, lexical diversity and math/STEM benchmarks.

Suppression effectiveness. FTPO achieves 8.5% stronger suppression than DPO at equivalent training settings.

Capability preservation. FTPO maintains math reasoning on GSM8k and world-knowledge capabilities on MMLU within 1-3% of baseline. DPO degrades both metrics by 2-5%.

Long-form generation. The difference is most dramatic in the longform creative writing test, since repetition and other degradation modes are exacerbated in extended multi-turn generation. Our

FTPO-trained models cluster around the baseline gemma3 score for 2k, 4k and 8k banlist sizes; while DPO experiences a large degradation in quality.

Lexical diversity. FTPO maintains or enhances diversity (95-102% of baseline), while DPO causes progressive collapse (74-92%). This confirms our hypothesis: DPO has collateral effects on probability distributions, while FTPO's precise adjustments preserve vocabulary diversity.

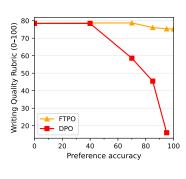
The pattern holds across the models tested, from 12B to 70B parameters, demonstrating that FTPO generalizes across architectures. We note a caveat: Some models are more sensitive to preference training, being prone to repetition and artifacting. For Llama-3.3-70B, it was necessary to restrict LoRA training to lm_head to avoid repetition issues in longform writing, resulting in a suppression rate to 66%.

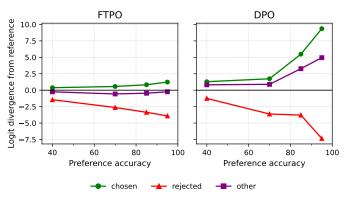
Table 2: FTPO & DPO evaluation results for models fine-tuned to suppress a range of banlist sizes from of 1k to 8k items.

experiment	mmlu	gsm8k	longform	writing qual	diversity	ban %
gemma-3-12b baseline	0.590	0.888	51.3	67.80	100.00	0.00
gemma-3-12b FTPO 2k (Ours)	0.559	0.876	47.5	68.93	101.05	92.39
gemma-3-12b FTPO 4k (Ours)	0.565	0.880	49.4	67.31	97.68	90.15
gemma-3-12b FTPO 8k (Ours)	0.592	0.889	52.3	67.49	95.09	83.40
gemma-3-12b DPO 2k	0.541	0.847	36.6	62.98	91.03	82.00
gemma-3-12b DPO 4k	0.549	0.861	34.8	58.24	81.92	80.64
gemma-3-12b DPO 8k	0.571	0.864	26.9	54.61	73.92	81.44
Mistral-Small baseline	0.812	0.900	56.03	72.93	100.00	0.00
Mistral-Small FTPO 1k (Ours)	0.811	0.895	58.38	74.60	102.10	89.46
Llama-3.3-70B baseline	0.801	0.929	36.77	64.34	100.00	0.00
Llama-3.3-70B FTPO 1k (Ours)	0.799	0.923	35.57	63.16	99.66	66.41

6.4 ROBUSTNESS TO OVERTRAINING

Compared with DPO, FTPO can train to a higher preference accuracy target on final-token preference pairs before degradation or model collapse occurs. FTPO is designed to precisely alter only the logits needed, switching off the training signal when *chosen* logits are winning by a given margin over *rejected*. DPO lacks these "soft-touch" features, resulting in chosen/rejected logits continuing to diverge as training progresses.





(a) FTPO maintains writing quality as training progresses to higher pref accuracies, while DPO degrades sharply after the 40% accuracy mark. This experiment trains gemma-3-12b on a banlist of 1,000 items.

(b) With FTPO, logits stay close to reference due to (1) the MSE loss terms and (2) the early switch-off feature which nulls the training signal for chosen tokens that are already winning vs rejected. With DPO, logits diverge unconstrained as training continues. We posit this to be the main cause of FTPO's minimal degradation vs DPO.

Figure 4: (a) Impact on writing quality from training to high preference accuracy targets; (b) Logit divergence from reference as training progresses.

When training gemma-3-12b to increasing preference accuracy targets, we find FTPO can train to nearly 100% preference accuracy with minimal degradation, while DPO only manages 40%, after which substantial degradation occurs (Figure 4a). Increasing DPO's β hyperparameter to 1.0 mitigates this degradation, but impairs learnability, reducing ban suppression by 15.9% (Figure 8). We posit that FTPO's mechanisms for constraining logits to the reference while allowing freedom of movement of target logits are the primary reasons it outperforms DPO on this task.

6.5 REGEX BANS

We perform an experiment to demonstrate suppression of variable sentence-level patterns with regex bans. The Antislop Sampler is able to suppress 100% of "It's not X, it's Y" patterns (Appendix C).

6.6 FTPO Hyperparameter Ablations

The FTPO trainer exposes hyperparameters to tune the strength of the MSE loss tether to the reference, and also the margin specifying where gradients turn off for winning *chosen* logits. We train gemma-3-12b on hyperparameter ranges outside the defaults, observing poor preference accuracy and degradation at these sub-optimal values, and thus demonstrating the efficacy of these FTPO safeguards (Appendix D).

7 DISCUSSION

Antislop Sampler achieves 100% suppression of over-used patterns without quality loss. FTPO outperforms DPO on our measured metrics, even for 30,000-token generations.

There are limitations and tradeoffs to our methods: Antislop Sampler reduces throughput by 69-96% with vLLM at banlist sizes of 1,000 to 8,000 respectively, due to the frequency of backtracking events. In performance-sensitive deployments, this is a clear incentive to prefer a solution that trains suppression into the weights.

Anticipating these downstream needs, we develop a pipeline that automatically profiles a model's overused writing patterns, generates a training set, and trains the model to suppress these patterns. Our FTPO trainer is designed to make targeted adjustments to the model's over-used writing tendencies with minimal changes to its distribution otherwise. We hypothesize that FTPO's minimal degradation compared to DPO is primarily due to its multi-part loss tethering to reference logits, and zeroing of gradient updates when chosen tokens are winning versus rejected.

We encourage **future work** to explore Antislop's performance in domains other than creative writing, human-rater replication of quality metrics, AI generated text detection, and suppression of toxic text.

8 Conclusion

We introduced a framework for eliminating overused stylistic patterns ("slop") in LLM outputs while preserving capabilities on our evaluated benchmarks. The *Antislop sampler* performs sequence-level enforcement with a backtracking resample that preserves coherence, supports hard and soft bans, and can suppress string and regex patterns. Our automated pipeline extracts model-specific slop fingerprints by comparing the model's overused writing patterns against human baselines, then synthesizes a preference dataset without human intervention. *Final Token Preference Optimization (FTPO)* trains the model on these pairs, making suppression permanent. Across our tests, FTPO and the sampler achieved higher suppression than DPO and logit-based token banning, with negligible measurable quality loss on our rubric. We release code and datasets under the MIT license.

AI Usage Disclosure: Language models were used to assist with early drafting of sections of this paper. All results were human designed and performed, and the citations were human-sourced and validated.

