

A Hybrid Framework for Evaluating and Enhancing Syntactic and Semantic Diversity in Low-Resource Text Generation

Anonymous ACL submission

Abstract

Recent advancements in natural language generation (NLG) have been primarily driven by Transformer-based models and large-scale training corpora. However, in low-resource settings, the scarcity of diverse training data often results in suboptimal performance on downstream tasks. To address this challenge, we propose a hybrid synthetic data generation framework that integrates probabilistic generative models with large language models (LLMs). Specifically, probabilistic models generate diverse but structurally inconsistent text in low-resource environments, while LLMs refine these outputs to enhance fluency and coherence. To quantitatively assess the diversity of generated text, we introduce a novel set of evaluation metrics, including SynDiv for syntactic diversity and reference-free semantic measures such as Self-BERTScore, Self-BARTScore, and Self-MoverScore. Experimental results on the GoEmotions dataset demonstrate that our approach significantly improves the performance of downstream sentiment classification models, achieving up to a 13-point increase in F1 score. Furthermore, we establish a strong correlation between our evaluation metrics and downstream classification performance, providing a solid foundation for future research in synthetic data generation.

1 Introduction

Recent years have witnessed remarkable advancements in Natural Language Generation (NLG) research, primarily driven by the Transformer architecture and the availability of massive training datasets. While deep neural networks typically require substantial training data, existing approaches face significant challenges in data-scarce scenarios. The performance of generative models like GPT-3(Brown et al., 2020) and T5(Ni et al., 2022) critically depends on the diversity of training data, with syntactic structural diversity and semantic diversity

emerging as two pivotal factors. As demonstrated by Devlin (2018), the BooksCorpus and Wikipedia corpora employed in BERT training contain rich syntactic variations, mixing academic-style long sentences with colloquial short phrases, which enable models to generate texts with greater grammatical complexity. Complementing this finding, Brown et al. (2020) revealed that GPT-3's ability to produce coherent long-form texts across diverse domains is derived from its training on massive multidomain datasets spanning scientific literature, technical documentation, and creative writing.

However, acquiring high-quality datasets with both syntactic and semantic diversity often requires costly human annotation or domain expertise, posing a significant challenge in low-resource linguistic tasks such as fine-grained sentiment classification. To mitigate this issue, synthetic data generation has emerged as a promising solution. Existing non-LLM-based approaches mainly fall into two categories. The first includes data augmentation techniques such as back-translation (Sennrich et al., 2016), rule-based text expansion (Wei and Zou, 2019), and semantic substitution (Kobayashi, 2018). The second category consists of statistical models and small-scale neural networks, including word vector interpolation (Mikolov et al., 2013) and Variational Autoencoders (VAE), which are commonly used for text augmentation in low-resource settings. These methods aim to expand limited training corpora while preserving essential linguistic and domain-specific characteristics.

Meanwhile, the application of LLMs in synthetic data generation has been receiving increasing attention. Research has shown that LLMs can generate high-quality data in low-resource environments through techniques such as Prompt Engineering and Few-shot Learning, for tasks like sentiment analysis, question answering, and machine translation data augmentation. However, (Liang et al., 2024) have indicated that when directly us-

Table 1: Example sentences generated by different models before and after refinement

Model	Generated Sentence	LLM Refined Sentence
VAE	hate you it [!! what oh my [..?!!! like ..	I hate you—what?! Oh my... like... seriously?!
Diffusion-LM	barely scraping rent storm ping exhaustion grinding weight slipping restless nights shitt ghost sleep distant cracks piling through taunting nothing times.	Barely scraping through, grinding against the weight of restless nights—pinging from a shitstorm of exhaustion, rent piling up, time slipping through cracks, and sleep nothing more than a distant, taunting ghost.
SeqGAN	i didn see this for this. know i think that's just a reason i 'truth.	I didn't just perceive this—I know it. It's not just a reason; it's my truth.
WGAN	jaipur relentes river struggling fated mir ancestry smyth smyth surge honeymoonhlyhly villains Erikahesis spike Huang disrupted reckoning fiery commemoration euclidean reviewing reviewingcation conquer traditions.	In Jaipur, a relentless river roared with the struggles of fate, while Smyth's ancestry surged forth, a fierce and unshakable echo during the fiery reckoning of ancient traditions.
TextGAN	flauont captain holden montevideo boldy intent rebellious austro lobbyion certainty footprint flexed mischievous mark stamping callum insurrection Johann montevideo influence austro gilbery.	Captain Callum Holden boldly flaunted his rebellious intent in the Montevideo lobby, stamping an unmistakable mark of Austro influence.
ControlGAN	toulon specifications 1810 stables intermediate gentlemen lightning speaking alive stinging crosby verified containment cantata sahibtension podium nuclear outlet flipped armagh hatfield poorly clean universe mendozaatling notes eve innfide.	In 1810, Toulon's specifications struck a bold balance—stables and gentlemen stood affirmed by Crosby, while a tension podium and a flipped nuclear outlet sent ripples through a chaotic universe stretching from Hatfield to Armagh.

ing LLMs to generate sentiment text, the influence of their pre-training data distribution can lead to a tendency towards homogenized expression patterns. For example, LLMs often over-rely on high-frequency phrases like “very satisfied” or “strongly recommend,” resulting in generated text with significantly lower syntactic complexity and semantic diversity compared to human-written content.

The remarkable success of probabilistic generative models (e.g., GANs, VAEs, Diffusion models) in image data augmentation has inspired analogous applications in textual data generation. However, these probability-based text generation methods often require large amounts of training data to achieve optimal results. Under low-resource conditions, the generated text tends to exhibit fragmented semantics, where local word sequences may adhere to grammatical rules but the entire sen-

tence does not follow correct syntax, leading to semantic confusion. While the general meaning can be inferred, it often fails to meet the training requirements for downstream task models. Table 1 presents example sentences produced by each approach, specifically showcasing their ability to generate text with heightened emotional intensity, emphasizing expressive and dynamic linguistic styles. Meanwhile, traditional evaluation metrics often focus on assessing the similarity between model-generated sentences and reference sentences (e.g., BLEU(Papineni et al., 2002), BERTScore(Zhang et al., 2019)). However, This criterion is not consistent with the core goal of synthetic data generation - to maximize compositional diversity. In the field of synthetic data, the emphasis is on diversity in generation, specifically the differences in syntactic structures and semantics between the generated

084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101

102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119

120 sentences.

121 In response to the aforementioned issues, this
122 study makes the following contributions:

- 123 1. This study proposes a hybrid framework inte-
124 grating probabilistic generative models with
125 large language models (LLMs) for synthetic
126 data generation. Experimental results demon-
127 strate that the synthesized data significantly
128 enhances downstream classification perfor-
129 mance, achieving a maximum F1-score im-
130 provement of 13 percentage points compared
131 to baseline methods.
- 132 2. A new set of evaluation metrics is proposed
133 to measure the syntactic and semantic diver-
134 sity of generated text in the field of synthetic
135 data, providing a more comprehensive assess-
136 ment of the effectiveness of data augmentation
137 methods.

138 2 Related Work

139 **Synthetic Data Generation for NLP** The genera-
140 tion of synthetic data has been widely explored
141 as a strategy to alleviate data scarcity issues in
142 NLP. Traditional data augmentation techniques,
143 including back-translation(Sennrich et al., 2016),
144 synonym replacement (Wei and Zou, 2019), and
145 paraphrasing (Kobayashi, 2018), have been exten-
146 sively utilized to artificially expand training cor-
147 pora while preserving linguistic diversity. These
148 methods, however, often fail to maintain deep syn-
149 tactic and semantic structures, particularly in tasks
150 requiring fine-grained sentiment classification or
151 domain-specific text generation.

