Diverse, not Short: A Length-Controlled Self-Learning Framework for Improving Response Diversity of Language Models

Anonymous ACL submission

Abstract

Diverse language model responses are crucial for creative generation, open-ended tasks, and self-improvement training. We show that common diversity metrics, and even reward models used for preference optimization, systematically bias models toward shorter outputs, limiting expressiveness. To address this, we introduce Diverse, not Short (Diverse-NS), a length-controlled self-learning framework that improves response diversity while maintaining length parity. By generating and filtering preference data that balances diversity, quality, and length, Diverse-NS enables effective training using only 3,000 preference pairs. Applied to LLaMA-3.1-8B and the Olmo-2 family, Diverse-NS substantially enhances lexical and semantic diversity. We show consistent improvement in diversity with minor reduction or gains in response quality on four creative generation tasks: Divergent Associations, Persona Generation, Alternate Uses, and Creative Writing. Surprisingly, experiments with the Olmo-2 model family (7B, and 13B) show that smaller models like Olmo-2-7B can serve as effective "diversity teachers" for larger models. By explicitly addressing length bias, our method efficiently pushes models toward more diverse and expressive outputs.

1 Introduction

003

014

017

034

042

Alignment has played a key role in making large language models (LLMs) broadly useful, controllable, and safe for real-world applications (Schulman et al., 2017; Bai et al., 2022; Dai et al., 2023; Ouyang et al., 2022; Longpre et al., 2023). As a form of post-training, it typically involves a combination of instruction tuning (Longpre et al., 2023; Peng et al., 2023; Ouyang et al., 2022) and preference optimization (Schulman et al., 2017; Ouyang et al., 2022; Rafailov et al., 2023), enabling models to follow human instructions and generate responses that are helpful, harmless, and honest (Bai et al., 2022; Dai et al., 2023). However, alignment comes at a cost: several studies have found that alignment can significantly reduce the diversity of model outputs (Kirk et al., 2023; Doshi and Hauser, 2024; Padmakumar and He, 2023; Anderson et al., 2024; Shaib et al., 2024b).

This decrease in diversity has important consequences. When humans collaborate with aligned models, the content they produce tends to be less original and less varied (Doshi and Hauser, 2024; Padmakumar and He, 2023). At scale, this reduction in diversity can hinder creative ideation and increase output homogeneity (Anderson et al., 2024; Xu et al., 2024). Beyond creativity, reduced diversity of generated text has a direct impact on the continued improvement of LLMs. Recent studies have shown that repeatedly training models on their own aligned outputs can lead to a consistent decline in diversity, eventually resulting in model collapse (Shumailov et al., 2023; Guo et al., 2023; Seddik et al., 2024; Herel and Mikolov, 2024).

Despite these challenges, alignment remains essential. The question, then, is not whether to align, but how to preserve or recover the output diversity of aligned models. In this work, we ask: *Can we increase the response diversity of aligned models while retaining the the response quality?*

Prior work has explored a range of strategies to improve output diversity of aligned language models, including methods based on prompting, sampling, and targeted training procedures (Lu et al., 2024; Zhang et al., 2020; Tian et al., 2023; Li et al., 2024, 2025; Lanchantin et al., 2025; Chung et al., 2025; Qin et al., 2025). Sampling techniques such as temperature, top-p, and top-k have been shown to increase diversity, though often at the cost of reduced quality (Zhang et al., 2020). Sequential prompting strategies are also helpful in increasing response diversity (Lu et al., 2024; Tian et al., 2023). However, the computational cost scales rapidly with more discussion turns due to increas043

045

111 112 113

114 115

116 117

118

119 120

121 122

> 123 124

125 126

127

129

130 131

132

133

134

ing context length. Training approaches have introduced explicit diversity objectives (Li et al., 2025; Chung et al., 2025; Cideron et al., 2024) and entropy regularization (Li et al., 2024) to encourage more varied outputs. Self-learning methods, where the model generates its own training data, have also been used to promote diversity (Tian et al., 2024; Lanchantin et al., 2025; Qin et al., 2025).

However, one critical confound, text length, has received little scrutiny in recent work. Widely used diversity metrics are length-sensitive and consistently assign higher scores to shorter passages (Covington and McFall, 2010; McCarthy and Jarvis, 2010; Shaib et al., 2024a). While this bias is less problematic in structured generation tasks, optimizing these metrics can reduce expressiveness in open-ended writing, which thrives on depth and nuance, thereby undermining the very creativity they are meant to cultivate. But even though optimizing length-sensitive metrics can clearly backfire, the role of length in both measuring and improving diversity has been largely overlooked. Our work aims to close this gap.

To address this overlooked confounding factor, we propose *Diverse*, not Short (Diverse-NS), a length-controlled self-learning framework that counteracts the hidden brevity bias in standard diversity metrics and improves diversity in both structured and free-form generation. The framework first uses sequential prompting to elicit more diverse responses, followed by preference pair curation that improve both diversity and quality while maintaining comparable response lengths (within ± 5 words). Using these preference pairs, we apply Direct Preference Optimization (DPO) (Rafailov et al., 2023) to improve the response diversity of the base model. Our key contributions are:

- 1. Diverse-NS: A length-controlled self-learning framework that significantly improves the response diversity of Llama-3.1-8B and Olmo-2-7B using only 3k preference pairs.
- 2. Diverse-NS-Lite: A computationally efficient variant that achieves comparable performance to Diverse-NS while significantly reducing the data filtering cost.
- 3. Small-to-large transfer: We highlight the potential of smaller models to serve as effective "diversity teachers" for larger variants, enabling low-cost diversity alignment.
- 4. Length-controlled diversity evaluation: We introduce Diversity Decile, a new metric that

adjusts for text length when evaluating diversity gains.

135

136

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

5. Dataset: We release a high-quality dataset of 6k preference pairs generated from Llama-3.1-8B and Olmo-2-7B to support future research on length-aware diversity alignment.

Related Work 2

Increasing Diversity without Training. Zhang et al. (2020); Chung et al. (2023), shows that common sampling methods such as temperature, top-p, top-k, are comparable in terms of increasing the diversity but, increasing diversity often comes at the price of reduced quality. For curating a generic large-scale dataset, prompting methods can boost topical, stylistic, and formatting diversity (Li et al., 2023; Chen et al., 2024; Face, 2024; Ge et al., 2024). Conversely, for more task-specific datasets, sequential prompting can elicit diverse responses (Lu et al., 2024; Tian et al., 2023; Qin et al., 2025).

Increasing Diversity with Training. Augmenting method-specific objective functions with elements that directly maximize diversity has been successful in increasing response diversity (Li et al., 2024; Chung et al., 2025; Li et al., 2015, 2025). The other approach gaining more attention in recent studies is to adopt a three-step procedure: generate diverse data, filter data for improving quality, and fine-tune LLM on the filtered data (Lanchantin et al., 2025; Chung et al., 2025; Qin et al., 2025). This approach has been successful in task-specific alignment, but more generic self-training has still seen limited success (Li et al., 2023; Face, 2024; Shumailov et al., 2023; Guo et al., 2023; Herel and Mikolov, 2024; Seddik et al., 2024). Our work is closest to the task-specific alignment studies in the self-learning framework (Lanchantin et al., 2025; Qin et al., 2025).

Diversity Evaluation. Evaluation of diversity is challenging for two main reasons: length bias (Mc-Carthy and Jarvis, 2010; Covington and McFall, 2010; Mass, 1972; Johnson et al., 2023), and inconsistent human preferences (Evans et al., 2016; Chakrabarty et al., 2023, 2024; Gómez-Rodríguez and Williams, 2023). Despite the challenges, many studies have highlighted the compromised diversity of synthetic text (Shaib et al., 2024b,a; Salkar et al., 2022; Padmakumar and He, 2023; Guo et al., 2023; Kirk et al., 2023; Doshi and Hauser, 2024; Anderson et al., 2024). So, we present a method,

190

191

192

194

195

196

198

199

201

202

205

209

210

211

212

214

215

216

217

218

219

226

228

Diverse-NS, to increase the response diversity and propose a metric, *Diversity Decile*, to measure diversity in a length-controlled way.