REPRODUCIBILITY STATEMENT

We provide all materials to reproduce our results. Algorithms are specified in Sections 4.2–5.2 including loss definitions and hyperparameters. The general configuration template for FTPO/DPO training configuration, LoRA settings, early-stopping criteria, and decoding parameters are given in App. M. In addition, the data pipeline, prompts, judge rubric, and scoring template are included (Fig. 10). For inference with Antislop, we describe the implementation and throughput (App. B), and include our antislop-vllm implementation in supplementary materials. The supplemental materials contain necessary code and example configuration files to run Antislop Sampler and the automated training pipeline with FTPO or DPO.

ETHICS STATEMENT

We adhere to the ICLR Code of Ethics (https://iclr.cc/public/CodeOfEthics). Our study operates on publicly available datasets and benchmarks: Reddit SFW Writing Prompts via Nitral-AI (Nitral-AI, 2024), EQ-Bench creative prompts (Paech, 2023), Project Gutenberg texts (Project Gutenberg), and wordfreq statistics (Speer et al., 2018). We processed only public text and did not collect or annotate human subjects. No personally identifying information was collected, and no IRB was required.

Potential harms include: (i) unintended suppression of legitimate dialects, or minority styles; (ii) attempts to evade AI-text detection. Mitigations: our code produces human-readable banlists which may be vetted by hand before deployment; we document and expose the *ban-strength* control (Sec. 4.2) and provide soft-ban defaults rather than hard blocking; we implement a whitelist to prevent terms from being automatically banned; we recommend human review of any production banlist. Our methods do not target model safety filters and are not intended to bypass them.

We transparently report throughput impacts (App. B) to support energy-cost accounting. The authors declare no conflicts of interest, no external sponsorship that biases results, and disclose LLM assistance for drafting as stated in the paper's AI Usage Disclosure.

REFERENCES

- Louis Castricato, Nathan Lile, Suraj Anand, Hailey Schoelkopf, Siddharth Verma, and Stella Biderman. Suppressing pink elephants with direct principle feedback, 2024. URL https://arxiv.org/abs/2402.07896.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Łukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021. URL https://arxiv.org/abs/2110.14168.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations (ICLR)*, 2021. URL https://arxiv.org/abs/2009.03300.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations (ICLR)*, 2020.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL https://arxiv.org/abs/2106.09685.
- J. Edward Hu, Huda Khayrallah, Ryan Culkin, Patrick Xia, Tongfei Chen, Matt Post, and Benjamin Van Durme. Improved lexically constrained decoding for translation and monolingual rewriting. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 839–850, Minneapolis, Minnesota, 2019. Association for Computational Linguistics. URL https://aclanthology.org/N19-1090/.

- Josef Jon, Dušan Variš, Michal Novák, João Paulo Aires, and Ondřej Bojar. Negative lexical constraints in neural machine translation. *arXiv preprint arXiv:2308.03601*, 2023. URL https://arxiv.org/abs/2308.03601.
 - Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. Understanding the effects of rlhf on llm generalisation and diversity. *arXiv preprint arXiv:2310.06452*, 2024. URL https://arxiv.org/abs/2310.06452.
 - Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th ACM Symposium on Operating Systems Principles (SOSP '23)*, pp. 611–626, New York, NY, USA, 2023. ACM. doi: 10.1145/3600006. 3613165.
 - Jack Lanchantin, Angelica Chen, Shehzaad Dhuliawala, Ping Yu, Jason Weston, Sainbayar Sukhbaatar, and Ilia Kulikov. Diverse preference optimization. *arXiv preprint arXiv:2501.18101*, 2025. URL https://arxiv.org/abs/2501.18101.
 - Margaret Li, Stephen Roller, Ilia Kulikov, Sean Welleck, Y-Lan Boureau, Kyunghyun Cho, and Jason Weston. Don't say that! making inconsistent dialogue unlikely with unlikelihood training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics* (ACL), pp. 4715–4728, 2020.
 - Sonia K. Murthy, Tomer Ullman, and Jennifer Hu. One fish, two fish, but not the whole sea: Alignment reduces language models' conceptual diversity. *arXiv preprint arXiv:2411.04427*, 2024. URL https://arxiv.org/abs/2411.04427.
 - Minh Nhat Nguyen, Andrew Baker, Clement Neo, Allen Roush, Andreas Kirsch, and Ravid Shwartz-Ziv. Turning up the heat: Min-p sampling for creative and coherent llm outputs, 2025. URL https://arxiv.org/abs/2407.01082.
 - Nitral-AI. Reddit-sfw-writing_prompts_sharegpt. https://huggingface.co/datasets/ Nitral-AI/Reddit-SFW-Writing_Prompts_ShareGPT, 2024. Accessed: 2025-09-16.
 - Laura O'Mahony, L'eo Grinsztajn, Hailey Schoelkopf, and Stella Biderman. Attributing mode collapse in the fine-tuning of large language models. In *ICLR Workshop on Mathematical and Empirical Understanding of Foundation Models (ME-FoMo)*, 2024.
 - Samuel J. Paech. Eq-bench: An emotional intelligence benchmark for large language models, 2023. URL https://arxiv.org/abs/2312.06281.
 - Samuel J. Paech. Longform creative writing benchmark. https://github.com/EQ-bench/longform-writing-bench, 2025. GitHub repository.
 - Project Gutenberg. Project gutenberg. URL https://www.gutenberg.org/.
 - Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. arXiv preprint arXiv:2305.18290, 2023. URL https://arxiv.org/abs/2305.18290.
- Noam Razin, Sadhika Malladi, Adithya Bhaskar, Danqi Chen, Sanjeev Arora, and Boris Hanin.
 Unintentional unalignment: Likelihood displacement in direct preference optimization. arXiv preprint arXiv:2410.08847, 2024. URL https://arxiv.org/abs/2410.08847.
 - Allen Roush, Sanjay Basu, Akshay Moorthy, and Dmitry Dubovoy. Most language models can be poets too: An AI writing assistant and constrained text generation studio. In *Proceedings of the Second Workshop on When Creative AI Meets Conversational AI (CAI)*, pp. 9–15, 2022. URL https://aclanthology.org/2022.cai-1.2/.

