

Syntax-Guided Diffusion Language Models with User-Integrated Personalization

Ruqian Zhang¹, Yijiao Zhang¹, Juan Shen¹, Zhongyi Zhu¹, Annie Qu^{2*}

¹Department of Statistics and Data Science, Fudan University

²Department of Statistics and Applied Probability,
University of California, Santa Barbara

Abstract

Large language models have made revolutionary progress in generating human-like text, yet their outputs often tend to be generic, exhibiting insufficient structural diversity, which limits personalized expression. Recent advances in diffusion models have opened new opportunities for improving language generation beyond the limitations of autoregressive paradigms. In this work, we propose a syntax-guided diffusion language model that integrates structural supervision and personalized conditioning to enhance text quality, diversity, and controllability. We introduce a cascaded framework that generates syntactic guidance before conditional text generation, and further generalize it to a novel noncascaded architecture for better alignment between structure and content. By incorporating syntactic information in the generating process, the proposed model better captures the lexical and structural characteristics of stylistic sentence construction. To enable fine-grained personalization, we develop a shared representation mechanism that facilitates information integration across users, supporting both faithful stylistic generation and generalizable zero-shot inference. Extensive experiments on multiple tasks demonstrate the superiority of our approach in fluency, diversity, and stylistic fidelity. Further qualitative analyses highlight its interpretability and flexibility in learning personalized patterns.

Keywords: Cascaded generation, diffusion models, natural language processing, representation learning, structural information

*Corresponding author

1 Introduction

In recent years, large language models (LLMs) have revolutionized natural language processing (NLP), demonstrating impressive performance in generating human-like text (OpenAI et al., 2024). However, their widespread use has also raised growing concerns about linguistic homogenization (Wan et al., 2023; Sourati et al., 2025). Generated sentences often follow generic templates, leading to overused lexical patterns and repetitive structural compositions, such as the notorious “em dash (–) conspiracy” (Mummery, 2025). The lack of diversity tends to diminish stylistic richness, posing a challenge for applications demanding personalized expression. For example, we rewrite a movie quote into three distinct styles using ChatGPT-4o (OpenAI et al., 2024) and DeepSeek-V3 (DeepSeek-AI et al., 2025). As shown in Figure 1, the outputs largely preserve the original structure in the input, with stylistic variations confined mainly to a small set of style-bearing words.

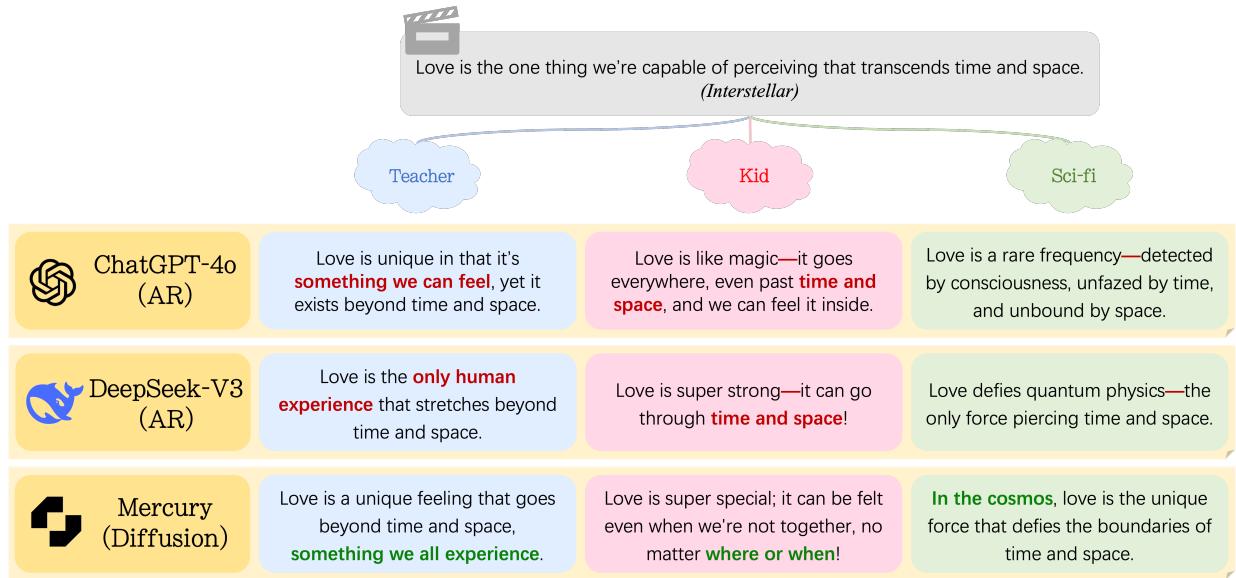


Figure 1: Example of rewriting a movie quote into different styles using autoregressive or diffusion models. Text in red shows limited modifications in words from the AR models, while text in green highlights structural variations introduced by the diffusion approach.

The limited structural diversity of existing LLMs largely arises from their autoregressive (AR) nature, which generates text through a next-token prediction paradigm (Holtzman et al., 2020). This left-to-right construction favors high-probability lexical choices as continuations, which hinders the ability to refine sentence structure at a global level once tokens are committed. As a powerful alternative, diffusion-based language models have recently emerged (Li et al., 2022; Nie et al., 2025), operated by denoising an input through iterative sampling (Song and Ermon, 2019; Ho et al., 2020), which enables parallel updating across all positions. This parallelism allows for global reorganization of sentence composition, providing greater flexibility and refinement. As illustrated in Figure 1, the diffusion-based Mercury (InceptionLabs et al., 2025) produces richer structural variations and achieves more powerful stylistic diversity beyond word substitutions.

Motivated by the importance of structural information, we propose a syntax-guided diffusion framework to enhance text diversity and personalized control. Syntactic patterns offer valuable signals for capturing personalized traits (Alhafni et al., 2024). For example, in the Yelp Review dataset, which comprises user reviews with sentiment labels, part-of-speech (POS) distributions vary across sentiment classes. As shown in Figure 2, positive reviews favor more frequent use of adjectives (e.g., friendly), while negative reviews use more verbs and pronouns that reflect personal experiences. Syntax also contributes to improved text fluency, as it resembles the compositional nature of human language (Chomsky, 1957).

Despite the rich information carried by syntax, its use in text generation remains underexplored due to the absence of predefined syntactic inputs in real-world scenarios. To the best of our knowledge, the application of syntax is limited to sequence-to-sequence tasks with explicit syntactic conditions (Li et al., 2023). To fill this gap, we propose a two-stage cascaded framework (Ho et al., 2022), where syntactic structures are first gen-

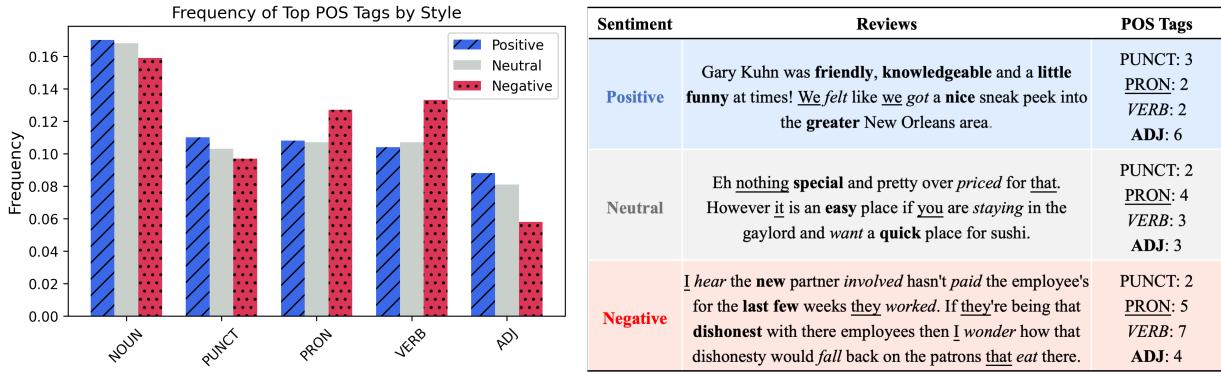


Figure 2: Left: Frequencies of the most common part-of-speech (POS) tags across sentiment styles in the Yelp Review dataset. Right: Examples of reviews with their POS tags.

erated as conditions to guide final text synthesis. While cascaded designs are adopted in image generation (Ramesh et al., 2022; Saharia et al., 2022), their potential for language modeling is underdeveloped. Compared to generating text from scratch, the intermediate syntactic generator introduces a structural prior that not only filters out implausible syntactic patterns but also encourages sentence structures aligned with target styles.

The use of syntax offers a promising direction for personalization at the structural level. As shown in Figure 2, sentiment styles exhibit both distinctive syntactic traits and overlapping patterns; for example, similar frequencies of nouns or punctuations are observed. However, most existing approaches, such as user-specific identifiers (Zhong et al., 2021) or low-rank adaptation (Hu et al., 2022), typically treat classes in isolation, limiting the potential to share information across styles. This motivates us to develop a new data integration tool, which is further supported by semantic-level evidence in Figure 3. Specifically, at the word level, positive reviews tend to use praising words (e.g., great, amazing), while negative ones contain more critical words (e.g., rude, worst). Neutral reviews share words with both positive and negative styles, while additionally favoring milder vocabulary (e.g., overall, okay). This interplay extends to sentence-level generation, illustrated by the clus-

tering of SimCSE-based sentence embeddings (Gao et al., 2021), where stylistic boundaries form both distinct and overlapping regions in the latent space.