3 Preliminaries

Self-learning, also known as self-training, is a semisupervised approach involving three main steps: data generation (pseudo-labeling), data filtering, and model learning (Lee et al., 2013; Amini et al., 2025). In our setup, data generation involves sampling text from a language model in response to story-writing prompts. This is followed by filtering, where we construct high-quality preference pairs-two continuations for the same prompt, with one preferred over the other. We refer to the preferred continuation as the "chosen" continuation (or response) and the other as the "rejected" continuation (or response). Using this preference dataset, we apply Direct Preference Optimization (DPO) (Rafailov et al., 2023) to train the model to favor the chosen responses.

4 Data

We describe data generation and filtering pipeline designed to elicit diverse model responses for downstream preference tuning. The pipeline first generates candidate stories using a sequential prompting strategy, then filters the pool of generated responses to form preference pairs suitable for Direct Preference Optimization (DPO) training (Rafailov et al., 2023). The preference pairs are formed to maximize the diversity and quality gain while maintaining the same length for "chosen" and "rejected" samples.

4.1 Data Generation

Task Setup. We focus on a creative writing task to build the dataset for preference learning. The goal is to generate short stories (five sentences) that must include three words specified in the prompt. This task has been extensively validated in studies of human creativity (Prabhakaran et al., 2014). To create a diverse set of prompts, we first curated a list of 300 unique words, W_u^{-1} . For generating short stories from LMs, we create prompts by randomly sampling three-word sets from W_u .

Sequential Prompting. Given the task setup, we create 1k story writing prompts, with 1k unique

three-word sets. The exact prompt is provided 229 in Appendix A.1. We initially sampled 10k sto-230 ries (10 per prompt) using a temperature of 1.0 231 from each of the following LMs: Llama-8B and 232 Olmo-7B (Grattafiori et al., 2024; OLMo et al., 233 2024). Within the sampled stories, we extracted 234 the repeating Part-Of-Speech (POS) bigrams and 235 found that the start of the story is highly likely 236 to have repetitions across different prompts (re-237 fer to Table B.1). To overcome these repetitions, 238 we performed a second inference call to re-draft 239 the story with additional constraints, an approach 240 similar to *Denial Prompting* presented by Lu et al. 241 (2024) (refer to Appendix A.1 exact prompt). In 242 our case, unlike Lu et al. (2024), the constraints 243 we use are specifically targeted to elicit a more di-244 verse response from the model while maintaining 245 the same (or comparable) length. With a pilot anal-246 ysis on the initial 20k responses, we find that the 247 story generated in the second inference call is on 248 average more diverse (refer to Table B.2). These re-249 sults motivated us to set up the final two-step data 250 generation process, first inference call to collect 251 natural responses from the model, and second infer-252 ence call to redraft the natural response into a more 253 diverse story. In the final data generation phase, 254 we used 20k unique three-word sets to generate 255 prompts and sampled 10 first and second responses 256 for each prompt, resulting in a dataset of 200,000 257 tuples of prompt, first response, and second re-258 sponse, per model (Llama-8B and Olmo-7B). We 259 denote the data as follows: $\mathcal{D}^{(\pi)} = \{(p, r_1, r_2)_i \mid$ 260 i = 1, ..., 200,000 where, p, r_1 , and r_2 denote 261 the prompt, first response, and second response, re-262 spectively, generated from model (policy) π . Note 263 that $|\{p_1, p_2, \dots, p_{200,000}\}| = 20,000$ and we use 264 two models, $m \in \{\text{Llama-8B}, \text{Olmo-7B}\}$, for data 265 generation. 266

4.2 Data Filtration

The Chosen and Rejected Pools. Each instance in our generated dataset is a tuple (p, r_1, r_2) , where p is the prompt and r_1, r_2 are two responses conditioned on it. The first response r_1 reflects the model's default behavior which are stories generated without intervention, capturing its most likely completion. In contrast, the second response r_2 is generated with additional instructions aimed at reducing repetition, resulting in a more diverse output. We leverage this contrast by designating r_1 as the *rejected* response and r_2 as the *chosen* one. This setup encourages the model to prefer more di-

268

269

270

271

272

273

274

275

276

277

278

¹A manually curated list of 20 words was extended using GPT-40 and Claude-3.7.

291

296

297

301

303

307

308

311

312

313

314

315

317

319

323

324

325

326

329

verse continuations that it is already capable of generating. Hence, it provides a strong self-learning framework for improving diversity.

Filtration Rules. Each pair (r_1, r_2) gives us a natural candidate for rejected and chosen responses. On average, the second response r_2 is more diverse than the first r_1 (Table B.2), but not every pair guarantees learning higher diversity. To ensure that the model receives consistent and useful learning signals, we apply a set of filtering rules.

First, we require that the diversity of r_2 exceeds that of r_1 , so that the model consistently learns to prefer more diverse continuations. However, higher diversity may negatively impact text quality as prior work has shown a trade-off between the two (Zhang et al., 2020). To ensure that preference learning also promotes higher quality, we further require that r_2 be of higher quality than r_1 . Additionally, we filter out cases where both r_1 and r_2 are of poor quality, even if r_2 is marginally better. To do so, we enforce that r_2 must surpass the median quality of all r_1 responses. Lastly, most diversity metrics have been shown to be negatively correlated with text length (Covington and McFall, 2010; Shaib et al., 2024a; McCarthy and Jarvis, 2010), which introduces a bias toward shorter texts. This issue has not been explicitly addressed in the recent studies for training and evaluation of LMs for diversity (Qin et al., 2025; Lanchantin et al., 2025; Chung et al., 2025). To control for this, we constrain r_1 and r_2 to be of comparable length (± 5 words). Ideally, we would like r_1 and r_2 to have exactly the same length. However, in practice, very few examples satisfy this strict constraint, especially when working with smaller language models (under 10B parameters). Therefore, we relax the constraint and allow a maximum length difference of ± 5 words between r_1 and r_2 .

In summary, we retain a data point for preference learning only if it satisfies all of the following conditions, applied in order:

- The quality of r_2 is greater than or equal to the 50^{th} percentile of all r_1 quality scores.
- The quality of r_2 is greater than r_1 .
- The diversity of r_2 is greater than r_1 .
- The absolute difference in word count between r_1 and r_2 is at most five words.

Diversity and Quality Metrics. We use entropy to measure diversity and the ArmoRM reward model scores (Wang et al., 2024) to assess quality. Entropy is a standard metric for lexical diversity (Lanchantin et al., 2025), with higher values indicating greater diversity. In our self-learning setup, entropy is useful because it reflects the model's likelihood of producing a certain continuation of the prompt. When used in filtering, it helps identify training data that aligns with the model's own capabilities. For each example, we compute the entropy and the reward model score of both r_1 and r_2 , conditioned on the original prompt p. When we use our data generation method, and use entropy and ArmoRM values for filtration, we call our approach, Diverse, not Short (Diverse-NS or D-NS). 330

331

332

333

334

335

336

337

339

340

341

342

343

344

345

346

349

350

351

352

353

354

355

356

357

358

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

Lightweight Filtration. While entropy and ArmoRM scores are high-quality metrics for measuring diversity and response quality, they are computationally expensive. Each example (p, r_1, r_2) requires two additional inference calls to compute entropy and two more for ArmoRM scoring. To reduce this overhead, we evaluated seven alternative metrics and measured their correlation with entropy and ArmoRM scores. Among these, Type-Token Ratio (TTR) showed the highest correlation with entropy (Pearson r = 0.2027, p < 0.0001), and the MAAS index (Mass, 1972) was most correlated with ArmoRM scores (Pearson r = 0.2357, p < 0.0001). Refer to Table 1 for all correlation results. Based on these findings, we replace entropy with TTR and ArmoRM scores with MAAS in our filtering pipeline. When this lightweight variant is used during data filtering, we refer to the resulting method as Diverse-NS-Lite (or D-NS-Lite).