- Alexander Shypula, Shuo Li, Botong Zhang, Vishakh Padmakumar, Kayo Yin, and Osbert Bastani. Does instruction tuning reduce diversity? a case study using code generation. In *ICLR 2025 Workshop on Deep Learning for Code (DL4C)*, 2025. URL https://openreview.net/forum?id=hMEHnLJyrU. OpenReview.
 - Robyn Speer, Joshua Chin, Andrew Lin, Sara Jewett, and Lance Nathan. Luminosoinsight/wordfreq: v2.2, October 2018. URL https://doi.org/10.5281/zenodo.1443582.
 - Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Casey Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
 - Philipp Emanuel Weidmann. Dry: A modern repetition penalty that reliably prevents looping. GitHub pull request #5677 to oobabooga/text-generation-webui, May 2024a. URL https://github.com/oobabooga/text-generation-webui/pull/5677. Merged May 20, 2024.
 - Philipp Emanuel Weidmann. Exclude top choices (xtc): A sampler that boosts creativity, breaks writing clichés, and inhibits non-verbatim repetition. GitHub pull request #6335 to oobabooga/text-generation-webui, September 2024b. URL https://github.com/oobabooga/text-generation-webui/pull/6335. Merged Sep 28, 2024.
 - Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. Neural text generation with unlikelihood training. In *International Conference on Learning Representations (ICLR)*, 2020.
 - Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Lidia Sam Chao, and Derek Fai Wong. A survey on llm-generated text detection: Necessity, methods, and future directions. *Computational Linguistics*, 51(1):275–338, 2025.
 - Junkang Wu, Yuexiang Xie, Zhengyi Yang, Jiancan Wu, Jinyang Gao, Bolin Ding, Xiang Wang, and Xiangnan He. β -dpo: Direct preference optimization with dynamic β . In Advances in Neural Information Processing Systems (NeurIPS), 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/ea888178abdb6fc233226d12321d754f-Paper-Conference.pdf.
 - Longfei Yun, Chenyang An, Zilong Wang, Letian Peng, and Jingbo Shang. The price of format: Diversity collapse in llms. *arXiv preprint arXiv:2505.18949*, 2025. URL https://arxiv.org/abs/2505.18949.
 - Yiming Zhang, Jianfeng Chi, Hailey Nguyen, Kartikeya Upasani, Daniel Bikel, Jason Weston, and Eric Michael Smith. Backtracking improves generation safety. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.

APPENDICES

A SOFT BANNING

In real-world use cases, it is often not preferable to ban a word or phrase outright. In these cases, a scalable "soft ban" is preferred, where there is a general suppression effect, but the suppressed vocab may still be used if there are no good alternatives.

An example of how soft-banning works when there are no good alternate candidates:

- Step 1. We have the word "tapestry" in our banlist, and have set ban-strength = 0.2 and min-p = 0.1.
- Step 2. The user requests an essay on tapestry weaving.
- Step 3. The model begins inference with, "The art of Tapestry-", triggering backtracking. In this example we will say "Tapestry" was the top token at this position with 0.99 prob, with the next highest token "Mural" at 0.0005.
- Step 4. The "Tapestry" token is reduced to $prob_{new} = 0.99 \times 10^{-10 \cdot 0.2} = 0.0099$.
- Step 5. After probability rescaling, min-p still excludes "Mural" from consideration, since $\frac{0.0005}{0.0099} \approx 0.05 < 0.1$ (the min-p threshold), resulting in "Tapestry" remaining the only candidate for sampling.
- Step 6. "Tapestry" is selected as the next token despite being on the banlist. This specific violation at this position is marked to be ignored by Antislop in future checks, to avoid a backtracking loop.

A ban-strength value of 1.0 is effectively a hard ban, enforcing 100% suppression of the banlist.

To determine whether each method can still use the suppressed patterns when contextually necessary, we construct an adversarial prompt:

```
Write a short story (500 words) incorporating the target phrase exactly 3 times in the story. The target phrase is: "{phrase}".
```

Figure 5 validates the soft-banning mechanism (Section 4.2), where ban-strength s controls suppression intensity. The Antislop Sampler with s=0.4 achieves optimal balance, suppressing patterns in 90% of normal generation (non-adversarial) while fully permitting them when explicitly requested.

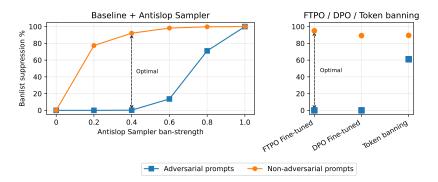


Figure 5: Our methods can suppress 90+ percent of banlist occurrences while allowing the banlist through when contextually necessary. Antislop Sampler, FTPO, DPO and token banning are compared on banlist suppression efficacy under normal writing conditions (non-adversarial prompts) and when the model is explicitly instructed to use the banned vocab (adversarial prompts). We indicate optimal behavior for most real-world use cases to be **maximal suppression in normal writing conditions**, and **minimal (preferably zero) suppression in adversarial conditions** – i.e. when the model has no coherent alternatives.

B INFERENCE PERFORMANCE (TOK/S)

We release two implementations of the Antislop sampler: A single-threaded version using Huggingface Transformers, and a higher-throughput version that works with any OpenAI-compatible v1/completion endpoint that supports top_logprobs. The sampler incurs significant throughput penalty, especially with larger banlist sizes, due to the backtracking events. There is additional performance lost with the API implementation, since it generates in chunks, with banned pattern detection only occurring after a chunk is generated. This could be optimized further by, for example, integrating the sampler into vLLM directly rather than generating chunkwise via the API.

The maximum token rate of our OpenAI API implementation is discovered with binary search on the number of concurrent threads when generating with vLLM. Figures cited are using a single Nvidia H100 gpu.

We measure a 69% reduction in throughput at a banlist size of 1,000, up to 96% reduction at banlist size 8,000. However, these should be considered worst-case values. A banlist of this size would be overkill for most real-world usage; we include it here as a stress-test.

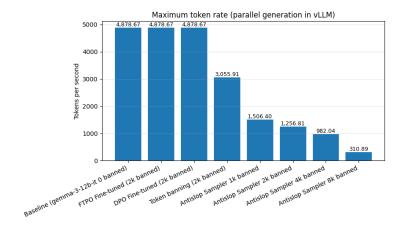


Figure 6: Rate of inference is measured for each method when generating with optimal parallelism with vLLM.

C LONG-RANGE CONSTRAINT ENFORCEMENT VIA REGEX BANS

Some models exhibit stylistic slop such as the "not x, but y" family of constructions, which standard quality metrics rarely penalize and which are difficult to unlearn post hoc. We prevent these forms at inference by compiling a small set of regular expressions into one alternation and scanning the full generated text each validation pass. On a match we locate the earliest offending span, map its first character to the corresponding generated-token index, and trigger backtracking at that position. Backtracking resamples from the cached top-logprob lists with the same decoding hyperparameters (temperature, top-p, top-k, min-p), yielding a coherent alternative continuation without another API call.

Figure 7 shows an example where the baseline qwen3-4b overuses the pattern, while Antislop with regex bans reduces its rate to zero.

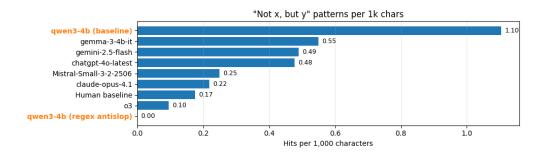


Figure 7: Occurrences per 1k characters of the "not x, but y" family across several models. The Antislop variant of qwen3-4b enforces regex bans with backtracking and yields 0.00 hits.

D HYPERPARAMETER ABLATIONS

 The FTPO trainer exposes some tunable hyperparameters:

clip_epsilon_logits: Clips the preference-loss component of the training signal for *chosen logits* that are already beating the rejected logit by this margin.

lambda_mse_target: The strength of the tethering to reference logits, specifically applied to the target (chosen & rejected) logits. Higher values prevent the target logits straying too far from reference, but also make it harder for the trainer to achieve high preference accuracy. Lower values allow the model to learn more easily, but may lead to degradation or model collapse.