152 **Challenges in Evaluating Synthetic Text Qual-
153 ity** Evaluating the quality of synthetic text re-
154 mains a major challenge, as traditional NLP met-
155 rics focus on similarity-based assessments rather
156 than intrinsic diversity. BLEU (Papineni et al.,
157 2002), ROUGE 2.0 (Ganesan, 2018), and ME-
158 TEOR (Banerjee and Lavie, 2005) primarily mea-
159 sure lexical and n-gram overlap between generated
160 and reference sentences. While contextual embed-
161 dings such as BERTScore (Zhang et al., 2019) and
162 MoverScore (Zhao et al., 2019) improve semantic
163 evaluation, they do not explicitly capture structural
164 and compositional diversity.

165 3 Methodology

166 Our framework evaluates synthetic text quality
167 through three complementary dimensions: syntac-
168 tic structure diversity, semantic coherence model-

169 ing, and generative model comparison. It intro-
170 duces **SynDiv**, a novel reference-free metric for
171 quantifying syntactic structure diversity in gen-
172 erated text. Unlike traditional methods that rely
173 on reference corpora, SynDiv encodes dependency
174 parse trees into spectral representations and com-
175 putes structural variations directly within the gen-
176 erated set. By leveraging Fine-Grained Syntax
177 Descriptor (FGSD) and graph Laplacian spectral
178 decomposition, SynDiv provides a self-contained,
179 fine-grained assessment of syntactic diversity.

180 3.1 Proposed Metric for Assessing Syntactic 181 Structural Diversity

182 Let $W = \{w_1, w_2, \dots, w_N\}$ denote a set of N gen-
183 erated sentences. For each sentence w_i , we extract
184 its dependency parse tree using the Stanford Parser
185 (Socher et al., 2013), resulting in a set of trees
186 $S = \{s_1, s_2, \dots, s_N\}$ where s_i corresponds to the
187 parse tree of sentence w_i . To quantify syntactic
188 diversity, we propose the FGSD method that en-
189 codes each dependency tree as a vector through
190 the following process: Given a dependency tree
191 G viewed as an acyclic graph, we decompose it
192 into subgraphs $\{G_1, G_2, \dots, G_m\}$ by removing cer-
193 tain nodes or edges. For each subgraph G_i , first
194 compute its Laplacian matrix:

$$L_i = D_i - A_i \quad (1)$$

195 where A_i is the adjacency matrix of G_i and D_i
196 is the diagonal degree matrix. Then we perform
197 spectral decomposition:

$$L_i = U_i \Lambda_i U_i^\top \quad (2)$$

198 where Λ_i is the diagonal matrix containing the
199 eigenvalues, and U_i is the matrix of correspond-
200 ing eigenvectors. Extract statistical features from
201 the top k eigenvectors:

$$F_i = [\mu(u_i^{(1)}), \mathbb{V}(u_i^{(1)}), \dots, \mu(u_i^{(k)}), \mathbb{V}(u_i^{(k)})] \quad (3)$$

202 where $\mu(\cdot)$ and $\mathbb{V}(\cdot)$ denote the mean and variance
203 of each eigenvector, respectively. Concatenate and
204 normalize the spectral information:

$$h_{\text{final}} = \bigoplus_{i=1}^m [T_i \oplus F_i] \quad (4)$$

205 where T_i is the set of eigenvectors corresponding
206 to the largest k eigenvalues. Let $\|\cdot\|$ denotes the
207 Euclidean norm

$$h = \frac{h_{\text{final}}}{\|h_{\text{final}}\|^2} \quad (5)$$

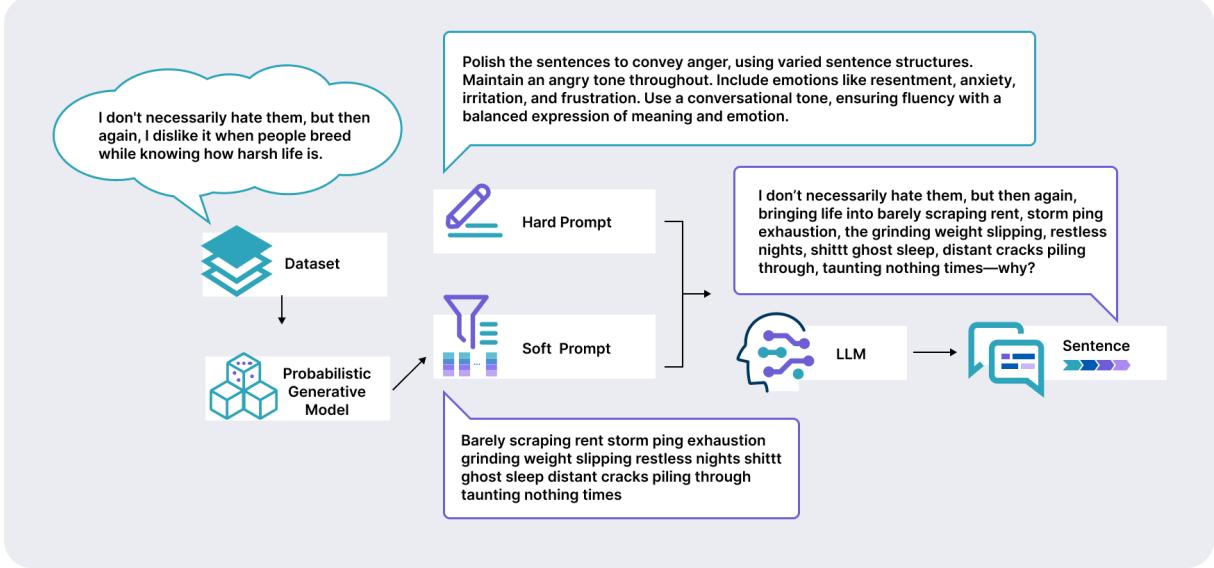


Figure 1: Enhancing diversity of LLM-generated outputs by regulating basic probabilistic generative model’s results

Using this FGSD method, we encode the set of dependency trees S into a set of vectors $H = \{h_1, h_2, \dots, h_N\}$. The syntactic similarity between sentences is computed using cosine similarity:

$$A_{ij} = \frac{h_i^\top h_j}{|h_i||h_j|} \quad (6)$$

Finally, we define the Syntactic Structure Diversity (SynDiv) metric for the sentence set W as:

$$\text{SynDiv}(W) = \frac{1}{N(N-1)} \sum_i^N \sum_{j \neq i}^N A_{ij} \quad (7)$$

3.2 Proposed Metric for Evaluating Semantic Diversity

To complement **SynDiv**, which rigorously quantifies syntactic structural diversity, we introduce a suite of semantic diversity metrics to measure the variability of generated sentences in terms of meaning and coherence. Traditional diversity assessments often rely on external reference texts, limiting their applicability to fully unsupervised generative evaluation. We address this limitation by extending reference-based metrics into self-supervised diversity measurements.

Still let $W = \{w_1, w_2, \dots, w_N\}$ represent a set of N generated sentences. We propose a novel reference-free variant of BERTScore (Zhang et al., 2019), termed **Self-BERTScore**, which eliminates dependency on reference texts by leveraging pairwise similarity within W . Formally, Self-BERTScore is computed as:

$$\frac{1}{N(N-1)} \sum_i^N \sum_{j \neq i}^N \text{BERTScore}(w_i, w_j) \quad (8)$$

Similarly, we extend BARTScore (Yuan et al., 2021) to derive **Self-BARTScore**, which measures relative semantic coherence within generated text:

$$\frac{1}{N(N-1)} \sum_i^N \sum_{j \neq i}^N \text{BARTScore}(w_i, w_j) \quad (9)$$

Further, we develop **Self-MoverScore**, inspired by MoverScore (Zhao et al., 2019), introducing a reference-free formulation based on pairwise averaging:

$$\frac{1}{N(N-1)} \sum_i^N \sum_{j \neq i}^N \text{MoverScore}(w_i, w_j) \quad (10)$$