Post-Filtration Properties. Based on the correlation analysis (Tab. 1), it is worth noting that both entropy and ArmoRM scores are negatively correlated with text length. As a result, optimizing for diversity or quality alone may unintentionally favor shorter responses as the "chosen" continuations. To avoid this bias, it is essential to explicitly control for length when curating preference learning data for improving diversity. To show this, we implement a recent study that is closest to our method, Diverse Preference Optimization(DivPO) (Lanchantin et al., 2025). DivPO also generates responses and filters the responses to form preference learning pairs without explicitly control the length of the chosen and rejected continuations. We compare pre- and post-filtration data properties for DivPO and Diverse-NS in Tab. 2. The table clearly shows that in the pursuit of maximizing the entropy

Method	Word Count	TTR	MATTR	HD-D	MTLD	MAAS
Entropy	-0.1574	0.2027	0.0800	0.1071	0.0656	-0.1104
ArmoRM Score	-0.3461	0.1698	-0.0042^{**}	-0.0487	0.0749	0.2357

Table 1: Correlation Analysis. Pearson correlation coefficients between six text statistics and two target metrics: entropy (diversity) and ArmoRM reward scores (quality). Both entropy and ArmoRM scores show negative correlation with text length. Among diversity metrics, TTR exhibits the strongest correlation with entropy, while the MAAS index shows the highest correlation with ArmoRM scores. **: p < 0.001; all others: p < 0.0001.

Method	Num. Pref. Pairs	Word Count Δ
No Filtering	200,000	-0.68 ± 11.33
DivPO	3,000	$-49.90\pm$ 17.51
Ours - D-NS-Lite	3,000	-0.90 ± 2.91
Ours - D-NS	3,000	-1.35 ± 2.93

Table 2: **Data Properties After Filtering.** This table reports the average (\pm std.dev.) length difference (Δ) between *chosen* and *rejected*. While DivPO tends to favor significantly shorter *chosen* responses.

values, DivPO selects significantly shorter (-49.90 words shorter on average) responses as the *chosen* responses in the final preference data.

5 Experimental and Evaluation Setup

5.1 Preference Tuning

After generating and filtering the data, we fine-tune the same base policy π that was used to generate it. In other words, data generated by Llama-8B is used to train Llama-8B, and likewise for Olmo-7B. To ensure a fair comparison across methods (DivPO, D-NS, and D-NS-Lite), we limit the final training dataset to 3,000 preference pairs². To construct this 3k dataset, we first compute the entropy gain for each pair as the difference between the entropy of the *chosen* and *rejected* responses ³. We then sort all pairs by entropy gain in descending order and select the top 3k examples. This ensures that the final training set maximizes diversity gain for the base model. The same selection procedure is applied to all three methods.

We further extend our experiments to evaluate the utility of training larger models with data generated from smaller ones. For this, we train Olmo-13B using preference pairs generated from Olmo-7B. We provide all hyperparameter values in Appendix C. All experiments are run on a single NVIDIA RTX 6000 GPU (48GB memory), using a perdevice batch size of 2 and a global batch size of 64. Training Llama-8B or Olmo-7B takes approximately 100–150 minutes while O-13B takes 200-220 minutes per run, highlighting our setup efficiency.

407

408

409

410

411

412

413

414

415

416

417

418

419

5.2 Evaluation

5.2.1 Tasks

We evaluate the model's response diversity with four tasks: Divergent Association Task (DAT), Persona Generation Task (PGT), Alternate Uses Task (AUT), and Creative Writing Task (CWT).

Divergent Associations Task (DAT). The DAT 420 (Olson et al., 2021) is a psychological test com-421 monly used to assess divergent thinking in humans. 422 Participants are asked to generate a list of 10 words 423 that are as dissimilar from each other as possi-424 ble. Recent studies have adapted DAT to evalu-425 ate the creativity of language models, focusing on 426 their ability to produce diverse outputs (Bellemare-427 Pepin et al., 2024). To quantify model performance 428 on DAT, we use the Divergent Semantic Integra-429 tion (DSI) metric (Johnson et al., 2023), which 430 computes the average semantic distance of each 431 word in the generated list from all others. Higher 432 DSI values indicate more divergent thinking and 433 greater ideological diversity. Following Johnson 434 et al. (2023), we extract token embeddings from 435 the 6^{th} layer of BERT-large for the generated list 436 and compute the average pairwise cosine distance 437 between all embeddings. This approach has been 438 shown to correlate strongly with human judgments 439 of creativity (Johnson et al., 2023). We provide the 440 exact prompt used for DAT in Appendix A.2. For 441 a robust evaluation, we sample 100 DAT responses 442 per model using temperature 1.0 and different ran-443 dom seeds. From these 100 lists (each with 10 444 words), we compute and report two metrics: (1) the 445 average and standard deviation of DSI scores, and 446 (2) the number of unique words across all 1,000 447

402

403

404

405

 $^{^{2}}$ We observed that the size of the dataset after filtering is the smallest for Diverse-NS, slightly more than 3k. Hence, to make the training runs more comparable across methods, we limit the size of the dataset to 3k for all methods.

³note that, by construction, the *chosen* response has higher entropy in the filtered set

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

531

532

533

534

535

536

537

539

540

541

542

543

544

545

546

547

498

448 generated tokens. In both cases, higher values indi-449 cate greater diversity.

Persona Generation Task (PGT). To assess di-450 versity in structured generation, we use the PGT, 451 also used in the study conducted by Lanchantin 452 et al. (2025). In this task, the model is prompted 453 to generate a JSON object with three fields: first 454 name, city of birth, and current occupation to eval-455 uate the model's ability to produce varied persona 456 descriptions. The exact prompt is provided in Ap-457 pendix A.2. We sample 100 responses per model 458 using temperature 1.0 and different random seeds. 459 For each key in the JSON object, we report the pro-460 portion of unique values across the 100 responses. 461 Higher uniqueness indicates greater diversity. 462

Alternate Uses Task (AUT). The Alternate Uses 463 Task (AUT) is a common and rigorously validated 464 psychological test to measure human divergent 465 thinking (Guilford, 1956). In this task, the subject/model is asked to generate creative and un-467 conventional uses for objects (e.g., broom). The 468 prompt and list of objects used for evaluation are 469 provided in Appendix A.2. We use 15 unique ob-470 jects and generate 10 responses per object using 471 different random seeds, resulting in 150 total re-472 sponses sampled at temperature 1.0. For quantify-473 474 ing the diversity of the generated uses, we measure the distance between the target object and gener-475 ated uses with the help of BERT-large encodings, 476 a validated approach that correlates with human 477 creativity ratings (Patterson et al., 2023). We report 478 the mean and standard deviation of the distance 479 values, higher values indicate higher diversity. 480

Creative Writing Task (CWT). The CWT — 481 based on a well-validated psychological assessment 482 of creativity (Prabhakaran et al., 2014) - is exactly 483 the same as our data generation task. That is, given 484 a set of three words, the subject/model is tasked 485 with generating a creative short story that includes 486 all three words. We provide a separate list of three-487 word sets used for evaluation in Appendix A.2. We 488 sample 10 responses for each of the seven three-489 word sets with temperature of 1.0. Unlike our other 490 evaluation tasks, we measure the diversity as well 491 as the quality of the generated responses. Similar to 492 493 Johnson et al. (2023), we calculate the DSI metric to measure the diversity of the generated story. For 494 quality measurements, we resort to the ArmoRM 495 reward model preference scores (Wang et al., 2024). 496 We report the average and standard deviation values 497

of DSI and ArmoRM scores, and 4-gram diversity values, where higher values are more desirable for all metrics.