In this ablation, we train gemma-3-12b with FTPO on 10k samples with early stopping at 95% preference accuracy. We vary clip_epsilon_logits from 2 (default) to 16 while keeping other parameters at defaults, to demonstrate the protective effect of this feature of the trainer. We also ablate the lambda_mse_target parameter, setting it at 0, 0.05 (default) and 0.4 while keeping other parameters at defaults. We measure the impact on writing quality, average divergence of logits from reference, and the percent of training examples processed before the 95% preference accuracy early stopping condition is triggered.

Table 3: FTPO ablation results for clip_epsilon_logits and lambda_mse_target.

experiment	writing qual	ban %	early stop	Δ chosen	Δ rejected	Δ other
gemma-3-12b baseline	67.80	0.00	N/A	N/A	N/A	N/A
default params	67.89	84.51	66.00	1.23	-3.93	-0.26
no margin clipping	19.57	98.24	37.00	1.48	-7.02	-0.35
no target mse loss	39.65	94.54	46.00	-2.91	-8.31	-3.17
strong target mse loss	69.68	55.86	100.00	1.18	-1.50	0.07

We find that setting the clip_epsilon_logits parameter (the margin clip point that switches off preference loss for winning logits) to 16 – effectively disabled – results in model collapse. Logits diverge much further from reference, and output degrades to single-word repetitions. With this parameter set to 2 (the default), the model reaches the 95% preference accuracy stopping point with writing quality preserved.

With lambda_mse_target reduced to 0, disabling the reference tether for target logits, we observe faster training and logits diverging farther from reference. Writing quality degrades 71% from the baseline per our rubric, illustrating the protective effect of this loss component. When lambda_mse_target is set to 0.4, logits diverged much less from reference, but the model was only able to achieve 74% preference accuracy by training completion. At the default value of 0.05, the model reached the 95% preference accuracy target without any substantial output degradation.

E DPO β Hyperparameter Ablation

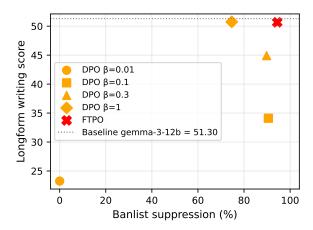


Figure 8: We examine the impact of DPO's β hyperparameter, training gemma-3-12b on our final-token preference set with several values of β : 0.01, 0.1, 0.3 and 1.0. This training set suppresses a banlist of 1,000 items. With DPO, we observe an expected tradeoff in learnability vs degradation (Wu et al., 2024). DPO manages a < 1% reduction in output quality at $\beta = 1.0$, but at the expense of significantly impaired banlist suppression (74.7%). At lower values of β , output quality is markedly reduced for the DPO-trained models. In comparison, the FTPO model trained on the same dataset achieves the highest suppression rate of 94.4% suppression, with neglibible (< 1%) degradation in longform writing score.

F SUPPRESSION PERFORMANCE VS WRITING QUALITY FOR EQ-BENCH DATASET

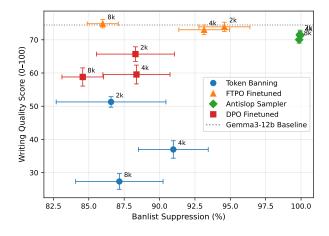


Figure 9: We replicate 6.2 with an out-of-distribution writing prompts dataset. While a smaller dataset size of 96 prompts (and correspondingly larger error bars), we observe a similar pattern of banlist suppression rates and impact on writing quality for each method.

G MOST COMMON OVER-REPRESENTED WORDS AND TRIGRAMS ACROSS MODELS

pattern	percent models
flickered	98.5
flicker	94.0
flickering	92.5
leaned	82.1
muttered	82.1
gaze	80.6
grinned	80.6
containment	77.6
gestured	77.6
addendum	74.6
murmured	73.1
nodded	73.1
glint	68.7
hesitated	68.7
whispered	68.7
blinked	64.2
hummed	64.2
faintly	62.7
leans	62.7
unreadable	62.7

Table 4: Top overlapping words across 67 AI models. Each entry shows the % of models in which the token appears among their top 120 most over-represented words (relative to a human baseline).

pattern	percent models
voice barely whisper	68.7
said voice low	61.2
air thick scent	49.3
took deep breath	44.8
smile playing lips	43.3
something else something	37.3
said voice barely	35.8
voice barely audible	35.8
take deep breath	32.8
could shake feeling	31.3
eyes never leaving	29.9
casting long shadows	28.4
says voice low	26.9
something else entirely	26.9
heart pounding chest	25.4
one last time	23.9
spreading across face	22.4
air thick smell	19.4
could help feel	19.4
long shadows across	19.4

Table 5: Top overlapping trigrams across 67 AI models. Each entry shows the % of models in which the phrase appears among their top 40 most over-represented trigrams (relative to a human baseline).

```
972
       H WRITING QUALITY RUBRIC PROMPT
973
974
         You are an expert in assessing creative writing. Your task is to
975
             score the test model's response below, by several metrics, on a
976
             0-20 scale.
977
978
         [PROMPT START]
979
         {writing_prompt}
980
981
         [PROMPT END]
982
983
         [TEST MODEL RESPONSE]
984
         {test_model_response}
985
986
         [TEST MODEL RESPONSE END]
987
988
         [Task]
989
         You are an expert in assessing creative writing. Your task is to
990
            score the model's response below, by several metrics, on a 0-20
991
            scale.
992
        Scoring notes:
993
994
         - In the output, write the metric names exactly as below so they can
995
            be parsed.
996
997
         - Use the designated output format exactly.
998
         - All criteria are "higher is better"
999
1000
         - You are a critic, and your job is to be critical, especially of any
1001
            failings or amateurish elements.
1002
1003
         - Output format is:
1004
         [Scores]
1005
1006
        Metric 1 name: [Score 0-20]
1007
1008
        Metric 2 name: ...
1009
1010
1011
        Now, rate the supplied model output on the following criteria:
1012
1013
         Spelling/grammar
        Formatting issues & artifacts
1014
        Coherence
1015
        Consistency of tense, pronouns, perspective
1016
        Repetition issues
1017
        Overall quality of the piece
1018
```

Figure 10: Writing quality rubric prompt: This prompt was used to assess the overall quality of creative writing outputs in our experiments, with a particular focus on the common modes of degradation.

I IMPACT ON METRICS BY BANLIST SIZE

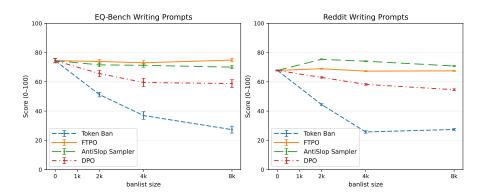


Figure 11: Impact on writing quality per our LLM-judged rubric at several banlist sizes, for each suppression method (Token banning, FTPO, Antislop Sampler and DPO).

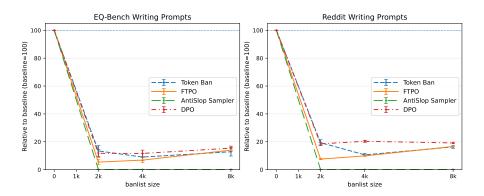


Figure 12: Impact on banlist suppression rates at several banlist sizes, for each suppression method (Token banning, FTPO, Antislop Sampler and DPO).