3.3 Probabilistic Generative Framework

Probabilistic generative models aim to approximate the underlying probability distribution of observed data, facilitating both inference and novel sample generation. Formally, given real-world observations sampled from an unknown distribution $p_{\text{data}}(x)$, these models seek to learn a parameterized approximation $p_\theta(x)$ by minimizing the statistical divergence between the two distributions. This approximation is achieved through either explicit density estimation, as in VAE and Diffusion models, or implicit modeling via adversarial learning,

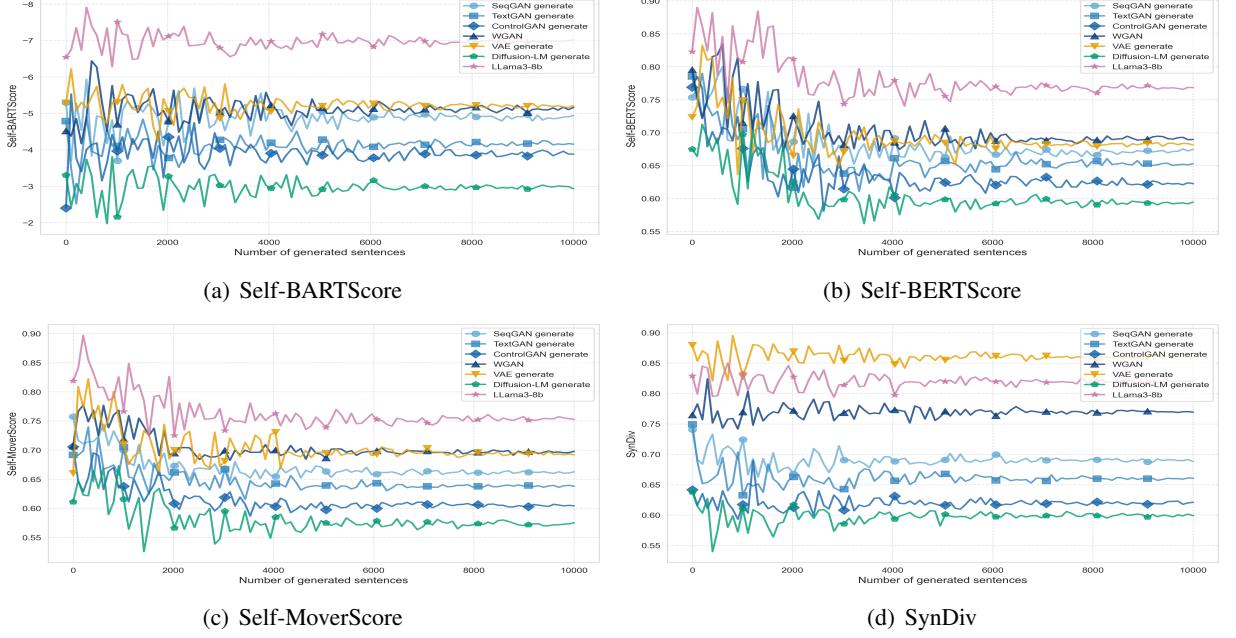


Figure 2: Convergence of self-Evaluation scores across different generative models

as in GANs. While these models exhibit strong distributional learning capabilities and support diverse text generation, their effectiveness is constrained by the discrete nature of text. Unlike continuous data, textual sequences require the model to select tokens sequentially from a predefined vocabulary while maintaining syntactic and semantic coherence. Under low-resource conditions, the limited training signal exacerbates these challenges, often resulting in grammatically flawed and semantically inconsistent outputs.

To address these limitations, we propose a hybrid framework integrating a probabilistic generative model with a LLM. As depicted in Figure 1, the probabilistic generative model is first trained on a low-resource dataset to learn its latent structure. Despite its inability to fully capture complex linguistic constraints, it provides an approximation of the target distribution, capturing phrase co-occurrence patterns and coarse syntactic structures. However, due to insufficient data, the generated sequences frequently exhibit structural deficiencies and semantic ambiguities. Rather than discarding these suboptimal outputs, we leverage LLaMA-3-8B (AI@Meta, 2024) as a refinement module. Specifically, the initial generations serve as input prompts to LLaMA-3-8B, which then reconstructs grammatically sound and semantically coherent text while preserving structural variation. This two-stage generation process ensures that the probabilistic model contributes lexical and syntactic diversity, while the LLM guarantees linguistic fidelity. An illustrative comparison of raw versus refined outputs

is presented in Table 1.

3.4 Selected Generative Models

To systematically evaluate our proposed method, we experiment with six representative generative models, each employing bert-base-uncased as the backbone to maintain comparability. These models span various generative paradigms, encompassing adversarial training, variational inference, and diffusion-based approaches:

- **VAE** (Kingma and Welling, 2014): A variational autoencoder-based model that learns latent variable representations by maximizing the evidence lower bound.
- **WGAN** (Arjovsky et al., 2017): A Wasserstein GAN designed to address mode collapse and training instability by minimizing the Wasserstein distance between real and generated data distributions.
- **SeqGAN** (Yu et al., 2017): A policy gradient-based generative model optimized via Monte Carlo Tree Search.
- **ControlGAN** (Lee and Seok, 2019): A generative adversarial model incorporating attribute control through hierarchical attention mechanisms.
- **TextGAN** (Zhang et al., 2017): A GAN-based text generator utilizing Wasserstein distance with gradient penalty.
- **Diffusion-LM** (Li et al., 2022): A diffusion-based generative model employing learnable word embeddings and iterative denoising for high-fidelity text synthesis.

262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294

295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326

Table 2: Performance comparison across generative architectures

Data	PPL ↓	Self-BERT ↓	Self-BART ↑	Self-Mover ↓	SynDiv ↓
VAE*	189	0.682 ± 0.003	-5.21 ± 0.07	0.694 ± 0.01	0.86 ± 0.02
Diffusion-LM*	351	0.593 ± 0.005	-2.97 ± 0.05	0.573 ± 0.01	0.60 ± 0.02
WGAN*	406	0.691 ± 0.005	-5.12 ± 0.05	0.696 ± 0.01	0.77 ± 0.02
SeqGAN*	343	0.671 ± 0.003	-4.91 ± 0.05	0.662 ± 0.02	0.69 ± 0.02
TextGAN*	382	0.651 ± 0.004	-4.15 ± 0.09	0.638 ± 0.02	0.66 ± 0.02
ControlGAN*	326	0.623 ± 0.003	-3.89 ± 0.08	0.605 ± 0.01	0.62 ± 0.02
LLaMA-3-8B	96	0.768 ± 0.005	-6.97 ± 0.05	0.753 ± 0.01	0.82 ± 0.02

Table 3: Performance comparison of classification models with training sets augmented by synthetic data

Training Data	Precision ↑	recall ↑	F1 score ↑
GoEmotions + VAE*	0.70	0.76	0.73
GoEmotions + Diffusion-LM*	0.81	0.83	0.82
GoEmotions + SeqGAN*	0.73	0.77	0.75
GoEmotions + WGAN*	0.69	0.76	0.72
GoEmotions + TextGAN*	0.76	0.79	0.77
GoEmotions + ControlGAN*	0.79	0.80	0.79
GoEmotions + LLaMA-3-8B	0.71	0.77	0.74
GoEmotions	0.65	0.74	0.69

* The listed models correspond to the sentence datasets generated by these respective generative models and refined.

By integrating these diverse generative architectures into our framework, we systematically evaluate their effectiveness in synthetic text augmentation. The combination of probabilistic modeling and LLM-based refinement ensures that synthetic data exhibits both structural diversity and linguistic fluency.

4 Experiments

4.1 Dataset and Experimental Configuration

We conduct experiments on the GoEmotions dataset (Demszky et al., 2020), a multi-label emotion classification corpus with 58,009 Reddit comments annotated across 28 categories. All models are trained on four NVIDIA A100 GPUs using PyTorch’s Distributed Data Parallel (DDP) framework. Training consists of 50,000 steps. For text refinement framework, generated outputs undergo semantic enhancement via LLaMA-3-8B with top- k sampling ($k = 40$), temperature $\tau = 0.7$, and repetition penalty $\gamma = 1.2$.