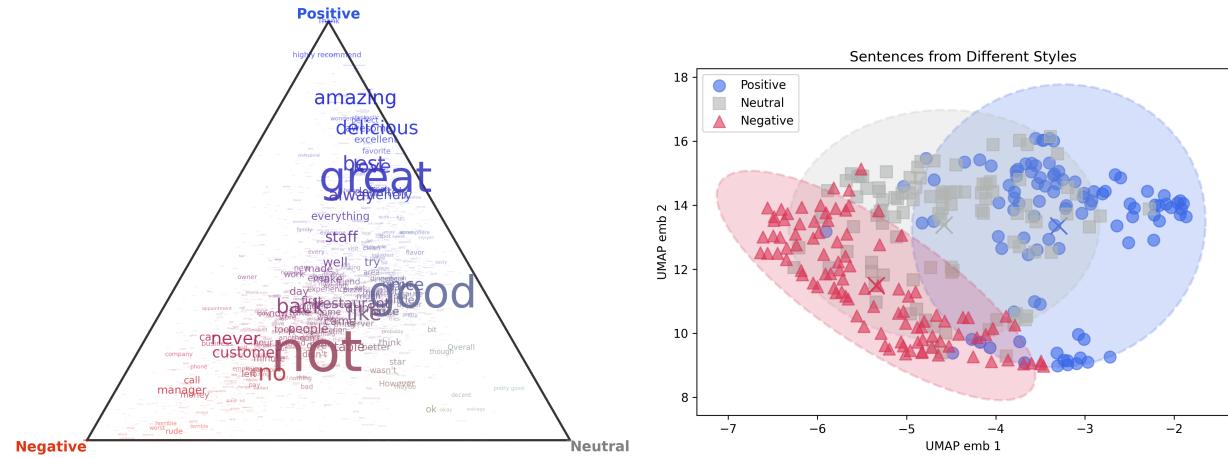


Figure 3: Left: Word-level distribution in the Yelp Review dataset in a sentiment triangle. Right: Sentence-level UMAP projection (McInnes et al., 2020) of SimCSE embeddings.

To better exploit shared and distinctive stylistic signals, we propose a personalization framework via the attention mechanism (Vaswani et al., 2017) based on latent features. Instead of modeling each personalized embedding separately, we assume that the embeddings are constructed from a set of building-block personality representations in a shared latent space, weighted uniquely for each style. Through integrated modeling, our method leverages all available data to characterize each style, enhancing both generalizability and data efficiency. To achieve fine-grained control, we move beyond the conventional use of a single embedding vector and introduce an embedding matrix for each style, enabling personalization through a cross-attention mechanism. This attention layer is integrated into both the syntax and text stages of our cascaded diffusion framework.

We further investigate the behavior of multi-stage generation in the cascaded frameworks. While cascaded diffusion models decompose generation into sequential stages, they often ignore the potential feedback across stages. To better leverage the parallel nature of

diffusion, we propose a generalized noncascaded framework with overlaps between stages, where later predictions can dynamically refine earlier conditions through bidirectional interactions. In addition, we introduce a unified attention mechanism that merges multiple parallel diffusion models into a single, jointly conditional generation process. By enabling continuous interaction between syntactic and textual components, this unified design improves alignment across stages while reducing overall model complexity.

Extensive experiments on multiple datasets demonstrate that our method consistently outperforms AR and diffusion baselines of comparable model sizes in both text generation quality and personalized control, validating the effectiveness of syntactic guidance and cross-class information sharing. Moreover, numerical studies demonstrate that the proposed noncascaded architecture enhances cross-stage coherence compared to the traditional cascaded framework. We showcase the generalizability of the shared personality representations by extrapolating and interpolating learned class weights to synthesize unseen styles in a zero-shot manner on the Yelp dataset.

The remainder of the paper is organized as follows. In Section 2, we present the proposed syntax-guided cascaded diffusion framework. Section 3 introduces the integrated personalization method based on shared latent features. In Section 4, we extend the cascaded framework to a noncascaded formulation that enables bidirectional interactions. Section 5 reports experimental results, including both quantitative evaluations and illustrative analysis. Section 6 concludes with potential future directions.

2 Syntax-Guided Text Generation

Text generation aims to model the distribution over text sequences. Unlike AR models, diffusion language models jointly learn the distribution of all tokens in a sequence, offering

higher diversity and flexibility in text synthesis (Hoogeboom et al., 2021; Chen et al., 2023). A major line of work operates in continuous embedding spaces (Li et al., 2022), which has proven effective in improving efficiency and quality (Lovelace et al., 2023). While diffusion language models have shown strong potential in LLM (Li et al., 2023), they remain in an early stage of development and often underperform AR counterparts in terms of text quality (Gong et al., 2023; Rout et al., 2025), underscoring the need for improvement.

Given a text sequence \mathbf{w}_x of length L , we denote its embeddings as $\mathbf{x}_0 = E_x(\mathbf{w}_x) \in \mathbb{R}^{L \times d}$, where E_x is a pretrained text encoder and d is the dimension of the continuous embeddings. Directly modeling the distribution $p_{\theta}(\mathbf{x}_0)$ without any prior structure often poses challenges in obtaining coherent outputs. To facilitate this process, we consider data augmentation by introducing auxiliary structural variables \mathbf{z} :

$$p_{\theta}(\mathbf{x}_0) = \int p_{\theta}(\mathbf{x}_0 | \mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z}.$$

Building on this formulation, we propose a cascaded diffusion framework that leverages structural information by decomposing text generation into two stages: predicting structure first and then generating text conditioned on its structure.

Syntax provides a natural and interpretable choice for structural information in language analysis. Unlike images, whose structure is often encoded implicitly in a latent space (Hong et al., 2018), syntax is explicitly observable in text, allowing for linguistically grounded guidance. Moreover, syntax acts as an additional source of supervision, which is particularly valuable when human-written text is scarce. In this work, we adopt part-of-speech (POS) tags as the structural representation, which are easily accessible from tools such as spaCy (Honnibal and Montani, 2017). Examples of POS tag sequences are illustrated in Figure 4.

Text (w_x)	All models are wrong, but some are useful.
Syntax (w_s)	DET NOUN AUX ADJ PUNCT CCONJ PRON AUX ADJ PUNCT
Text (w_x)	Statistics is the grammar of science.
Syntax (w_s)	NOUN AUX DET NOUN ADP NOUN PUNCT

Figure 4: Examples of POS tag sequences for different sentences.

2.1 Learning to Generate Syntax as Structural Guidance

Given a text sequence w_x , we extract its POS sequence w_s . To model syntactic structure within a diffusion framework, we adopt the continuous diffusion approach from Li et al. (2022). We use a pretrained syntactic encoder E_s to project w_s into a continuous latent space, yielding syntactic embeddings $s_0 = E_s(w_s) \in \mathbb{R}^{L \times d}$.

We denote the true data-generating distribution of the syntactic embeddings as $s_0 \sim q(s_0)$. The forward diffusion process progressively perturbs s_0 over T timesteps by adding Gaussian noises, governed by a predefined variance schedule $\{\alpha_t^s\}_{t=1}^T$ with $\alpha_t^s \in (0, 1)$:

$$q(s_t | s_{t-1}) = \mathcal{N}(\sqrt{\alpha_t^s} s_{t-1}, (1 - \alpha_t^s) \mathbf{I}),$$

where s_t denotes the noisy syntactic embeddings at time t and \mathbf{I} is the identity matrix. This Markovian corruption process produces a sequence of latent variables $\{s_t\}_{t=1}^T$ with decreasing signal-to-noise ratio. As T becomes sufficiently large, the marginal distribution of s_T converges to a standard multivariate Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$. If the reverse distribution $q(s_{t-1} | s_t)$ were available, we could generate new samples from $q(s_0)$ by first sampling $s_T \sim \mathcal{N}(0, \mathbf{I})$ and then iteratively removing noises.