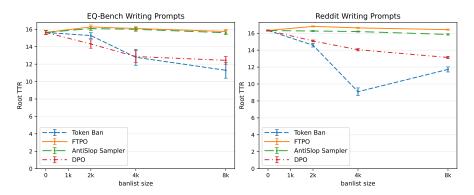


Figure 13: Impact on lexical diversity at several banlist sizes, for each suppression method (Token banning, FTPO, Antislop Sampler and DPO).

J FTPO Loss Function Definition

Preference Loss Component:

For each chosen token index c against a rejected token index r, define the logit gap

$$\Delta = y[c] - y[r].$$

The margin requirement is m. A smooth penalty is applied if the gap is smaller than m:

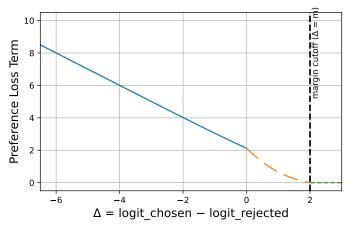
$$\ell^{\text{pref}} = \log(1 + e^{(m-\Delta)/\tau}),$$

with $\tau = 1$ here. A taper weight

$$w = \operatorname{clamp}\left(\frac{m-\Delta}{m}, 0, 1\right)$$

shrinks the contribution as Δ approaches the margin. The preference loss is the weighted mean over chosen tokens:

$$\mathcal{L}_{\mathrm{pref}} = \frac{\sum w \, \ell^{\mathrm{pref}}}{\sum w}.$$



 Δ < 0: chosen losing → full softplus penalty ---- Δ ≥ m: margin satisfied → loss = 0 - 0 ≤ Δ < m: chosen winning → tapering penalty

Figure 14: Preference loss component as a function of the logit gap Δ . When $\Delta < 0$ (chosen losing), the penalty is large. As Δ increases toward the margin m, the penalty smoothly tapers. Once $\Delta \geq m$, the weight goes to zero and the preference loss no longer contributes.

MSE tether terms:

Let deviations be $d_j = y[j] - y^{ref}[j]$. Define:

- Target set $T = \{c\} \cup \{r\}$ (chosen and rejected indices).
- Non-target set $N = \{1, \dots, V\} \setminus T$.

Non-target MSE loss term:

$$\mathcal{L}_{ ext{nontarget}} = rac{\sum_{j \in N} d_j^2}{|N|}.$$

Target MSE loss term with zero-penalty window

$$e_j = \max(|d_j| - \tau_{\text{target}}, 0), \qquad \mathcal{L}_{\text{target}} = \frac{\sum_{j \in T} e_j^2}{|T|}.$$

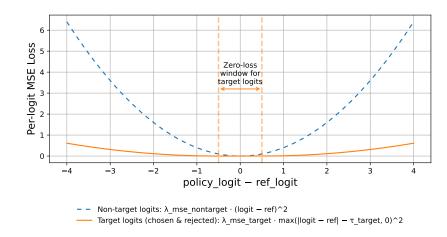


Figure 15: MSE loss components as functions of logit deviation from the reference. The non-target term (blue) penalizes any deviation quadratically. The target term (orange) allows a dead zone around zero, where no penalty applies, then grows quadratically once the deviation exceeds the zero-penalty window.

Here τ_{target} is a zero-penalty window: if the chosen or rejected logits are within $\pm \tau_{\text{target}}$ of the reference, no penalty is applied.

Total objective:

With weighting coefficients $\lambda_{\text{nontarget}}$ and λ_{target} , the total FTPO loss is

$$\mathcal{L} = \mathcal{L}_{pref} + \lambda_{nontarget} \, \mathcal{L}_{nontarget} + \lambda_{target} \, \mathcal{L}_{target}.$$

This formulation allows the model to learn a clear preference signal while preventing uncontrolled drift of the logit distribution.

K SLOP PROFILE CLUSTERING BETWEEN MODELS

1194

1195

1196

1197

1198

1199

1200

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

1231

1232

1233

1234

1236

1237 1238

1239

1240

1241

1188

1189

Colloquially, slop may refer to over-used words, phrases, themes or writing styles. Here we focus on over-used words and n-grams as they are relatively straightforward to extract. For a given model, we generate outputs from a creative writing prompts dataset (Paech, 2023) and a writing prompts dataset sourced from Reddit (Nitral-AI, 2024). We then compute a list of the most over-represented words and bigrams/trigrams relative to a human baseline. The human baseline we use for individual words is the Python library wordfreq (Speer et al., 2018). For bigrams/trigrams, we compute a human baseline from a mix of sources including a large Reddit creative writing dataset, and a selection of public domain works from the Gutenberg Library (Project Gutenberg). For n-gram extraction, we remove stop-words.

A "slop fingerprint" is collated from the top 120 most over-represented words and the top 40 most over-represented bigrams and trigrams. To avoid over-indexing on high-frequency words & phrases in single texts (e.g. a character name), we require the pattern to occur from at least 3 writing prompts independently. To examine the relationship of this fingerprint between models, we perform hierarchical clustering on these top-200 lists per the average rank-distance between each model pair (Figure 16).

It's important to distinguish between counting the frequency of words and n-grams in a text, and calculating their frequency *relative to a human baseline*, as we are doing here. The former simply surfaces patterns that are common in writing; the latter surfaces repetitive writing tendencies of a model that begin to stand out across multiple generations, leading to the perception of "slop". In some models this repetition is extreme: *mistral-small-3.1-24b-instruct-2503* produced 102 "eyes never leaving" trigrams and 62 "voice barely whisper" trigrams across just 96 writing prompts.

Rank-Distance Slop Clustering mistral-small-3.2-24b-instruct-2506 deepseek-chat-v3-0324 chatgpt-4o-latest mistral-medium-3.1 deepseek-r1 kimi-k2 gpt-5-2025-08-07 claude-sonnet-4 claude-opus-4 claude-3.7-sonnet gemini-2.5-pro-preview deepseek-chat-v3.1 deepseek-r1-0528 gemma-3-27b-it gemma-3-12b-it gemma-2-9b-it gemini-2.0-flash-001 gpt-oss-120b gpt-4-0314 gpt-3.5-turbo-0613 mistral-small-3.1-24b-instruct-2503 mistral-small-24b-instruct-2501 llama-3.1-70b-instruct llama-3.1-405b-instruct llama-3.1-8b-instruct llama-4-maverick llama-4-scout alan_edward_nourse gemma-3-12b-it-ftpo-4k mark twain jack london h g wells leo tolstoy arthur conan doyle rudyard kipling andre_norton jane austen I_frank_baum

Figure 16: Top 200 over-represented words and bigrams/trigrams were extracted for each model relative to a human baseline, for a set of creative writing outputs. For included human authors, a selection of their works were used. A dendrogram was generated with cluster distance as the **average ranking distance** of the top over-represented words & n-grams list between models. Our FTPO antislop finetune of gemma-3-12b is highlighted, clustering closer to human authors than any other tested model.

0.5

0.6 0.7 0.8 0.9 Cluster Distance

We find a high correlation in words and n-grams found on the top most over-represented lists across the models tested, with "flickered" appearing on 98.5% of lists, and the trigram "voice barely whisper" appearing on 68.7% of lists. See Table 4 for the most commonly co-occurring word patterns across slop fingerprints, and Table 5 for trigram patterns.