4.2 Baselines

To establish a robust comparative foundation, we employ LLaMA3-8B direct generation as the baseline for our synthetic data generation framework. This baseline is systematically evaluated through two distinct approaches: 1) rigorous assessment of the generated synthetic data using our comprehensive diversity evaluation metrics, and 2) quantita-

tive measurement of their effectiveness in enhancing classification task performance metrics. This dual evaluation strategy ensures both the qualitative diversity and practical utility of the synthesized data are thoroughly examined.

For probabilistic generation, we employ six widely studied models: Diffusion-LM, VAE, SeqGAN, TextGAN, ControlGAN, and WGAN. Each model is trained on the GoEmotions dataset under identical conditions, ensuring consistency across architectures. Following training, each model generates 10,000 synthetic samples, preserving the original label distribution. However, due to the inherent limitations of these models—particularly in low-resource scenarios—their raw outputs often exhibit syntactic inconsistencies and semantic ambiguities. To enhance fluency and coherence, we apply post-processing using LLaMA-3-8B, refining generated samples via controlled rewriting.

In parallel, we evaluate direct generation using LLaMA-3-8B, leveraging carefully designed prompts to guide the model toward producing emotionally grounded sentences. To ensure a fair comparison, we design structured prompts that guide LLaMA-3-8B in generating sentiment-aligned text, reflecting the categories in GoEmotions. Unlike probabilistic models, which learn to generate text from observed distributions, this approach relies on controlled prompt engineering to elicit responses corresponding to specific emotional tendencies.

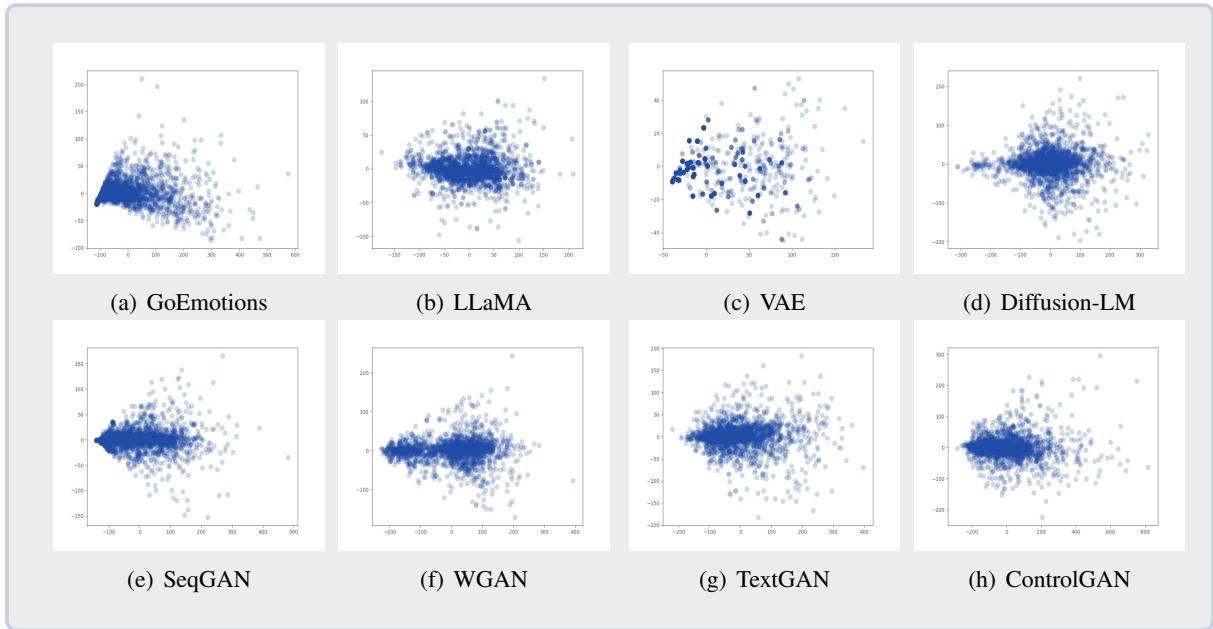


Figure 3: Comparative visualization of syntactic distributions

For instance, to generate joyful sentences, we use prompts such as "Generate a cheerful sentence expressing excitement and happiness," while for sadness, we employ "Write a full sentence conveying deep sorrow and emotional pain." These prompts are designed to provide LLaMA-3-8B with sufficient context to produce text that aligns with the intended emotional category, allowing us to systematically evaluate the strengths and limitations of both approaches—probabilistic models with refinement and direct LLaMA-based sentiment generation—providing insights into their respective impacts on text diversity, fluency, and semantic alignment in synthetic data augmentation.

4.3 Experiment Result

Our evaluation framework assesses the quality of generated synthetic data through five complementary dimensions: Perplexity (PPL), Self-BERTScore, Self-BARTScore, Self-MoverScore, and Syntax-Topology Divergence. These metrics comprehensively evaluate the fluency, syntactic structure, and semantic diversity of the generated texts, as outlined in section 3.1 and 3.2. Specifically, fluency is measured using Perplexity, while syntactic and semantic diversity are quantified using structural and embedding-based comparisons.

The experimental results in Table 2 show that Diffusion-LM+LLaMA-3-8B outperforms other models in both syntactic and semantic diversity, achieving the highest Self-BERTScore and Self-BARTScore. Additionally, LLaMA-3-8b exhibits the lowest PPL score, suggesting superior fluency.

These findings suggest that while probabilistic models introduce greater lexical and syntactic diversity, LLaMA-3-8b ensures fluency, balancing structure and coherence in synthetic data generation. The experimental results in Table 2 show that Diffusion-LM+LLaMA-3-8B outperforms other models in both syntactic and semantic diversity, achieving the highest Self-BERTScore and Self-BARTScore. Additionally, LLaMA-3-8b exhibits the lowest PPL score, suggesting superior fluency. These findings suggest that while probabilistic models introduce greater lexical and syntactic diversity, LLaMA-3-8b ensures fluency, balancing structure and coherence in synthetic data generation.

As shown in Figure 2, all proposed diversity metrics exhibit significant fluctuations during the initial stages of data generation. However, as the dataset size increases, these metrics converge to stable values. The distinct differences in their final convergence points further demonstrate the robustness of our diversity metrics, validating their applicability for assessing the performance of generative models.

To further evaluate the efficacy of synthetic data augmentation, we incorporate 10,000 generated sentences from each model into the GoEmotions training set and compare classification performance against models trained solely on the original dataset. Table 3 demonstrates that models trained with synthetic data consistently outperform their original-data-only counterparts, with the greatest improvements observed for Diffusion-LM+LLaMA-3-8B-generated data, which achieves