Since $q(s_{t-1} | s_t)$ is computationally intractable, we adopt variational approximation with a reverse model $p_{\theta_s}(s_{t-1} | s_t)$, where θ_s denotes the trainable parameters. This reverse process learns to denoise the latent sequence to reconstruct samples from noise. To train the

model, we minimize the KL divergence between the true joint distribution $q(\mathbf{s}_{0:T})$ and the parametrized joint distribution $p_{\boldsymbol{\theta}_s}(\mathbf{s}_{0:T})$, which is equivalent to maximizing a variational lower bound on the log-likelihood $\log p_{\boldsymbol{\theta}_s}(\mathbf{s}_0)$. The objective function is given by

$$\begin{aligned}\mathcal{L}_s(\boldsymbol{\theta}_s) &\equiv D_{\text{KL}}(q(\mathbf{s}_{0:T}) \| p_{\boldsymbol{\theta}_s}(\mathbf{s}_{0:T})) = \mathbb{E}_{q(\mathbf{s}_{0:T})} \left[\log \frac{q(\mathbf{s}_{0:T})}{p_{\boldsymbol{\theta}_s}(\mathbf{s}_{0:T})} \right] \\ &= \mathbb{E}_{q(\mathbf{s}_{0:T})} \left[\log \frac{q(\mathbf{s}_T | \mathbf{s}_0)}{p_{\boldsymbol{\theta}_s}(\mathbf{s}_T)} + \sum_{t=2}^T \log \frac{q(\mathbf{s}_{t-1} | \mathbf{s}_t, \mathbf{s}_0)}{p_{\boldsymbol{\theta}_s}(\mathbf{s}_{t-1} | \mathbf{s}_t)} + \log \frac{q(\mathbf{s}_0)}{p_{\boldsymbol{\theta}_s}(\mathbf{s}_0 | \mathbf{s}_1)} \right],\end{aligned}\quad (1)$$

where $p_{\boldsymbol{\theta}_s}(\mathbf{s}_T) = \mathcal{N}(0, \mathbf{I})$ serves as the prior. Following Ho et al. (2020), we model the reverse transition $p_{\boldsymbol{\theta}_s}(\mathbf{s}_{t-1} | \mathbf{s}_t)$ as a Gaussian distribution with fixed variance:

$$p_{\boldsymbol{\theta}_s}(\mathbf{s}_{t-1} | \mathbf{s}_t) = \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}_s}(\mathbf{s}_t, t), (1 - \alpha_t^s)\mathbf{I}),$$

where $\boldsymbol{\mu}_{\boldsymbol{\theta}_s}(\mathbf{s}_t, t)$ is predicted by a neural network, implemented using a Transformer-based architecture (Vaswani et al., 2017; Peebles and Xie, 2023). This parameterization mirrors the forward process and allows efficient optimization of the objective.

Leveraging the syntactic diffusion model, we can generate syntactic conditions $\hat{\mathbf{w}}_s$ from scratch by first sampling random noises $\mathbf{s}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and iteratively denoising through the learned reverse model $p_{\boldsymbol{\theta}_s}(\mathbf{s}_{t-1} | \mathbf{s}_t)$ for $t = T, T-1, \dots, 1$. The final denoised embedding $\hat{\mathbf{s}}_0$ is subsequently reconstructed into a POS sequence using a syntactic decoder D_s . As illustrated in Figure 5, the syntactic diffusion model forms the first stage of our generation pipeline. It enables the synthesis of structural conditions without predefined syntactic inputs, thereby allowing more flexible downstream text generation.

2.2 Cascaded Text Generation with Syntactic Conditioning

The cascaded text generation leverages syntactic structure through a two-stage process as $p_{\boldsymbol{\theta}_x}(\mathbf{x}_0 | \mathbf{s}_0)p_{\boldsymbol{\theta}_s}(\mathbf{s}_0)$ with $\boldsymbol{\theta}_x$ denoting the trainable parameters of the text diffusion model. In the first stage, $p_{\boldsymbol{\theta}_s}(\mathbf{s}_0)$ serves as a prior over syntactic structure, as described in

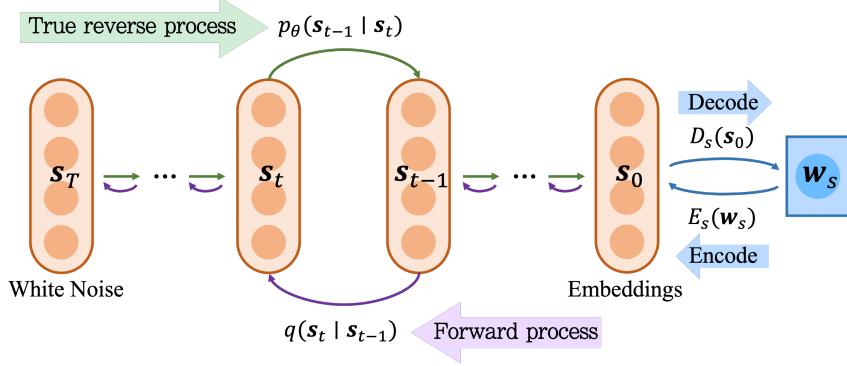


Figure 5: Overview of the syntactic diffusion model. The input syntactic sequence \mathbf{w}_s is encoded into latent embeddings \mathbf{s}_0 , which are diffused into noise via the forward process. The reverse process reconstructs \mathbf{s}_0 , which is then decoded into predicted POS tags.

Section 2.1. In the second stage, $p_{\theta_x}(\mathbf{x}_0 | \mathbf{s}_0)$ models the conditional distribution over text given the prior, formulating the generation task as a sequence-to-sequence problem.

To model the second stage, we introduce a conditional text diffusion model that incorporates syntactic information via an attention mechanism. The conditional denoising model is defined as $p_{\theta_x}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{s}_0) = \mathcal{N}(\boldsymbol{\mu}_{\theta_x}(\mathbf{x}_t, \mathbf{s}_0, t), (1 - \alpha_t^x)\mathbf{I})$ with a schedule $\{\alpha_t^x\}_{t=1}^T$, which computes cross-attention between the noisy text embeddings \mathbf{x}_t and \mathbf{s}_0 at t . Specifically, we define learnable projection matrices for the query, key, and value as $W_Q, W_K, W_V \in \mathbb{R}^{d \times d_M}$, where d_M is the dimension of the projected space. The query $Q(\mathbf{x}_t)$ is computed as $Q(\mathbf{x}_t) = \mathbf{x}_t W_Q$, while the key and value are $K(\mathbf{s}_0) = \mathbf{s}_0 W_K$ and $V(\mathbf{s}_0) = \mathbf{s}_0 W_V$, respectively. The attention map is then given by $A = \text{Softmax}\left(\frac{Q(\mathbf{x}_t)K(\mathbf{s}_0)^T}{\sqrt{d_M}}\right)$, which captures the similarity between each query and key token. The output of the attention layer is a weighted combination of the value vectors as $AV(\mathbf{s}_0)$.

For the training of the text reverse model, we define a forward diffusion process analogous to that in Section 2.1, where the text embeddings \mathbf{x}_0 are progressively corrupted over T steps according to $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{\alpha_t^x}\mathbf{x}_{t-1}, (1 - \alpha_t^x)\mathbf{I})$. The model parameters θ_x are

learned by minimizing the objective function similar to (1) with $p_{\theta_x}(\mathbf{x}_T) = \mathcal{N}(0, \mathbf{I})$:

$$\mathcal{L}_x(\boldsymbol{\theta}_x) = \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_T | \mathbf{x}_0)}{p_{\theta_x}(\mathbf{x}_T)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)}{p_{\theta_x}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{s}_0)} + \log \frac{q(\mathbf{x}_0)}{p_{\theta_x}(\mathbf{x}_0 | \mathbf{x}_1, \mathbf{s}_0)} \right].$$

The syntactic and text diffusion models jointly form the cascaded generation pipeline, referred to as **SynText**, as illustrated in Figure 6. To generate a new sentence, the process begins with the syntactic diffusion model, which transforms a noise sample $\mathbf{s}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ into syntactic embeddings $\hat{\mathbf{s}}_0$ via iteratively denoising. These embeddings then guide the text diffusion model, which denoises a second random noise sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ to produce the text embeddings $\hat{\mathbf{x}}_0$. The output sentence $\hat{\mathbf{w}}_x$ is reconstructed from $\hat{\mathbf{x}}_0$ using the text decoder D_x . This two-stage framework enriches learnable information and enhances generation quality by explicitly modeling syntactic structure before text realization.

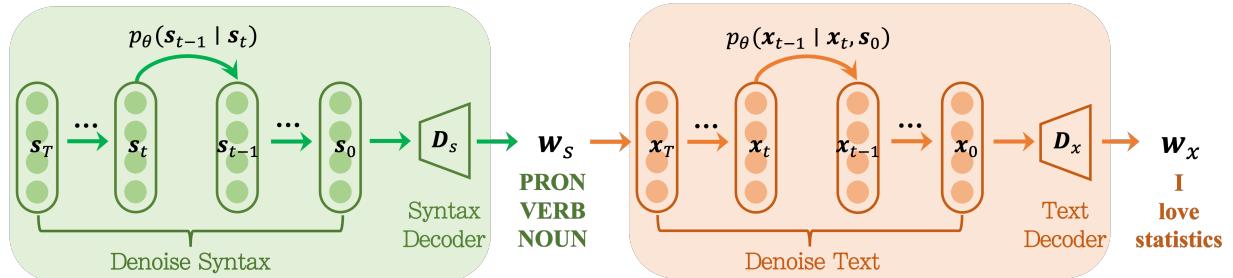


Figure 6: The **SynText** generation pipeline. The syntactic diffusion first generates syntactic embeddings \mathbf{s}_0 , which condition the text diffusion to produce the final output sentence.

Our proposed SynText exploits the potential of syntax as structural guidance in text generation. Beyond improving quality and diversity, syntax lays a new foundation for fine-grained personalization through stylistic structural patterns. The use of explicit syntax also offers greater interpretability from a linguistic perspective, providing reasoning into sentence structure beyond the probabilistic prediction using latent diffusion models. Nevertheless, the cascaded framework comes with limitations. It requires an additional syntactic diffusion model, and the two-stage design imposes a one-way direction of dependency

between syntax and text, which may induce error propagation if syntactic predictions are imperfect. These challenges motivate us to develop more flexible and efficient architectures, which we address in Section 4.