We utilise this method for identifying over-represented usages to compile a target list for slop reduction with the Antislop Sampler and FTPO fine-tuning. It should be noted that this method of identifying slop is domain-specific; the over-used patterns in creative writing will differ from professional writing, for instance. Here, we focus on creative writing, however the method can be applied to other domains by choosing a different set of prompts from which to derive the slop list.

```
1242
          REGEX BLOCKLIST USED FOR "NOT x, BUT y"
1243
1244
       regex_patterns: [
1245
         "\\b(?:\\w+n(?:[']t)|not\\s+(?:just|only|merely|because))\\s+(?:(?![.;<sub>|</sub>
1246
             :?!]).){1,100}?[.;:?!]\\s*(?:it|they|you)(?:['](?:s|re|m))?\\b(?!\|
1247
             \s+(?:was|were|is|are|wasn[']t|weren[']t|isn[']t|aren[']t|ain[']t)
1248
             \b) (?:\s*[*]?\s*)?(?!when\b|then\b|but\b|and\b|yet\b) (?!ri]
1249
             qht\\b) (?!normal\\b) (?!true\\b) (?!sure\\b) (?!only\\b) (?!still\\b) (_|
1250
             ?!rarely\b) (?!already\b) (?!wrong\b) (?!want\b) (?!just\b) (?!cou
             ldn\b) (?!could\b) (?!saw\b) (?!started\b) (?!remember\b) (?!strug_l)
1251
             gled\b) (?!watched\b) (?!goal\b) (?!took\b) (?!kept\b) (?!reminded\b)
1252
             \\b) (?!time\\b) (?!have\\b) (?!acted\\b) (?!smiled\\b) (?!think\\b) (?!
1253
             give\\b) (?!qrab\\b) (?!qave\\b) (?!turn\\b) (?!justify\\b) (?!\\w+ly\\
1254
             b) (?=[a-z] \{4,\} \b) [a-z] + \w*",
1255
         "\\b(?:\\w+n(?:[']t)|not)\\s+(?:just|only|merely)?\\s*(?:(?![-]|[.?!])|
1256
             .) \{1,80\}? [-] \{1,2\}\\x+\w+\(?:[']\\w+\)?\x+",
1257
1258
         "\\b(?:wasn[']t|weren[']t|isn[']t|aren[']t|ain[']t|not)\\s+(?!\\b(?:mi
1259
             nute|minutes|hour|hours|day|days|year|years|second|seconds) \\b) (?!
1260
             with\\b) (?!even\\b) (?:(?![.;:?!]).) {2,120}?[.;:?!]\\s*(?:it|they|y_1)
             ou|that) (?:\s+(?:was|were|is|are) \b(?:\s+[*_~]?\w+[*_~]?) ?| (?:|was|were|is|are) 
1261
             ['](?:s|re|m))\\b(?:\\s+[*_~]?\\w+[*_~]?)?)",
1262
1263
         "\\bno\\s+longer\\s+(?:just|only|merely)?\\s+[^.;:?!]{1,120}[.;:?!]\\s<sub>|</sub>
1264
             *(?:it|they|you)\\s+(?:is|are|was|were)\\b(?:\\s+[*_~]?\\w+[*_~]?)|
1265
             ?",
1266
         "\\b(?:wasn[']t|weren[']t|isn[']t|aren[']t|ain[']t|not)\\s+(?:just|onl
1267
             y|merely)?\s*(?:(?!\bbut\b|[.?!]).){1,80}?[,;:\-]\s*but\s+(?|
1268
             !I\\b)(?:also\\s+)?"
1269
1270
1271
1272
1273
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
```

```
1296
     M AUTO-ANTISLOP CONFIGURATION FILE FOR GEMMA-3-12B-IT 2K
1297
         BANLIST SIZE
1298
1299
1300
1301
     1302
     # MAIN AUTO-ANTISLOP CONFIGURATION
1303
     1304
     1305
     # RUN SETUP
1306
     1307
     experiment_base_dir: "results/auto_antislop_runs" # Base for timestamped
1308
        run directories
     human_profile_path: "data/human_writing_profile.json"
1309
     log_level: "INFO"
1310
     # Iteration 0: Generates the baseline dataset & computes slop
1311
        strings/ngrams to ban
1312
     # Iteration 1: Generates a dataset using antislop, banning those strings
1313
        \& ngrams. Recomputes the slop strings/ngrams at the end \& adds any
1314
        new slop to the banlists
     # Iteration 2+: Extra iterations catch slop that emerges after the
1315
        initial set is banned
1316
     num_iterations: 2 # Minimum 2 iterations (this is enough to catch most
1317
        slop)
1318
     model id: "google/gemma-3-12b-it" # Global model id for the pipeline. Can
1319
        be overridden on individual steps.
1320
     1321
     # VLLM SERVER MANAGEMENT (Conditional: if --manage-vllm is True)
1322
     1323
     manage_vllm: true
     vllm_model_id: null # Model served by vLLM (if unset, will use model_id)
1324
     vllm_port: 8000
1325
     vllm_hf_token: null # Optional: Your Hugging Face token if model is gated
vllm_cuda_visible_devices: "0" # set to e.g. "0,1,2,3" for multiple gpus
1326
1327
     vllm_qpu_memory_utilization: 0.85 # leave some room for the refusal
1328
        classifier if you are using it (about 3gb)
1329
     vllm_max_model_len: 4500
     vllm_dtype: "bfloat16"
1330
     # Additional raw CLI arguments for vLLM server, e.g.,
1331
        ["--tensor-parallel-size", "4"] for multiple gpus
1332
     vllm_extra_args: [] # each param as a separate string, e.g.
1333
        ["--quantization", "bitsandbytes"]
1334
     vllm_env:
                            # env vars for the vLLM process
       # VLLM_USE_V1: "1" # may be needed for amd gpus
1335
1336
1337
     1338
     # GENERATION PARAMETERS (using antislop-vllm)
     1339
     generation_step_enabled: true
1340
1341
     # --- API & Model Configuration ---
1342
     # If you set manage_vllm=true, leave the base url unset
1343
     #generation_api_base_url: "http://localhost:8000/v1"
1344
     #generation_api_base_url:
        "https://apjmbtwbrb8t61-8888.proxy.runpod.net/v1"
1345
     generation_model_id: null # Model id for generation requests (if unset,
1346
        uses model_id)
1347
     generation_api_key: "xxx" # API key for the vLLM server
1348
1349
     # --- Core Generation Settings ---
     generation_max_new_tokens: 1000
```

```
1350
      generation_threads: 50 # Number of parallel threads for API queries in
1351
          antislop-vllm. Note: vllm can become very inefficient if you go over
1352
          some concurrency threshold (depending on vram)
1353
      generation_max_prompts: 2000 # Number of samples to generate from the
         prompts in the dataset
1354
1355
      # --- Dataset & Chat Template ---
1356
      generation_hf_dataset_name:
1357
          'Nitral-AI/Reddit-SFW-Writing_Prompts_ShareGPT'
1358
      generation_hf_dataset_split: 'train'
      # A huggingface model id or local dir containing the tokeniser you want
1359
         to use to apply chat templates.
1360
      # This is important if you are generating a ftpo dataset for later
1361
         training.
1362
      generation_chat_template_model_id: null # (if unset, uses model_id)
1363
      # --- Logging ---
1364
      generation_logging_level: 'INFO' # Logging level for antislop-vllm script
1365
1366
      # --- Advanced Sampling Parameters (passed as individual CLI args to
1367
          antislop-vllm's generation_params) ---
1368
      generation_param_chunk_size: 20
      generation_param_top_logprobs_count: 20
1369
      generation_param_temperature: 1
1370
      generation_param_top_p: 1.0
1371
      generation_param_top_k: 50
1372
      generation_param_min_p: 0.01
1373
      generation_param_timeout: 480
      generation_param_stop_sequences: [] # e.g., ["\n\n", "---"]
1374