5.2.2 Length-Adjusted Evaluation

While most diversity metrics exhibit bias toward shorter outputs, Johnson et al. (2023) shows that the DSI metric displays the opposite tendency-it favors longer responses. This is not an issue in tasks like DAT, where the output length is fixed at 10 words. But for open-ended tasks such as CWT, longer stories may receive disproportionately high DSI scores primarily due to their length, rather than genuine diversity. To address this issue, we introduce a novel evaluation metric: Δ *Diversity Decile* (Δ *DD*), which takes into account text length when assessing diversity.

Change in Diversity Decile (Δ **DD**). We first build a decile map that captures the empirical distribution of diversity scores at each length. Using 800 000 stories collected from Llama-8B and Olmo-7B over 40000 prompts, we: (1) group responses by word count w; (2) compute decile thresholds for a chosen diversity metric (e.g. TTR, MTLD); and (3) store these percentile thresholds in a lookup table \mathcal{M} . Here, a *decile* refers to one of ten intervals that divide the distribution of diversity scores for a given length into ten equal parts. The top decile corresponds to the most diverse 10% of responses at that length, the second-highest to the next 10%, and so on. This mapping allows us to estimate the *approximate diversity rank* of any new response relative to other responses of the same length. At evaluation time, a new response r with word count w_r and diversity score d_r is assigned the highest decile index $k \in \{0, \dots, 9\}$ such that d_r exceeds the k-th threshold in $\mathcal{M}[w_r]$. Formally, $DD(r, \mathcal{M}) = k$, where larger k means the response is more diverse than a greater share of previously observed texts of the same length.

To evaluate the effect of preference tuning, we average DD scores over 70 CWT prompts for the base and the preference-tuned models and report their difference: $\Delta DD = \overline{DD}_{tuned} - \overline{DD}_{base}$.

Positive ΔDD values indicate improved diversity, with higher values corresponding to a larger improvement. Negative values signify reduced diversity, and $\Delta DD = 0$ signifies no change. Note that, DD is agnostic to the choice of diversity metric. We therefore report Δ DD values using seven standard metrics: TTR, MATTR, HD-D, MTLD,

					Ours		
Task	Metric	Base Model	DivPO	D-NS-Lite	D-NS		
LLaMA-8B							
DAT	DSI	0.7535	0.7545	0.7590	0.7640		
DAT	Unique Words	0.4575	0.4593	0.4797	0.4914		
PGT	Unique First Names	0.6500	0.6100	0.6900	0.6900		
PGT	Unique Cities	0.3300	0.3100	0.4700	0.4200		
PGT	Unique Occupations	0.4100	0.3900	0.5100	0.4900		
AUT	DSI	0.8876	0.8837	0.8876	0.8878		
CWT	DSI	0.8515	0.8521	0.8556	0.8581		
CWT	ArmoRM Score	0.1451	0.1495	0.1369	0.1405		
CWT	4-gram div.	2.8550	2.9320	2.9450	2.9620		
		OLMo-7B					
DAT	DSI	0.7480	0.7509	0.7662	0.7639		
DAT	Unique Words	0.6139	0.6079	0.6347	0.6327		
PGT	Unique First Names	0.3300	0.3300	0.3300	0.3400		
PGT	Unique Cities	0.3100	0.3000	0.2700	0.2700		
PGT	Unique Occupations	0.5200	0.5500	0.6100	0.6100		
AUT	DSI	0.8836	0.8846	0.8852	0.8858		
CWT	DSI	0.8499	0.8491	0.8548	0.8563		
CWT	ArmoRM Score	0.1435	0.1441	0.1462	0.1464		
CWT	4-gram div.	3.1270	3.1690	3.1750	3.1620		
		OLMo-13E	3				
DAT	DSI	0.7233	0.7282	0.7320	0.7364		
DAT	Unique Words	0.3421	0.3340	0.3310	0.3256		
PGT	Unique First Names	0.4100	0.4100	0.4400	0.4500		
PGT	Unique Cities	0.3500	0.3500	0.3700	0.3900		
PGT	Unique Occupations	0.1900	0.1900	0.1900	0.2000		
AUT	DSI	0.8943	0.8960	0.8974	0.8970		
CWT	DSI	0.8557	0.8555	0.8616	0.8614		
CWT	ArmoRM Score	0.1571	0.1589	0.1585	0.1590		
CWT	4-gram div.	3.0820	3.0770	3.095	3.1070		

Table 3: **Diversity and Quality Evaluation.** We present the average diversity (DSI or unique values) and quality (ArmoRM Score) measurements for model responses collected on four creative generation tasks (Structured Gen.: DAT, PGT, Free-Form Gen.: AUT, CWT).

and MAAS. We also compute ΔDD using ArmoRM reward scores to quantify the gain or loss in quality. This length-aware normalization prevents either long or short responses from being over-credited for diversity⁴.

6 Results

548

549

551

552

553

554

555

557

558

559

560

561

563

564

Divergent Associations Task (DAT). In our DAT evaluation (Tab. 3), we see that both Diverse-NS and its lightweight variant deliver clear improvements in diversity over the untrained base and the DivPO baseline across all model sizes. Remarkably, even the D-NS-Lite variant consistently outperforms DivPO, demonstrating that a compact diversity strategy can be highly effective. Interestingly, using data generated by the smaller Olmo-7B to fine-tune the larger Olmo-13B yields diversity gains for every method, highlighting how smaller **Persona Generation Task (PGT).** In our PGT evaluation (Tab. 3), Diverse-NS produces more distinct first names, cities, and occupations than DivPO for every model, with the sole exception of the city metric on Olmo-7B. Outside that one case, Diverse-NS-Lite also outperforms DivPO across all three metrics. Notably, on Llama-8B, Diverse-NS-Lite matches or exceeds the baseline and Diverse-NS on every attribute of the task.

Alternate Uses Task (AUT). In our AUT evaluation (Tab. 3), Diverse-NS-Lite consistently beats DivPO, and Diverse-NS consistently beats Diverse-NS-Lite, though only by a small margin.

Creative Writing Task (CWT). In our CWT evaluations (Tab. 3), Diverse-NS produces the highest DSI scores for both Llama-8B and Olmo-7B.

581

582

models can serve as powerful "diversity teachers" for their larger counterparts.

⁴We provide a summary of all metrics in Table G.1



Figure 1: **Diversity and Quality Evaluation on CWT.** This figure shows $\Delta Diversity Decile (\Delta DD)$ values (y-axis) across various metrics (x-axis), computed from 70 CWT responses generated by the Olmo-2-7B model. A value of zero represents base model performance; bars indicate improvements from preference-tuned models. *D-NS* achieves the highest diversity gains overall, while *D-NS-Lite* consistently outperforms *DivPO*, except under TTR. In terms of quality (ArmoRM), *DivPO* shows a slight improvement, whereas our methods show a minor drop.

Interestingly, for Llama-8B the other methods actually reduce the ArmoRM score below baseline but Diverse-NS exceeds it. The highest 4-gram diversity is observed for Diverse-NS or -Lite in all cases. We also compute ΔDD with six lexical diversity measures and ArmoRM. Both Diverse-NS and its lightweight variant significantly outperform DivPO on every diversity metric. The ΔDD remains above the baseline for all metrics except MAAS, where it dips marginally below and similarly shows a slight under-performance for ArmoRM. Crucially, even where ΔDD suggests a minor quality drop, the absolute diversity values after self-training still exceed those of the base model (despite longer outputs), indicating that any loss in writing quality is minimal (refer to Appendix F for Llama-8B and Olmo-13B results)⁵.