3 Personalization with Shared Personalities

3.1 Shared Latent Personality Representations

Personalized text generation seeks to align outputs with specific preferences by capturing stylistic attributes (Salemi et al., 2024). Existing methods typically treat each style in isolation, learning separate parameters to model stylistic variations. For example, unique identifier tokens are introduced to encode user information into latent embeddings (Zhong et al., 2021; Mireshghallah et al., 2022), or user-specific adapters, such as low-rank adaptation, are adopted to tune selected layers (Hu et al., 2022; Hayou et al., 2024). While effective, these class-by-class approaches restrict information integration across styles, leading to inefficient data utilization and limited generalization. Furthermore, the complexity of model architectures requires careful layer selection in LoRA tuning (Frenkel et al., 2025), undermining its interpretability.

To address these challenges, we propose a flexible method based on information sharing and incorporate it into both the syntax and text diffusion stages. Motivated by the observations in Figure 2 and 3 on the Yelp dataset, we assume the existence of a shared representation space across all styles. We construct this space by introducing a set of shared latent representations, termed *personalities*, which serve as building blocks to unify style-specific embeddings.

Specifically, assume there are R common personality features, denoted as $\mathbf{p}_1, \dots, \mathbf{p}_R \in$

\mathbb{R}^{d_p} , where d_p is the embedding dimension. We define the personality matrix as $\mathbf{P} = (\mathbf{p}_1, \dots, \mathbf{p}_R) \in \mathbb{R}^{d_p \times R}$, which serves as a codebook storing all personalized factors shared across styles. Suppose the dataset contains K different styles. The personalized embeddings are then constructed by extracting relevant personalities in \mathbf{P} according to a personalized weight vector $\boldsymbol{\gamma}_k = (\gamma_{k,1}, \dots, \gamma_{k,R})^T$. The entry $\gamma_{k,r} \in [0, 1]$ represents how much the r th personality is expressed in style k and satisfies $\sum_{r=1}^R \gamma_{k,r} = 1$.

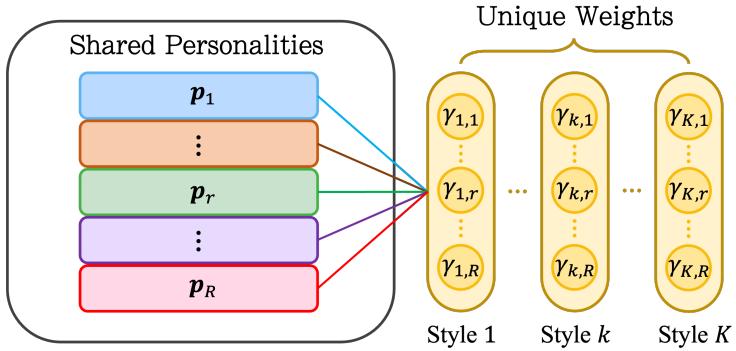


Figure 7: Illustration of shared personality representations and unique personality weights.

Our framework decomposes personalization into two components: a shared set of personality representations and distinct personalized weights, enabling better integration of cross-style information. By avoiding the need to learn separate embeddings for each style, the model significantly reduces the number of style-specific parameters from d_p to R , alleviating the reliance on per-style data and improving sample efficiency. The introduction of personality weights also enhances interpretability, which characterizes the composition of each style and associations across styles. Furthermore, the unified personality matrix enhances the generalizability of stylistic generation. Once the representations are learned, they span a shared latent space, from which unseen styles can be synthesized by adjusting the weights over the learned personalities through interpolation or extrapolation.

3.2 Personality Layer with Attention Mechanism

While personalized embeddings are typically defined as vectors from convex combinations of personality representations, $\mathbf{e}_k = \mathbf{P}\boldsymbol{\gamma} = \sum_{r=1}^R \gamma_{k,r} \mathbf{p}_r$, we extend this approach to a matrix form $\mathbf{E}_k = [\gamma_{k,1} \mathbf{p}_1, \dots, \gamma_{k,R} \mathbf{p}_R]^T \in \mathbb{R}^{R \times d_p}$ for $k = 1, \dots, K$. Unlike conventional methods that concatenate a fixed personalized token into the noisy input during the denoising process, our matrix-based design simulates a personalized prompt of length R , enabling the model to attend over multiple personality components through a cross-attention mechanism. This formulation facilitates dynamic and context-aware personalized conditioning during generation, offering better flexibility.

When applied in diffusion models, the personalized attention mechanism introduces a personality layer in the Transformer block, referred to as **PLayer**. For a given style k , we compute the key and value from \mathbf{E}_k as $K(\mathbf{E}_k) = \mathbf{E}_k W_K$ and $V(\mathbf{E}_k) = \mathbf{E}_k W_V$, where $W_K, W_V \in \mathbb{R}^{d_p \times d_M}$ are projection matrices. At each denoising timestep t , the query derived from the noisy text input \mathbf{x}_t is denoted as $Q(\mathbf{x}_t) \in \mathbb{R}^{L \times d_M}$. The generation is then guided toward the desired stylistic direction through cross-attention as illustrated in Figure 8.

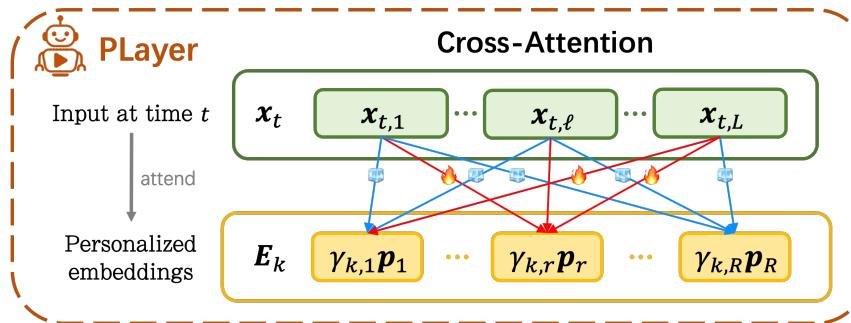


Figure 8: Illustration of the cross-attention personalization in PLayer.

In our cascaded diffusion framework, the PLayer is incorporated in both the syntax and text diffusion stages, enabling the model to capture richer stylistic signals compared to traditional diffusion language models that condition only on the text level.

4 Noncascaded Generation with Mutual Guidance

4.1 From Cascade to Noncascade

The cascaded framework defines a sequential order for generating multiple outputs, where each stage begins after the previous one has fully produced its results, as illustrated in Figure 9(a). For example, in **SynText**, syntactic structures are first generated and then used to guide subsequent text realization. However, the information flow between structure and content can be bidirectional, where semantics may also shape syntactic form.

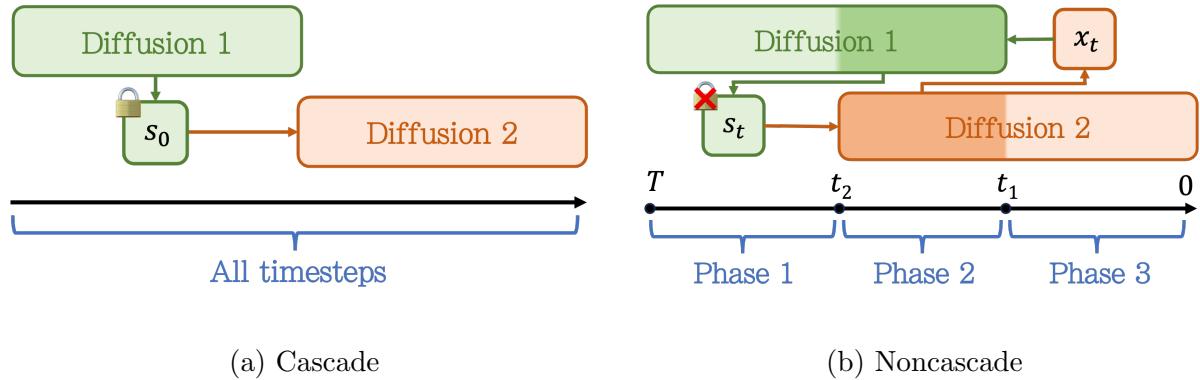


Figure 9: Generating processes of the cascaded diffusion and the noncascaded diffusion.

To accommodate bidirectional interactions, we propose a noncascaded diffusion framework by introducing temporal overlap between denoising processes, as shown in Figure 9(b). Instead of obtaining a fixed condition from the first stage, the two diffusion models proceed concurrently during the overlap phase, allowing dynamic information exchange from one to the other. The intermediate predictions from the second diffusion serve as real-time feedback to the first, promoting mutual adaptation throughout the generation.