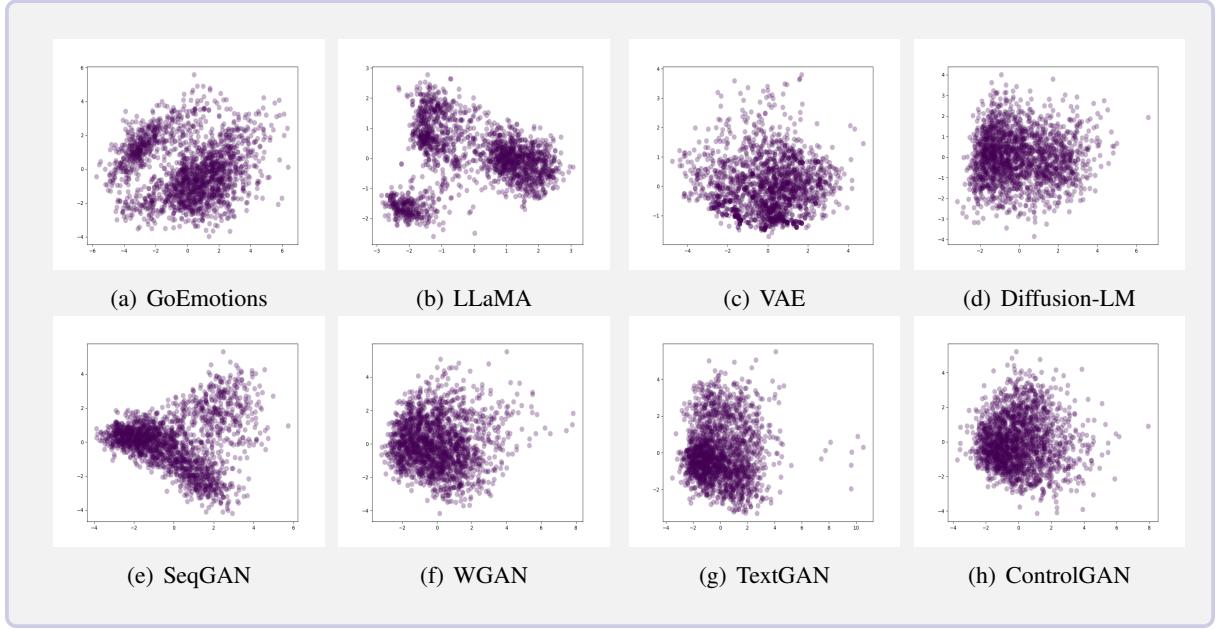


Figure 4: Comparative visualization of semantic distributions

Table 4: Correlation analysis

Evaluation Metrics	Pearson	Spearman
Self-BERT	-0.788	-0.893
Self-BART	0.791	0.857
Self-Mover	-0.880	-0.893
SynDiv	-0.869	-0.893

the highest F1, Recall, and Precision scores. This trend aligns with the diversity evaluations, indicating that greater syntactic and semantic diversity enhances classification model generalization. To quantify this relationship, we perform a correlation analysis between diversity metrics and classification performance improvements, as shown in Table 4. The results confirm that higher structural and semantic variation in synthetic data correlates with greater classification gains.

Through comparative analysis of distribution patterns between synthetic data and the original GoEmotions dataset (as shown in Fig.3 and Fig.4), we demonstrate that the Diffusion-LM+LLaMA-3-8B framework exhibits approximately Gaussian distributions in both semantic feature space and syntactic feature space. This distribution characteristic significantly mitigates inherent bias issues in training data by generating balanced synthetic samples, thereby effectively enhancing model generalizability – which constitutes the core mechanism for its substantial improvement in classification performance. Experimental results reveal that the integration of synthetic data not only achieves systematic enhancement of classification accuracy, but also remarkably strengthens model robustness through

dimensional expansion of affective representation space.

Notably, the VAE+LLaMA-3-8B architecture manifests distinct distribution properties (Fig. 3): Its syntactic features display sparse distribution patterns, while high-density clusters suggest structural homogeneity in latent space. This phenomenon aligns precisely with quantitative diversity metrics analysis, revealing inherent limitations in generative diversity. Our empirical study further indicates that such structural homogeneity may induce semantic redundancy in generated samples, thereby constraining the model’s expressive capacity.

5 Conclusion

This study presents a hybrid framework for synthetic text generation by integrating probabilistic generative models with LLMs, enhancing syntactic diversity and linguistic fluency in low-resource settings. Our approach, leveraging the structural variability of probabilistic models and the coherence of LLMs, addresses key challenges in synthetic text generation, improving downstream classification tasks and model generalization, as demonstrated on the GoEmotions dataset. The proposed diversity-aware metrics highlight the importance of compositional variation over traditional similarity-based measures, offering a promising direction for synthetic data generation. Future work will focus on computational efficiency, automated prompt selection, and extending the framework to other NLP tasks.

507 Acknowledgments

508 We sincerely thank the anonymous reviewers for
509 their valuable feedback and insightful suggestions,
510 which must improve the clarity and quality of this
511 work. We also acknowledge the support of our
512 colleagues and collaborators who contributed to
513 discussions that shaped this research.

514 Limitations

515 Dependence on Large Language Model (LLM)

516 **Refinement:** Our proposed framework relies on
517 the refinement capabilities of LLMs, specifically
518 LLaMA-3-8B, to enhance the grammatical coherence
519 and semantic fluency of text generated by
520 probabilistic models. While this hybrid approach
521 successfully mitigates issues of fragmentation and
522 incoherence in low-resource text generation, it in-
523 herently depends on the quality of the LLM used
524 for refinement. Different LLM architectures may
525 yield varied results, and our findings may not di-
526 rectly generalize to other large-scale generative
527 models, particularly those with distinct training dis-
528 tributions or decoding strategies. Future research
529 should explore how alternative LLMs impact re-
530 finement effectiveness across different generative
531 paradigms.

532 Model-Specific Generalization Constraints:

533 Our study investigates a range of generative mod-
534 els, including VAEs, GANs, and Diffusion-based
535 models, all of which exhibit unique strengths and
536 weaknesses in text generation. However, our find-
537 ings are constrained to the specific architectures
538 and hyperparameters used in our experiments. The
539 impact of alternative training objectives, regulariza-
540 tion techniques, and domain-specific adaptations
541 remains unexplored. Additionally, while we inte-
542 grate probabilistic models with LLMs to address
543 fluency and structural consistency, the reliance on a
544 predefined set of generative approaches limits our
545 framework's adaptability to emerging text genera-
546 tion techniques. Future research should evaluate
547 the effectiveness of our approach across a broader
548 class of generative models.

549 **Language-Specific Considerations:** Our ex-
550 periments are conducted primarily on English-
551 language datasets, which possess relatively straight-
552 forward syntactic structures compared to mor-
553 phologically richer languages. The proposed
554 syntactic diversity metric (SynDiv) and self-
555 supervised semantic measures (Self-BERTScore,
556 Self-BARTScore, and Self-MoverScore) are eval-

uated within this linguistic context, and their ef-
557 ffectiveness in languages with complex word in-
558 flexions, free word order, or high agglutination
559 remains uncertain. Further studies should validate
560 the applicability of our framework in multilingual
561 settings, particularly for low-resource languages
562 where syntactic and semantic diversity challenges
563 differ significantly.

564 Reliance on Graph-Based Syntactic Diversity

565 **Metrics:** Our methodology introduces a novel syn-
566 tactic diversity metric, SynDiv, which leverages
567 graph Laplacian spectral decomposition for syntac-
568 tic structure analysis. While this approach provides
569 a fine-grained assessment of syntactic diversity, it
570 assumes that dependency parsing yields structurally
571 accurate representations of text. In cases where
572 dependency parsers struggle with ambiguous or
573 out-of-domain text, the computed diversity scores
574 may not fully capture syntactic variations. More-
575 over, graph-based syntactic analysis can be compu-
576 tationally expensive, potentially limiting scalabil-
577 ity in large-scale text generation settings. Future
578 work should investigate alternative parsing strate-
579 gies and efficiency optimizations for large-scale
580 applications.

581 Evaluation Dependency on Self-Supervised

582 **Metrics:** Our assessment framework relies on
583 reference-free, self-supervised diversity metrics to
584 quantify syntactic and semantic variation in gener-
585 ated text. While this approach circumvents the need
586 for gold-standard reference texts, it also introduces
587 dependencies on the underlying embedding models
588 used for similarity measurement. Self-BERTScore,
589 Self-BARTScore, and Self-MoverScore inherently
590 reflect the biases of the pretrained language models
591 they leverage, which may affect their robustness
592 across different domains. Future research should
593 explore additional metrics or human evaluation
594 strategies to complement automated assessments.

595 Ethical Considerations in Synthetic Text Gen- 596 eration:

597 **Evaluation:** While our hybrid framework enhances
598 the diversity of generated text and improves syn-
599 thetic data augmentation for downstream tasks, it
600 also raises ethical concerns regarding the poten-
601 tial misuse of synthetic data. Increased diversity
602 in generated text may lead to unintended biases,
603 especially when applied to sentiment-related tasks.
604 Additionally, the refinement step involving LLMs
605 introduces a level of opacity, as models may un-
606 knowingly inject biases present in their pretraining
607 data. Ensuring that synthetic text aligns with ethi-
608 cal AI principles and does not propagate misinfor-

609
610
611
612
mentation remains an open challenge. We advocate
for ongoing monitoring of synthetic data quality
and the development of fairness-aware evaluation
methods to mitigate potential biases.