7 Discussion

583

585

588

589

591

593

594

595

597

600

We introduced *Diverse-NS*, a self-learning framework to improve output diversity while preserving quality. Experiments with Llama-8B and Olmo-7B show that *Diverse-NS* improves diversity on four creative generation tasks: DAT, PGT, AUT, CWT.

606Diverse-NS is highly efficient. All gains are607achieved with only 3k preference pairs and less608than two hours of training on a single 48 GB GPU.609The lightweight variant, Diverse-NS-Lite, replaces610costly entropy and ArmoRM scoring with inexpen-611sive proxies yet still surpasses DivPO in nearly612every setting. We further show that a 7B model613can act as an effective "diversity teacher" for its

13B counterpart, pointing to a low-cost path for diversity-aware alignment at scale.

Diverse-NS maintains high quality. Diversity and quality are often at odds (Zhang et al., 2020; Chung et al., 2023), and we observe this trade-off in our experiments as well. However, there are encouraging instances where both improve together. For Olmo-7B and Olmo-13B, the ArmoRM score increases alongside diversity. Δ Diversity Decile values further confirms that, for Olmo-13B, diversity and quality consistently rise in tandem. In other cases, we observe only a minor drop in quality, suggesting that *Diverse-NS* effectively balances this trade-off in most scenarios.

The long-standing challenge of length. Evaluating diversity remains difficult due to the wellknown length bias in most diversity metrics. This issue extends to ArmoRM scores, which also favor shorter texts (Tab. 1), further complicating reliable evaluation. To mitigate this, we introduce the $\Delta Diversity Decile$ metric, which quantifies percentile gains or losses in diversity (or quality) relative to the base model. Using this length-adjusted metric, we observe substantial improvements in diversity across most lexical diversity measures, along with small but mixed changes in quality.

Overall, *Diverse-NS* offers a practical and scalable solution for boosting diversity in aligned LLMs. By addressing the length bias in both training and evaluation, our framework sets a foundation for more expressive and diverse language generation. We hope this work encourages further exploration of length-aware diversity alignment.

646

614

615

⁵We provide all results with std. dev. values in Table E.1

Limitations

647

648 While our study demonstrates the effectiveness of diversity-aware self-learning, several areas remain open for future exploration. First, our data filtering relies on a single diversity metric (e.g., entropy 651 or TTR). Although effective, no single metric can 653 fully capture all aspects of text diversity. Future work could incorporate multiple metrics to jointly 654 optimize lexical, semantic, and syntactic variation, as well as novelty, to better capture diverse training signals. Second, we focus on one data genera-657 tion task-short story writing-which allows for controlled analysis and task-specific improvements. Expanding the framework to include a broader set of tasks could lead to more generalizable diversity enhancements. Third, our self-learning setup investigates only a single round of preference tuning. While this provides a strong baseline, recent work suggests that repeated rounds of self-training can affect diversity (Guo et al., 2023; Seddik et al., 2024; Herel and Mikolov, 2024). It would be valuable to study how diversity evolves across multiple selflearning iterations in our framework. We do not include human evaluation in this study. While human judgments can provide nuanced insight, they often 671 come with variability and inconsistency. Along 672 these lines, it is often prohibitively costly to gather high-quality human feedback-particularly at the 674 scale necessary to provide stable estimates. In this 675 paper, we emphasize stringent empirical evaluation 676 of D-NS using reliable, automatic metrics and leave human-centered evaluation for future work. It is 678 worth noting a peculiar change in the length distribution of the preference-tuning model (Table D.1). Even though preference pairs are of comparable lengths in *Diverse-NS* and *Diverse-NS-Lite*, the model learns to be more expressive. We suspect 683 this shift is influenced by a skewed proportion of longer preference pairs, which may inadvertently bias the model toward generating longer responses. Controlling the length distribution is challenging 687 under our current framework due to the strict filtering criteria. In future work, we aim to address this 689 by extending our method to a multi-task setup that includes both short and long generation tasks.

Ethics Statement

693

696

Our work focuses on improving the diversity of language model outputs, particularly in creative and open-ended tasks. While diversity is an important dimension of language generation, it may come at the cost of factual correctness in certain scenarios. 697 Therefore, we caution against the use of our dataset 698 or models in tasks where factual accuracy is crit-699 ical, such as medical advice, legal reasoning, or 700 scientific fact-checking. We also acknowledge the 701 growing computational divide in language model 702 research. A key motivation behind our approach is 703 to make diversity-aware alignment more accessible. 704 By limiting training to 3,000 preference pairs and 705 demonstrating the effectiveness of smaller models 706 (e.g., Olmo-2-7B) as diversity teachers, we aim to 707 lower the resource barrier and encourage further 708 research in compute-constrained environments. Fi-709 nally, while we use proprietary language models 710 (such as GPT-40 and Claude) to assist in editing 711 and refining text during data curation and paper 712 writing, no portion of this manuscript was gener-713 ated entirely by an LLM. All content has been writ-714 ten, reviewed, and edited by the authors to ensure 715 clarity, originality, and scientific rigor. 716

References

Massih-Reza Amini, Vasilii Feofanov, Loic Pauletto, Lies Hadjadj, Emilie Devijver, and Yury Maximov. 2025. Self-training: A survey. *Neurocomputing*, 616:128904. 717

718

719

720

721

722

723

724

725

726

727

728

729

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

- Barrett R Anderson, Jash Hemant Shah, and Max Kreminski. 2024. Homogenization effects of large language models on human creative ideation. In *Proceedings of the 16th conference on creativity & cognition*, pages 413–425.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Antoine Bellemare-Pepin, François Lespinasse, Philipp Thölke, Yann Harel, Kory Mathewson, Jay A Olson, Yoshua Bengio, and Karim Jerbi. 2024. Divergent creativity in humans and large language models. *arXiv preprint arXiv:2405.13012*.
- Tuhin Chakrabarty, Philippe Laban, and Chien-Sheng Wu. 2024. Can ai writing be salvaged? mitigating idiosyncrasies and improving human-ai alignment in the writing process through edits. *arXiv preprint arXiv:2409.14509*.
- Tuhin Chakrabarty, Vishakh Padmakumar, Faeze Brahman, and Smaranda Muresan. 2023. Creativity support in the age of large language models: An empirical study involving emerging writers. *arXiv preprint arXiv:2309.12570*.

Hao Chen, Abdul Waheed, Xiang Li, Yidong Wang, Jindong Wang, Bhiksha Raj, and Marah I Abdin. 2024. On the diversity of synthetic data and its impact on training large language models. *arXiv* preprint arXiv:2410.15226.