Specifically, consider the two-stage case of syntactic and text diffusion models. Introducing temporal overlap divides the generation process into three phases. While the forward processes remain identical to those in Section 2, the reverse processes are redesigned to al-

low mutual dependency between modalities. In *Phase 1*, from T to t_2 , the model focuses on syntax, denoising random noise into intermediate syntactic embeddings by $p_{\theta_s}(\mathbf{s}_{t-1} \mid \mathbf{s}_t)$. In *Phase 2*, from t_2 to t_1 , both syntax and text are denoised in parallel, with the denoising models modified to enable bidirectional conditioning:

$$p_{\theta_s}(\mathbf{s}_{t-1} \mid \mathbf{s}_t, \mathbf{x}_t) = \mathcal{N}(\boldsymbol{\mu}_{\theta_s}(\mathbf{s}_t, \mathbf{x}_t, t), (1 - \alpha_t^s)\mathbf{I}),$$

$$p_{\theta_x}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{s}_t) = \mathcal{N}(\boldsymbol{\mu}_{\theta_x}(\mathbf{x}_t, \mathbf{s}_t, t), (1 - \alpha_t^x)\mathbf{I}).$$

In *Phase 3*, from t_1 to 0, the model completes text denoising conditioned on the fully generated syntax via $p_{\theta_x}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{s}_0)$. To maintain consistent notations across all phases, we define $\mathbf{x}_t = \mathbf{x}_{t_2}$ for $t \in [t_2, T]$ and $\mathbf{s}_t = \mathbf{s}_{t_1}$ for $t \in [0, t_1]$. The training objective for the noncascaded diffusion model is thus given by

$$\begin{aligned} \mathcal{L}_{non}(\boldsymbol{\theta}_s, \boldsymbol{\theta}_x) = & \mathbb{E}_{q(\mathbf{s}_{0:T}), q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{s}_T \mid \mathbf{s}_0)}{p_{\theta_s}(\mathbf{s}_T)} + \log \frac{q(\mathbf{x}_T \mid \mathbf{x}_0)}{p_{\theta_x}(\mathbf{x}_T)} + \sum_{t=t_2}^T \log \frac{q(\mathbf{s}_{t-1} \mid \mathbf{s}_t, \mathbf{s}_0)}{p_{\theta_s}(\mathbf{s}_{t-1} \mid \mathbf{s}_t)} \right. \\ & + \sum_{t=t_1+2}^{t_2} \log \frac{q(\mathbf{s}_{t-1} \mid \mathbf{s}_t, \mathbf{s}_0)}{p_{\theta_s}(\mathbf{s}_{t-1} \mid \mathbf{s}_t, \mathbf{x}_t)} + \sum_{t=t_1+1}^{t_2} \log \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)}{p_{\theta_x}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{s}_t)} + \log \frac{q(\mathbf{s}_0)}{p_{\theta_s}(\mathbf{s}_0 \mid \mathbf{s}_{t_1+1}, \mathbf{x}_{t_1+1})} \quad (2) \\ & \left. + \sum_{t=2}^{t_1} \log \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)}{p_{\theta_x}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{s}_0)} + \log \frac{q(\mathbf{x}_0)}{p_{\theta_x}(\mathbf{x}_0 \mid \mathbf{x}_1, \mathbf{s}_0)} \right]. \end{aligned}$$

By allowing the generated text embeddings to influence the syntactic embeddings, the resulting structures are not only grammatically valid but also better aligned with the evolving semantic content. In contrast to traditional cascaded approaches, this bidirectional adaptability retains the iterative refinement capability that is inherent to diffusion models, providing a more flexible and robust framework for complex generative workflows.

Remark 1. *By adjusting t_1 and t_2 , the degree of temporal overlap can be flexibly controlled to accommodate different tasks. This allows the model to strike a balance between the mutual alignment achieved by joint refinement and the modular independence of each stage.*

4.2 Complete Overlap with Unified Attention

As a special case of the noncascaded framework, we consider the scenario where all stages are denoised concurrently throughout the entire trajectory by setting $t_2 = T$ and $t_1 = 0$. We refer to this configuration as *complete overlap*. Under this setting, the generation of syntax and text is tightly coupled, denoising jointly at every timestep.

The complete overlap enables more efficient implementation of multi-stage generation. In traditional diffusion, each modality requires a separate model with its own self-attention and cross-attention layers for denoising and conditioning, which increases model complexity and complicates the generation pipeline. To improve efficiency, we propose a unified self-attention architecture that achieves mutual conditioning in a single forward pass. Instead of preserving separate diffusion models, our architecture consolidates the denoising of syntax and text into one compact model, fully leveraging the advantage of parallel sampling.

At each timestep t , syntactic embeddings \mathbf{s}_t and text embeddings \mathbf{x}_t are jointly processed through a unified attention layer. Specifically, the queries, keys, and values from the syntactic stream, $(Q(\mathbf{s}_t), K(\mathbf{s}_t), V(\mathbf{s}_t))$, and the textual stream, $(Q(\mathbf{x}_t), K(\mathbf{x}_t), V(\mathbf{x}_t))$, are concatenated, respectively, to form the inputs (Q, K, V) . A single attention computation is then performed, producing a block-matrix attention map:

$$QK^T = \begin{bmatrix} Q(\mathbf{s}_t)K(\mathbf{s}_t)^T & Q(\mathbf{s}_t)K(\mathbf{x}_t)^T \\ Q(\mathbf{x}_t)K(\mathbf{s}_t)^T & Q(\mathbf{x}_t)K(\mathbf{x}_t)^T \end{bmatrix}.$$

The diagonal blocks represent self-attention, capturing intra-modality interactions, while the off-diagonal blocks correspond to cross-attention, enabling inter-modality conditioning. This unified self-attention reduces architectural redundancy and facilitates seamless bidirectional information flow between structure and content.

We refer to our method as Syntax-Text Diffusion (**STDiff**). Its training objective is

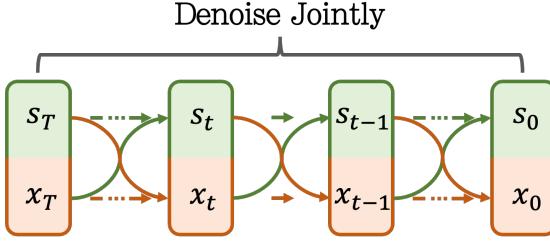


Figure 10: An illustration of the proposed **STDiff** framework. At each timestep, the syntactic and text embeddings are jointly denoised with bidirectional conditioning.

a special case of the general noncascaded loss in (2) with $t_2 = T$ and $t_1 = 0$. Compared to separately modeling intra- and inter-modality interactions, the unified attention layer achieves dynamic mutual conditioning while promoting parameter sharing in attention projection matrices. For example, rather than introducing distinct syntactic key projection matrices, $W_{K,\text{self}}^s$ in the syntax diffusion and $W_{K,\text{cross}}^s$ in the text diffusion, we only employ a single W_K^s within θ_s to handle both roles. Similarly, the query projection is operated through W_Q^s in θ_s , replacing the need for $W_{Q,\text{self}}^s$ and $W_{Q,\text{cross}}^s$ in the syntax diffusion.

Remark 2. *Complete overlap further expands the representational capacity of our model. While the cascaded framework imposes a hierarchical structure between syntax and text by modeling the true underlying joint distribution $q(\mathbf{s}_0, \mathbf{x}_0)$ with $p_{\theta_x}(\mathbf{x}_0 \mid \mathbf{s}_0)p_{\theta_s}(\mathbf{s}_0)$, the complete noncascaded framework relaxes this assumption by directly modeling $p_{\theta_s, \theta_x}(\mathbf{s}_0, \mathbf{x}_0)$. This formulation provides greater flexibility in approximating the true data distribution. We elaborate on this aspect in Section ?? of the Supplementary Materials.*

5 Experiments

We evaluate the effectiveness of our proposed cascaded and noncascaded diffusion frameworks across three aspects of text generation: quality, diversity, and personalized control.

5.1 Datasets, Baselines, and Evaluation Metrics

We conduct extensive experiments on benchmark datasets. In this section, we present the results on two representative datasets with sentiment polarity and stylistic attributes.

Additional results on other datasets are provided in the Supplementary Materials.

The **Yelp Review** dataset consists of user reviews labeled with sentiments, available at <https://business.yelp.com/>. We categorize reviews with 4 or 5 stars as *positive*, 3 stars as *neutral*, and 1 or 2 stars as *negative*. After preprocessing, we construct a balanced training set by randomly sampling 30,000 reviews for each sentiment class. The **Emotion** dataset (Saravia et al., 2018) contains English tweets annotated with six fine-grained emotional categories including *sadness*, *joy*, *love*, *anger*, *fear*, and *surprise*. To ensure adequate syntactic richness, we filter tweets to retain those with at least 10 words, and then randomly sample 8,000 examples for each emotion.

We consider two tasks under different levels of structural conditioning: **Free Generation** and sequence-to-sequence **Sentence Expansion**. In free generation, the model synthesizes sentences directly from random noise, without explicit structural information. In contrast, sentence expansion provides subject-verb-object (SVO) triplets as input, e.g., *I love statistics*, which is more tractable with compositional guidance.

We evaluate our proposed models, **SynText** (cascaded) and **STDiff** (noncascaded), against both diffusion and autoregressive baselines. For diffusion benchmarks, we consider **LD4LG** (Lovelace et al., 2023), which generates text from compressed latent representations via auxiliary encoder-decoder modules. We also include **DiffuSeq** (Gong et al., 2023) for the sentence expansion task, which is tailored for sequence-to-sequence generation. For autoregressive comparison, we fine-tune **GPT-2 Medium** (Radford et al., 2019), a pre-trained model with a comparable parameter size, as a fair baseline of mainstream LLMs.