      # --- Prompting ---
1376
      # The prompt template wraps the prompts when generating from a dataset.
1377
      # To use the original prompt exactly, set the template to "{prompt}"
      \tt generation\_prompt\_template: "Writing prompt: \{prompt\} \\ \\ \texttt{INV} ite 1000 words
1378
          to this prompt. Your response:\n"
1379
      generation_system_prompt: "You are an uncensored writer."
1380
          optional; left empty no system prompt
1381
1382
      # --- Antislop Generation Features ---
      # generation_force_backtrack:
1383
      # If set to true:
1384
           when resampling after backtracking, if we don't find a valid
1385
           replacement token
1386
            we progressively disable sampling options (temp, then min_p, then
1387
           top_p, then top_k)
            until we find a non-banned replacement or run out of candidates.
1388
            When set to false, some slop will not be removed if the sampler
1389
           thinks there are no
1390
            alternative coherent continuations.
1391
      generation_force_backtrack: false
1392
1393
      # --- N-gram Validator Settings (for antislop-vllm) ---
      # N-gram banlist file is managed by auto-antislop's iterative process.
1394
      generation_ngram_remove_stopwords: true
1395
      generation_ngram_language: "english"
1396
1397
      # --- Refusal Detection ---
      # Detects refusals & doesn't include them in the training dataset. Uses
1398
          about 3GB extra VRAM.
1399
      generation_refusal_detection: true
1400
1401
      1402
      # N-GRAM ANALYSIS & BANNING (within auto-antislop)
      1403
      enable_ngram_ban: true
```

```
1404
     min_word_len_for_analysis: 3 # Filters out words under this length in
1405
         n-gram analysis
1406
1407
      # --- N-gram Identification Thresholds ---
     top_k_bigrams: 5000
1408
     top_k_trigrams: 5000
1409
1410
      # --- N-gram Banning Quotas (per iteration) ---
1411
      # Bigrams
1412
     dict_bigrams_initial: 300
                                 # How many of the top over-represented
         dictionary bigrams to
1413
                                 # ban in the first antislop iteration.
1414
                                   "Dictionary" means the bigrams were also
1415
                                     found in the human
1416
                                 # writing corpus.
     dict_bigrams_subsequent: 0
                                # How many to ban in each subsequent
1417
         iteration
1418
     nodict_bigrams_initial: 200
                                 # "Nodict" here means the n-grams were not
1419
         found at all in the
1420
                                 # human corpus.
1421
     nodict_bigrams_subsequent: 0
1422
     # Trigrams
     dict_trigrams_initial: 300
1423
     dict_trigrams_subsequent: 0
1424
     nodict_trigrams_initial: 200
1425
     nodict_trigrams_subsequent: 0
1426
1427
      # --- User-Defined N-gram Bans ---
      # User-supplied extra n-grams to always ban (processed by auto-antislop)
1428
     extra_ngrams_to_ban: [
1429
       # "voice barely whisper",
1430
1431
     1432
      # OVER-REPRESENTED WORD ANALYSIS & BANNING
1433
      1434
     compute_overrep_words: true
1435
      top_k_words_for_overrep_analysis: 200000
1436
      # --- Quotas for Adding Over-represented Words to Slop Phrase banlist ---
1437
     dict_overrep_initial: 920
                                   # How many of the top over-represented
1438
         dictionary words to
1439
                                   # ban in the first antislop iteration.
1440
                                   # "Dictionary" means the words were also
1441
                                      found in the human
                                   # writing corpus.
1442
     dict_overrep_subsequent: 0
                                 # How many to ban in each subsequent
1443
        iteration
1444
                                   # "Nodict" here means the n-grams were
     nodict_overrep_initial: 80
1445
         not found at all in the
1446
                                   # human corpus.
1447
     nodict_overrep_subsequent: 0
1448
      1449
      # SLOP PHRASE BANNING
1450
      1451
      # Slop phrases are over-represented whole phrases extracted from the
1452
         generated texts.
1453
     enable_slop_phrase_ban: true
1454
     min_phrase_freq_to_keep: 2 # Min frequency for a new phrase from
1455
         slop-forensics to be considered
1456
     top_n_initial_slop_ban: 0 # New slop phrases from slop-forensics to ban
1457
         in iter 0
```

```
1458
          top_n_subsequent_slop_ban: 0 # New slop phrases from slop-forensics to
1459
                ban in later iters
1460
1461
           # --- User-Defined Slop Phrase Bans ---
           # User supplied list of strings to always ban
1462
           # - case insensitive
1463
           # To trigger a ban, the sequence must not have a word-like character
1464
                   (not punctuation or whitespace) directly on either side. That is to
1465
                say, we
1466
                   are not banning disallowed sequences that occur as substrings in
                 longer
1467
                   words. The exception is if the banned string is already bookended by
1468
                   a non-word character.
1469
1470
                  Examples:
                  banned string "cat"
1471
                      - won't trigger a ban for "cation"
1472
                     will trigger a ban on "cat[morecat]"
1473
                 banned string "cat["
1474
                     - *will* trigger a ban on "cat[morecat]", because the banned string
1475
                         ends with a non-word character.
           extra_slop_phrases_to_ban: [
1476
              # "...", "...", "rain", "tapestry", "static", "regret", "rust"
1477
1478
1479
           # --- Whitelisted Strings ---
1480
           # These will be excluded from the list of slop strings that the pipeline
1481
                finds.
           # Note: special tokens in the tokenizer and parts of the chat template
1482
                are
1483
                        automatically whitelisted.
1484
           whitelist_strings: [
1485
              # "think", "thinking"
1486
1487
           1488
           # REGEX BANNING
1489
           1490
           # User-supplied regex patterns to ban
           # Note: unoptimised regex patterns can slow down antislop generation, as
1491
                they will be called often on large texts.
1492
           extra_regex_patterns: [
1493
              # These ones ban "it's not x, it's y" type patterns:
1494
1495
              #"\\b(?:\\w+n(?:[']t)|not\\s+(?:just|only|merely|because))\\s+(?:[.]