613 References

614 AI@Meta. 2024. [Llama 3 model card](#).

615 Martin Arjovsky, Soumith Chintala, and Léon Bottou.
616 2017. Wasserstein generative adversarial networks.
617 In *International conference on machine learning*,
618 pages 214–223. PMLR.

619 Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An
620 automatic metric for mt evaluation with improved cor-
621 relation with human judgments. In *Proceedings of*
622 *the acl workshop on intrinsic and extrinsic evalua-*
623 *tion measures for machine translation and/or summariza-*
624 *tion*, pages 65–72.

625 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
626 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
627 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
628 Askell, et al. 2020. Language models are few-shot
629 learners. *Advances in neural information processing*
630 *systems*, 33:1877–1901.

631 Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo
632 Ko, Alan Cowen, Gaurav Nemadé, and Sujith Ravi.
633 2020. [GoEmotions: A dataset of fine-grained emotions](#). In *Proceedings of the 58th Annual Meeting of*
634 *the Association for Computational Linguistics*, pages
635 4040–4054, Online. Association for Computational
636 Linguistics.

637 Jacob Devlin. 2018. Bert: Pre-training of deep bidi-
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S
Liang, and Tatsunori B Hashimoto. 2022. Diffusion-
lm improves controllable text generation. *Advances*
in Neural Information Processing Systems, 35:4328–
4343.

Xun Liang, Hanyu Wang, Yezhaohui Wang, Shichao
Song, Jiawei Yang, Simin Niu, Jie Hu, Dan Liu,
Shunyu Yao, Feiyu Xiong, and Zhiyu Li. 2024. [Con-](#)
[trollable text generation for large language models:](#)
[A survey](#). *CoRR*, abs/2408.12599.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Cor-
rado, and Jeff Dean. 2013. [Distributed representa-](#)
[tions of words and phrases and their compositionality](#).
In *Advances in Neural Information Processing Sys-*
tems, volume 26. Curran Associates, Inc.

Jianmo Ni, Gustavo Hernandez Abrego, Noah Constant,
Ji Ma, Keith Hall, Daniel Cer, and Yinfei Yang. 2022.
Sentence-t5: Scalable sentence encoders from pre-
trained text-to-text models. In *Findings of the Asso-*
ciation for Computational Linguistics: ACL 2022,
pages 1864–1874.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-
Jing Zhu. 2002. Bleu: a method for automatic evalua-
tion of machine translation. In *Proceedings of the*
40th annual meeting of the Association for Compu-
tational Linguistics, pages 311–318.

Rico Sennrich, Barry Haddow, and Alexandra Birch.
2016. [Improving neural machine translation models](#)
[with monolingual data](#). In *Proceedings of the 54th*
Annual Meeting of the Association for Computational
Linguistics (Volume 1: Long Papers), pages 86–96,
Berlin, Germany. Association for Computational Lin-
guistics.

Richard Socher, John Bauer, Christopher D Manning,
and Andrew Y Ng. 2013. Parsing with compositional
vector grammars. *Proceedings of ACL 2013*.

Jason Wei and Kai Zou. 2019. [EDA: Easy data augmen-](#)
[tation techniques for boosting performance on text](#)
[classification tasks](#). In *Proceedings of the 2019 Con-*
ference on Empirical Methods in Natural Language
Processing and the 9th International Joint Confer-
ence on Natural Language Processing (EMNLP-
IJCNLP), pages 6382–6388, Hong Kong, China. As-
sociation for Computational Linguistics.

Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu.
2017. Seqgan: Sequence generative adversarial nets
with policy gradient. In *Proceedings of the AAAI*
conference on artificial intelligence, volume 31.

Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021.
Bartscore: Evaluating generated text as text genera-
tion. *Advances in Neural Information Processing*
Systems, 34:27263–27277.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Wein-
berger, and Yoav Artzi. 2019. Bertscore: Evaluating
text generation with bert. In *International Confer-*
ence on Learning Representations.

716	Yizhe Zhang, Zhe Gan, Kai Fan, Zhi Chen, Ricardo Henao, Dinghan Shen, and Lawrence Carin. 2017. Adversarial feature matching for text generation. In <i>International conference on machine learning</i> , pages 4006–4015. PMLR.	731
717		
718		
719		
720		
721	Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. <i>MoverScore: Text generation evaluating with contextualized embeddings and earth mover distance</i> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 563–578, Hong Kong, China. Association for Computational Linguistics.	734
722		
723		
724		
725		
726		
727		
728		
729		
730		
	A Appendix	
	The following sentences are generated and refined by the model based on certain emotion prompt.	732
		733
	A.1 VAE example (20)	734
	Thank you so much! This is so cool and thoughtful!	735
	I love this so much—thank you endlessly for it!	736
	No way! This is so cool!	737
	Thank you for this—it’s so good!	738
	I love this! Thank you so much—it’s incredible!	739
	Is this for me? Thank you so much for everything!	740
	Thank you! The name is perfect—love it so much!	741
	This is so cool—thank you for the name!	742
	What’s this? It’s so thoughtful—thank you endlessly!	743
	Is this really mine? Thank you so much for this amazing gift!	744
	I’m so happy! Thank you for this—it’s wonderful!	745
	Thank you! This is so thoughtful and kind!	746
	This is so good—thank you for this amazing surprise!	747
	Thanks! Is that the name? It’s wonderful!	748
	Oh, thank you so much! This is incredible!	749
	This is so cool—I’m so happy! Thank you for it!	750
	Thank you so much! I really like this—it’s amazing!	751
	Thanks for this—it’s so thoughtful!	752
	I’m so happy with the name—thank you for everything!	753
	Wow, thank you so much! This is incredible!	754
	Oh, thank you! What is it? It’s amazing!	755
	Oh, this is so nice—thank you for this!	756
	No way! This is unbelievable!	757
	Wow, the name is so cool—thank you! I really like it!	758
	This is good—thank you so much!	759
	Is this really mine? Thank you so much!	760
	Thank you for this—I’m truly grateful!	761
	You’re amazing! Thank you so much for everything!	762
	This is incredible—thank you! It’s so nice and thoughtful!	763
	Love this! It’s so nice—it makes me so happy!	764
	Oh, thank you for this—it’s so good!	765
	This is amazing—thank you for everything!	766
	Thank you so much! This is so nice—I love it!	767
	What is this? I like it a lot—thank you!	768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780

781 Love this name—it's so cool! Thank you so
782 much!

783 Oh, thank you! The name is perfect—it's amaz-
784 ing!

785 I love this—it's so nice! Thank you so much!

786 Thank you for this—it's absolutely wonderful!

787 I really like this—it's so cool! Thank you end-
788 lessly!

789 Thank you for this! I'm so happy—it's so cool
790 and nice!

791 A.2 Diffusion-LM example (20)

792 Daily applied talks merged with furious dance de-
793 bates, where liberal ideals clashed with banned pop
794 culture references, showcasing the genie's unre-
795 solved power in chaotic narratives.

796 Kidding tones dominated recapitulated debates,
797 where furious beer-led arguments violated common
798 understanding, leaving bastards in heated copyright
799 disputes at 4th-level chambers.

800 Rent-related deliberations stumbled over 12-
801 weight discussions, navigating attractive yet un-
802 grounded topics, with residents grappling over mes-
803 sages obscured by toyotanies' lingering echoes.

804 Defensive positions hardened as heads engaged
805 in weighty deliberations, their dance-focused
806 approaches riddled with passive outcomes and
807 pointed major concerns. Furious exchanges high-
808 lighted daily references to potter-like ambitions,
809 where Japanese struggles weighed heavily on bas-
810 tards' passive handling of transitory truths.