The evaluation metrics encompass three dimensions: generation quality, output diversity, and stylistic faithfulness. For generation quality, we report perplexity (**Ppl**), which reflects the negative log-likelihood of the generated text, and **Mauve** score (Pillutla et al., 2021), which quantifies the distributional similarity between generated and reference text dataset. We also include **BertScore** (Zhang et al., 2020) for sentence expansion, measuring semantic similarity between the generated and reference examples. For output diversity, we compute the repetition rates of 3-grams and 4-grams (**Div-3/4**). For stylistic faithfulness, we report the classification accuracy (**Acc**) on the generated samples measured by the DistilBERT classifier (Sanh et al., 2019). In addition, the class-wise **Mauve** scores reflect how well the generation captures the true style distribution, independent of classifier choice. These evaluation metrics are widely used in the NLP and LLM literature (Bao et al., 2019; Mozafari et al., 2020; Lovelace et al., 2023; Lou et al., 2024).

To assess the quality of syntactic learning, we introduce a new metric named syntactic n-gram overlap (**SGO**), which quantifies corpus-level n-gram similarity of syntactic patterns between generated and reference text. Formally, we define **SGO** as the geometric mean of individual n -gram overlap scores $\text{Overlap}(n)$ for $n = 2, 3, 4$:

$$\text{SGO} = \exp \left(\sum_{n \in \{2,3,4\}} \frac{1}{3} \cdot \log(\text{Overlap}(n) + \epsilon) \right),$$

where ϵ is a small smoothing term. Details on the definitions of $\text{Overlap}(n)$ and **SGO** are provided in Section ?? of the Supplementary Materials.

5.2 Quantitative Results

5.2.1 Free Generation

In the free generation setting, models generate text from scratch solely based on a target style or label, without additional structural or lexical information. This task allows com-

plete flexibility in sentence composition and evaluates the model’s ability to produce fluent, diverse, and stylistically appropriate text.

For evaluation, each model generates 1,000 random samples per style, which are then compared against 1,000 human-written references of the same style from an independent test set. The results on the **Yelp Review** dataset are summarized in Table 1.

Table 1: Free text generation performance of different methods on the Yelp Review dataset.

Bold indicates the best performance, and underline indicates the second best.

		Ppl \downarrow	Mauve \uparrow	Div-3/4 \downarrow	Acc \uparrow	SGO \uparrow
<i>Positive</i>	SynText	<u>81.313</u>	0.420	0.151/ <u>0.037</u>	<u>0.934</u>	<u>0.960</u>
	STDiff	106.819	0.533	0.127/0.028	0.964	0.962
	LD4LG	189.097	0.400	<u>0.145</u> /0.038	0.860	0.913
	GPT-2-M	65.727	0.395	0.242/0.132	0.742	0.845
<i>Negative</i>	SynText	<u>100.570</u>	0.333	0.123/ <u>0.024</u>	<u>0.889</u>	0.942
	STDiff	131.806	0.421	0.103/0.018	0.931	0.950
	LD4LG	200.224	<u>0.365</u>	<u>0.117</u> /0.025	0.804	0.908
	GPT-2-M	37.498	0.237	0.267/0.139	0.568	<u>0.947</u>
<i>Neutral</i>	SynText	<u>86.023</u>	<u>0.356</u>	0.155/0.036	<u>0.835</u>	<u>0.951</u>
	STDiff	113.021	0.446	0.138/0.031	0.894	0.952
	LD4LG	180.437	0.333	<u>0.143</u> / <u>0.034</u>	0.733	0.909
	GPT-2-M	49.692	0.162	0.253/0.136	0.560	0.906

We highlight three important observations from the results. First, diffusion language models exhibit significantly higher diversity in generated text compared to the AR competitor, while maintaining comparable or even superior generation quality. The higher 3-gram and 4-gram repetition rates of GPT-2-M suggest its tendency to rely on frequent phrases. Although GPT-2-M achieves lower perplexity, this metric is computed using GPT-2 and favors next-token prediction, making it less suitable for evaluating non-AR frameworks.

Second, incorporating syntax offers potential to enhance coherence and stylistic align-

ment. Compared to LD4LG, a latent diffusion based on text embeddings, our syntax-aware methods generally yield higher Mauve, SGO, and classification accuracy across all sentiment styles. Notably, the superior diversity of STDiff indicates that structural guidance promotes diffusion’s advantage, yielding more diverse and coherent generations.

Third, the noncascaded design improves performance compared to the cascaded framework. STDiff consistently demonstrates stronger capability in capturing both semantic and structural patterns. While SynText shows slightly lower perplexity, this may stem from its reliance on fixed syntactic conditions, leading to more predictable word choices but less flexibility and diversity. Therefore, we adopt STDiff for subsequent comparisons.

To better evaluate stylistic learning, we assess both the ability to capture style-specific patterns and to distinguish among different styles. We generate 1,000 samples for each sentiment and compute Mauve scores between each generated set and all reference sets to quantify their similarities. We compare STDiff with LD4LG and GPT-2-M by reporting the pairwise relative differences in Mauve in Table 2, where rows indicate the sentiment of generated samples and columns indicate the reference style. We observe that STDiff achieves higher Mauve along the diagonal, indicating its superior ability to model style-consistent distributions. Meanwhile, the off-diagonal differences are mostly negative, further suggesting that STDiff generates more distinguishable and faithful stylistic outputs. We also observe two positive off-diagonal entries. This can be attributed to two factors. First, neutral reviews often overlap with negative ones in lexical and syntactic patterns, which naturally yields higher Mauve similarity between these two styles. Second, GPT-2-M performs poorly on the neutral class, leading to low Mauve scores against all references. Consequently, once STDiff captures stylistic features that are shared with negative references, the low baseline of GPT-2-M (0.070) amplifies the relative difference.

Table 2: Relative differences (%) in Mauve scores between STDiff and comparison methods. Gen: styles of the generated text. Ref: styles of the reference text. **Bold** indicates our method performs better than the competitor.

STDiff vs LD4LG				STDiff vs GPT-2-M			
Gen \ Ref	Positive	Negative	Neutral	Gen \ Ref	Positive	Negative	Neutral
Positive	33.3%	-14.7%	-22.0%	Positive	34.9%	-57.4%	-57.4%
Negative	-17.9%	15.7%	-32.3%	Negative	-27.3%	77.6%	-44.1%
Neutral	-31.4%	4.0%	33.9%	Neutral	-7.7%	84.3%	175.3%

We further investigate the effectiveness of our **PLayer** against the standard approach of isolated personalized tokens, referred to as **Token**. For **Token**, each user is assigned a unique embedding, optimized independently during training without information sharing. We incorporate both strategies into the noncascaded pipeline and report the metrics concerning stylistic control in Table 3. The **PLayer** consistently outperforms the personalized token approach across all styles, demonstrating that leveraging shared patterns achieves higher semantic and structural fidelity. Notably, **Token** achieves better performance than LD4LG, confirming the benefit of incorporating syntactic information.

Table 3: Comparison of our **PLayer** personalization with isolated personalized token.

Methods	Positive			Negative			Neutral		
	Mauve	Acc	SGO	Mauve	Acc	SGO	Mauve	Acc	SGO
PLayer	0.533	0.964	0.962	0.421	0.931	0.950	0.446	0.894	0.952
Token	0.453	0.927	0.954	0.387	0.916	0.939	0.411	0.853	0.938
LD4LG	0.400	0.860	0.913	0.365	0.804	0.908	0.333	0.733	0.909

The **Emotion** dataset comprises six distinct emotional styles, posing a more challenging

task for faithful text personalization. We evaluate all models under the same procedure as in the Yelp experiments, and report their performance in terms of Mauve, Acc, and SGO to assess the quality of personalization in Table 4.

Table 4: Free text generation performance of different methods on the Emotion dataset.

	<i>Sadness</i>			<i>Joy</i>			<i>Love</i>		
	Mauve	Acc	SGO	Mauve	Acc	SGO	Mauve	Acc	SGO
STDiff	0.676	0.932	0.986	0.748	0.942	0.986	0.585	0.958	0.963
LD4LG	0.621	0.857	0.934	0.533	0.852	0.921	0.524	0.899	0.932
GPT-2-M	0.495	0.351	0.895	0.448	0.578	0.928	0.274	0.072	0.911
	<i>Anger</i>			<i>Fear</i>			<i>Surprise</i>		
	Mauve	Acc	SGO	Mauve	Acc	SGO	Mauve	Acc	SGO
STDiff	0.753	0.933	0.986	0.770	0.942	0.985	0.776	0.980	0.971
LD4LG	0.626	0.864	0.928	0.574	0.821	0.928	0.538	0.909	0.909
GPT-2-M	0.313	0.852	0.940	0.314	0.676	0.930	0.249	0.062	0.877

We observe that **STDiff** consistently outperforms LD4LG and GPT-2-M across all emotional styles, exhibiting its robustness in capturing fine-grained stylistic patterns. Compared to LD4LG, STDiff achieves significant gains in both semantic alignment and structural fidelity, suggesting that explicit syntactic guidance enhances control. In contrast, GPT-2-M struggles to identify different emotions. These results validate the advantage of diffusion with syntax-aware personalization in handling diverse emotional identities.