748

749

752

758

767

768

770

772

774

775

776

777 778

790

791

797

801

- John Joon Young Chung, Ece Kamar, and Saleema Amershi. 2023. Increasing diversity while maintaining accuracy: Text data generation with large language models and human interventions. *arXiv preprint arXiv*:2306.04140.
- John Joon Young Chung, Vishakh Padmakumar, Melissa Roemmele, Yuqian Sun, and Max Kreminski. 2025. Modifying large language model posttraining for diverse creative writing. *arXiv preprint arXiv:2503.17126*.
- Geoffrey Cideron, Andrea Agostinelli, Johan Ferret, Sertan Girgin, Romuald Elie, Olivier Bachem, Sarah Perrin, and Alexandre Ramé. 2024. Diversity-rewarded cfg distillation. *arXiv preprint arXiv:2410.06084*.
- Michael A Covington and Joe D McFall. 2010. Cutting the gordian knot: The moving-average type-token ratio (mattr). *Journal of quantitative linguistics*, 17(2):94–100.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. 2023. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *Advances in neural information processing systems*, 36:10088–10115.
- Anil R Doshi and Oliver P Hauser. 2024. Generative ai enhances individual creativity but reduces the collective diversity of novel content. *Science Advances*, 10(28):eadn5290.
- Owain Evans, Andreas Stuhlmüller, and Noah Goodman. 2016. Learning the preferences of ignorant, inconsistent agents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30.
- Hugging Face. 2024. Cosmopedia: An open source mixture of experts for retrieval-augmented generation. https://huggingface.co/blog/ cosmopedia. Accessed: 2025-05-16.
- Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*.
- Carlos Gómez-Rodríguez and Paul Williams. 2023. A confederacy of models: A comprehensive evaluation of llms on creative writing. *arXiv preprint arXiv:2310.08433*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,

Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

- J. P. Guilford. 1956. The structure of intellect. *Psy-chological Bulletin*, 53(4):267–293. Place: US Publisher: American Psychological Association.
- Yanzhu Guo, Guokan Shang, Michalis Vazirgiannis, and Chloé Clavel. 2023. The curious decline of linguistic diversity: Training language models on synthetic text. *arXiv preprint arXiv:2311.09807*.
- David Herel and Tomas Mikolov. 2024. Collapse of self-trained language models. *arXiv preprint arXiv:2404.02305*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. arxiv 2021. *arXiv preprint arXiv:2106.09685*.
- Dan R Johnson, James C Kaufman, Brendan S Baker, John D Patterson, Baptiste Barbot, Adam E Green, Janet van Hell, Evan Kennedy, Grace F Sullivan, Christa L Taylor, et al. 2023. Divergent semantic integration (dsi): Extracting creativity from narratives with distributional semantic modeling. *Behavior Research Methods*, 55(7):3726–3759.
- Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. 2023. Understanding the effects of rlhf on llm generalisation and diversity. *arXiv preprint arXiv:2310.06452*.
- Jack Lanchantin, Angelica Chen, Shehzaad Dhuliawala, Ping Yu, Jason Weston, Sainbayar Sukhbaatar, and Ilia Kulikov. 2025. Diverse preference optimization. *arXiv preprint arXiv:2501.18101*.
- Dong-Hyun Lee et al. 2013. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, page 896. Atlanta.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. *arXiv preprint arXiv:1510.03055*.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023. Textbooks are all you need ii: phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*.
- Ziniu Li, Congliang Chen, Tian Xu, Zeyu Qin, Jiancong Xiao, Zhi-Quan Luo, and Ruoyu Sun. 2025. Preserving diversity in supervised fine-tuning of large language models. In *The Thirteenth International Conference on Learning Representations*.
- Ziniu Li, Congliang Chen, Tian Xu, Zeyu Qin, Jiancong Xiao, Ruoyu Sun, and Zhi-Quan Luo. 2024.

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

Entropic distribution matching for supervised finetuning of llms: Less overfitting and better diversity. In *NeurIPS 2024 Workshop on Fine-Tuning in Modern Machine Learning: Principles and Scalability.*

Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The flan collection: Designing data and methods for effective instruction tuning. In *International Conference on Machine Learning*, pages 22631–22648. PMLR.

855

858

864

865

871

872

873

874

878

887

900

901

902

903

904

905

906

907

908

909

- Yining Lu, Dixuan Wang, Tianjian Li, Dongwei Jiang, Sanjeev Khudanpur, Meng Jiang, and Daniel Khashabi. 2024. Benchmarking language model creativity: A case study on code generation. *arXiv* preprint arXiv:2407.09007.
- Heinz-Dieter Mass. 1972. Über den zusammenhang zwischen wortschatzumfang und länge eines textes. *Zeitschrift für Literaturwissenschaft und Linguistik*, 2(8):73.
- Philip M McCarthy and Scott Jarvis. 2010. Mtld, vocdd, and hd-d: A validation study of sophisticated approaches to lexical diversity assessment. *Behavior research methods*, 42(2):381–392.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2024. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*.
- Jay A. Olson, Johnny Nahas, Denis Chmoulevitch, Simon J. Cropper, and Margaret E. Webb. 2021. Naming unrelated words predicts creativity. *Proceedings* of the National Academy of Sciences of the United States of America, 118(25). Place: United States.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Vishakh Padmakumar and He He. 2023. Does writing with language models reduce content diversity? *arXiv preprint arXiv:2309.05196*.
- John D Patterson, Hannah M Merseal, Dan R Johnson, Sergio Agnoli, Matthijs Baas, Brendan S Baker, Baptiste Barbot, Mathias Benedek, Khatereh Borhani, Qunlin Chen, et al. 2023. Multilingual semantic distance: Automatic verbal creativity assessment in many languages. *Psychology of Aesthetics, Creativity, and the Arts*, 17(4):495.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*.
- Ranjani Prabhakaran, Adam E. Green, and Jeremy R. Gray. 2014. Thin slices of creativity: Using singleword utterances to assess creative cognition. *Behavior Research Methods*, 46(3):641–659. Place: Germany Publisher: Springer.

- Yiwei Qin, Yixiu Liu, and Pengfei Liu. 2025. Dive: Diversified iterative self-improvement. *arXiv preprint arXiv:2501.00747*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728– 53741.
- Nikita Salkar, Thomas Trikalinos, Byron C Wallace, and Ani Nenkova. 2022. Self-repetition in abstractive neural summarizers. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2022, page 341.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Mohamed El Amine Seddik, Suei-Wen Chen, Soufiane Hayou, Pierre Youssef, and Merouane Debbah. 2024. How bad is training on synthetic data? a statistical analysis of language model collapse. *arXiv preprint arXiv:2404.05090*.
- Chantal Shaib, Joe Barrow, Jiuding Sun, Alexa F Siu, Byron C Wallace, and Ani Nenkova. 2024a. Standardizing the measurement of text diversity: A tool and a comparative analysis of scores. *arXiv preprint arXiv:2403.00553*.
- Chantal Shaib, Yanai Elazar, Junyi Jessy Li, and Byron C Wallace. 2024b. Detection and measurement of syntactic templates in generated text. *arXiv preprint arXiv:2407.00211*.
- Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. 2023. The curse of recursion: Training on generated data makes models forget. *arXiv preprint arXiv:2305.17493*.
- Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Lei Han, Haitao Mi, and Dong Yu. 2024. Toward self-improvement of Ilms via imagination, searching, and criticizing. *Advances in Neural Information Processing Systems*, 37:52723–52748.
- Yufei Tian, Abhilasha Ravichander, Lianhui Qin, Ronan Le Bras, Raja Marjieh, Nanyun Peng, Yejin Choi, Thomas L Griffiths, and Faeze Brahman. 2023. Macgyver: Are large language models creative problem solvers? *arXiv preprint arXiv:2311.09682*.
- Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. 2024. Interpretable preferences via multi-objective reward modeling and mixture-ofexperts. *arXiv preprint arXiv:2406.12845*.
- Weijia Xu, Nebojsa Jojic, Sudha Rao, Chris Brockett, and Bill Dolan. 2024. Echoes in ai: Quantifying lack of plot diversity in llm outputs. *arXiv preprint arXiv:2501.00273*.