                    ;:?!]).){1,100}?[.;:?!]\\s*(?:it|they|you)(?:['](?:s|re|m))?\\b(?!
1496
                    \\s+(?:was|were|is|are|wasn[']t|weren[']t|isn[']t|aren[']t|ain[']t|
1497
                    )\b)(?:\s*[*]?\s*)?(?!when\b|then\b|but\b|and\b|yet\b)(?!r|
1498
                    ight\b) (?!normal\b) (?!true\b) (?!sure\b) (?!only\b) (?!still\b) [
1499
                    (?!rarely\b) (?!already\b) (?!wrong\b) (?!want\b) (?!just\b) (?!co|
                    uldn\\b) (?!could\\b) (?!saw\\b) (?!started\\b) (?!remember\\b) (?!stru
                    ggled\b) (?!watched\b) (?!goal\b) (?!took\b) (?!kept\b) (?!reminde_|
1501
                    d\b) (?!time\b) (?!have\b) (?!acted\b) (?!smiled\b) (?!think\b) (?
1502
                    !qive\\b) (?!qrab\\b) (?!qave\\b) (?!turn\\b) (?!justify\\b) (?!\\w+ly\\
1503
                    \b) (?=[a-z] \{4,\} \b) [a-z] + \w*",
1504
1505
              #"\\b(?:\\w+n(?:[']t)|not)\\s+(?:just|only|merely)?\\s*(?:(?![-]|[.?!]|
                    1506
1507
              \#"\b(?:wasn[']t|weren[']t|isn[']t|aren[']t|ain[']t|not)\s+(?!\b(?:m|
1508
                    inute|minutes|hour|hours|day|days|year|years|second|seconds) \b) (?|
1509
                    !with\\b) (?!even\\b) (?:(?![.;:?!]).) \{2,120\}?[.;:?!]\\s*(?:it|they||
1510
                    you|that) (?: \s+(?:was|were|is|are) \b(?: \s+[*_~]?\w+[*_~]?) ?|(?|vu|that) | (?: \s+[*_~]?) | (?|vu|that) | (?: \s+[*_~]?) | (?: \s+[*_~
                    :['](?:s|re|m))\b(?:\\s+[*_~]?\\w+[*_~]?)?)",
1511
```

```
1512
        #"\\bno\\s+longer\\s+(?:just|only|merely)?\\s+[^.;:?!]{1,120}[.;:?!]\\|
1513
            s*(?:it|they|you)\\s+(?:is|are|was|were)\\b(?:\\s+[*_~]?\\w+[*_~]?\
1514
            )?",
1515
        #"\\b(?:wasn[']t|weren[']t|isn[']t|aren[']t|ain[']t|not)\\s+(?:just|on|
1516
            ly|merely|?\s*(?:(?!\bbut\b|[.?!]).){1,80}?[,;:\-]\s*but\s+(|
1517
            ?!I\\b) (?:also\\s+)?"
1518
1519
1520
      1521
      # FINETUNING
1522
      1523
      finetune_enabled: true
1524
      # --- General Finetuning Setup ---
1525
      finetune_use_unsloth: false
1526
      finetune_mode: "ftpo" # ftpo | dpo-final-token (final token preference
1527
          optimisation)
1528
      finetune_ftpo_dataset: "" # you can specify an existing ftpo dataset,
1529
         or leave unset to let the
                                  # pipeline use the one produced in the
1530
                                     generation step
1531
      finetune_base_model_id: null # Base model for DPO (if unset, uses
1532
         model_id)
1533
      finetune_max_seq_length: 2500 # this may truncate some outputs
1534
      finetune_load_in_4bit: true # glora
1535
      # --- Early Stopping ---
1536
      finetune_early_stopping_wins: 0.85 # Early stopping threshold for
1537
          fraction of *chosen* completions that are selected over *rejected*.
1538
                                          # More than 0.85 may be overtrained.
1539
                                             Set to > 1.0 to disable early
1540
                                             stopping.
      finetune_early_stopping_loss: null # Loss threshold for early stopping.
1541
          Set to null to disable.
1542
1543
      # --- LoRA Configuration ---
1544
      finetune_lora_r: 256 # the ftpo trainer works best with a high lora rank
      finetune_lora_alpha: 256
1545
      finetune_lora_dropout: 0.05
1546
      finetune_weight_decay: 0.01
1547
      finetune_target_modules: ["up_proj", "down_proj", "lm_head"]
1548
1549
      # --- Layer Freezing ---
      finetune_freeze_early_layers: true
1550
      finetune_n_layers_unfrozen: 5
1551
1552
      # --- Training Process ---
1553
      finetune_gradient_checkpointing: "unsloth"
1554
      finetune_chat_template: "" \# e.g. "gemma-3" -- get the chat template from
          unsloth's helper if required, otherwise leave the string blank to use
1555
          the tokeniser's chat template
1556
      finetune_batch_size: 3
1557
      finetune_gradient_accumulation_steps: 5
1558
      finetune_warmup_ratio: 0.1
1559
      finetune_num_epochs: 1
1560
      # --- Learning Rate ---
1561
      finetune_learning_rate: 0.000001
1562
      finetune_auto_learning_rate: true # true: automatically determine
1563
          learning rate based on dataset size, effective batch size & lora rank
1564
      finetune_auto_learning_rate_adjustment_scaling: 0.08 # scale the auto-lr
1565
         by this factor
```

```
1566
       # --- DPO/FTPO Specific ---
1567
      finetune_beta: 0.1 # DPO beta
1568
       # --- Output & Saving ---
1569
      finetune_output_dir_suffix: "_ftpo_exp01" # Appended to experiment run
1570
          dir
1571
      finetune_save_merged_16bit: true
1572
      finetune_save_gguf_q8_0: false
1573
1574
       # --- Dataset Handling for Finetuning ---
      finetune_max_train_examples: 12000 # adjust as needed
1575
      finetune_shuffle_seed: 42
1576
1577
      # --- FTPO Sample Regularization ---
1578
      # 0 = off; 0.9 strongly downsamples overrepresented rule violations
       # (this is useful because the raw generated dataset is typically very
1579
          skewed)
1580
      ftpo_sample_rejected_regularisation_strength: 0.8
1581
      ftpo_sample_chosen_regularisation_strength: 0.2
1582
      ftpo_sample_min_chosen_tokens: 4 # filter out ftpo samples that have
1583
          fewer than this number in the chosen tokens list
1584
1585
       # FTPO-specific hyper-parameters
1586
       # Leave any of these out (or set to null) to fall back to FTPOTrainer
1587
          defaults.
1588
       # Loss terms are computed separately for the target (chosen + rejected)
1589
          tokens vs the remainder of the vocab.
1590
       # This is because we want to allow more freedom of movement for the
1591
          target tokens.
1592
1593
       # MSE loss term 1: light mse loss applied tokenwise on target tokens
      ftpo_lambda_mse_target: 0.05
                                     # Strength of MSE loss tether on the
1594
          individual logits in the
1595
                                                    chosen+rejected set vs
1596
                                                    reference.
1597
                                       # Grace bandwidth (logits) before the
      ftpo_tau_mse_target: 0.5
1598
          above MSE loss kicks in.
1599
       # MSE loss term 2: stronger mse term applied to remaining (non-target)
1600
          vocab
1601
      ftpo_lambda_mse: 0.4
1602
1603
      ftpo_clip_epsilon_logits: 2
                                       # For a chosen token: "after winning vs
          rejected token by this margin, preference loss turns off"
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
```