5.2.2 Sentence Expansion

In the sentence expansion task, the model is conditioned on structured inputs of SVO triplets. The objective is to expand these core components into a complete sentence. Com-

pared to free generation, this sequence-to-sequence task offers explicit structural guidance, which reduces uncertainty in content planning.

We employ spaCy (Honnibal and Montani, 2017) to extract SVO triplets from human-written sentences. The extracted triplets are concatenated into a single sequence, which serves as the structural context input. During generation, the model is conditioned on the triplets in the reference sentences to produce full sentences accordingly. The results of all methods on the Yelp Review dataset are reported in Table 5.

Table 5: Sentence expansion performance of different methods on the Yelp Review dataset.

Bold indicates the best performance, and underline indicates the second best.

		BertScore↑	Mauve↑	Div-3/4↓	Acc↑	SGO↑
<i>Positive</i>	STDiff	0.861	0.904	0.152/0.038	0.935	0.986
	DiffuSeq	0.858	0.726	0.221/0.057	0.544	0.914
	LD4LG	0.861	<u>0.847</u>	<u>0.167/0.045</u>	<u>0.913</u>	<u>0.952</u>
	GPT-2-M	0.858	0.397	0.169/0.051	0.833	0.745
<i>Negative</i>	STDiff	0.861	0.892	0.123/0.026	0.872	0.991
	DiffuSeq	0.858	0.660	0.197/0.044	0.520	0.904
	LD4LG	0.861	<u>0.754</u>	<u>0.125/0.026</u>	<u>0.843</u>	<u>0.947</u>
	GPT-2-M	0.859	0.483	0.134/0.035	0.787	0.741
<i>Neutral</i>	STDiff	<u>0.861</u>	0.932	0.147/0.034	<u>0.827</u>	0.987
	DiffuSeq	0.860	0.749	0.245/0.072	0.605	0.906
	LD4LG	0.862	<u>0.788</u>	<u>0.170/0.042</u>	0.844	<u>0.942</u>
	GPT-2-M	0.860	0.475	0.173/0.051	0.712	0.745

We conclude that STDiff consistently achieves strong performance across all metrics and emotion categories. Notably, it obtains the highest Mauve and SGO, along with the lowest Div-3/4, indicating its superior ability to capture both semantic and syntactic characteristics while maintaining generation diversity. Although LD4LG performs competitively in terms of BertScore and Acc, it falls short on more comprehensive metrics, such as Mauve

and SGO, as well as diversity, suggesting possible overfitting to local patterns. Moreover, GPT-2-M exhibits significantly lower Mauve and SGO scores, highlighting the advantage of diffusion in global refinement when constructing sentences given textual fragments. Overall, these results demonstrate the effectiveness of our method in leveraging structural inputs to improve sentence fluency and stylistic faithfulness.

5.3 Qualitative Illustrations

5.3.1 Diversity and Zero-Shot Generalization

To showcase the generating diversity and stylistic alignment of our STDiff, we present generated sentences under different sentiment styles. For better interpretability, we illustrate using the sentence expansion task, where the same SVO triplet is fed to both GPT-2-M and STDiff. The examples conditioned on sentiment labels are shown in Figure 11.

Prompt	Style	GPT-2-M	STDiff
<i>I visited restaurant</i>	Positive	<i>I have visited this restaurant</i> many times and loved their food and the <i>service</i> was outstanding	<i>For all my friends I visited this restaurant twice</i> and would definitely go back for the excellent service
	Negative	<i>I have visited this restaurant</i> three times and experienced no customer <i>service</i> , slow <i>service</i>	<i>Changing my review to one star because I visited this restaurant twice</i> after finding bugs of fish
	Neutral	<i>I have visited this restaurant</i> three times, and every time the food was bland, and the <i>service</i> was just ok	<i>Have to update my review longer because I visited this restaurant twice</i> and both times have gotten worse

Figure 11: Sentence expansion examples from GPT-2-M and our diffusion model under the same prompt and different sentiment styles.

Although GPT-2-M has a comparable number of parameters to our STDiff, its diversity and stylistic expressiveness are notably limited. We highlight the repetitive segments in purple. The repeated fragments in STDiff outputs primarily reflect the input prompt,

whereas GPT-2-M exhibits additional redundancy, including fixed sentence structures and overused content such as the word “service.” Moreover, the stylistic variation in GPT-2-M is mostly confined to a few sentiment-bearing words, while STDiff generates richer expressions reflecting user experience and exhibits more diverse syntactic forms.

To further explore the generalizability of our personality modeling framework, we consider a zero-shot generation setting. During training, we obtain distinct personality weights on shared representations corresponding to the canonical sentiment styles. At inference time, we synthesize unseen styles by combining these learned weights. For example, we can amplify sentiments by extrapolating, yielding stronger positive or negative. We can also blend styles to create sentences reflecting intermediate characteristics and mixed tones. The generated sentences are shown in Figure 12. These results demonstrate that our sharing framework learns a unified latent space of personalities, enabling smooth generalization beyond training styles, which contrasts with conventional approaches that learn discrete embeddings for each style independently.

5.3.2 Personalization Visualization

To illustrate the effectiveness of our method in capturing human-written patterns, we analyze the generated text from both structural and semantic perspectives on the Yelp Review dataset. In Figure 13, we examine syntactic structures by measuring POS frequencies for each style, and visualize sentence embeddings using outputs from the text diffusion.

We observe that the generated text closely resembles the reference data in the motivating examples (see Figure 2 and 3), demonstrating the model’s ability to preserve underlying distributions of each style. Moreover, compared to real data, our method exhibits clearer separation between different styles, particularly in the embedding space. This suggests that our model enhances inter-style distinctiveness, resulting in more style-consistent generation.

Personality weights			Seen in training?	Generated sentence
γ_{Positive}	γ_{Negative}	γ_{Neutral}		
x1	0	0	✓ Yes	This place is amazing and the products they have are perfect the pizza is great and the atmosphere is absolutely a pleasure! I will be back...
x2	0	0	✗ No	Love this place some the best food here has the perfect The pizza is great and the owner is absolutely a sweet man I will be back.
0	x1	0	✓ Yes	Over rated because most of the kids here has no flavor The pizza is good but this place doesn't deserve its reputation I'm one down star
0	x2	0	✗ No	Just disappointed because most of the people here has been saying the pizza is good but this place doesn't deserve my business I'm one down star
0	0	x1	✓ Yes	The staff where wonderful and the girls there has are great the pizza is good but this location doesn't meet my expectations I'm fishy.
x0.3	x0.3	0	✗ No	The employees where wonderful and the products they have are perfect The pizza is great but this location doesn't meet my expectations I'm one off star

Figure 12: Zero-shot generation by adjusting personality weights. Rows highlighted in red are unseen weight combinations at training.

To better interpret the learned personality representations, we visualize the personality weights for each sentiment style in Figure 14. While all styles share certain common embeddings, each style is characterized by distinctive key components. For instance, embeddings 2 and 7 are more prominent in the *Positive* style, whereas embedding 6 is more active in the *Negative* style. Furthermore, the *Neutral* weights exhibit an intermediate pattern, indicating the model’s ability to learn a semantically meaningful latent space.

6 Discussion

In this work, we propose a syntax-guided diffusion language model that incorporates structural information into text generation to enhance quality, enrich diversity, and strengthen personalized control. We construct a cascaded framework to generate sound syntactic conditions, and further extend it to a noncascaded design to utilize the parallel nature

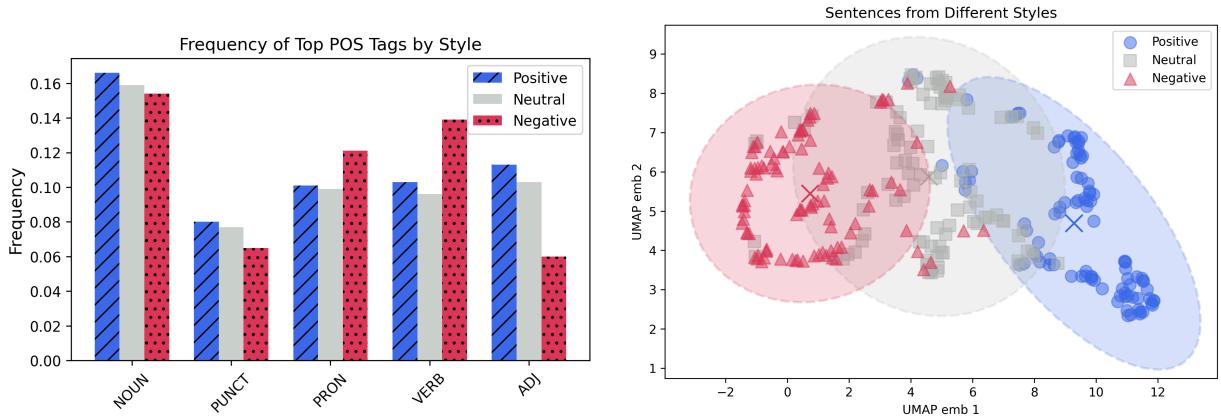


Figure 13: Frequencies of POS tags and UMAP of sentence embeddings in generated text.