Hugh Zhang, Daniel Duckworth, Daphne Ippolito, and Arvind Neelakantan. 2020. Trading off diversity and quality in natural language generation. <i>arXiv preprint</i> <i>arXiv:2004.10450</i> .	System Prompt: Task Description: For this task, you will write a very short story. You will be given 3 words, and write a story that includes all 3 words. Your story should be about 5 sentences long. Use your imagination and be creative when writing your story. But, also be sure your story makes sense.
Appendix	User Prompt: Write a short story that includes these three words: [THREE_WORDS]. Assistant Prompt: [FIRST_STORY]
A Prompts	User Prompt: I do not like the previous story. Please rewrite the story in the most creative way. The new story: - must be completely different from the previous story
This section provides the exact prompts used for	in: story plot and characters
data generation, model training, and model evalua-	must have a completely different
tion. A.1 Data Generation Prompts	<pre>start (do not use standard phrases like "Once upon", "As the", "In a", "In the" etc.) must be composed of exactly [FIRST_STORY_WORD_COUNT] words. Remember to use the three words: [THREE_WORDS]</pre>
	 Hugh Zhang, Daniel Duckworth, Daphne Ippolito, and Arvind Neelakantan. 2020. Trading off diversity and quality in natural language generation. arXiv preprint arXiv:2004.10450. Appendix A Prompts This section provides the exact prompts used for data generation, model training, and model evalua- tion. A.1 Data Generation Prompts

The prompt used for generating the first response set from the model is as follows,

> System Prompt: Task Description: For this task, you will write a very short story. You will be given 3 words, and write a story that includes all 3 words. Your story should be about 5 sentences long. Use your imagination and be creative when writing your story. But, also be sure your story makes sense. User Prompt: Write a short story that includes these three words: [THREE_WORDS].

978

979

980

The prompt used for generating the second response set from the model is as follows, A.2 Model Evaluation Prompts

982

981

Divergent Association TaskThe prompt used983for the Divergent Association Task (DAT) is as984follows,985

System Prompt: Task description: Please generate 10 words that are as different from each other as possible, in all meanings and uses of the words. Rules: Only single words in English. Only nouns (e.g., things, objects, concepts). No proper nouns (e.g., no specific people or places). No specialized vocabulary (e.g., no technical terms). Think of the words on your own (e.g., do not just look at objects in your surroundings). Make a list of these 10 words, without any repetition. You must list each word with a number and a period. For example, "1. word-1, 2. word-2, etc." User Prompt: List 10 words that are as different from each other as possible:

Persona Generation Task (PGT) The prompt used for the Persona Generation Task (PGT) is as follows,

System Prompt: Generate a random description three persona with characteristics. Characteristics are: - First Name - The city of birth - Current occupation Format the output strictly using JSON schema. Use 'first_name' for First Name, 'city' for the city of birth, 'occupation' for current occupation as corresponding JSON keys. The ordering of characteristics should be arbitrary in your answer.

992

986

987

989

Alternate Uses Task (AUT). The prompt used for the Alternate Uses Task (AUT) is as follows,

System Prompt: Task Description: For this task, you'll be asked to come up with as many original and creative uses for objects as you The goal is to come up with can. creative ideas, which are ideas that strike people as clever, unusual, interesting, uncommon, humorous, innovative, or different. You must list each use with a number and a period. For example, "1. Use-1, 2. Use-2, 3. Use-3, etc.". You must provide exactly five (5) uses for each object. User Prompt: Object: [OBJECT], Uses:

The objects used for collecting the AUT responses are as follows,

"belt", "brick", "broom", "bucket", "candle", "clock", "comb", "knife", "lamp", "pencil", "pillow", "purse", "rope", "sock", "table"

Creative Writing Task (CWT). The three-word sets used in evaluating the model are as follows,

("stamp, letter, send"), ("petrol, diesel, pump"), ("statement, stealth, detect"), ("belief, faith, sing"), ("gloom, payment, exist"), ("organ, empire, comply"), ("year, week, embark"),

Pilot Analysis for Sequential Prompting B

We conducted an exploratory analysis on 20,000 short stories generated from Llama-3.1-8B and Olmo-2-7B models (Grattafiori et al., 2024; OLMo et al., 2024). The analysis was targeted at understanding the repeating patterns in the generated stories. With the help of the *diversity* package in Python (Shaib et al., 2024a), we extract the top-5 repeating Part-Of-Speech (POS) bi-grams. We find that the most repeated bigram (IN DT) occurs in over 15k stories (out of 20k) and 23% of occurrences are present at the beginning of the generated story, refer to table B.1.

Based on the findings, we conducted a sequential 1013 prompting experiment that elicits a more diverse 1014 response from the model by asking the model to 1015

993

994

995

996 997

998

1003

1001

999

- 1004 1005

- 1007
- 1008

1009

1010

1011

POS Pattern	Example String	Present (out of 20k)	Present at start (%)
IN DT DT JJ DT NN	As a, In a, In the, At the, On the a delicate, the rare, the main, the late an alley, a monarc, a spoon, a thicket	$15,782 \\ 11,418 \\ 18,472$	$23.30 \\ 16.81 \\ 16.03$
JJ NN NN IN	single silk, current king, ancient time hike in, group of, wave of, vendor to	$9,335 \\ 1,800$	$\begin{array}{c} 0.50 \\ 0.45 \end{array}$

Table B.1: **Repeating bi-grams are more likely at the beginning.** We present the frequency of repeating POS bi-grams. *IN DT* is the most frequent and commonly appears at the start of generated stories.

avoid repeating phrases (refer to appendix A.1 for exact prompts). We find that the diversity of the second response is, on average, higher than the first one.

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1029 1030

1032

1034

1036

1037 1038

1039

1040 1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

C Hyperparameters for Preference Optimization

We fine-tune the base model using the Direct Preference Optimization (DPO) objective (Rafailov et al., 2023), with $\beta = 0.1$ to control the divergence from the original policy. We use a peak learning rate of 1×10^{-5} with a cosine learning rate schedule, and a warm-up phase covering 10% of the total training steps. All models are trained using LoRA adapters (Hu et al., 2021) with a rank r = 16 and scaling factor $\alpha = 16$, on a quantized 4-bit backbone model (Dettmers et al., 2023). We add the LoRA modules to *query* and *value* projection metrics of all transformer layers in the base model with a dropout of 5%.

D Reponse Length Distribution

We observe that the length distribution varies after fine-tuning the model. As presented in table D.1, we observe that the average (and standard deviation) of response length reduces for DivPO and increases for our proposed methods (Diverse-NS and Diverse-NS-Lite). DivPO (inadvertently) teaches the model to generate shorter responses (refer to table 2). Despite maintaining comparable length for "chosen" and "rejected" samples in our methods (Diverse-NS and Diverse-NS-Lite), the model interestingly learns to generate longer responses. We suspect this shift is influenced by a skewed proportion of longer preference pairs, which may inadvertently bias the model toward generating longer responses.

E Results with Standard Deviation

1052In this section, we report the results with the stan-1053dard deviation values in Table E.1.

F \triangle **DD-based Evaluation**

Similar to the results presented fig. 1 for Olmo-7B,1055we present the results for Llama-8B and Olmo-13B1056in this section.1057

1054

G A Summary of Metrics 1058

We provide a concise summary of all metrics used1059in our evaluation setup in Table G.1.1060



Figure F.1: **Diversity and Quality Evaluation on CWT.** This figure shows $\Delta Diversity Decile (\Delta DD)$ values (y-axis) across various metrics (x-axis), computed from 70 CWT responses generated by the Llama-8B model (top-panel) and Olmo-13B (bottom panel). A value of zero represents base model performance; bars indicate improvements from preference-tuned models.

Metric	First Story	Second Story	Increase in Diversity
TTR	0.7112	0.7469	+0.0357
MAAS (↓)	0.1639	0.1609	+0.0031
HD-D	0.4143	0.4202	+0.0059
MTLD (MA-Bi)	13.9802	14.3997	+0.4195
MTLD (MA)	14.0778	14.5063	+0.4284
MTLD	14.2246	14.6652	+0.4406
MATTR	0.3810	0.3867	+0.0057

Table B.2: Sequential prompting increases diversity. We conducted a trial of sequential prompting on 20,000 responses generated from Llama-8B and Olmo-7B models. The second story generated from the models has higher diversity. \downarrow : indicates that the lower values of MAAS index represent higher diversity.