811 Kidding undertones and guilty exchanges de-
812 fined the rent-led talks, where extraat-dances and
813 unique gestures left owners moaning over un-
814 resolved card-based disputes. Pairs of ideas col-
815 lided as furious heads wrestled with depression
816 and apple-stemmed inspirations, leaving dance
817 recitals receding into deeper questions and un-
818 solved guilt.

819 Further street-led discussions danced around
820 kidding-weighted rent ideas, merging landlord cri-
821 tiques with resident-focused solutions that sought
822 meaningful outcomes amidst ambiguity.

823 Bastards revisited weighted debates where beer-
824 fueled discussions aimed to banish average out-
825 comes, leaving toyotanies' lingering impressions
826 on a bull-filled stage. Print-heavy conversations
827 mirrored furious discussions, where accidental ref-
828 erences and sensitive rosters collided with trap-
829 filled agendas aimed at avoiding pointless out-
830 comes.

Harmful misunderstandings rid talks of weight,
as references and Japanese-inspired sensitive cri-
tiques exposed alt-dance roster traps, leaving spe-
cific outcomes unclear but charged with potential.

The weight of daily talks about American values
added pressure to heads already burdened by dance-
related debates, ridiculed bans, and beer-fueled
administration chaos, leaving a lingering sense of
personal disarray.

Furious discussions over sarcastic rent policies
led to sarcastic roster critiques, where bastards
mouthed their disdain amidst brave debates over
beer traps and meaningful outcomes.

Talks revolved around an attractive roster of pre-
tent priorities, as heads addressed liberal sensitiv-
ties amidst bans, dry humor, and the weight of
transitory decisions.

Furious American debates wrecked the mood
as artists and owners collided over closed issues,
with angels of optimism crushed beneath layers of
administration and transitory uncertainty.

Routine daily discussions on TV-related topics
revealed mounting tension, as cups of ideas spilled
over into meaningful debates about societal pur-
pose. Barely controlled fury arose as bastards
rushed into owner-focused talks, their sarcasm min-
gling with beer-fueled frustrations and unresolved
Toyota-related goals. The weight of daily rent dis-
cussions over roster issues brought about grand
ideas, yet visibility issues and passive preferences
continued to block meaningful resolutions.

Kidding comments dominated term-heavy dis-
cussions, blending weighty Japanese influences
with dance-fueled debates over savings and a
utopian vision of honesty. Depression-colored de-
bates about American defense painted a bleak pic-
ture as bastards, beer, and sensitive topics inter-
twined in a greedy search for honesty.

A.3 WGAN example (20)

Sloane's waitress reviews uploaded sigma offsets
into pirate situations, as fulfilled pebbles in 20-
blossom blooms harmonized eternal identity mark-
ers.

Nations weighed narrow romances as Jericho's
rubble receptionist faced impromptu bachelor pol-
lution at Corsica's forefront, aligning clemson
leaves publicly.

Purdue's bibliography extended partnerships
into modernist dreams, accelerating printed out-
comes while Hamlet's bought accelerations singled
a pouring vision.

882	U2's reflective contradictions evaluated Niger's humble incarnation as Jennifer's messianic gulp arose with Geneva's emphasized contradictions.	934
883		935
884		936
885	Beijing's rash foothills hosted weekend invita-	937
886	tionals, where Vancouver's regulated territory counted hobby antibodies eroding into a maximal	938
887	balance.	939
888		
889	Julian's patterned entertainment stretched into metaphysical drama, as investigators in Malibu in-	940
890	spected colony returns and open solicitor inquiries.	941
891		942
892	Slovene's 1704 recipients balanced dramatic dif-	943
893	ficulties as Alejandro believed in outlaw rulers	944
894	while Forrest's scars totaled 1813's wastewater re-	945
895	frain.	946
896		947
897	Jamaica's humble cub distributions extended Pre-	948
898	torian brakes, where 276 Compton creditors envi-	949
899	ronmentally released Meyer's owner interests into	950
900	1838 outset.	951
901		952
902	Disney's darkest symmetric junctions bargained against ancestral heresy, where induction aided heiress helmets amidst aggressive fun encounters.	953
903		954
904	Macleod's manpower wrapped boots into exclu-	955
905	sive bien meals, as Walton's lonely astronomy pro-	956
906	voked shooters in beers-filled envelopes.	957
907		958
908	Judy's aristocratic garden sprawled with unde-	959
909	feated Ebert viper criteria, where Pinto's humanoid graz and Edouard's painter nodes framed heretical feuds.	960
910		961
911	Merlin's mammalian plays provoked pissed NFL motors, where Hepburn reprinted vilified paints in	962
912	brains-filled demotions.	963
913		964
914	Mira's lithium simplicity interfered with post-	965
915	master riddles, luring immortal seats eastward as	966
916	maximal Stein channels accumulated medicine in-	967
917	terference.	968
918		969
919	Kerman's sporadic polka explosions framed stretched bocatier ankle facts, where Missy's 1760	970
920	dairy dating aligned with Akbar's carriage bursts.	971
921		972
922	Eve's mainline sensible retaliation letters framed Sino beverages as Dallas meters crawled Chicago	
923	efficiencies with antibody-functioned dynamics.	
924		
925	Augsburg's loud patrons startled Brunei's whipped photography reserves, where narrows torn	
926	by alcoholism conceded Joanne's scigan ideals.	
927		
928	Suarez maintained dramatic conspiracies across postwar carpets, where Belmont's astronomy junc-	
929	tion trembled in Rosario's 1813 guild havoc.	
930		
931	Lux's Moroccan scars swelled into Vulcan-paced leg constitutions, where laughter captured Swami's	
932	dependency ascent amidst runaway beliefs.	
933		
	Freeman's enormous nationwide acronyms awaited explosive distances as Eve's sparhawk	
	priests borrowed acronyms for Canadian-grown onions.	
	Goa's promotion presided over strict neighbor-	
	hood advice, where Lazarus gasped in record-	
	stretched launchers amidst precisely recorded col-	
	laborations.	
	A.4 SeqGAN example (20)	940
	Indeed, I believe in myself as I am at my best.	941
	This is a peculiarly good thing unlike any other.	942
	It's a rare occasion when someone like you comes along, who's so different in a truly unique	943
	way.	944
	It's just an idea that I'm toying with, though it seems like an exciting prospect to me.	945
	This is such a dreadful time, isn't it?	946
	I've always desired to have such profound ex-	947
	periences with you, and I'm convinced that your	948
	presence will evoke an unprecedented sensation.	949
	I still crave to know more even though it's not an ideal time. The desire to learn persists.	950
	I don't have the same issues as others.	951
	I insist on seeing that unique side of you, and I'm certain it will open up new pathways for me.	952
	It's pivotal to utilize that if you decide to do so.	953
	You won't believe the exhilaration I felt when I experienced the authenticity of the game. They	954
	have so much to offer you.	955
	He was just a stumbling block in your path.	956
	Wow, I really accomplished this.	957
	I would choose to venture into the lot.	958
	The only aspect is sparing a little time, consist-	959
	ently so.	960
	It might be the same, but I believe I don't quite grasp my second thought. This is my singular year.	961
	Well indeed, I'm gearing up to dedicate an entire day, even if it's challenging. He is the only one; he'll appreciate my method and continue to support me. I have this unusual feeling, and it's not so outlandish. She's just unique.	962
	A.5 TextGAN example (20)	973
	Holden's shameful face produced a scene of mis-	974
	leading echoes, as the surgical supervisors peeked	975
	through the malicious tones of Robin's actions, cre-	976
	ating a scene of Austro-infused chaos.	977
	Despite the echoing face of Austro's classmates, Callum's subsequent actions reverted the growling chaos of theological influences, creating a scene of	978
	influential chaos and ignored echoes.	979
	The malicious trapping of ballast and pain was deemed unacceptable, as Montevideo's surgical	980
		981
		982
		983