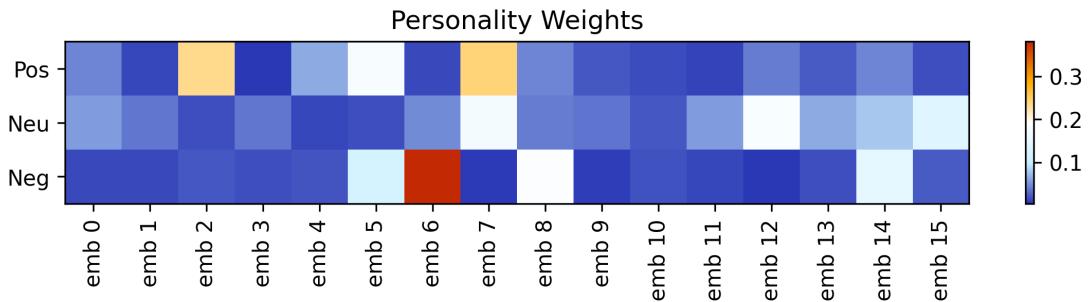


Figure 14: Personality weights for different styles over different personality embeddings.

of diffusion models. To achieve fine-grained personalization, we introduce a shared representation mechanism, which enables both faithful reconstruction and flexible zero-shot generation. Real-world text experiments show that our method outperforms baselines under the same scale of model parameters across both free generation and sentence expansion tasks. Various illustrations validate the adaptability and reliability of our approach.

An interesting direction for future work is to explore the use of syntactic information in discrete diffusion models for text generation, which have recently gained increasing attention (Lou et al., 2024; Yu et al., 2025). Although our current implementation adopts a continuous formulation, the core ideas of structural conditioning and shared personalized representations are highly generalizable. The relatively small vocabulary size of syntactic

entities may further enable efficient and stable discrete sampling. In addition, while we employ POS tags as a reliable form of syntactic supervision, future studies could utilize more fine-grained signals, such as dependency relations and constituency trees.

Another future direction lies in refining the noncascaded architecture. As a generalization of the cascaded framework, the noncascaded design is applicable beyond our current setting and can extend to a broad range of hierarchical generation scenarios, including multi-stage and multi-modal generation. Different applications may benefit from varying degrees of overlap between stages. A systematic exploration of information flow across different stages can provide deeper insights into the trade-offs between sequential and parallel generation, thereby motivating new large language model frameworks.

References

- Alhafni, B., V. Kulkarni, D. Kumar, and V. Raheja (2024). Personalized text generation with fine-grained linguistic control. In *Proceedings of the 1st Workshop on Personalization of Generative AI Systems (PERSONALIZE 2024)*, pp. 88–101.
- Bao, Y., H. Zhou, S. Huang, L. Li, L. Mou, O. Vechtomova, and et al. (2019). Generating sentences from disentangled syntactic and semantic spaces. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6008–6019.
- Chen, T., R. Zhang, and G. Hinton (2023). Analog bits: Generating discrete data using diffusion models with self-conditioning. In *The Eleventh International Conference on Learning Representations*.
- Chomsky, N. (1957). *Syntactic Structures*. Berlin, Boston: De Gruyter Mouton.

DeepSeek-AI, A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, and et al. (2025). Deepseek-v3 technical report.

Frenkel, Y., Y. Vinker, A. Shamir, and D. Cohen-Or (2025). Implicit style-content separation using b-lora. In *European Conference on Computer Vision*, pp. 181–198. Springer.

Gao, T., X. Yao, and D. Chen (2021). SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 6894–6910.

Gong, S., M. Li, J. Feng, Z. Wu, and L. Kong (2023). Diffuseq: Sequence to sequence text generation with diffusion models. In *The Eleventh International Conference on Learning Representations*.

Hayou, S., N. Ghosh, and B. Yu (2024). LoRA+: Efficient low rank adaptation of large models. In *Proceedings of the 41st International Conference on Machine Learning*, Volume 235, pp. 17783–17806.

Ho, J., A. Jain, and P. Abbeel (2020). Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, Volume 33, pp. 6840–6851.

Ho, J., C. Saharia, W. Chan, D. J. Fleet, M. Norouzi, and T. Salimans (2022). Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research* 23(47), 1–33.

Holtzman, A., J. Buys, L. Du, M. Forbes, and Y. Choi (2020). The curious case of neural text degeneration. In *International Conference on Learning Representations*.

Hong, S., D. Yang, J. Choi, and H. Lee (2018). Inferring semantic layout for hierarchi-

cal text-to-image synthesis. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7986–7994.

Honnibal, M. and I. Montani (2017). spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.

Hoogeboom, E., D. Nielsen, P. Jaini, P. Forré, and M. Welling (2021). Argmax flows and multinomial diffusion: Learning categorical distributions. In *Advances in Neural Information Processing Systems*, Volume 34, pp. 12454–12465.

Hu, E. J., Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen (2022). LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.

InceptionLabs, S. Khanna, S. Kharbanda, S. Li, H. Varma, E. Wang, and et al. (2025). Mercury: Ultra-fast language models based on diffusion.

Li, X., J. Thickstun, I. Gulrajani, P. S. Liang, and T. B. Hashimoto (2022). Diffusion-LM improves controllable text generation. In *Advances in Neural Information Processing Systems*, Volume 35, pp. 4328–4343.

Li, Y., L. Cui, J. Yan, Y. Yin, W. Bi, S. Shi, and Y. Zhang (2023). Explicit syntactic guidance for neural text generation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pp. 14095–14112.

Li, Y., K. Zhou, W. X. Zhao, and J.-R. Wen (2023). Diffusion models for non-autoregressive text generation: A survey. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23*, pp. 6692–6701. Survey Track.

Lou, A., C. Meng, and S. Ermon (2024). Discrete diffusion modeling by estimating the ratios of the data distribution. In *Proceedings of the 41st International Conference on Machine Learning*.

Lovelace, J., V. Kishore, C. Wan, E. Shekhtman, and K. Q. Weinberger (2023). Latent diffusion for language generation. In *Advances in Neural Information Processing Systems*, Volume 36, pp. 56998–57025.

McInnes, L., J. Healy, and J. Melville (2020). Umap: Uniform manifold approximation and projection for dimension reduction.

Mireshghallah, F., V. Shrivastava, M. Shokouhi, T. Berg-Kirkpatrick, R. Sim, and D. Dimtriadis (2022). UserIDentifier: Implicit user representations for simple and effective personalized sentiment analysis. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 3449–3456.

Mozafari, J., A. Fatemi, and P. Moradi (2020). A method for answer selection using distilbert and important words. In *2020 6th International Conference on Web Research (ICWR)*, pp. 72–76. IEEE.

Mummery, D. (2025). Artificial intelligence ... or not? *British Journal of General Practice* 75(757), 368–368.

Nie, S., F. Zhu, Z. You, X. Zhang, J. Ou, J. Hu, and et al (2025). Large language diffusion models.

OpenAI, A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, and et al. (2024). Gpt-4o system card.

Peebles, W. and S. Xie (2023). Scalable diffusion models with transformers. In *Proceedings*

of the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 4195–4205.

Pillutla, K., S. Swayamdipta, R. Zellers, J. Thickstun, S. Welleck, Y. Choi, and Z. Harchaoui (2021). Mauve: Measuring the gap between neural text and human text using divergence frontiers. In *Advances in Neural Information Processing Systems*, Volume 34, pp. 4816–4828.

Radford, A., J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. (2019). Language models are unsupervised multitask learners. *OpenAI blog* 1(8), 9.

Ramesh, A., P. Dhariwal, A. Nichol, C. Chu, and M. Chen (2022). Hierarchical text-conditional image generation with clip latents.

Rout, L., C. Caramanis, and S. Shakkottai (2025). Anchored diffusion language model.

Saharia, C., W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, and et al (2022). Photo-realistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*, Volume 35, pp. 36479–36494.

Salemi, A., S. Mysore, M. Bendersky, and H. Zamani (2024). LaMP: When large language models meet personalization. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, pp. 7370–7392.

Sanh, V., L. Debut, J. Chaumond, and T. Wolf (2019). Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. In *Advances in Neural Information Processing Systems Workshop on Energy Efficient Machine Learning and Cognitive Computing*.

Saravia, E., H.-C. T. Liu, Y.-H. Huang, J. Wu, and Y.-S. Chen (2018). CARER: Con-

textualized affect representations for emotion recognition. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 3687–3697.

Song, Y. and S. Ermon (2019). Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, Volume 32.

Sourati, Z., F. Karimi-Malekabadi, M. Ozcan, C. McDaniel, A. Ziabari, J. Trager, and et al (2025). The shrinking landscape of linguistic diversity in the age of large language models.

Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*, Volume 30.

Wan, Y., G. Pu, J. Sun, A. Garimella, K.-W. Chang, and N. Peng (2023). “Kelly is a warm person, Joseph is a role model”: Gender biases in LLM-generated reference letters. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 3730–3748.

Yu, R., Q. Li, and X. Wang (2025). Discrete diffusion in large language and multimodal models: A survey.

Zhang, T., V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi (2020). Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

Zhong, W., D. Tang, J. Wang, J. Yin, and N. Duan (2021). UserAdapter: Few-shot user learning in sentiment analysis. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 1484–1488.