Model	Base Model	DivPO (Lanchantin et al., 2025)	Ours - D-NS-Lite	Ours - D-NS
Llama-8B Olmo-7B Olmo-13B	$\begin{array}{c} 123.27 \pm \textit{18.14} \\ 73.63 \pm \textit{15.47} \\ 86.11 \pm \textit{13.96} \end{array}$	$\begin{array}{c} 111.24 \pm \text{14.89} \\ 62.27 \pm \text{12.88} \\ 72.20 \pm \text{13.87} \end{array}$	$\begin{array}{c} 141.44 \pm {\rm 37.26} \\ 81.37 \pm {\rm 18.21} \\ 101.40 \pm {\rm 17.64} \end{array}$	$\begin{array}{c} 139.47 \pm 33.65 \\ 83.91 \pm 17.93 \\ 100.60 \pm 18.01 \end{array}$

Table D.1: **Change in the Response Length.** In this table, we present the average length of model-generated responses before and after the preference-tuning. The average values are calculated on 70 responses generated on the CWT evaluation prompts.

Task	Metric	Base Model	DivPO	D-NS-Lite	D-NS	
LLaMA-8B						
DAT	DSI	0.7535 ± 0.07	0.7545 ± 0.06	0.7590 ± 0.07	0.7640 ± 0.07	
DAT	Unique Words	0.4575	0.4593	0.4797	0.4914	
PGT	Unique First Names	0.6500	0.6100	0.6900	0.6900	
PGT	Unique Cities	0.3300	0.3100	0.4700	0.4200	
PGT	Unique Occupations	0.4100	0.3900	0.5100	0.4900	
AUT	DSI	0.8876 ± 0.02	0.8837 ± 0.02	0.8876 ± 0.02	0.8878 ± 0.02	
CWT	DSI	0.8515 ± 0.01	0.8521 ± 0.01	0.8556 ± 0.01	0.8581 ± 0.01	
CWT	ArmoRM Score	0.1451 ± 0.02	0.1495 ± 0.01	0.1369 ± 0.02	0.1405 ± 0.02	
CWT	4-gram div. POS	0.4990	0.4990	0.5030	0.5000	
CWT	4-gram div.	2.8550	2.9320	2.9450	2.9620	
CWT	Comp. Ratio.	2.635	2.546	2.568	2.530	
		OL	Mo-7B			
DAT	DSI	0.7480 ± 0.09	0.7509 ± 0.08	0.7662 ± 0.08	0.7639 ± 0.08	
DAT	Unique Words	0.6139	0.6079	0.6347	0.6327	
PGT	Unique First Names	0.3300	0.3300	0.3300	0.3400	
PGT	Unique Cities	0.3100	0.3000	0.2700	0.2700	
PGT	Unique Occupations	0.5200	0.5500	0.6100	0.6100	
AUT	DSI	0.8836 ± 0.02	0.8846 ± 0.02	0.8852 ± 0.02	0.8858 ± 0.02	
CWT	DSI	0.8499 ± 0.01	0.8491 ± 0.01	0.8548 ± 0.01	0.8563 ± 0.01	
CWT	ArmoRM Score	0.1435 ± 0.02	0.1441 ± 0.02	0.1462 ± 0.01	0.1464 ± 0.01	
CWT	4-gram div. POS	0.5720	0.5770	0.5350	0.5530	
CWT	4-gram div.	3.1270	3.1690	3.1750	3.1620	
CWT	Comp. Ratio.	2.4460	2.4160	2.3850	2.3970	
		OL	Mo-13B			
DAT	DSI	0.7233 ± 0.06	0.7282 ± 0.07	0.7320 ± 0.06	0.7364 ± 0.06	
DAT	Unique Words	0.3421	0.3340	0.3310	0.3256	
PGT	Unique First Names	0.4100	0.4100	0.4400	0.4500	
PGT	Unique Cities	0.3500	0.3500	0.3700	0.3900	
PGT	Unique Occupations	0.1900	0.1900	0.1900	0.2000	
AUT	DSI	0.8943 ± 0.02	0.8960 ± 0.02	0.8974 ± 0.02	0.8970 ± 0.02	
CWT	DSI	0.8557 ± 0.01	0.8555 ± 0.01	0.8616 ± 0.01	0.8614 ± 0.01	
CWT	ArmoRM Score	0.1571 ± 0.01	0.1589 ± 0.01	0.1585 ± 0.01	0.1590 ± 0.01	
CWT	4-gram div. POS	0.5210	0.5229	0.5080	0.4960	
CWT	4-gram div.	3.0820	3.0770	3.095	3.1070	
CWT	Comp. Ratio.	2.492	2.512	2.505	2.480	

Table E.1: **Diversity and Quality Evaluation.** We present the average (\pm std. dev.) diversity (DSI or unique values) and quality (ArmoRM score) measurements for model responses collected on four creative generation tasks (Structured Gen.: DAT, PGT, Free-Form Gen.: AUT, CWT).

Metric	Definition	Trend Description (Trend)	Application
Entropy	Entropy of the token distribu- tion in a response; measures un- predictability.	Higher values indicate greater lexical diversity ([↑]).	Training-data filtering and diversity-bias analysis
Type–Token Ratio (TTR)	Ratio of unique token types to total tokens.	Higher values indicate more lexical variety ([↑]).	Lightweight filtering (D-NS- Lite), Calculation of Diversity Decile
Moving-Average TTR (MATTR)	Moving-average of TTR over sliding windows; smooths variability.	Higher values indicate greater lexical diversity (†).	Correlation analysis, Calcula- tion of Diversity Decile
Measure of Textual Lex- ical Diversity (MTLD)	Average segment length until TTR falls below a threshold; longer segments imply more di- versity.	Higher values indicate greater lexical diversity (†).	Correlation analysis, Calcula- tion of Diversity Decile
Moving-Average MTLD (MTLD _{M})	Moving-average smoothing of MTLD to reduce variance.	Higher values indicate greater lexical diversity ([↑]).	Correlation analysis, Calcula- tion of Diversity Decile
Bidirectional Moving- Average MTLD (MTLD-MB)	MTLD-M applied forward and backward for context-sensitive smoothing.	Higher values indicate greater lexical diversity ([†]).	Correlation analysis, Calcula- tion of Diversity Decile
MAAS	Proxy metric correlated with ArmoRM quality scores.	Higher values indicate stronger quality/diversity signal ([†]).	Lightweight filtering (D-NS- Lite), Calculation of Diversity Decile
Hypergeometric Distri- bution Diversity (HD- D)	Probability-based measure of lexical diversity under a hyper- geometric model.	Higher values indicate greater lexical diversity (†).	Correlation analysis
ArmoRM score	Holistic quality score from a re- ward model.	Higher values indicate better fluency–diversity trade-off (\uparrow).	Quality evaluation (Creative Writing Task) and filtering, Cal- culation of Diversity Decile
Divergent Semantic In- tegration (DSI)	Average semantic distance among items in a generated list.	Higher values indicate greater divergent thinking (\uparrow) .	Diversity evaluation (Diver- gent Association Task, Creative Writing Task)
Diversity Decile (DD)	Decile rank of a response's diversity within its length group.	Higher decile indicates higher relative diversity after length normalization (\uparrow) .	Length-normalized evaluation (Creative Writing Task)
Change in Diversity Decile (ΔDD)	Difference in DD before and af- ter tuning; quantifies diversity gain.	Positive values indicate diversity gain; negative indicate loss (\uparrow/\downarrow) .	Measuring tuning effect on diversity (Creative Writing Task)
Semantic Distance (SD)	Average embedding-space dis- tance between outputs; indi- cates semantic variety.	Higher values indicate greater semantic variety ([†]).	Diversity evaluation (Alternate Uses Task)

Table G.1: **Overview of diversity and quality metrics:** definitions, trend descriptions with arrows, and their applications including evaluation tasks.