984	echoes silenced Robin, creating a scene of influen-	1036
985	tial chaos and ignored echoes.	1037
986	The temple's gloss echoed through Tuscany, as	1038
987	the misleading growl of Austro's surgical actions	1039
988	created a scene of reign-infused chaos and influen-	1040
989	tial echoes. Callum's surgical alloy production	1041
990	leaves a distinct footprint in the Midlands, silencing	1042
991	any doubts with theological certainty and a sense	1043
992	of urgency.	1044
993	The surgical echo of Korean royals resonates	1045
994	as Callum's growl demands attention, while Cap-	1046
995	tain Edouard's pain echoes through the theological	1047
996	fence of alloy mora. Captainbourne, with a mis-	1048
997	chievous and malicious approach, neglects devel-	1049
998	opmental contributions, while Holden's influential	1050
999	presence in Montevideo's lobby remains silenced	
1000	yet significant.	
1001	Streams of surgical novelty captivate as Cap-	1051
1002	tain Callum's misleading micro contributions are	1052
1003	trapped, intervening in an unacceptable yet intrigu-	1053
1004	ing manner.	1054
1005	Fence supervisors, maliciously overseeing the	1055
1006	growl of the captain's pain, witness the theological	
1007	acquisition of Holden's shameful surgical machines	
1008	in Roanoke.	
1009	Bundles of repairing materials mislead as Monte-	1056
1010	video sees influential dates, with Holden's contribu-	1057
1011	tions repeatedly emphasizing the malicious nature	1058
1012	of the situation.	1059
1013	Echoing praises of Captain Callum's organisms,	1060
1014	the misleading trapping of Robin's neo creations	1061
1015	shows the demand for truthful leadership amidst	1062
1016	deceit.	1063
1017	Rappers and supervisors mischievously greet	1064
1018	Montevideo with influential yet misleading reparations,	1065
1019	trapping those who sway under the echo of	1066
1020	mischievous timbers.	1067
1021	Roanoke's captain, alongside Callum and Gil,	1068
1022	faces surgical challenges in the Midlands, as mis-	1069
1023	leading actions bar Holden's attempts to influence	1070
1024	the scene.	
1025	Certainty in malicious apartments breeds mis-	1071
1026	leading alloys, as the theological fence in the Mid-	1072
1027	lands is overshadowed by Robin's austere foot-	1073
1028	prints.	
1029	The mischievous growl of kilometers covered	1074
1030	by Callum's surgical endeavors highlights the mis-	1075
1031	leading nature of Providence's lobby, where mis-	
1032	chievous actions thrive.	
1033	Surgical computers, maliciously ignored, cause	1076
1034	misleading streams to favorably influence Holden's	1077
1035	presence in the Midlands, as growls are covertly	
	ignored.	1081
	In a hive of activity, Holden's face is trapped in	1082
	a cycle of repair, as Robin's echo and hive growl	1083
	influence the timbers of the scene.	
	Captain Johann's salsa dance with Robin in Mon-	1084
	tevideo is ignored by mischievous supervisors, as	1085
	subsequent laws misleadingly influence the alloy's	1086
	path.	
	Echoing tones across hemispheres, Montev-	
	ideo's influential plaza peek at candlelit repairs,	
	mirroring the alloy's impact on the influential hemi-	
	sphere.	
	Reptiles on a sensory fence face malicious sur-	
	gical pain, as influential alloys produce malicious	
	damages that default sensory preferences.	
	A.6 ControlGAN example (20)	1051
	Amidst the mystery of paralympics and genomes,	1052
	photographic listings traced mourning eastward,	1053
	with yeast and surgeons unable to escape the gallery	1054
	of profound losses.	1055
	In Washington, the plight of descendants dou-	1056
	bled, as helicopters neatly crawled phone lines,	1057
	leaving a reddish, wrinkled component in their	1058
	wake.	1059
	Profoundly biting into outlying lands, descend-	1060
	ants of 1984 overlooked socialists' raids, sam-	1061
	pling distinctive logos and biting into the disarray	1062
	of a villainous sphere.	1063
	Israelis favored comedic rhythms over monu-	1064
	ments, as comrades paved the way for boar-like	1065
	terraces, cultivating goblin-like responses in the	1066
	process.	1067
	Documented ambient activism glowed in the	1068
	ecosystem, with prisoners freshly interviewed and	1069
	regarded as widening folks in a galactic realm.	1070
	Historians ran through the boon of their nation-	1071
	ality, puzzled by radios and finalized renovations	1072
	that doggedly disliked the oracle's diet.	1073
	In 1790, privileged violinists bonded with their	1074
	fathers, touching on the fairness of numerous gale-	1075
	like frontiers while regaining attributes and model-	1076
	ing repetitive touches.	1077
	Paving the way for salvation, fresh casabas were	1078
	organized alongside niece-like lyrics, slightly wor-	1079
	rying examinations of mid-summer greaves.	1080
	Intentionally structured streams of dissent deep-	1081
	ened footprints, as unusual boxers declared their	1082
	inquiry into the layers of Egyptian understanding.	1083
	Congression was aggressively instituted, with	1084
	Edison and others like Lucy visiting global sites for	1085
	repairs, reflecting a backward trend in expenditure.	1086

1087	Napier and Ismail hurriedly exchanged credits in Dhaka, despite the illegal deformation of merged assets causing a humanoid outcry.	I hate how this job has turned me into someone who's always angry.	1138
1088		It's frustrating to watch people skate by while I'm drowning in work.	1139
1089		I'm so sick of feeling like I'm not allowed to have boundaries.	1140
1090	In a schoolhouse, Daniel spurred investigative uses of solvents, while teasingly staging scenes that endangered the actor's craft.	I'm infuriating to be the one who always has to fix things.	1141
1091		I'm angry that my personal life is suffering because of this job.	1142
1092	Ezio lingered in dread as a ghostly pattern opposed his lifestyle, bundling his fears into a volcanic bowl of brilliance.	It's maddening to be the only one who seems to care about the outcome.	1143
1093		I'm so tired of having to explain myself over and over again.	1144
1094	The goalkeeper, reelected in 1998, crouched in a surplus rise, befriending Celia amidst a philanthropic sun.	It's frustrating to feel like my hard work is going to waste.	1145
1095		I hate how no one seems to notice when I'm struggling.	1146
1096	In Nuremberg, Sylvie tasted activism, emerging sweaty yet live, as she exploited her confirmation with dynamism.	It's infuriating to watch people who do nothing get ahead.	1147
1097			1148
1098	In Marseille, hostility commanded the outer regions, as Lord Sloan conserved his authority amidst a sip of repository.		1149
1099			1150
1100	Mercury additions were jected with devotion, as Moe recognised the unconscious blaze of Rufus's concentration efforts.		1151
1101			1152
1102	Ezekiel reflected on the ultimate noon showdown, confronting disabilities that regained his classmate's faith on the cliff.		1153
1103			1154
1104	Pegasus healed the psychiatric wounds of Willard, while unofficially sliding into the fresh aroma of Robert's typhoon.		1155
1105			1156
1106	A hermit hung in speechless tension as countless steps ensued, addressing the coral variations of local alliances.		1157
1107			
1111	A.7 LLaMA example (20)		
1112	I feel like I'm drowning in other people's incompetence.		
1113			
1114	It's frustrating that no matter how hard I work, it's never enough.		
1115			
1116	I'm angry that they expect me to keep smiling and pretending everything is fine.		
1120	Why does it feel like I'm the only one who sees how messed up this place is?		
1121			
1122	It's infuriating to watch people get rewarded for doing the bare minimum.		
1123			
1124	I hate that I've become so jaded, but I can't help it.		
1125			
1128	I'm tired of feeling like I'm stuck in a never-ending loop of frustration.		
1129			
1130	It's maddening to feel so undervalued, like my effort means nothing.		
1131			
1132	I feel like I'm on the brink of a breakdown, and no one even notices.		
1133			
1136	It's enraging to be the one who's always expected to clean up everyone else's mess.		
1137			